



UNIVERSITÄT
LEIPZIG

Automatisierte Analyse von Impedanzspektren mittels konstruktivistischen maschinellen Lernens

Von der Fakultät für Mathematik und Informatik
der Universität Leipzig angenommene

DISSERTATION

zur Erlangung des akademischen Grades

DOCTOR RERUM NATURALIUM
(Dr. rer. nat)

im Fachgebiet Informatik

vorgelegt von

Dipl.-Inform. Thomas Schmid

geboren am 30. Januar 1982 in Schwäbisch Hall

Die Annahme der Dissertation wurde empfohlen von

1. Prof. Dr. Martin Bogdan (Universität Leipzig)
& Prof. Dr. Dorothee Günzel (Charité Berlin)
2. Prof. Dr. Günther Palm (Universität Ulm)

Die Verleihung des akademischen Grades erfolgt mit Bestehen
der Verteidigung am 11. Juni 2018 mit dem Gesamtprädikat Magna cum laude.

Inhaltsverzeichnis

Vorwort	v
Werdegang	vi
Selbstständigkeitserklärung	vi
Zusammenfassung	vii
I. MOTIVATION UND STAND DER TECHNIK	1
1. Einleitung	3
2. Menschliches Lernen	9
2.1. Konstruktivismus	9
2.2. Paradigmen der Psychologie	10
2.3. Modelle und Modellbildung	11
2.4. Kognitive Ebenen	12
2.5. Intersubjektivität	13
3. Maschinelles Lernen	15
3.1. Überwachung versus Selbstorganisation	15
3.2. Das neurophysiologische Paradigma	17
3.3. Multiperspektivisches Lernen	18
3.4. Auswahl und Extraktion von Features	20
3.5. Evaluation und Bewertungsmaßstäbe	21
4. Epithel-Analyse mittels Impedanzspektroskopie	23
4.1. Wechselstrom und Widerstand	23
4.2. Impedanzspektroskopie in der Praxis	24
4.3. Struktur und Funktion von Epithelien	26
4.4. Epitheliale Impedanzanalyse	27
II. METHODIK UND ERGEBNISSE	29
5. Konstruktives maschinelles Lernen	31
5.1. Grundannahmen	31
5.2. Pragmatisch definierte maschinelle Modelle	32
5.2.1. Begriffsklärung	32
5.2.2. Verwandtschaft, Vererbung und Hierarchien	33
5.2.3. Implementierungsaspekte	34
5.3. Konstruktivistisches maschinelles Lernen	35
5.3.1. Ablauf	35
5.3.2. Implementierungsaspekte	37
5.4. Konstruktion maschineller Modelle	38
5.4.1. Konzeptuelles Wissen	39

5.4.2.	Prozedurales Wissen	40
5.5.	Rekonstruktion maschineller Modelle	41
5.5.1.	Lernverfahren	41
5.5.2.	Komplexitätsreduktion	42
5.5.3.	Modellselektion	43
5.6.	Dekonstruktion maschineller Modelle	45
5.6.1.	ΣZ -Dekonstruktion	46
5.6.2.	$T\Sigma$ -Dekonstruktion	46
5.6.3.	TZ -Dekonstruktion	47
5.6.4.	Vollständige Dekonstruktion	47
6.	Simulation von Impedanzspektroskopie-Messungen	49
6.1.	Ersatzschaltkreise für epitheliales Gewebe	49
6.2.	Eigenschaften epithelialer Zelllinien und ihrer funktionalen Zustände	51
6.3.	Berechnung einer theoretischen Impedanz	53
6.4.	Modellierung von Messfehlern	54
6.5.	Systematische Synthetisierung von Datensätzen	56
6.6.	Abgleich modellierter und gemessener Daten	56
7.	Konzeptuelles Wissen für Impedanzspektren	59
7.1.	Exploration einer konzeptuellen Wissensdomäne	60
7.1.1.	Explorationsverlauf	60
7.1.2.	Charakterisierung der explorierten Wissensdomäne	62
7.2.	Adaption konzeptuellen Wissens	65
7.2.1.	Adaptionsverlauf	66
7.2.2.	Anwendung auf Testdaten	67
8.	Prozedurales Wissen für Impedanzspektren	69
8.1.	Exploration einer prozeduralen Wissensdomäne	70
8.1.1.	Explorationsverlauf	70
8.1.2.	Charakterisierung der explorierten Wissensdomäne	72
8.2.	Adaption prozeduralen Wissens	75
8.2.1.	Adaptionsverlauf	75
8.2.2.	Anwendung auf Testdaten	77
III.	DISKUSSION UND SCHLUSSFOLGERUNGEN	79
9.	Epitheliale Impedanzanalyse	81
9.1.	Modellierung impedanzspektroskopischer Messungen	81
9.2.	Unterscheidung von Zelllinien	83
9.3.	Quantifizierung der epithelialen Kapazität	84
9.4.	Generalisierbarkeit des Analyseverfahrens	85
10.	Implementierungs- und Verfahrensaspekte	87
10.1.	Automatisierung von unüberwachtem und überwachtem Lernen	87
10.2.	Abstraktion und Wissens Ebenen	89
10.3.	Modellverwaltung	90
10.4.	Parallelisierung	92

11. Chancen und Risiken konstruktivistischen Lernens	95
11.1. Umgang mit Mehrdeutigkeiten	95
11.2. Nachvollziehbarkeit	97
11.3. Modellierungslücken	98
11.4. Fazit und Ausblick	100
ANHANG	104
A. Implementierung und Konfiguration	105
A.1. Unüberwachte Lernverfahren	105
A.2. Überwachte Lernverfahren	107
A.3. Ablaufsteuerung	108
A.4. Konfigurationsreferenz	110
B. Modellierung epithelialer Zellkulturen	115
B.1. Berechnung der transepithelialen Impedanz Z^T	116
B.2. Abgleich modellierter und gemessener Daten	119
B.2.1. Methode M1: Diskrete Approximation	119
B.2.2. Methode M2: Cole-Cole-Fit	119
B.2.3. Methode M3: Maschinelle Lernverfahren	120
B.2.4. Ähnlichkeitsdiagramme	120
B.3. HT-29/B6	121
B.3.1. Kontrollbedingungen	121
B.3.2. Manipulation mit Nystatin	122
B.3.3. Manipulation mit EGTA	123
B.3.4. Manipulation mit EGTA und Nystatin	124
B.4. IPEC-J2	125
B.4.1. Kontrollbedingungen	125
B.4.2. Manipulation mit Nystatin	126
B.4.3. Manipulation mit EGTA	127
B.4.4. Manipulation mit EGTA und Nystatin	128
B.5. MDCK I	129
B.5.1. Kontrollbedingungen	129
B.5.2. Manipulation mit Nystatin	130
B.5.3. Manipulation mit EGTA	131
C. Charakterisierung der Datenbasis	133
C.1. Input-Features	133
C.2. Zielparameter	135
C.3. Metadaten	136
D. Charakterisierung des adaptierten Wissens	137
D.1. Konzeptuelles Wissen	137
D.2. Prozedurales Wissen	140
Literaturverzeichnis	143
Abbildungsverzeichnis	160
Tabellenverzeichnis	162

VORWORT

Diese Arbeit ist am Institut für Informatik der Universität Leipzig und am Department of Biomedical Engineering der University of Alabama at Birmingham entstanden. Ihre Fragestellungen gingen aus einem anwendungsbezogenen Kooperationsprojekt mit der Charité Berlin hervor, aus dem sich am Ende grundsätzliche Fragen des maschinellen Lernens ergaben. Insgesamt stellt diese Arbeit somit das Ergebnis einer problemorientierten Grundlagenforschung dar.

Mein besonderer Dank gilt Prof. Dr. Martin Bogdan, der meine Forschung bereits seit meiner Diplomarbeit fachlich begleitet und unterstützt hat. Er weckte nicht nur mein Interesse an künstlichen neuronalen Netzen und biomedizinischen Anwendungen, sondern war auch bereit, die Entwicklung eines Ansatzes maschinellen Lernens zu fördern, der von etablierten Pfaden abweicht. Seine Betreuung und Unterstützung machten diese Arbeit überhaupt erst möglich.

Ebenso dankbar für ihre Betreuung bin ich Prof. Dr. Dorothee Günzel vom Institut für Klinische Physiologie der Charité Berlin. Sie hat mich bereits während meines Studiums für die computergestützte Analyse von Epithelien begeistert und anschließend ermutigt, dieses Thema im Rahmen einer Promotion wissenschaftlich weiterzuverfolgen. Ihre langjährige kompetente und engagierte Beratung hat maßgeblich zum Gelingen dieser Arbeit beigetragen.

Diese Arbeit zu verfassen wäre mir jedoch nicht möglich gewesen, hätte ich neben ideeller nicht auch materielle Unterstützung erfahren. Ich möchte insbesondere der FAZIT-Stiftung Gemeinnützige Verlagsgesellschaft danken, die meine Forschungen zwei Jahre lang mit einem Promotionsstipendium unterstützt und damit deren Beginn ermöglicht hat. Ebenso bin ich der deutsch-amerikanischen Fulbright-Kommission dankbar, die mir mit ihrer Förderung einen spannenden und erfolgreichen achtmonatigen Forschungsaufenthalt in den USA ermöglicht hat.

Daneben habe ich regelmäßig Unterstützung und Ermutigung durch weitere Wegbegleiter erfahren. So möchte ich insbesondere Prof. Dr. Yuhua Song danken, die mich im Rahmen meines Aufenthalts an der University of Alabama at Birmingham betreut hat. Meinen Kollegen am Lehrstuhl Technische Informatik danke ich für die langjährige gute Zusammenarbeit, insbesondere Dr. Karim el-Laithy, Jörn Hoffmann und Karin Wenzel. Meinen beiden guten Freunden Hans und Carsten danke ich für das sorgfältige und kritische Gegenlesen der einzelnen Kapitel.

Ein gern zitiertes Sprichwort besagt, in der Ruhe liege die Kraft. Doch Ruhe braucht Zeit, und Zeit ist ein kostbares Gut. Von aller Unterstützung, die ich bei der Arbeit an dieser Dissertation erfahren habe, bin ich daher nicht zuletzt für das Verständnis und die Geduld dankbar, die mein privates Umfeld hierfür aufgebracht hat. Meinen Eltern möchte ich darüber hinaus für die langjährige moralische Unterstützung sowie die Ermöglichung meines Studiums danken.

WERDEGANG

Ich wurde am 30. Januar 1982 in Schwäbisch Hall (Baden-Württemberg) geboren. Kindheit und Jugend verbrachte ich im benachbarten Künzelsau, wo ich 2001 die Allgemeine Hochschulreife erworben habe. Ab 2002 studierte ich in Tübingen, Gaborone (Botswana) und Berlin Bioinformatik. Meine Diplomarbeit mit dem Titel “Using an Artificial Neural Network to Determine Electric Properties of Epithelia” fertigte ich am Institut für Klinische Physiologie der Charité Berlin an. Von 2010 bis 2013 forschte ich als Stipendiat an der Universität Leipzig und der University of Alabama at Birmingham (USA). Seither bin ich als wissenschaftlicher Mitarbeiter am Institut für Informatik der Universität Leipzig tätig.

SELBSTSTÄNDIGKEITSERKLÄRUNG

Hiermit erkläre ich, die vorliegende Dissertation selbstständig und ohne unzulässige fremde Hilfe angefertigt zu haben. Ich habe keine anderen als die angeführten Quellen und Hilfsmittel benutzt und sämtliche Textstellen, die wörtlich oder sinngemäß aus veröffentlichten oder unveröffentlichten Schriften entnommen wurden, und alle Angaben, die auf mündlichen Auskünften beruhen, als solche kenntlich gemacht. Ebenfalls sind alle von anderen Personen bereitgestellten Materialien oder erbrachten Dienstleistungen als solche gekennzeichnet.

Leipzig, den 15. Dezember 2017

ZUSAMMENFASSUNG

Empirische Wissenschaften wie Biologie oder Psychologie konstituieren sich in ihrem Kern aus einer Menge vorläufiger Vermutungen über ihren Untersuchungsgegenstand. Ihre Konzepte und Gesetzmäßigkeiten stellen Verallgemeinerungen zurückliegender Beobachtungen dar, die als gültig angenommen werden, solange es niemandem gelingt, diese zu widerlegen. In den empirischen Naturwissenschaften bilden Messinstrumente wie die Impedanzspektroskopie, deren Ergebnisse unter anderem in Materialwissenschaften oder in der Biomedizin genutzt werden, die Grundlage zur Erstellung von Hypothesen. Eine wissenschaftliche Analyse erfordert jedoch nicht nur das Suchen nach Gesetzmäßigkeiten, sondern ebenso ein Suchen nach Widersprüchen und Alternativen. Hypothesen nicht nur aufzustellen, sondern auch zu hinterfragen, gilt dabei als genuin menschliche Fähigkeit. Zwar werden zur Hypothesenbildung aus empirischen Daten häufig maschinelle Lernverfahren genutzt, doch die Bewertung solcher Hypothesen bleibt bislang ebenso dem Menschen vorbehalten wie das Suchen nach Widersprüchen und Alternativen.

Um diese menschliche Fähigkeit nachzubilden, schlägt die vorliegende Arbeit eine Strategie maschinellen Lernens vor, die sowohl am Leitbild eines kritischen Rationalismus als auch an Prinzipien konstruktivistischer Lerntheorien ausgerichtet ist. Im Gegensatz zu etablierten maschinellen Lernverfahren sehen konstruktivistische Lerntheorien nicht nur ein unüberwachtes oder überwachtes Lernen vor, sondern auch ein Lernen mittels Zweifel. Um einen solchen Lernprozess operationalisieren und automatisieren zu können, werden maschinell erlernte Zusammenhänge hier als Modelle im Sinne der Allgemeinen Modelltheorie nach Herbert Stachowiak interpretiert. Die damit verbundene Definition pragmatischer Eigenschaften als Metadaten erlaubt nicht nur die Selektion zu erlernender Daten aus einem gegebenen Datensatz, sondern auch das Erzeugen und Identifizieren von Beziehungen zwischen Modellen. Dadurch wird es möglich, konkurrierende Modelle für einen gegebenen Datensatz zu unterscheiden und deren Kohärenz zu überprüfen. Insbesondere können so Mehrdeutigkeiten mittels Modell-Metadaten erkannt werden.

Chancen und Risiken eines solchen Ansatzes werden hier anhand automatisierter Analysen impedanzspektroskopischer Messungen aufgezeigt, wie sie in physiologischen Untersuchungen an Epithelien erhoben werden. Da in empirischen Messungen naturgemäß nur Näherungswerte für die Ziel-Messgröße bestimmt werden können, wird das Verhalten von Epithelien hier detailliert modelliert und daraus synthetisierte Impedanzspektren als Grundlage von Analysen mittels konstruktivistischen maschinellen Lernens verwendet. Diese Analysen erfolgen in einem ersten Schritt in Form eines selbstständigen Explorierens eines Teils der Impedanzspektren, welches in einer hierarchisch geordneten Menge von Modellen resultiert. Anschließend werden diese Modelle zur Adaption konkreter Anwendungen genutzt. Als Beispiel für eine Klassifikationsanwendung werden Modelle adaptiert, die eine verlässliche Zuordnung eines Impedanzspektrums zu der zugrundeliegenden Zelllinie erlauben. Als Beispiel für eine Regressionsanwendung werden Modelle adaptiert, die eine Quantifizierung der epithelialen Kapazität erlauben. In beiden Anwendungen identifiziert das konstruktivistische maschinelle Lernen selbstständig die Grenzen der Gültigkeit der von ihm aufgestellten Hypothesen und liefert dadurch eine differenzierte und für Menschen nachvollziehbare Interpretation der analysierten Daten.

Teil I.

**MOTIVATION UND
STAND DER TECHNIK**

*Dieser Schematismus unseres Verstandes,
in Ansehung der Erscheinungen und ihrer bloßen Form,
ist eine verborgene Kunst in den Tiefen der menschlichen Seele,
deren wahre Handgriffe wir der Natur schwerlich jemals abraten,
und sie unverdeckt vor Augen legen werden.*

— Immanuel Kant (1724-1805)



Abbildung 1.1.: Gemeinfreie fotografische Reproduktion des Gemäldes "Blüte und Verwesung". Je nach Perspektive sind darin entweder ein Liebespaar beim Picknick oder ein kieferloser menschlicher Schädel erkennbar. Aus geringer Entfernung (oben) ist in der Regel das Liebespaar leichter zu erkennen, aus größerer Entfernung (rechts) dagegen eher der menschliche Schädel. Ursprung und Künstler dieses Werkes aus dem 19. Jahrhundert gelten als unbekannt.





1

Einleitung

Maschinen und Computerprogramme werden zunehmend als intelligent bezeichnet. Nicht nur Massenmedien, auch Wissenschaftler sprechen heute in Publikationen und Projektanträgen regelmäßig von intelligenten Systemen. Gestützt wird dieses Narrativ vor allem durch publikumswirksame Ergebnisse in eng definierten, anwendungsbezogenen Aufgabenstellungen wie etwa dem Brettspiel Go [240]. Generellere Herausforderungen wie etwa den Turing-Test [253] hingegen konnte bislang noch kein künstliches System zweifelsfrei meistern [267]. Vielmehr gilt für viele lernende Systeme nach wie vor: Je komplexer die Aufgabe, desto unpräziser die Ergebnisse. Hinzu kommt, dass auch vergleichsweise einfache kognitive Fähigkeiten wie das Erkennen von Katzen in Bildern meist nur mit erheblichem technischen Aufwand realisierbar sind [149].

Begründet wird diese Divergenz zwischen dem weit verbreiteten Sprachgebrauch und den bisherigen wissenschaftlichen Befunden mit unterschiedlichsten Erklärungen. In vielen Studien werden ein Mangel an ausreichenden Daten oder eine zu hohe Komplexität der zu lösenden Aufgabe angeführt. Entsprechend wird häufig darauf verwiesen oder zumindest suggeriert, dass die jeweils zu lösende Aufgabe lösbar werden würde, wenn nur die Menge der Daten bzw. die Rechnerressourcen groß genug wären. Angesichts der Tatsache, dass dieses Ressourcenargument bereits seit mehreren Jahrzehnten angeführt wird (vgl. z.B. [154]), spricht wenig dafür, dass dies alleine ausreichend sein könnte. So ist zwar mittlerweile etwa die Architektur des visuellen Kortexes bereits in Form integrierter Schaltkreise imitierbar [169], doch komplexe kognitive Funktionen können damit nicht realisiert werden.

Ein ebenfalls häufig formulierter Erklärungsversuch ist der Verweis auf möglicherweise mehrdeutige Daten. Im Falle solcher Mehrdeutigkeiten, so heißt es dann, sei keine eindeutige, deterministisch bestimmbare Lösung möglich. Untermalt wird dies teils durch Vergleiche mit optischen Täuschungen. So genannte Kippbilder etwa führen zu einer Wahrnehmung, die Psychologen als multistabil bezeichnen [153]: Die aktuelle Wahrnehmung kann spontan durch eine andere ersetzt werden; Abb. 1.2 zeigt einige bekannte Beispiele. Zu Bedenken ist beim Vergleich von Datensätzen mit Kippbildern allerdings die Rolle des Menschen als aktiver Betrachter. Denn Psychologen bezeichnen hier ausdrücklich die Wahrnehmung des Betrachters als veränderlich – und nicht etwa den Stimulus bzw. das Kippbild [132].

Unbestreitbar ist, dass mehrdeutige Sinneswahrnehmungen eine Herausforderung sind, der nicht nur Maschinen, sondern auch Menschen regelmäßig ausgesetzt sind. Dies betrifft nicht

1. Einleitung

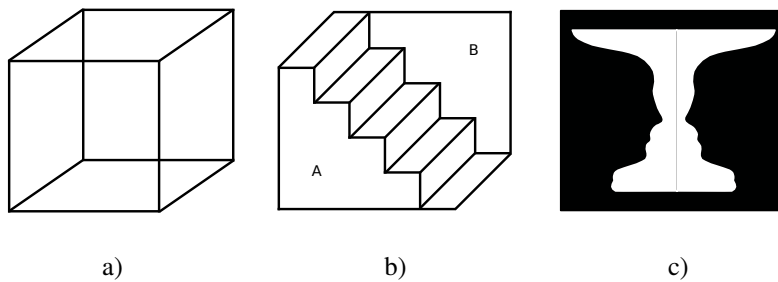


Abbildung 1.2.: Bekannte Kippbilder aus psychologischen Studien. a) Der Necker-Würfel, ursprünglich Louis Albert Necker (1786-1793) zugeschrieben, kann als dreidimensionale Darstellung sowohl aus der Perspektive “von unten nach oben” als auch “von oben nach unten” wahrgenommen werden (Abbildung nach [178]). b) Ähnlich wie der Necker-Würfel kann auch die von Heinrich Georg Friedrich Schröder (1810-1885) beschriebene Schröder-Treppe entweder als von links oben abwärts gerichtet (A) oder als auf dem Kopf stehend und links unten beginnend (B) wahrgenommen werden (Abbildung nach [234]). c) Die Rubinische Vase, ursprünglich Edgar J. Rubin (1886-1951) zugeschrieben, kann sowohl als weiße Vase auf schwarzem Grund als auch als zwei schwarze, gegenüberliegende Gesichtsprofile auf weißem Grund wahrgenommen werden (Abbildung nach [220]).

nur visuelle, sondern ebenso akustische Wahrnehmungen, etwa bei der Oktavbestimmung von Tönen oder der Lokalisation von Sprachquellen [23]. In der Kunst stellt Mehrdeutigkeit sogar ein elementares und seit Jahrhunderten verwendetes Gestaltungsmerkmal dar. Optische Illusionen etwa wurden in zahlreichen bekannten Werken verwendet [35], etwa dem Gemälde “Blüthe und Verwesung” aus dem 19. Jahrhundert (Abb. 1.1). Aber auch literarische Kippfiguren, die verbale und narrative Ambiguität aufweisen, sind der Fachliteratur bekannt [212].

Das Phänomen der Mehrdeutigkeit verdeutlicht, dass eine gegebene Menge an Information nicht nur von unterschiedlichen Menschen, sondern auch vom selben Menschen zu unterschiedlichen Zeitpunkten unterschiedlich interpretiert wird. Und es wirft Fragen auf: Kann der Mensch eine eindeutige, deterministische Lösung finden, wo es eine Maschine nicht kann? Lassen sich unter solchen Voraussetzungen überhaupt Maschinen mit menschenähnlichen Fähigkeiten schaffen? Und falls ja, auf welchem Weg? Wie kann eine computergestützte Informationsverarbeitung an höheren kognitiven Funktionen ausgerichtet werden? Und welche Ansätze sollten lernende Systeme dann verfolgen, um Mehrdeutigkeiten angemessen zu verarbeiten?

Einer der in der Informatik am breitesten akzeptierten Ansätze zur Entwicklung intelligenter Systeme ist die Imitation neurophysiologischer Strukturen. Seit in den 40er Jahren des vergangenen Jahrhunderts das erste Mal so genannte künstliche Neuronen vorgeschlagen wurden [167], sind diese in immer neuen Varianten als Grundbaustein für eine am Menschen orientierte Intelligenz propagiert worden. Kern dieses Paradigmas sind stark vereinfachte Annahmen über neurophysiologische Prozesse menschlicher Neuronen (Abb. 1.3). Auch neuere Versuche, diese Prozesse beispielsweise durch Nachahmung der so genannten Aktionspotentiale [31] zeitlich oder durch Nachahmung der Wirkung neuronaler Botenstoffe [59] molekular realistischer nachzubilden, basieren auf diesem Prinzip. Netze aus solchen modellierten Neuronen, so die Annahme, könnten dann ähnliche Funktionen realisieren wie Netze aus biologischen Neuronen.

Das zentrale Problem dieses neuroinspirierten Ansatzes ist, dass eine Orientierung am menschlichen Gehirn bislang nur auf der Ebene einzelner Zellen sowie auf der makroskopischen Ebene möglich ist. Denn nur auf diesen Ebenen ist das Vorbild in einem ausreichenden Maß verstanden. Fast alle dazwischen ablaufenden Vorgänge, insbesondere wie Zellverbände interagieren und

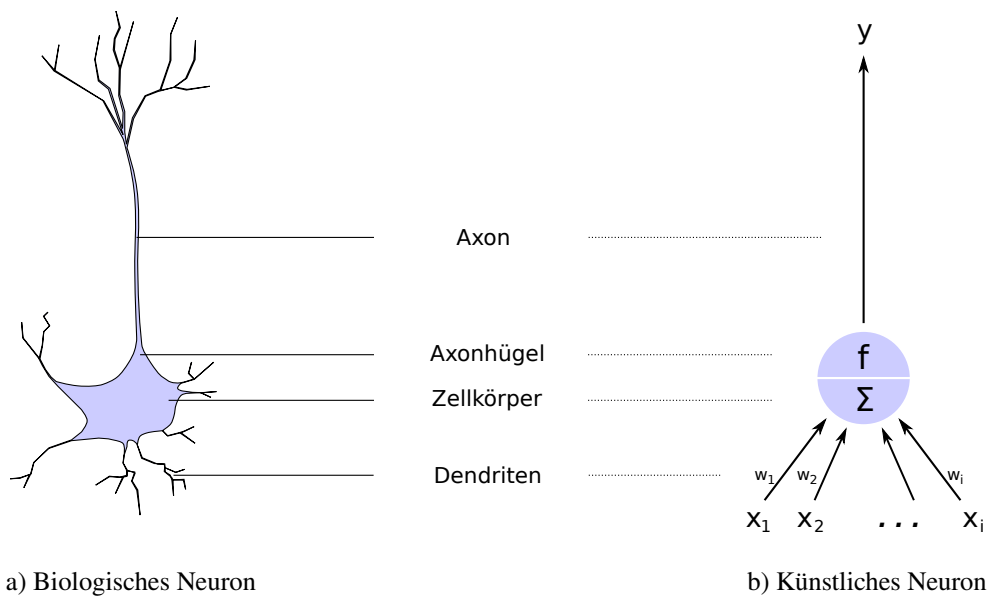


Abbildung 1.3.: Grundbaustein des neurophysiologischen Paradigmas. a) Schematische Darstellung eines prototypischen Motoneurons. Dendritische Ausläufer nehmen Signale anderer Neuronen auf, die über die Zellmembran zum Zellkörper weitergeleitet werden. Durch Überlagerung ergibt sich am Axonhügel ein gemeinsames Signal, das bei Überschreiten eines Schwellenwerts durch einen Axon genannten Ausläufer an andere Neuronen weitergeleitet wird. b) Schematische Darstellung der Funktion $y = f(\sum_{j=1}^i w_j x_j)$, die ein einfaches künstliches Neuron realisiert. x_j repräsentiert einen dendritischen Eingang, die Summation aller Eingänge die Signalintegration und die Aktivierungsfunktion f die Signalreaktion. f ist meist eine nichtlineare Funktion. (Abbildung nach Cheng und Titterington [40])

dadurch höhere kognitive Funktionen ermöglichen, sind dagegen noch weitgehend unverstanden. Gerade so zentrale menschliche Fähigkeiten wie das Abstrahieren von Eigenschaften oder Funktionen lassen sich daher auf diesem Weg bislang weder erklären noch auf Basis eines neurophysiologischen Ansatzes technisch realisieren. Zwar haben Informatiker immer wieder Lernverfahren vorgeschlagen, die nach deren Vorstellung Strukturen und Prozesse des Gehirns nachbilden; etwa im internationalen Großprojekt Human Brain Project [106], das von der Europäischen Union mit rund einer Milliarde Euro gefördert wird. Doch ein Nachweis für die Übereinstimmung mit Funktionsweise und Leistungsfähigkeit des Gehirns steht aktuell noch aus.

Ein anderes zentrales Hindernis auf dem Weg zu Maschinen, die sich mit menschlichen Leistungen messen können, stellt die Annahme eines naiven Realismus für Rechnersysteme dar. Ziel zahlreicher Informatikanwendungen ist das Extrahieren impliziter, zuvor unbekannter und potentiell nützlicher Information aus Daten. Dies wird häufig als *Knowledge Discovery* bezeichnet [68] und gilt als eines der wichtigsten Ziele moderner Informationsverarbeitung¹. Dieses angestrebte Ziel deckt sich mit dem, was man beim Menschen als Erkenntnis bezeichnet. Während jedoch in vielen Knowledge-Discovery-Projekten implizit angenommen wird, Systeme könnten die Welt so erkennen, wie sie ist, haben Erkenntnistheoretiker die Idee eines solchen naiven Realismus weitgehend verworfen. Der aus der moderneren Schule des Konstruktivismus stammende Gedanke, dass Realität erst durch den Betrachter bzw. das wahrnehmende System geschaffen wird, wird dagegen algorithmisch nicht berücksichtigt. Auch die grundlegende Frage, was auf diesem Weg überhaupt erkannt werden kann, wird in den meisten Arbeiten nicht gestellt.

¹Für den US-amerikanischen Mathematiker Richard W. Hamming ist die Gewinnung von Erkenntnis mutmaßlich sogar die einzige Daseinsberechtigung der Fachrichtung Informatik [151].

1. Einleitung

Ebenso unüblich ist die Definition von Kriterien, nach denen entschieden wird, ob von Maschinen nominell erfolgreich gelernte Zusammenhänge überhaupt relevant bzw. nutzbar sind. Denn solche Kriterien müssen insbesondere in gesellschaftlich relevanten Anwendungen nach menschlichen Maßstäben erklär- und überprüfbar sein. Während ein Mangel an Nachvollziehbarkeit in industriellen Anwendungen wie Maschinensteuerung oder Objekterkennung eher selten als problematisch bewertet wird, stellen mangelnde Interpretierbarkeit und ein computergestützter naiver Realismus für die wissenschaftliche Suche nach neuen, bislang unbekanntem Zusammenhänge in Datensätzen eine zentrale Herausforderung dar. Forscher versuchen daher in den vergangenen Jahren zunehmend, Interpretationsmöglichkeiten für die vorherrschenden überwachten Lernverfahren zu schaffen (z.B. [258, 100, 226]). Für die Entwicklung einer "Explainable Artificial Intelligence" stellte die US-amerikanische Forschungsagentur DARPA im Jahr 2016 millionenschwere Förderungen in Aussicht (Projektausschreibung DARPA-BAA-16-53).

Ein potentiell Bindeglied, um moderne erkenntnistheoretische Prinzipien mit psychologischen Befunden in Einklang zu bringen, stellt das Konzept so genannter mentaler Modelle dar. Psychologen verstehen darunter ganz allgemein die Repräsentation eines Gegenstandes oder eines Prozesses im Bewusstsein eines Menschen [48]. Dies stellt zunächst nur einen grundsätzlichen Erklärungsansatz dar. Dieser deckt sich allerdings – insbesondere hinsichtlich einer explizit beschränkten Gültigkeit von Modellen für einzelne Individuen – unmittelbar mit konstruktivistischen Annahmen. Nimmt man nun an, dass solche Modelle hierarchisch geordnet und kombinierbar sind, könnten sie in ihrer Gesamtheit die Grundlage menschlicher Kognition darstellen. Der Philosoph Herbert Stachowiak, der diesen Gedanken mit seiner Allgemeinen Modelltheorie (AMT) formalisiert hat, kommt sogar zu dem Schluss, dass jegliche menschliche Erkenntnis eine Erkenntnis in Modellen sein müsse [243]. Mit dieser Theorie, dem Lebenwerk Stachowiaks, existiert ein Konzept, mit dem sich nicht nur erkenntnistheoretische und psychologische Überlegungen integrieren lassen, sondern dessen Prinzipien sich auch zur Umsetzung mit informatischen Methoden eignen. Soweit bekannt, wurde dieser Ansatz jedoch bislang nicht verfolgt.

Sowohl Psychologen als auch Informatiker argumentieren häufig, die Implementierung und Evaluation einer Theorie menschlichen Denkens in Form eines Computerprogramms könne deren Richtigkeit bestätigen. Für das großangelegte Human Brain Project etwa wurde dies von dessen Initiator sogar als zentrale Existenzberechtigung formuliert [82]. Einen solch umfassenden Anspruch verfolgt die vorliegende Arbeit ausdrücklich nicht. Ziel ist es vielmehr, die angewandte Datenanalyse weiterzuentwickeln und sowohl leistungsfähiger als auch transparenter zu machen. Dazu werden zwar menschliche Denkstrukturen maschinell nachgebildet, anders als im klassischen neurophysiologischen Ansatz wird hier jedoch von Modellen und nicht von Neuronen als Elementarbestandteilen des Denkens ausgegangen. Zentrale Zielsetzungen des vorliegenden konstruktivistischen Ansatzes sind die Nachbildung von Intersubjektivität und Abstraktion mithilfe von Modellen. Dies geschieht auf Basis von Stachowiaks Allgemeiner Modelltheorie und unter Berücksichtigung relevanter erkenntnistheoretischer und psychologischer Paradigmen. Zur Implementierung werden etablierte Verfahren des maschinellen Lernens sowie in empirischen Disziplinen bewährte statistische Methoden kombiniert. Insgesamt wird damit ein neuartiges Framework für ein konstruktivistisches maschinelles Lernen entworfen.

Die Evaluation eines solchen grundlegend neuen Frameworks stellt eine besondere Herausforderung dar. Einerseits müssen dafür umfangreiche Datensätze zur Verfügung stehen, die nicht nur Mehrdeutigkeit aufweisen, sondern denen gleichzeitig belegbare kausale Zusammenhänge zugrunde liegen. Andererseits muss für diese Datensätze auch eine detaillierte Kontrolle und exakte Überprüfung von Vorhersagen des Verfahrens möglich sein. Messtechnisch gewonnene Daten können diese Forderung nach einem überprüfbaren wahren Zielwert allerdings grundsätzlich nicht erfüllen, da der wahre Wert einer Messung als rein ideeller Wert nicht exakt bekannt sein kann (vgl. Deutsche Industrienorm DIN 1319-1). Evaluationen müssen daher in diesem Fall zwingend

anhand synthetischer, komplexer und mehrdeutiger Datensätze erfolgen.

Ein gut verstandenes und mittels Berechnung gut nachbildbares Instrument empirischer Untersuchungen sind impedanzspektroskopische Messungen. Diese frequenzabhängigen Wechselstrommessungen lassen sich nicht nur für spezifische Schaltkreise und gerätespezifische Messfehler exakt modellieren und berechnen, sondern bei Annahme von Ersatzschaltkreisen für Proben aus Biologie und Materialwissenschaften auch für diese modellieren. Aufgrund der physikalischen Eigenschaften der Untersuchungsobjekte sowie der jeweiligen Fragestellung entstehen darüber hinaus häufig mehrdeutig interpretierbare Messergebnisse, deren Auflösung Forscher in der Praxis vor erhebliche Probleme stellen kann. Beispielhaft für diese Problematik ist die Analyse impedanzspektroskopischer Messungen an Epithelien. Screenings oder diagnostische Untersuchungen an diesen Körperzellen sind zwar technisch schnell und effizient durchführbar, in der Auswertung jedoch häufig mit teils schwer auflösbaren Mehrdeutigkeiten verbunden. Aufgrund der praktischen Relevanz solcher Messungen in der klinischen Physiologie ist von einer Auflösung dieser Mehrdeutigkeiten eine erhebliche Verbesserung der Diagnostik zu erwarten. Im Folgenden dienen daher synthetisierte impedanzspektroskopische Daten für gegebene Zelllinien als zentrales Untersuchungsobjekt zur Evaluation des Frameworks.

1. Einleitung



2

Menschliches Lernen

Die Frage, wie Menschen denken und lernen, beschäftigt mehr als nur eine wissenschaftliche Disziplin. Nicht nur die Psychologie, die seit einigen Jahrzehnten eine breite Schnittstelle zu den Naturwissenschaften aufweist, hat Antworten darauf entwickelt. Auch innerhalb von Philosophie und Erziehungswissenschaften haben sich im Laufe ihrer mehrere tausend Jahre währenden Tradition zahlreiche, teils konträre Denkschulen entwickelt. Selbst bei einer engen Beschränkung auf diese drei Disziplinen würde deren vollständige Abbildung den Rahmen der vorliegenden Arbeit bei weitem sprengen. Die folgende Darstellung beschränkt sich daher auf die jeweils wichtigsten Paradigmen und Konzepte menschlichen Lernens.

2.1. Konstruktivismus

Erkenntnis, das philosophische Äquivalent zum Informatikbegriff Knowledge Discovery, beginnt dort, wo bloße Kenntnis und Wiedergabe bereits bekannter und allgemein anerkannter Fakten enden. Doch was kann der Mensch überhaupt erkennen? Spätestens seit Platon und Aristoteles sind Erklärungsversuche zu dieser Frage überliefert. Die konkurrierenden Denkschulen lassen sich dabei im wesentlichen nach zwei grundsätzlichen Positionen unterscheiden: entweder es wird angenommen, dass der Mensch in seiner eigenen Wahrnehmung eine rein passiv-rezipierende Rolle ("Naiver Realismus") oder eine aktiv-gestaltende Rolle einnimmt.

Die erstere Position – also die Vorstellung, dass die Dinge im Wesentlichen so sind, wie sie dem Menschen erscheinen – wurde in der Philosophie im Verlauf des 20. Jahrhundert weitgehend verworfen. Vereinzelt unterscheiden Autoren diesen klassischen "naiven Realismus" noch von einem wissenschaftlichen Realismus, den sie zumindest bei reinen Naturbeobachtungen für zulässig halten (z.B. [245]). Spätestens seit den Arbeiten von Karl Popper gilt allerdings auch die als Positivismus bekannte Vorstellung als logisch widerlegt, dass sich Erkenntnisse beweisen oder verifizieren lassen [192]. Die daraus hervorgegangene Schule des kritischen Rationalismus postuliert folglich, dass jede Erkenntnis hypothetisch und vorläufig bleiben muss [193].

Noch konsequenter als vom kritischen Rationalismus wird eine absolute Erkennbarkeit der Realität vom so genannten Konstruktivismus abgelehnt. Unter diesem Begriff haben sich gleich mehrere Strömungen dem Gedanken verschrieben, dass Erkenntnis primär durch aktives Zutun des

2. Menschliches Lernen

Betrachters entsteht. Lernen ist demnach gleichbedeutend mit der Erschaffung einer individuellen Repräsentation der Welt und wird nicht durch eine absolute Wirklichkeit, sondern durch "unser Verhalten, Denken und Handeln" bestimmt [224, S. 220]. Wissen selbst besteht nach Richards und von Glasersfeld in der "Konstruktion und Aufrechterhaltung von Invarianzen" [208, S. 210]¹. Während jedoch ein radikaler Konstruktivismus eine Übereinstimmung von individueller Wahrnehmung und Realität grundsätzlich verneint und damit eine stets rein subjektive Wahrnehmung postuliert, geht der Sozialkonstruktivismus von einer grundsätzlich durch soziale Interaktion geschaffenen gemeinschaftlichen Wahrnehmung aus.

Die Ideen des Sozialkonstruktivismus wurden auch in den Erziehungswissenschaften aufgegriffen und erlangten so schnell praktische Bedeutung. Als Wegbereiter einer sozialkonstruktivistischen Didaktik gilt der russische Psychologe Vygotsky, der erstmals die Bedeutung der sozialen Interaktion für die Entwicklung von Kindern hervorhob [260]. Im angelsächsischen Raum sind konstruktivistische Lerntheorien heute in Pädagogik und Didaktik vorherrschend [66]. In Deutschland hat insbesondere die von Kersten Reich postulierte systemisch-konstruktivistische oder interaktionistisch-konstruktivistische Didaktik großen Einfluss auf die schulische und außerschulische Bildung. Reich postuliert drei grundsätzliche Lernprozesse: *Konstruktion* von neuem individuellem Wissen bzw. neuen Fähigkeiten, *Rekonstruktion* von Wissen bzw. Fähigkeiten anderer und *Dekonstruktion* eigenen Wissens bzw. eigener Fähigkeiten [203].

2.2. Paradigmen der Psychologie

Bei der ebenfalls am menschlichen Lernen interessierten Disziplin Psychologie handelt es sich im Gegensatz zur Philosophie um eine primär empirische Wissenschaft. Dies hat unter anderem zur Folge, dass sich nicht mit der grundsätzlichen Frage nach Wahrheit bzw. danach, was der Mensch überhaupt erkennen kann, beschäftigt wird. Vielmehr steht die Feststellung und – soweit möglich – Erklärung von Wahrnehmung und Verhalten des Individuums im Vordergrund.

Die heute dominierende Schule ist die Neuropsychologie bzw. Psychobiologie. Ihr Begründer Donald Hebb postulierte Mitte des 20. Jahrhunderts die Existenz eines direkten Zusammenhangs zwischen Verhalten und neuronalen Prozessen. In der Konsequenz müssten sich also Ergebnisse aus psychologischen Studien mit Ergebnissen aus neurophysiologischen Studien korrelieren lassen [94, S. xii]. Hebb's Ideen haben zur einer Ausdifferenzierung der modernen Psychologie in eine eigenständige interdisziplinäre Neurowissenschaft geführt (vgl. z.B. [17]). Parallel dazu haben diese Ideen in Form von neuronalen Netzen Eingang in maschinelle Lernverfahren gefunden.

Anders als etwa in der Physik existiert in der Psychologie jedoch traditionell kein zentrales Paradigma, das von allen konkurrierenden Denkschulen akzeptiert wird. Während etwa Voluntarismus, Stukturalismus und Funktionalismus das menschliche Bewusstsein in den Mittelpunkt und unbewusste Vorgänge ins Abseits stellten, kehrte Freuds Psychoanalyse dies um und stellte das Unbewusste in den Mittelpunkt. Der auf John B. Watson zurückgehende Behaviorismus wiederum postuliert, dass Bewusstsein sich grundsätzlich nicht aus einer Außenperspektive erforschen lässt und aus dieser auch nicht erforscht werden sollte. Behavioristische Untersuchungen beschränken sich daher ausschließlich auf beobacht- und messbares Verhalten. Bekannte Vertreter waren Ivan Pavlov und Burrhus Frederic Skinner, deren vielbeachtete Tierexperimente [188, 109] noch heute Bestandteil zahlreicher Lehrbücher sind.

Trotz dieser offensichtlichen Zersplitterung der Psychologie lassen sich mehrere weitverbreitete Paradigmen identifizieren. In ihrer Gesamtheit wird sie daher meist als "multiparadigmatische" oder teils auch als "präparadigmatische" Wissenschaft bezeichnet. Hergenhahn etwa unterscheidet

¹Wörtlich heißt es dort: "'Bestätigt sich eine anfängliche und vorläufige Interpretation in wechselnden Kontexten, so wird dieses bewährte Wissen zum festen Bestandteil der erlebten Wirklichkeit, es bildet sich ein Begriff heraus. Notwendig ist also immer die Wiederholbarkeit eines Erlebnisses, um Invarianten überhaupt bilden zu können."

im einfachsten Fall zwischen vier psychologischen Paradigmen [96]: einem funktionalistischen, einem assoziationalistischen, einem kognitiven und einem neurophysiologischen Paradigma. Das assoziationalistische Paradigma lässt sich dabei gegebenenfalls noch weiter in ein behavioristisches, ein psychoanalytisches und ein humanistisches Paradigma unterteilen [95].

2.3. Modelle und Modellbildung

Eines der Prinzipien funktionalistischer Psychologie ist die Annahme mentaler Modelle. Ursprünglich sollte so die motorische Steuerung des Menschen erklärt werden, etwa der Hand [257, 112]. Im weiteren Sinne werden mentale Modelle als "hypothetische Konstrukte" verstanden [265], die sich hierarchisch ordnen lassen [201] und einem Menschen in ihrer Gesamtheit zu Vorhersagen über die physische Umwelt dienen [180]. Es handelt sich jedoch nicht um statische Konstrukte, vielmehr unterliegen solche Modelle fortlaufender Veränderung, vergleichbar mit einem rekursiven Computerprogramm [180, S. 32f]. Als wichtigstes Gütekriterium gilt die Zweckmäßigkeit, als weniger wichtig gilt die Übereinstimmung mit physischen Gegebenheiten [173, S. 64ff].

Unabhängig davon halten auch Philosophen Modelle für ein wichtiges "Mittel der menschlichen Erkenntnisgewinnung" [123, S. 412], teils sogar für das einzige solche [243, S.56]. Während der Begriff des mentalen Modells in der Regel eher vage formuliert ist (z.B. [219]) existieren in der Logik und anderen Disziplinen eine Reihe expliziter und formaler Modellbegriffe. Obwohl diese sich teils stark unterscheiden, gilt es jedoch als mehrheitsfähig, dass es sich bei einem Modell um eine Abbildung handelt [276, S. 221ff]; dies impliziert insbesondere, dass Original und Abbild existieren und in einer bestimmten Relation zu einander stehen. Stellvertretend für viele weitere werden hier der (für die Informatik elementare) mathematisch-naturwissenschaftliche Modellbegriff, der kybernetische Modellbegriff sowie der Modellbegriff nach Stachowiak dargestellt.

Ausgehend von den Naturwissenschaften des 19. und 20. Jahrhunderts versteht man unter einer mathematischen Modellierung heute die Formulierung von Gesetzmäßigkeiten als mathematische Vorschriften. Diese sind, mit den Worten Heinrich Hertz', nicht als absolute Gesetzmäßigkeiten zu verstehen, sondern als "innere Scheinbilder oder Symbole" [97, S. 1f]. Hertz weist darauf hin, dass in der Regel für einen Zusammenhang mehrere Modellierungen möglich sind und benennt als Auswahlkriterien *Zulässigkeit*, *Richtigkeit* und *Zweckmäßigkeit* [97, S. 2f]. Die wichtigste Form eines mathematischen Modells stellt die mathematische Abbildung oder Funktion dar. Sie beschreibt eine Beziehung oder Relation zwischen zwei Mengen, genannt Urbild und Bild, die jedem Urbild-Element genau ein Bild-Element zuordnet (vgl. z.B. [88, S. 93]).

Ein Makel des mathematischen Modellbegriffs ist, dass dieser nur Bild und Urbild kennt und keinen unmittelbaren Bezug zur menschlichen Erkenntnisgewinnung hat. Im Gegensatz dazu ist in der Kybernetik jedes Modell abhängig von einem Modellsubjekt. Nach Klaus, auf dessen Definition sich hier stellvertretend bezogen wird, gilt für ein Modell M mit Analogien zu einem Modelloriginal O [123, S. 413]: Es ist ein Modell für ein Subjekt S , "sofern informationelle Beziehungen zwischen S und M dazu beitragen können, Verhaltensweisen von S gegenüber O zu beeinflussen". Ein Modell ist also eine Art Instrument, instrumentiert von einem Nutzer für einen bestimmten Zweck; insbesondere tritt dessen Realitätsbezug hinter seiner Zweckmäßigkeit zurück [173, S. 64ff]. Diese Perspektive ist zwar nicht nur in den Ingenieurwissenschaften (aus der sie stammt [216]) weit verbreitet, allerdings auch umstritten (vgl. z.B. [10]).

Nicht explizit berücksichtigt ist im kybernetischen Modell der Prozess der Modellbildung. Insbesondere wird nicht zwischen Modellbildner und Modellanwender unterschieden, was Aussagen über zeitliche und intersubjektive Gültigkeit eines solchen Modells unmöglich macht. Diese Lücke schließt Herbert Stachowiaks Allgemeine Modelltheorie. Nach dieser zeichnet sich ein Modell nicht nur dadurch aus, dass es ein verkürztes Abbild eines Modelloriginals für ein bestimmtes Modellsubjekt ist [243]. Vielmehr sei für jedes Modell eine dreifache pragmatische Relativierung

2. Menschliches Lernen

erforderlich: "Eine pragmatisch vollständige Bestimmung des Modellbegriffs hat nicht nur die Frage zu berücksichtigen, *wovon* etwas Modell ist, sondern auch, *für wen*, *wann* und *wozu* bezüglich seiner spezifischen Funktionen es Modell ist." [243, S. 133]. Zu beachten ist, dass diese an Grenzen orientierte Definition auch komplexe Modellierungen nicht ausschließt. Stachowiak geht insbesondere davon aus, dass jedes Modell selbst wieder Original für andere Modelle sein kann und definiert auf dieser Grundlage Operationen auf Modellen.

2.4. Kognitive Ebenen

Neurophysiologen konnten bislang kein Konzept menschlichen Erkennens und Lernens entwickeln, das sowohl einfache als auch höhere kognitive Funktionen aus einer gesamtheitlichen Perspektive erklärt. In Pädagogik und Didaktik hingegen ist eine solche Perspektive nicht nur Forschungsgegenstand, sondern auch Arbeitsgrundlage. Beide Disziplinen zielen seit jeher darauf ab, "Lehr- und Lernprozesse so zu organisieren, daß Wissen angeeignet, Fähig- und Fertigkeiten erworben werden, daß meß- und beurteilbare Leistungen entstehen" [203, S. 73].

Analog zur Entwicklung der psychologischen Schulen hat sich auch in den Erziehungswissenschaften in den vergangenen Jahrzehnten die Vorstellung verändert, was Lernen ausmacht. Mitte des 20. Jahrhunderts lag der Tätigkeit von Lehrenden ein eher behavioristisches Menschenbild zugrunde [92], d.h. Lernen anhand von Vorbild und Nachahmung stellte die vorherrschende Lehrform dar. Auf Basis der so genannten Gestalttheorie entwickelt sich in der ersten Hälfte des 20. Jahrhunderts die Idee eines Lernens durch Einsicht, die jedoch erst mit dem Aufkommen der kognitiven Psychologie in den 1960ern größere Verbreitung fand (auch als "kognitive Wende" bezeichnet). Vorläufiger Höhepunkt dieser Entwicklung ist eine Didaktik, die im konstruktivistischen Sinn auf Partizipation, Pluralität und Lernen im Austausch mit der Umwelt setzt [225].

Eine solche konstruktivistische Didaktik nimmt nicht nur Zeit- und Subjektgebundenheit für Konstrukte an², sondern auch eine grundsätzlich hierarchische Ordnung von Lernzielen. Grundlage dafür sind vor allem die Arbeiten von Benjamin S. Bloom und Robert M. Gagné. Gagné unterscheidet fünf wesentliche Kategorien [73]: verbale Informationen, intellektuelle Fähigkeiten, kognitive Strategien, Haltungen, motorische Fähigkeiten. Bloom hingegen unterscheidet zunächst zwischen einer kognitiven, einer affektiven und einer psychomotorischen Domäne, die jeweils weiter unterteilt werden. Besonders bekannt ist Blooms Taxonomie für kognitive Lernziele [27], die aus verschiedenen Fachrichtungen heraus kritisiert [242, 198, 71] und später überarbeitet wurde [9]. Aber auch für die affektive Domäne, also die Verarbeitung von Sinneseindrücken, wurde eine ausführliche Taxonomie formuliert [136]. Ebenso für die psychomotorische Domäne, also die Steuerung und Koordination von Muskeln [241].

In der praktischen Pädagogik und Didaktik bildet heute die überarbeitete Bloomsche Taxonomie eine wichtige Arbeitsgrundlage (Abb. 2.1 gibt einen schematischen Überblick). Dies spiegelt sich etwa in bildungspolitischen Vorgaben wie dem Europäischen Qualifizierungsrahmen bzw. dem Deutschen Qualifizierungsrahmen unmittelbar wieder. Kernpunkt der überarbeiteten Bloomschen Taxonomie ist die hierarchische Ordnung kognitiver Prozesse in sechs Stufen [9]. Die unterste Stufe entspricht einem reinen Erinnern bzw. Wiedererkennen. Auf den beiden nächsthöheren Stufen stehen zunächst die Fähigkeit, zu verstehen bzw. Bedeutung zu erkennen, und danach das Anwenden eines erlernten Verfahrens in einer bestimmten Situation. Als höhere kognitive Prozesse folgen dann die Analyse und Bewertung von Gegenständen und Prozessen. Die höchste Stufe ist mit einer Synthese, also dem Schaffen eines neuen Produkts, erreicht.

²So heißt es bei Reich: "In der Behauptung der Konstruktion von Wirklichkeit ist eingeschlossen, daß solche Konstruktionen jeweils zeitgebunden sind, von den spezifischen Beobachtern und deren Verständigungsgemeinschaft abhängen, daß sie keine ewigen Wahrheiten festschreiben können [...]" [203]



Abbildung 2.1.: Schematische Darstellung der überarbeiteten Bloomschen Taxonomie (in Anlehnung an [9]). Lernziele werden in einer ersten Dimension, hier der y-Achse, hinsichtlich der kognitiven Prozessebenen Merken (Ebene 1), Verstehen (2), Anwenden (3), Analysieren (4), Bewerten (5) und Erstellen (6) unterschieden. In einer zweiten Dimension, hier auf der x-Achse, werden Lernziele nach Art des zu erlernenden Wissens eingeteilt. Unterschieden wird zwischen Wissen um Fakten, Konzepte, Prozeduren sowie metakognitives Wissen.

In einer zweiten Dimension wird auch zwischen vier Wissensdomänen unterschieden, in denen kognitive Prozesse stattfinden [9]. Die für die vorliegende Arbeit bedeutsamsten Wissensdomänen sind konzeptuales und prozedurales Wissen. Unter konzeptuellem Wissen sind dabei etwa Klassifikationen und Kategorien, Prinzipien und Generalisierungen sowie Theorien, Modelle und Strukturen zu verstehen [135]. Unter prozeduralem Wissen dagegen subjektbezogene Fähigkeiten oder Algorithmen, subjektbezogene Techniken und Methoden sowie Kriterien zur Auswahl geeigneter Methoden [135].

2.5. Intersubjektivität

Moderne Wissenschaftstheorien basieren im Wesentlichen wie der Konstruktivismus auf der Annahme, dass Menschen die Welt individuell wahrnehmen. Nach Luhmann etwa wäre sonst gar keine Erkenntnis möglich [162, S. 52]. Erst durch Unterschiede ist das Abstrahieren rein subjektiver Sinnstrukturen und "ein gemeinsam intendiertes Handeln gegenüber einer als intersubjektiv verstandenen Außenwelt" möglich [140, S.64]. Im Umkehrschluss hat Charles S. Peirce, Begründer der Denkschule des Pragmatismus, über die Erfahrung des Einzelnen geschrieben: "Wenn er etwas sieht, was andere nicht sehen können, nennen wir es Halluzination." [189]

Ob eine subjektive Wahrnehmung von anderen bestätigt werden kann, ist in empirischen Wissenschaften wie der Psychologie oder auch den Sozialwissenschaften eine zentrale methodische Frage. Zwar unterliegen dort quantitative Untersuchungen wie etwa Beobachtungsstudien, standardisierte Befragungen oder Inhaltsanalysen üblicherweise einem eng definierten Studiendesign. Durchgeführt werden diese jedoch immer von Menschen, wodurch die Bewertung bestimmter Sachverhalte je nach individueller Einschätzung variieren kann (vgl. [101]). Dies gilt insbesondere

2. Menschliches Lernen

für die Zu- bzw. Einordnung von Untersuchungsobjekten in nominale Kategorien.

Um sicherzugehen, dass individuelle Bewertungsmaßstäbe das Ergebnis nicht systematisch verfälschen, werden solche Aufgaben meist von zwei oder mehr Personen parallel durchgeführt. Als Faustregel gilt: Je mehr Bewerter, desto zuverlässiger die Ergebnisse [139]. Anhand des Grades der Übereinstimmung – auch als Inter-Kodierer-Reliabilität (engl. Intercoder Reliability oder Interrater Reliability) bezeichnet [246] – wird überprüft, ob die einzelnen subjektiven Bewertungen intersubjektiv nachvollziehbar und die erzielten Ergebnisse damit für Außenstehende relevant sind [195]. Gleichzeitig gilt eine schlechte Inter-Kodierer-Reliabilität auch als deutliches Indiz für schlecht definierte nominale Kategorien oder schlechtes Kodierer-Training [131].

Statistisch gesehen ist bei solchen Vergleichen eine vollständige Übereinstimmung mehrerer Kodierer ebenso unwahrscheinlich wie eine vollständig unterschiedliche Bewertung. Um den Grad der Inter-Kodierer-Reliabilität zu messen, wurden insbesondere für nominal-skalierte Variablen zahlreiche Reliabilitätskoeffizienten vorgeschlagen [194]. Nur wenige davon werden jedoch tatsächlich regelmäßig angewendet [160]. Neben einer simplen prozentualen Übereinstimmung [101] werden häufig zufallsberücksichtigende Maßzahlen [20, 81, 236, 43, 63] verwendet. Aufgrund ihrer statistischen Eigenschaften werden diese Reliabilitätskoeffizienten jedoch als ungeeignet für die Verwendung in empirischen Studien kritisiert [139]. Die Koeffizienten nach [20] und [236] haben zusätzlich den Nachteil, auf den Vergleich nur zweier Kodierer beschränkt zu sein.

Zwar ist in solchen Studien vor allem die Inter-Kodierer-Reliabilität von nominal-skalierten Variablen von Bedeutung, doch auch die von Variablen metrischer Skala kann bewertet, da er auf Korrelation beruh werden. Viele für nominal-skalierte Variablen entwickelten Koeffizienten sind hier jedoch nicht anwendbar. Eine der wenigen Ausnahmen ist Cronbach's α [50], der zwar häufig verwendet, jedoch aufgrund seine korrelationsbasierten Metrik insbesondere für Inhaltsanalysen als ungünstig gilt [105, 139]. Als einziger zufallskorrigierender Reliabilitätskoeffizient, der sowohl für nominale als auch metrischen Skalen und gleichzeitig auf eine beliebige Anzahl von Kodierern angewendet werden kann, gilt daher Krippendorff's α [138]. Dieses wurde wiederholt als Standardmaßzahl für Intercoder-Reliabilität vorgeschlagen [139, 93].



3

Maschinelles Lernen

Das so genannte maschinelle Lernen (engl. machine learning) gilt als Schnittstellen-Disziplin zwischen künstlicher Intelligenz und kognitiver Psychologie [148]. Als zentrale Fragestellung gilt, wie sich Computersysteme schaffen lassen, die sich durch Erfahrung automatisch verbessern, und welche fundamentalen Gesetze Lernprozessen zugrunde liegen [170]. Die verwendeten und erforschten Algorithmen stützen sich vor allem auf statistische Verfahren für Mustererkennung [26] und Induktion [127].

3.1. Überwachung versus Selbstorganisation

Kennzeichen lernender Systeme ist, dass das Lernen anhand von Beispielen bzw. einer endlichen Menge an Datenvektoren erfolgt. Wie genau gelernt wird, kann jedoch stark variieren. In den vergangenen Jahrzehnten wurden zahlreiche sehr unterschiedliche Verfahren vorgeschlagen, die sich zwischen einem weitgehend fremdbestimmten überwachten Lernen (engl. supervised learning) und einem weitgehend selbstorganisierten unüberwachten Lernen (engl. unsupervised learning) bewegen. Abb. 3.1 ordnet einige wichtige Algorithmen zwischen diesen Polen ein.

Als überwachtes Lernen bezeichnet man Verfahren, die eine Abbildung auf einen vorgegebenen Zielparameter erlernen sollen. Je nachdem, ob deren Bildbereich metrisch oder kategorial definiert ist, spricht man entweder von Regression oder Klassifikation. Sofern eine Abschätzung der Wahrscheinlichkeitsverteilung für Eingabe- und Ausgabevariablen erzeugt wird, wie etwa bei Naive Bayes [155] oder Hidden Markov Models [57], spricht man von generativen Verfahren. In der anwendungsorientierten Informatik hingegen sind diskriminative Verfahren wie logistische Regression [190], Entscheidungsbäume [158] oder künstliche neuronale Netze [141] populärer. Neuronale Netze gelten dabei als besonders flexibel, da diese theoretisch ohne jedes Vorwissen jede erdenkliche mathematische Funktion erlernen können [104]. Ebenfalls weitverbreitet sind Supportvektor-Maschinen [49], deren Einsatz allerdings die Annahme und Festlegung einer Kernelfunktion voraussetzt. Verglichen mit menschlichem Lernen ähnelt überwachtes Lernen einer pädagogisch-didaktischen Lernsituation mit einem Lehrer oder Trainer.

Analog zum Lernen ohne Lehrer oder Trainer wurden so genannte unüberwachte Lernverfahren entwickelt, die keinen vorgegebenen Zielparameter benötigen. Bei den meisten unüberwachten

3. Maschinelles Lernen

Methoden handelt es sich um Klassifikationsverfahren, mit denen Cluster in einem gegebenen Datensatz identifiziert werden sollen. Zu den ältesten und weit verbreitesten Algorithmen zählt etwa das k-means-Clustering [113]. Neben weiteren künstlichen Verfahren wie etwa dem hierarchischen Clustering [175] haben sich in den vergangenen Jahren aber auch biologisch inspirierte Algorithmen wie Kohonens selbstorganisierende Karte (engl. self-organizing map, kurz: SOM) [130] oder Grossbergs adaptive Resonanztheorie (engl. adaptive resonance theory, kurz: ART) [34] etabliert. Eine generelle Schwierigkeit von Clusteringverfahren ist, dass deren Ergebnisse durch externe Verfahren überprüft werden müssen [133]. Auch für Regressionsaufgaben wurden einige unüberwachte Verfahren vorgeschlagen (z.B. [36, 134]), die vor allem zur Dimensionsreduktion [152] eingesetzt werden.

Nicht alle Verfahren lassen sich eindeutig als überwacht oder unüberwacht bezeichnen. So lässt sich etwa ein Multilayer-Perzeptron, also ein überwachtes Verfahren, dazu verwenden, einen gegebenen Datensatz auf sich selbst abzubilden [98]. Wird anschließend die Ausgabeschicht eines solchen Autoencoders entfernt, so verbleibt ein Netz, das den Ursprungsdatensatz entsprechend der Anzahl der versteckten Neuronen auf einen neuen Datensatz niedriger Dimension abbildet. Dieses Prinzip ist beispielsweise eine wichtige Grundlage des so genannten Deep Learning [19]. Ein weiteres Beispiel für Zwischenformen des maschinellen Lernens sind Algorithmen des so genannten Semi-supervised Learning, bei denen überwachtes und unüberwachtes Lernen kombiniert und so nur für einen Teil der verwendeten Daten ein vorgegebener Zielwert benötigt wird [238]. Dies ermöglicht nicht nur die Analyse unvollständiger Datensätze, sondern erzielt in manchen Fällen sogar bessere Ergebnisse als klassische überwachte Lernverfahren (z.B. [13, 79, 62]). Allerdings müssen bei Algorithmen des Semi-supervised Learning zwingend im Voraus Annahmen über Verteilungsdichten getroffen werden. Im Falle ungünstiger Annahmen können die Ergebnisse deutlich schlechter als bei einem überwachten Lernverfahren sein [272]. Auch gibt es bislang nur wenige Hinweise darauf, dass im menschlichen Gehirn überwachtes und unüberwachtes Lernen in ähnlich enger Verknüpfung stattfinden [273].

Eine Sonderrolle nimmt das so genannte Reinforcement Learning ein, das sich vor allem in den Bereichen Robotik [124] und Adaptive Kontrolle [156] als hilfreich erwiesen hat. Das Verfahren erhält dabei zwar eine Rückmeldung zu jeder Vorhersage, jedoch nicht den exakten

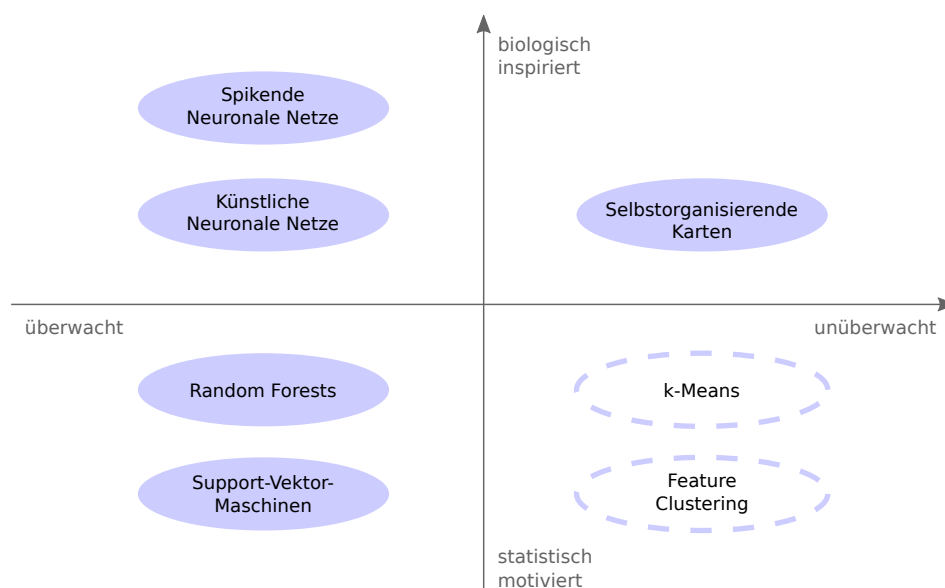


Abbildung 3.1.: Schematische Einordnung etablierter Lern- (Vollton) bzw. Cluster-Verfahren (gestrichelt).

Zielwert. Analog zum behavioristischen Lernparadigma wird stattdessen nur mitgeteilt, ob das Lernziel erreicht oder verfehlt wurde. Es konnte gezeigt werden, dass dieses Prinzip der Wirkung einzelner Neurotransmitter ähnelt [172, 67, 202]. Gleichzeitig nährten bereits in der Blütezeit der behavioristischen Psychologie manche seiner Vertreter Zweifel, dass allein mit einem solchen Mechanismus höhere kognitive Funktionen realisiert werden können [247, 248, 249]. In jüngerer Zeit wurde außerdem darauf hingewiesen, dass sich das Lernen mit Reinforcement-Verfahren im Vergleich zu typischen menschlichen Lernfortschritten als zu langsam darstellt [119].

3.2. Das neurophysiologische Paradigma

Obwohl zahlreiche rein statistische Verfahren zum maschinellen Lernen gerechnet werden, ist die Orientierung an neurophysiologischen Strukturen zweifellos ein zentrales Paradigma dieser Disziplin. Prototypisch dafür steht das Konzept künstlicher Neuronen, wie sie erstmals von McCulloch und Pitts vorgeschlagen wurden [167]. Für diese werden Integration und Weiterleitung elektrischer Signale, wie sie an Nervenzellen beobachtet wurden (siehe auch Abb. 3.2), durch eine Summations- und eine Aktivierungsfunktion abgebildet. Durch Verknüpfung des Ausgangs eines solchen künstlichen Neurons mit den Eingängen weiterer Neuronen können Netzwerke erzeugt werden, im einfachsten Fall sogenannte partiell veränderliche Perzeptronen [217]. Obwohl mehrere Variationen dieses Konzepts vorgeschlagen wurden (z.B. [103, 187]), werden in der Praxis am häufigsten mehrschichtige Feed-Forward-Netze in Kombination mit einem Backpropagation-Algorithmus verwendet (z.B. [99]).

Unberücksichtigt bleiben bei klassischen künstlichen neuronalen Netzen jedoch der Einfluss synaptischer Übertragungsmechanismen, die interzelluläre Kodierung von Signalen sowie zeitliche Effekte. So ist etwa bekannt, dass die Signale zwischen Nervenzellen frequenzmoduliert übertragen werden [24]. Diese Beobachtung wurde bei der Entwicklung so genannter spikender neuronaler Netze berücksichtigt, für die je nach zellulärem Vorbild individuelle Frequenz-Kodierungsmuster nachgebildet werden können [110]. Zwar konnte gezeigt werden, dass so die Funktionalität eines einzelnen Neurons bzw. die Komplexität der damit abgebbaren Funktion erhöht werden kann, allerdings steigt auch der erforderliche Rechenaufwand deutlich [163]. Weiter ist bekannt, dass Synapsen Signale nicht nur einfach weiterleiten, sondern modulieren, und dass diese Modulation veränderlich ist [51]. Um diese so genannte synaptische Plastizität – und damit komplexere Funktionen – nachzubilden, wurden in den vergangenen zwei Jahrzehnten mehrere konkurrierende

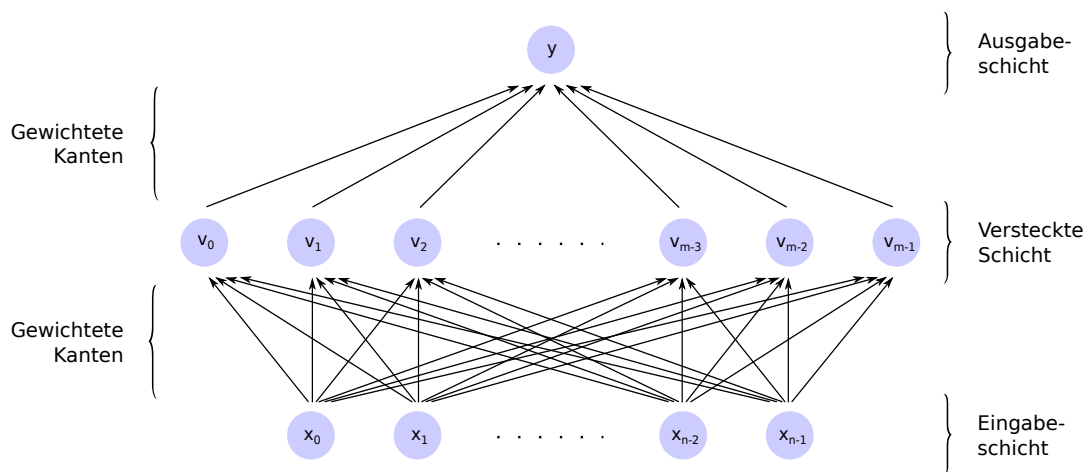


Abbildung 3.2.: Schematische Darstellung eines neuronalen Netzes (Abbildung in Anlehnung an [40]).

3. Maschinelles Lernen

Ansätze vorgeschlagen (z.B. [166, 150, 164]).

Während die meisten künstlichen neuronalen Netze als überwachte Lernverfahren konzipiert sind, gibt es nur wenige biologisch-inspirierte unüberwachte Lernverfahren. Das mit Abstand bekannteste ist die so genannte selbstorganisierende Karte, wie sie von Teuvo Kohonen vorgeschlagen wurde [128]. Deren Vorbild ist die Signalverarbeitung in der Gehirnrinde, dem so genannten Kortex, wo Signale eine räumlich verteilte Verarbeitung erfahren. Kohonen's Self-Organizing Feature Map, kurz SOM, bildet entsprechend ähnliche Eingabedaten auf ähnliche Neuronen ab. Insbesondere handelt es sich dabei um eine topologieerhaltende Abbildung. Obwohl ursprünglich kein unmittelbarer Bezug zu neurophysiologischen Mechanismen hergestellt wurde, lassen sich Analogien zwischen SOMs und der tatsächlichen Signalverarbeitung im Kortex herstellen [129].

Viele biologisch-inspirierte Verfahren, insbesondere die hier vorgestellten, ermöglichen nachweislich ein effizientes Lernen von Zusammenhängen. Gleichzeitig existieren zu diesen aber auch Alternativen, die nicht dem neurophysiologischen Paradigma folgen. In Klassifikationen liefern etwa Support-Vektor-Maschinen [46] oder Random Forests [30] ähnlich gute oder sogar besser Ergebnisse als neuronale Netze (z.B. [232, 233]); oft reichen auch bereits einfachere Verfahren wie Entscheidungsbäume oder der k-Nearest-Neighbor-Algorithmus [47] für gute Vorhersagen aus. Alle diese Verfahren lassen sich auch in Regressionsanwendungen einsetzen, in denen aber auch statistische Ansätze wie Multivariate Adaptive Regression Splines (MARS) ähnliche gute Ergebnisse liefern [52]. Als Alternative zu SOMs sind insbesondere k-Nearest-Neighbor-Clustering [161] oder hierarchisches Clustering [174] weitverbreitet. Insgesamt lässt sich daraus schließen, dass die Orientierung an neurophysiologischen Strukturen nicht zwingend erforderlich ist, um Modelle für Klassifikations- oder Regressionsaufgaben bilden zu können.

3.3. Multiperspektivisches Lernen

Während die ersten Verfahren maschinellen Lernens nur einen einzelnen lernenden Agenten vorsahen, wurde seit den 1970er Jahren vermehrt auch der parallele Einsatz mehrerer konkurrierender Agenten vorgeschlagen. Ziel ist es, mehrere unterschiedlich gute Ergebnisse zu einem verbesserten Gesamtergebnis zu kombinieren¹. Für solche so genannten Komitee- oder Ensemble-Verfahren werden häufig Entscheidungsbäume [145, 55] oder KNNs [90, 142] eingesetzt, aber auch Ensembles aus SVMs wurden bereits vorgeschlagen [121]. Abb. 3.3 illustriert diesen Ansatz schematisch für ein Ensemble aus neuronalen Netzen. Grundsätzlich ist es eher nachrangig, welches bzw. welche Lernverfahren eingesetzt werden; es ist sogar wahrscheinlich, dass die gleichzeitige Verwendung unterschiedlicher Verfahren das Ergebnis verbessert [144]. Entscheidender ist, wie sicher gestellt wird, dass sich die Lernergebnisse unterscheiden und auf welche Weise die Ergebnisse der Lernalgorithmen am Ende kombiniert werden [182].

Die naheliegendste Möglichkeit, bei gegebenem Lernverfahren und gegebenem Datensatz unterschiedliche Lernergebnisse zu erzeugen, ist eine Aufteilung des Datensatzes in k disjunkte Teildatensätze, die je einem Komitee-Mitglied zugeordnet werden. Prinzipiell entspricht dies einer k -fachen Kreuzvalidierung [76]. Allerdings kann ein solches Vorgehen dazu führen, dass sich die Teildatensätze – und entsprechend die Lernergebnisse – unerwünscht stark von einander unterscheiden [74]. Eine Alternative zu dieser Strategie stellt das so genannte bootstrap aggregating oder kurz bagging dar, das dem stochastischen Verfahren des Ziehens-mit-Zurücklegen entspricht und für jedes Komitee-Mitglied durchgeführt wird [29]. Dies erzeugt grundsätzlich eine Stichprobe von der Größe des Ausgangsdatsatzes, kann jedoch auch zur Steigerung der Effizienz auf eine Teilmenge begrenzt werden; dieses Vorgehen wird auch als subsample aggregating oder kurz subbagging bezeichnet.

¹Die mathematischen Wurzeln dieses Ansatzes reichen bis ins 18. Jahrhundert zurück. Eine gute historische Einordnung bietet [215].

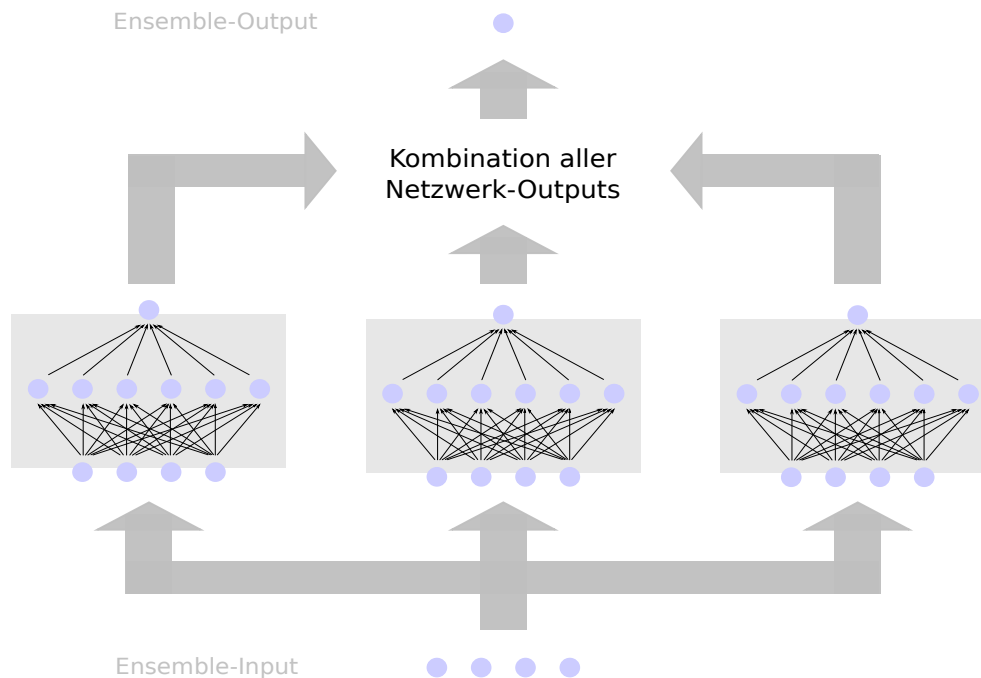


Abbildung 3.3.: Schematische Darstellung eines Ensembles aus drei neuronalen Netzen. Die Netze werden unabhängig voneinander mit dem Ensemble-Input trainiert und liefern im Idealfall divergierende individuelle Vorhersagen, die jeweils nach einer definierten Strategie zu einem einzigen Ensemble-Output kombiniert werden (Abbildung in Anlehnung an [159]).

Eine weitere Möglichkeit divergierende Ergebnisse zu erzielen, besteht darin, jedem Komitee-Mitglied nur eine Teilmenge der verfügbaren Features zuzuweisen. Bei manueller Zuweisung solcher Feature-Subsets können etwa mit neuronalen Netzen Verbesserungen erzielt werden [41]. Allerdings gibt es Hinweise darauf, dass dieses Prinzip nur bei einer großen Anzahl redundanter Features gut funktioniert [252]. Bekanntestes Beispiel für diese Strategie ist das Random-Forest-Verfahren nach Breiman [30]. Obwohl die Subsets dabei randomisiert erzeugt werden, konnte die hohe Treffsicherheit von Random Forests mittlerweile auch theoretisch begründet werden [25].

Auch um Ergebnisse mehrerer Lernalgorithmen zu kombinieren, wurden zahlreiche konkurrierende Ansätze vorgeschlagen. Die naheliegendste Möglichkeit besteht für Klassifikationsaufgaben in der einfachen Mehrheit [15]; für Regressionsaufgaben entsprechend im Mittelwert bzw. Median, auch Bragging genannt [33]. Bereits die Ergebnisse solcher einfacher Kombinationsverfahren sind in der Regel eindrucksvoll (vgl. z.B. [90]), was sich auch theoretisch begründen lässt [122]. Das populärste Verfahren hingegen, das so genannte Boosting, gewichtet jedes Komiteemitglied individuell [262]. Noch komplexer sind Verfahren, bei denen auch diese Gewichtung von Lernalgorithmen angepasst wird [111]. Einen Überblick über weitere Ansätze bietet [252].

Obwohl Ensemble- und Komitee-Verfahren explizit unterschiedliche Agenten vorsehen, können sie im Ergebnis keine intersubjektive Erkenntnis liefern. Ziel von Komitees ist ja gerade, für einen gegebenen Eingabevektor mehrere Ergebnisse zu erhalten, die nicht übereinstimmen. Somit ist wahrscheinlich, dass sich die durch die einzelnen Komitee-Mitglieder gebildeten Modelle deutlich unterscheiden. Während bei variierenden Datensätzen zumindest grundsätzlich denkbar ist, dass die zugehörigen Agenten weitgehend übereinstimmende Modelle erstellen, ist dies bei variierenden Features offensichtlich nicht der Fall. Auch die Kombination der Ergebnisse der Komitee-Mitglieder stellt keinen intersubjektiven Abgleich dar, da keine Evaluation und Bewertung anhand von Reliabilitätskoeffizienten erfolgt.

3.4. Auswahl und Extraktion von Features

Da empirische Datensätze häufig irrelevante oder redundante Variablen enthalten, spielen in der Praxis neben den eigentlichen maschinellen Lernverfahren auch Methoden zur Featureauswahl und -extraktion eine wichtige Rolle. Anwendungen der computergestützte Textverarbeitung, der Analyse von Geneexpressionsdaten oder aus der kombinatorischen Chemie etwa profitieren von einer solchen Vorverarbeitung erheblich [85], da diese in der Regel die Zahl der Eingangsdimensionen reduziert und so Rechen- und Speicherbedarf optimiert. Eng verwandt mit der automatisierten Auswahl von Features ist das Prinzip des Feature Ranking, bei dem eine automatisierte Bewertung jedes Features erfolgt. Zahlreiche Feature-Selection-Algorithmen nutzen intern ein Feature Ranking. Eingesetzt werden Feature Selection, Ranking und Extraction am häufigsten in überwachten Klassifikations- oder Regressionsaufgaben. Mit Modifikationen ist grundsätzlich aber auch ein Einsatz für unüberwachte Lernaufgaben möglich [3].

Eine grundsätzliche Frage beim Entwurf wie auch bei der Wahl eines zu verwendenden Feature-Selection-Verfahrens ist, ob die Eigenheiten des eigentlichen Lernalgorithmus berücksichtigt werden oder nicht. Im ersten Fall werden Teilmengen der verfügbaren Features mit dem eigentlichen Lernalgorithmus evaluiert. Solche auf überwachten Lernverfahren basierende Methoden werden in der Literatur als Wrapper bezeichnet und unterscheiden sich vor allem darin, ob sie deterministisch vorgehen (z.B. Sequential Forward Selection [264], Recursive Feature Elimination [86]) oder randomisiert (z.B. Simulated Annealing [1], genetische Algorithmen [227]). Im zweiten Fall werden zur Bewertung und Auswahl von Features statistische Maße herangezogen, etwa für Relevanz [18], Konsistenz [11] oder Korrelation [87]. Auch Abwandlungen des k-Nearest-Neighbors-Algorithmus wurden bereits zu diesem Zweck vorgeschlagen [177]. Solche vom verwendeten Lernalgorithmus unabhängigen Verfahren werden als Filter bezeichnet.

Wrapper liefern häufig Feature-Subsets, die nicht nur kleiner sind als mit Filtern bestimmte Subsets, sondern mit denen sich auch bessere Vorhersagen treffen lassen [126]. Allerdings müssen Filter im Gegensatz zu Wrappern nur einmal pro Datensatz ausgeführt werden, so dass diese insbesondere bei einer großen Anzahl von Features effizienter sind. Um die Vorteile beider Ansätze nutzen zu können, können diese kombiniert werden [64, 223]. So genannte Filter-Wrapper sortieren zunächst mittels eines Filters Features mit besonders geringem Nutzen aus; anschließend wird die verbleibende Feature-Teilmenge mit einer Wrapper-Methode evaluiert. Eine elegantere Lösung als die Anwendung eines expliziten Feature-Auswahlalgorithmus stellt es jedoch dar, Lernverfahren zu verwenden, die dies implizit als Teil des Lernverfahrens umsetzen. So ist etwa bei Entscheidungsbäumen nach erfolgreichem Training ablesbar, welche Feature-Teilmenge verwendet wurde. Ähnlich ist bei Random Forests ein Feature-Ranking ablesbar, und auch bei künstlichen neuronalen Netzen können etwa die Gewichte zwischen den Neuronen als Feature Ranking interpretiert werden [221].

Der rechnerische Aufwand und die Komplexität von Verfahren zur Extraktion neuer Features variiert sehr stark. Besteht der Eingabevektor beispielsweise aus einer Messreihe, können neue Features im einfachsten Fall durch Berechnung statistischer Kenngrößen bzw. Lagemaße erfolgen [232]. Eine weitere Möglichkeit zur Featureextraktion besteht darin, ähnliche Features zu clustern und die resultierenden Cluster anstelle der gruppierten Feature zu verwenden [259]. Seit langem bekannt und häufig eingesetzt ist auch die so genannte Hauptkomponentenanalyse (engl. principal component analysis, kurz PCA) [213] oder Faktoranalyse [116]. Neben diesen vergleichsweise intuitiven Verfahren wurden in den vergangenen beiden Jahrzehnten zahlreiche nichtlineare Verfahren vorgeschlagen, die jedoch in der Praxis kaum nennenswerte Verbesserungen brachten [256]. In jüngster Zeit sind vor allem neuronale Netze, die einen gegebenen Datensatz auf sich selbst abbilden, ein populäres und effizientes Werkzeug zur Dimensionsreduktion. Diese so genannten Autoencoder sind insbesondere im Kontext von Deep-Learning-Verfahren weit verbreitet [98].

3.5. Evaluation und Bewertungsmaßstäbe

Für die Bewertung des Lernerfolgs maschineller Lernverfahren existiert kein einheitliches Maß. Dies liegt nicht nur darin begründet, dass die verschiedenen beteiligten Disziplinen unabhängig voneinander entsprechende Kenngrößen entwickelt haben. Auch die unterschiedlichen Aufgabenstellungen machen die Definition eines einheitlichen Maßes unmöglich. So ist es beispielsweise aufgrund des unterschiedlichen Skalenniveaus des Zielparameters unvermeidbar, zwischen der Evaluierung von Klassifikations- und der Evaluierung von Regressionsaufgaben zu unterscheiden. Ebenso wenig können Bewertungsmaße für überwachtes Lernverfahren, die einen Zielwert voraussetzen, auf unüberwachtes Lernen ohne Zielwerte angewendet werden.

Zur Bewertung des Lernerfolgs überwachter Regressionsaufgaben stehen zahlreiche Kenngrößen zur Verfügung [239]. Im einfachsten Fall wird dazu die mittlere absolute Abweichung (engl. mean absolute error, MAE) oder analog zu Least-Squares-Methoden die Summe der Fehlerquadrate (engl. sum of squares error) genutzt. Gebräuchlich sind aber auch die Verwendung des Mittels der Fehlerquadrate (engl. mean squared error, MSE) bzw. die Wurzel daraus (engl. root mean square error, RMSE). Obwohl diese Maße in der Praxis weitverbreitet sind, haben sie den prinzipiellen Nachteil, dass sie für unterschiedliche Datensätze unterschiedliche Intervallgrenzen aufweisen. Zum Vergleich verschiedener durch Regressionsverfahren identifizierter Modellierungen eignen sich daher relative Maße besser, wie beispielsweise die mittlere absolute Abweichung im Verhältnis zum Zielwert.

Auch die Ergebnisse überwachter Klassifikationsaufgaben können unterschiedlich evaluiert werden. In klinischen oder psychologischen Kontexten ist vor allem das Kenngrößenpaar Spezifität und Sensitivität gebräuchlich bzw. abgeleitet davon Receiver-Operator-Kurven (engl. receiver operator curve, ROC) [206]. Ähnliche Aussagen sind aber auch mit Precision und Recall möglich [199]. Häufig soll die Evaluation anhand einer einzigen Kenngröße erfolgen. Dafür kann etwa die so genannte Accuracy verwendet werden, die die Summe der Hauptdiagonalen im Verhältnis zur Gesamtsumme darstellt. Diese leitet sich aus Spezifität und Sensitivität ab und kann auch aus einer so genannten Konfusionsmatrix bestimmt werden. Alternativ ist auch die Nutzung der durch eine ROC beschriebene Fläche (engl. area under curve, AUC) [28] möglich.

Kommen unüberwachte Lernformen zum Einsatz, so ist naturgemäß nur definiert, welche Eigenschaften das Ergebnis aufweisen soll. Ein Lernziel im Sinne eines eindeutigen Ideals, an dem sich das tatsächliche Lernergebnis messen ließe, existiert hingegen nicht. In der Regel bedeutet dies, dass zum Beispiel ein auf diesem Weg erhaltenes Clustering durch externe Verfahren verifiziert werden müssen [133]. Dies kann beispielweise erfolgen indem überprüft wird, ob die erhaltene Cluster-Zuordnung mithilfe eines überwachten Lernverfahrens zuverlässig erlernbar und somit unterscheidbar ist [232]. In realen Anwendungen kann jedoch auch dies stets nur ein erster Schritt der Verifikation sein. Für klinische Anwendungen etwa wird empfohlen, Ergebnisse zusätzlich anhand von Fachwissen und praktischen Erfahrungen abzugleichen [4].

3. Maschinelles Lernen



4

Epithel-Analyse mittels Impedanzspektroskopie

4.1. Wechselstrom und Widerstand

Eine Impedanz ist eine vor allem in der Elektrotechnik gebräuchliche physikalische Messgröße, die das Prinzip eines elektrischen Widerstands über Gleichstromkreise und den klassischen ohmschen Widerstand hinaus verallgemeinert. In einem Wechselstromkreis bzw. bei Anwendung von sinusförmiger Wechselspannung und sinusförmigem Wechselstrom gibt sie sowohl das Verhältnis dieser beiden Größen wie auch die Verschiebung der Phasenwinkel an. Der Wert einer Impedanz Z wird in Ohm angegeben.

Das Verhältnis von Wechselspannung zu Wechselstrom wird als Scheinwiderstand oder Magnitude r bzw. $|Z|$ bezeichnet. Die Phasenverschiebung wird mit dem Winkel ϕ angegeben, der Werte zwischen -90 und $+90$ Grad annehmen kann. Alternativ kann eine Impedanz auch komplexwertig angegeben werden, wobei der Realteil einem ohmschen Widerstand und der Imaginärteil einem induktiven (positiver Wert) bzw. kapazitiven (negativer Wert) Widerstand entspricht. Als alternative Repräsentation kommt gelegentlich auch die so genannte Admittanz zur Anwendung, die den komplexen Kehrwert einer Impedanz darstellt [186].

Der Wert einer Impedanz ist abhängig von der Frequenz der sinusförmigen Spannung und des sinusförmigen Stroms. Die Untersuchung der Änderung dieser Größe bei variierender Frequenz wird Impedanzspektroskopie genannt und erlaubt eine höhere Auflösung des elektrischen Verhaltens eines Untersuchungsobjekts als etwa der ohmsche Widerstand [58]. Typischerweise werden für eine solche Messung 50 bis 100 Frequenzen zwischen 1 Hz und 1 MHz verwendet [84], für manche Anwendungen werden allerdings auch Einzelfrequenzmessungen durchgeführt [146].

Eine hilfreiche Visualisierung impedanzspektroskopischer Messungen ist der so genannte Bode-Plot¹, bei dem Magnitude und Phasenverschiebung getrennt gegen den Logarithmus der Frequenz aufgetragen werden (siehe auch Abb. 4.1a). Dieser gilt als sehr aussagekräftig [58], ist jedoch nicht für jede Art Anwendung optimal. Eine auf der komplexwertigen Repräsentation basierende Alternative dazu stellt das Nyquist-Diagramm bzw. der Cole-Cole-Plot dar², bei

¹Benannt nach dem US-amerikanischen Elektrotechniker Hendrik Wade Bode (1905-1982), der diese Darstellung in den 1930er in den Bell Laboratories einführte.

²Das Nyquist-Diagramm wurde von dem Physiker Harry Nyquist (1889-1976), der Cole-Cole-Plot von dem Biophysiker Kenneth S. Cole (1900-1984) und dessen Bruder propagiert [45]. Beide Darstellungsformen gelten als äquivalent.

4. Epithel-Analyse mittels Impedanzspektroskopie

dem Real- und Imaginärteil einer Impedanz gegeneinander aufgetragen werden (Abb. 4.1b). Da dies deren Frequenzabhängigkeit nur implizit abbildet, wurde später auch eine dreidimensionale Variante vorgeschlagen [44], die aber in der Praxis seltener verwendet wird.

Wie bei jeder Messung kann auch bei einer Impedanzmessung nie der wahre Wert der zu bestimmenden Impedanz ermittelt werden³. Dies erklärt sich insbesondere durch den geräte-spezifischen Mess-Bias, den grundsätzlich jedes Messgerät aufweist. Um dem wahren Wert einer zu bestimmenden Impedanz näher zu kommen, wird daher empfohlen, den messgeräte-spezifischen Bias zu bestimmen und zur Bereinigung von Impedanz-Messergebnissen zu verwenden [185]. Eine weitere Quelle für Verzerrung ist der stochastische Fehler solcher Messungen, der mit steigender Frequenz zunimmt [185].

4.2. Impedanzspektroskopie in der Praxis

Das systematische Messen von Impedanzen wird in einer Reihe wissenschaftlicher Disziplinen als analytisches Hilfsmittel genutzt. In Physik und Materialwissenschaft sind die Charakterisierung von Festkörpern [200] bzw. deren Leitfähigkeitsmechanismen [38] typische Anwendungen, in der Elektrochemie die Charakterisierung von Batterien [222] und Brennstoffzellen [268] und in der Elektrotechnik die Charakterisierung von elektrischen Bauteilen wie Kondensatoren oder Spulen (z.B. [196]). Daneben wird die Impedanzspektroskopie auch in zahlreichen biomedizinischen Anwendungen eingesetzt wie etwa zur Körperfettbestimmung [147], zur Klassifizierung von Brustkrebs-Patienten [120] oder zur Untersuchung physiologischer und pathophysiologischer Charakteristika von Epithelien [84].

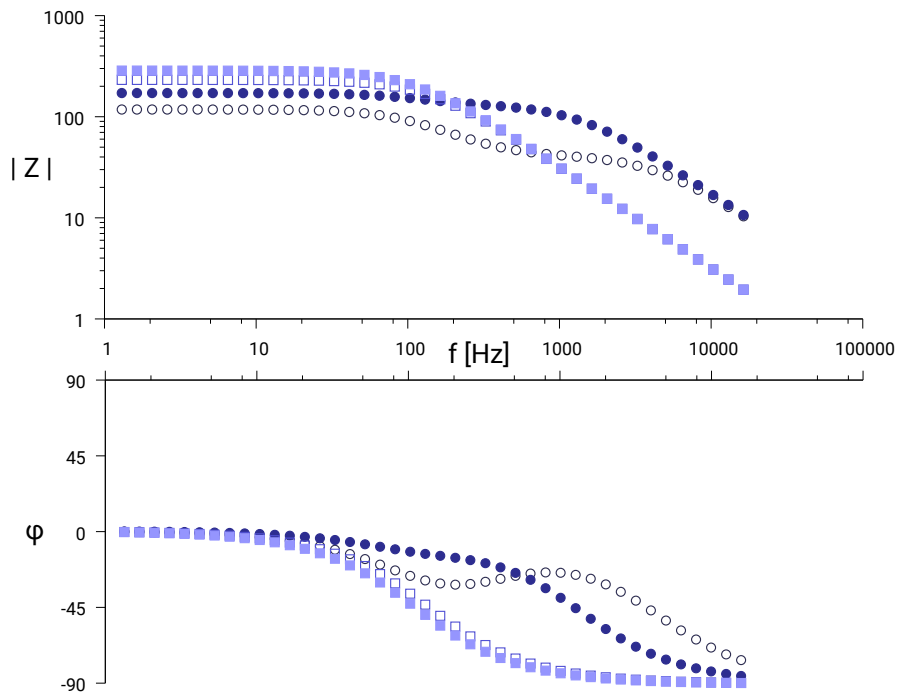
Um Rückschlüsse von Impedanzspektren auf einen Untersuchungsgegenstand ziehen zu können, werden für diesen typischerweise Ersatzschaltbilder angenommen. Bei homogenen Materialien lässt sich dies damit begründen, dass deren Widerstand umgekehrt proportional zu Länge und Fläche ist [165]. Bei inhomogenen Materialien wie Zellen und Gewebe ist dies zwar nicht der Fall, allerdings lässt sich empirisch eine Beziehung zwischen der Impedanz und dem Volumen des enthaltenen elektrolythaltigen Wassers beobachten [146]. Die Zusammensetzung angenommener Ersatzschaltbilder ist stark anwendungsabhängig und kann von einzelnen Bauteilen, über einfache RC-Glieder bis hin zu komplexen Schaltkreisen reichen [183].

Impedanzspektroskopie bietet mehrere Vorteile gegenüber anderen analytischen Methoden. Einerseits handelt es sich um eine vergleichsweise kostengünstige Technik, für die lediglich eine Impedanz-Schnittstelle und ein Messgerät zur Bestimmung des Frequenzgangs erforderlich ist [84]. Für biomedizinische Anwendungen ist darüber hinaus günstig, dass es sich um eine nicht-invasive Technik handelt [147]. Da somit Zellen und Zellproben unbeschädigt bleiben, sind Messwiederholungen möglich und das Risiko von Fehldiagnosen wird reduziert.

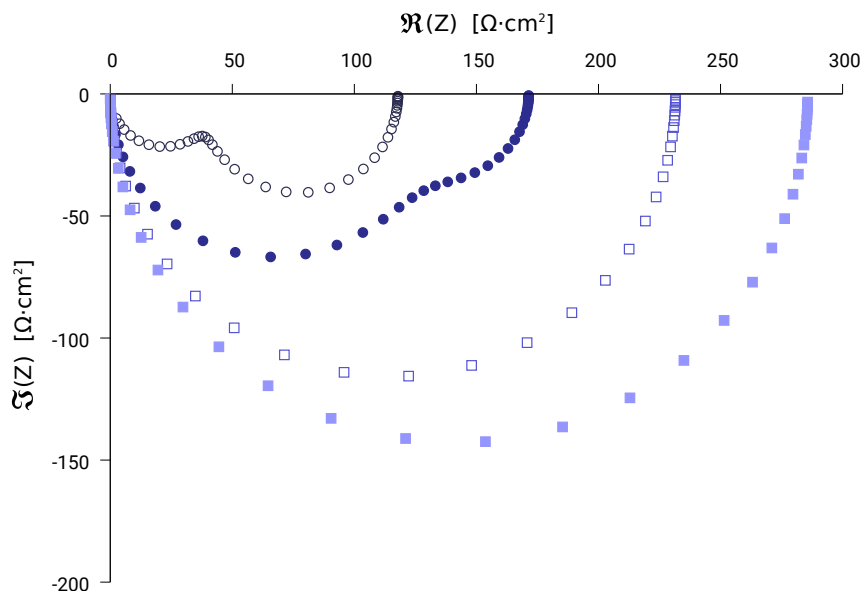
Eine Herausforderung jeder impedanzspektroskopischen Anwendung stellt die Ableitung eines Ersatzschaltbilds für das Untersuchungsobjekt dar. Einerseits ist häufig unklar, welche resistiven, kapazitiven oder induktiven Anteile repräsentiert sind. Andererseits kann eine gegebene Menge an Widerständen, Kondensatoren und Spulen auf mehrere verschiedene Weisen miteinander verknüpft oder auch in mehrere Teilelemente aufgeteilt sein. In der Konsequenz führt dies zu mehrdeutigen Messergebnissen, da diese von unterschiedlichen Ersatzschaltkreisen stammen können [165]; dies gilt insbesondere für symmetrische Schaltkreise, falls deren Größen symmetrisch vertauscht werden. Insgesamt bedeutet dies, dass impedanzspektroskopische Messungen grundsätzlich nur als Hilfsmessungen betrachtet und parallel ergänzende Analysen durchgeführt werden sollten, um eine gegebene Fragestellung beantworten zu können [184, S. xxiii].

³In der Messtechnik wird der wahre Wert einer Messung als rein ideeller Wert betrachtet, der grundsätzlich nicht exakt bekannt sein kann. Dieses Prinzip ist unter anderem in der Deutschen Industrienorm (DIN) 1319-1 festgeschrieben.

4.2. Impedanzspektroskopie in der Praxis



a) Bode-Plot



b) Nyquist-Diagramm

Abbildung 4.1.: Darstellungsformen für Impedanzspektren. Abgebildet sind vier überlappende modellierte Spektren mit je 42 Impedanzen Z . Die beiden RC-Glieder des Modellschaltkreises C weisen hier entweder sehr ähnliche (■, □) oder sehr unterschiedliche (●, ○) Zeitkonstanten auf. a) Darstellung als Bode-Plot, wobei Betrag $|Z|$ und Phasenwinkel ϕ getrennt gegen die Messfrequenz aufgetragen werden. b) Darstellung als Nyquist-Diagramm bzw. Cole-Cole-Plot, in dem Realteil \Re und Imaginärteil \Im gegeneinander aufgetragen werden.

4.3. Struktur und Funktion von Epithelien

Epithelien sind einer von vier Grundgewebetypen des menschlichen Körpers. Im Verbund bilden sie ein- oder mehrschichtige Epithelien, die das Körperinnere vom Körperäußeren trennen. Epithelien sind in jedem vielzelligen tierischen Organismus in großer Anzahl und vielen Variationen zu finden [2], etwa in Niere [274], Haut [69] oder Drüsengewebe [80, 60]. Als mucus-bildende Zellschichten findet man sie auch im Darm [118] oder in der Nase [91]. Neben einer Barrierefunktion haben Epithelschichten vor allem die Aufgabe, für einen geregelten Stoffaustausch zu sorgen [83]. Unterhalb eines Epithels schließt sich unmittelbar die deutlich dünnere Basalmembran und darunter subepitheliales Gewebe an, die gemeinsam eine weitere Barriere darstellen [254].

Obwohl sich Epithelien unterschiedlicher Gewebe teils stark unterscheiden, weisen sie einen grundsätzlich polaren Aufbau auf. Dies bedeutet, dass sich Zusammensetzung und Eigenschaften der apikalen und basolateralen Zellmembranen deutlich unterscheiden [89]. Insbesondere weisen beide Seiten eine unterschiedliche Häufigkeit spezifischer Membranproteinen auf, die Transportmechanismen realisieren. Die Aufrechterhaltung dieser molekularen Asymmetrie wird durch Zell-Zell-Verbindungen der so genannten Zonula occludens bedingt. Diese auch Tight Junctions genannten Strukturen bestehen aus Membranproteinen, die die gesamte Zelle umschließen und mit benachbarten Zellen eine enge Verbindung herstellen. Sie bewirkt unter anderem, dass Moleküle der einen Zellseite nicht frei auf die andere diffundieren können [266].

Eine noch wichtigere Aufgabe der Tight Junctions ist die Regulation des parazellulären Molekültransports, der ganz wesentlich die Durchlässigkeit des Epithels bestimmt. Denn während kleine anorganische Ionen die Zellzwischenräume in der Regel schnell passieren, ist sie für größere Moleküle nur sehr selektiv passierbar [214]. Insbesondere werden nicht alle Moleküle durchgelassen, und auch diejenigen, die durchgelassen werden, werden dies nicht zu jeder Zeit. Tight Junctions stellen damit einen mindestens ebenso selektiven und komplexen Transportmechanismus dar wie die Kanäle und Transporter der apikalen und basolateralen Membran [108].

Die Moleküldurchlässigkeit ist ein wesentliches Funktionsmerkmal eines Epithels. Sie wird in der Regel durch Quantifizierung des Austauschs von geladenen Teilchen bestimmt, also durch Bestimmung der elektrischen Leitfähigkeit bzw. des Widerstands. Ist die transzelluläre Leitfähigkeit größer als die parazelluläre, so spricht man von einem dichten Epithel (engl. tight epithelium) [42], im umgekehrten Fall von einem lecken Epithel (engl. leaky epithelium). Ausschlaggebend für die Dichte eines Epithels ist vor allem die Zusammensetzung der Tight Junctions, die im Darm und in der Niere unter anderem von der Position innerhalb des Organs abhängt.

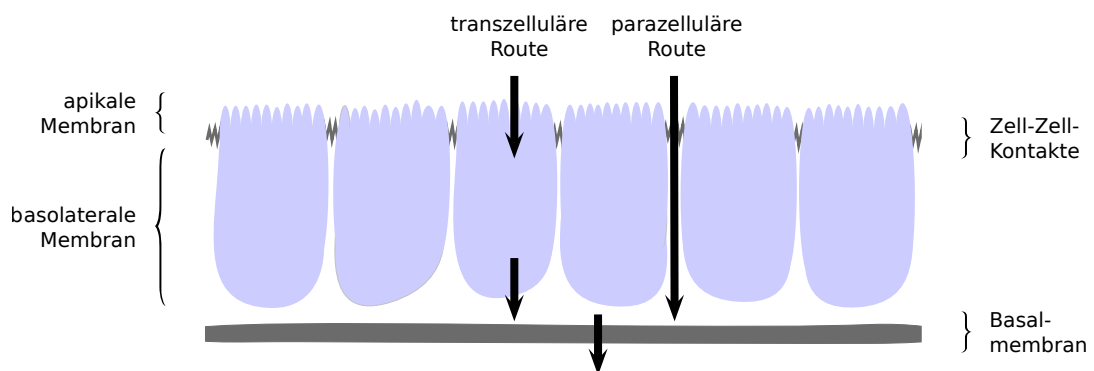


Abbildung 4.2.: Schematische Darstellung eines einschichtigen Epithels. Zell-Zell-Kontakte trennen die dem Äußeren (apikale) und die dem Inneren des Körpers (basolaterale) zugewandte Seite. Unterhalb des Epithels schließt sich die so genannte Basalmembran an.

Fehlfunktionen von Epithelien verursachen Störungen beim Austauschen von Stoffen zwischen Körper und Umwelt und damit teils schwere Erkrankungen. Ein bekanntes Beispiel einer solchen Fehlfunktion ist die genetisch bedingte Veränderung des apikalen Chlorid-Ionenkanaltyps CFTR, die zu Mukoviszidose bzw. zystischer Fibrose führt [114]. Eine besonders häufige Ursache für Fehlfunktionen von Epithelen sind bakterielle Infektionen, wie etwa im Fall einer Harnwegs-entzündung [8] oder einer durch bakterielle Toxine verursachten sekretorischen Diarrhöe [237]. Solche Infekte führen typischerweise zu einer erhöhten Durchlässigkeit des betroffenen Epithels für Wasser und andere darin gelöste Teilchen [244].

Neben Erkrankungen von Lunge oder Harnröhre stellen insbesondere Durchfallerkrankungen eine Herausforderung für den Menschen dar. Zwar sind diese in Industrieländern in der Regel gut behandelbar, weltweit gehören Durchfallerkrankungen jedoch zu den fünf häufigsten Todesursachen. Der Weltgesundheitsbehörde WHO zufolge sind Durchfallerkrankungen bei Kindern unter fünf Jahren sogar die zweithäufigste Todesursache [263]. Neben bakteriellen Infekten lösen auch chronisch entzündliche Darmerkrankungen wie Colitis ulcerosa oder Morbus Crohn schwere Diarrhöen aus [16]. Bei diesen Krankheiten führt eine über proinflammatorische Cytokine vermittelte erhöhte Durchlässigkeit zu Wasser- und Elektrolytverlust.

Die Erforschung epithelialer Fehlfunktionen wird anstelle von unmittelbar entnommenen Gewebeproben überwiegend an Zellkulturen durchgeführt. Dazu werden Epithelien in einem Nährmedium, also *in vitro*, kultiviert. Dies hat insbesondere den Vorteil, dass für Studien eine größere Menge an Untersuchungsobjekten zur Verfügung steht. Da dadurch unterschiedliche Studien an der gleichen Zelllinie durchgeführt werden können, erhöht dies außerdem die Vergleichbarkeit von Ergebnissen und erlaubt die Berücksichtigung früherer Ergebnisse. Als Modell für menschliche Darmepithelien dienen häufig die Zelllinien HT-29 und IPEC (z.B. [230]). Ein bekanntes Modell für menschliche Nierenepithelien ist die Zelllinie MDCK [168].

4.4. Epitheliale Impedanzanalyse

Dank der Erforschung ihrer elektrischen Eigenschaften hat das Verständnis von Epithelien und ihrer Funktion im vergangenen Jahrhundert große Fortschritte gemacht. Noch in den 1950er Jahren wurde etwa postuliert, es handle sich bei einem Epithel um eine geschlossene Zellschicht mit zwei Membranen [125]. Wenige Jahre später konnte gezeigt werden, dass diese Vorstellung nur als Vereinfachung und nur für bestimmte Epithel-Typen zulässig ist [53, 54]. Ein wichtiges Hilfsmittel stellen dabei impedanzspektroskopische Messungen dar, da diese nicht nur eine Bestimmung des Gesamtwiderstands, sondern auch detailliertere Aussagen über das Epithel erlauben. Um möglichst originalgetreue Bedingungen gewährleisten zu können, werden Epithelien dazu in einer so genannten Ussing-Kammer fixiert [157] (siehe auch Abb. 4.3). Dadurch ist es möglich, die Messungen in Gegenwart physiologischer Pufferlösungen durchzuführen.

Entscheidend für eine erfolgreiche Impedanzanalyse ist eine zweckmäßige Modellierung der Physiologie der untersuchten Zellen und deren Funktion. Für unbehandelte Zellkulturen und viele Gewebe erlaubt eine Methodik, die als Ein-Wege-Impedanzspektroskopie bezeichnet wird, zwischen der elektrischen Leitfähigkeit von Epithel und subepitheliale Gewebe zu unterscheiden [77]. Beide Parameter werden aus zu extrapolierenden Endpunkten impedanzspektroskopischer Messkurven abgeleitet. Untersuchungen von Patientenbiopsien mit dieser Methode haben entscheidend zum Verständnis der Veränderungen der Darmbarriere bei Morbus Crohn beigetragen [270]. Sofern die zu untersuchenden Messkurven keine Halbkreis-Form aufweisen, können diese Endpunkte allerdings nicht mit ausreichender Genauigkeit abgeschätzt werden [230].

Eine Weiterentwicklung stellt die sogenannte Zwei-Wege-Impedanzspektroskopie dar [143], deren Einsatz sich etwa in pharmakologischen Studien an Antidiarrhoika [235, 5] und so genannten Uptake-Enhancern [218] bewährt hat. Dieses relativ neue Verfahren erlaubt den Anteil von

4. Epithel-Analyse mittels Impedanzspektroskopie

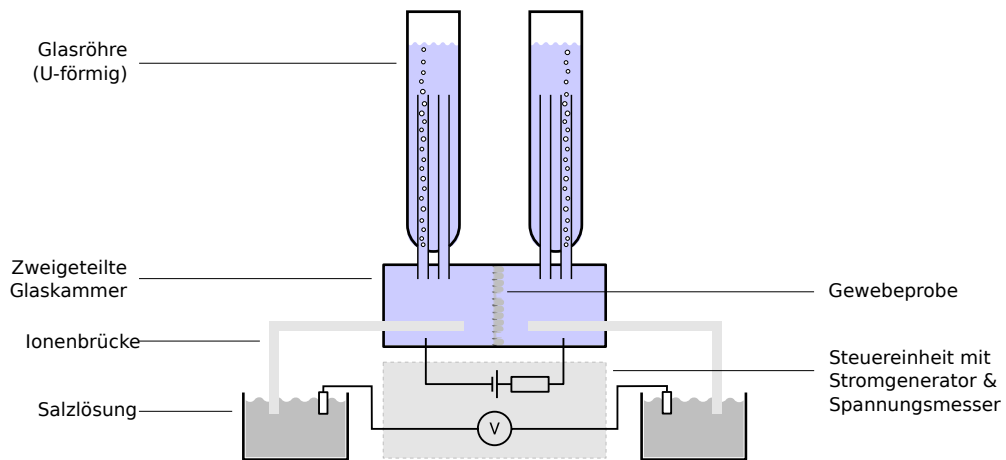


Abbildung 4.3.: Schematischer Aufbau einer Ussing-Kammer (Abbildung nach [157]).

transzellulärer und parazellulärer Leitfähigkeit an der Leitfähigkeit des Epithels zu bestimmen; dazu wird die parazelluläre Leitfähigkeit experimentell verändert und mithilfe von Konzentrationsmessungen eines Markermoleküls zusätzlich zur Wechselstromanalyse weitere Informationen gewonnen. Die Genauigkeit dieser abgeleiteten Parameter hängt jedoch stark von der Genauigkeit ab, mit der Werte für epitheliale und subepitheliale Leitfähigkeit bestimmt werden [230].

Da eine Unterscheidung von apikalen und basolateralen Membraneigenschaften durch eine gleichzeitige apparative Messung äußerst aufwendig ist, wurde dafür bereits in den 1990er Jahren ein computergestützter Ansatz vorgeschlagen [22]; eine Modifizierung impedanzspektroskopischer Messverfahren und die Anwendung nicht-linearer Fehlerquadrat-Minimierungen ermöglichte eine grobe Quantifizierung dieser Membraneigenschaften, allerdings bei Schätzfehlern von bis zu ± 20 Prozent. In neueren Ansätzen konnte gezeigt werden, dass die Komponenten eines Epithel-äquivalenten elektrischen Schaltkreises mit Algorithmen des maschinellen Lernens verlässlich aus impedanzspektroskopischen Messungen abschätzbar sind [229, 231, 230, 232, 233]. Für die Zellkulturlinien HT-29/B6 und IPEC-J2 konnte etwa detailliert gezeigt werden, dass epithelialer und subepithelialer Widerstand bei geeigneter Modellierung mit künstlichen neuronalen Netzen besser aus realen Messkurven abschätzbar sind als mit klassischen Verfahren [230].

Offen ist allerdings, wie die Anwendung maschineller Lernverfahren über einzelne Zelllinien und -kulturen hinaus verallgemeinert werden kann. Vorstellbar ist etwa, dass bestimmte funktionale Zustände unterschiedlicher Zelllinien Gemeinsamkeiten aufweisen und so gemeinsam betrachtet werden können. Ungeklärt ist außerdem, wie Lernverfahren mit mehrdeutigen Impedanzdaten sinnvoll und effizient umgehen können. Darüber hinaus stellen die bisher eingeführten Ansätze maschinellen Lernens in der epithelialen Impedanzanalyse prinzipiell nur Machbarkeitsstudien dar, deren Ergebnisse sich durch eine weitere Ausdifferenzierung der Verfahren optimieren lassen.

Teil II.

METHODIK UND ERGEBNISSE

*Wer A sagt, muss nicht B sagen.
Er kann auch erkennen, dass A falsch war.*

— Bertolt Brecht (1898-1956)



5

Konstruktivistisches maschinelles Lernen

Obwohl das maschinelle Lernen heute einen etablierten Forschungszweig darstellt, können existierende Verfahren in vielen Anwendungen nicht die gewünschte oder erforderliche kognitive Komplexität abbilden. Um menschlichen kognitiven Fähigkeiten näher zu kommen, wird hier die Orientierung an konstruktivistischen Lerntheorien sowie Stachowiaks Allgemeiner Modelltheorie vorgeschlagen. Insgesamt ergibt sich so ein neues Framework zur effektiveren Anwendung von Verfahren des maschinellen Lernens.

5.1. Grundannahmen

Auf Basis der im ersten Teil der Arbeit eingeführten Konzepte wird Lernen im Folgenden im Sinne dreier grundlegender Annahmen interpretiert:

1. **Die kleinste Sinneinheit einer kognitiven Funktion ist ein Modell.** Unter einem Modell wird hier ein pragmatisches Modell im Sinne Stachowiaks verstanden (siehe Kapitel 2.3). Dies schließt Abbildungen im mathematischen Sinn sowie deren Repräsentation oder Approximation in Form von maschinellen Lernverfahren ein. Die Zweckmäßigkeit eines Modells wird anhand der Genauigkeit bewertet, mit der es einen Zielparameter abbildet.
2. **Lernen erfolgt durch Konstruktion, Rekonstruktion oder Dekonstruktion von Modellen.** Die Konstruktion eines Modells entspricht im klassischen maschinellen Lernen weitgehend einem unüberwachten Lernen, die Rekonstruktion eines Modells einem überwachten Lernen. Ein Äquivalent zum Lernen durch Dekonstruktion ist für maschinelle Lernverfahren bislang weder konkretisiert noch implementiert worden.
3. **Abstraktion erfolgt durch Konstruktion von Modellen, die ausschließlich auf zuvor konstruierten Modellen basieren.** Zwischen Trainingsdaten und den Eingabefeatures eines abstrahierten maschinellen Modells muss folglich mindestens ein weiteres maschinelles Modell existieren. Ein abstrahiertes Modell basiert somit zwar nicht direkt auf den zu erlernenden Trainingsdaten, hängt jedoch indirekt davon ab.

5.2. Pragmatisch definierte maschinelle Modelle

5.2.1. Begriffsklärung

Im Unterschied zu Stachowiaks Modellbegriff beschränkt sich die vorliegende Arbeit auf die Betrachtung vektorbasierter einerseits und algorithmisch erzeugter Modelle andererseits. Dabei wird unterschieden zwischen Modellen mit und ohne explizit definierten pragmatischen Eigenschaften.

Vektorielle Modelle. In überwachten Lernverfahren setzt sich ein Lernvektor aus einem m -dimensionalen Inputvektor $I = (i_0, \dots, i_{m-1})$ und einem n -dimensionalen Outputvektor $O = (o_0, \dots, o_{n-1})$ zusammen, wobei für O oft $n = 1$ gilt. Da dadurch nicht nur Urbild- und Bildraum definiert sind, sondern implizit auch ein Abbildungszusammenhang zwischen I und O unterstellt wird, wird ein Lernvektor im Folgenden als (*vollständiges*) *vektorielles Modell* \mathcal{V} bezeichnet:

$$\mathcal{V} = (I, O) \quad (5.1)$$

$$= (i_0, \dots, i_{m-1}, o_0, \dots, o_{n-1}) \quad (5.2)$$

Ist einem gegebenen I kein O zugeordnet, etwa beim Trainingsvektor eines unüberwachten Lernverfahrens, wird I als *unvollständiges vektorielles Modell* bezeichnet.

Maschinelle Modelle. Wird mit einem maschinellen Lernverfahren eine mathematische Abbildung approximiert, so wird das Ergebnis dieser Approximation auf Basis einer Menge von j vektorieller Modelle als *maschinelles Modell* \mathcal{M} bezeichnet:

$$\mathcal{M} \sim \{\mathcal{V}_0, \dots, \mathcal{V}_{j-1}\} \quad (5.3)$$

Bei überwachten Lernverfahren erfolgt dies anhand einer endlichen Menge vollständiger, bei unüberwachten Verfahren anhand einer endlichen Menge unvollständiger vektorieller Modelle.

Pragmatische Eigenschaften. Das *Zeitintervall* T , innerhalb dessen ein Modell gültig ist, ist für ein vektorielles Modell \mathcal{V} durch den Zeitpunkt der Erhebung eindeutig definiert. Da hier eventuelle Fehlertoleranzen dieser Erhebung als vernachlässigbar angenommen werden, wird im folgenden insbesondere angenommen, dass gilt:

$$T_{\mathcal{V}} = \tau_{min} = \tau_{max} \quad (5.4)$$

Im Gegensatz dazu kann die zeitliche Gültigkeit eines maschinellen Modells \mathcal{M} grundsätzlich nur als hypothetisch angenommen und mittels hypothetischer Intervallgrenzen definiert werden:

$$T_{\mathcal{M}} = [\tau_{min}, \tau_{max}] \quad (5.5)$$

Der Urheber oder Betrachter eines Modells heißt im Folgenden *Modellsubjekt* σ . In naturwissenschaftlichen Erhebungen wird dies häufig ein Sensor oder ein Messgerät sein, bei Befragungen, Beobachtungsstudien oder Inhaltsanalysen typischerweise ein menschlicher Bewerter. Die Menge aller Modellsubjekte σ_i , für die ein gegebenes Modell \mathcal{M} gültig ist, heißt $\Sigma_{\mathcal{M}}$ und wird definiert als Teilmenge der (unendlichen) Menge Σ aller Subjekte:

$$\Sigma_{\mathcal{M}} \subset \Sigma \quad (5.6)$$

Ein Zielparameter eines Modells \mathcal{M} heißt im Folgenden *Modellzweck* ζ . Die Menge aller Modellzwecke ζ_i von \mathcal{M} heißt $Z_{\mathcal{M}}$ und wird definiert als Teilmenge der (unendlichen) Menge Z aller Modellzwecke:

$$Z_{\mathcal{M}} \subset Z \quad (5.7)$$

Ein vollständiges vektorielles Modell \mathcal{V} heißt *pragmatisch definiertes vektorielles Modell* \mathcal{V}^* , falls dessen pragmatische Eigenschaften explizit definiert sind. Dann folgt:

$$\mathcal{V}^* = (\mathcal{V}, T_{\mathcal{V}}, \Sigma_{\mathcal{V}}, Z_{\mathcal{V}}) \quad (5.8)$$

5.2. Pragmatisch definierte maschinelle Modelle

Ein maschinelles Modell \mathcal{M} , für das die pragmatischen Eigenschaften $T_{\mathcal{M}}, \Sigma_{\mathcal{M}}, Z_{\mathcal{M}}$ definiert sind, heißt *pragmatisch definiertes maschinelles Modell* \mathcal{M}^* . Dann folgt:

$$\mathcal{M}^* = (\mathcal{M}, T_{\mathcal{M}}, \Sigma_{\mathcal{M}}, Z_{\mathcal{M}}) \quad (5.9)$$

Ein pragmatisch definiertes maschinelles Modell \mathcal{M}^* mit $|\Sigma_{\mathcal{M}^*}| > 1$, dessen Subjektmenge also mehr als ein Element enthält, heißt *intersubjektives maschinelles Modell*.

5.2.2. Verwandtschaft, Vererbung und Hierarchien

Durch Analyse eines gegebenen Datensatzes, also einer endlichen Menge an vektoriellen Modellen, mittels konstruktivistischen maschinellen Lernens soll eine endliche Menge intersubjektiver maschineller Modelle identifiziert werden. Weiter soll auf dieser Menge eine Ordnung definiert sein, die die Modelle anhand von Verwandtschaftsbeziehungen hierarchisch ordnet.

Verwandtschaft. Zwei Modelle \mathcal{M}_1 und \mathcal{M}_2 (bzw. \mathcal{V}_1 und \mathcal{V}_2) heißen pragmatisch verwandt, wenn eine Übereinstimmung von mindestens zwei pragmatischen Eigenschaften gegeben ist. Es lassen sich folglich vier Verwandtschaftsgrade unterscheiden:

1. vollständig mit $T_{\mathcal{M}_1} = T_{\mathcal{M}_2}, \Sigma_{\mathcal{M}_1} = \Sigma_{\mathcal{M}_2}, Z_{\mathcal{M}_1} = Z_{\mathcal{M}_2}$.
2. subjektiv-intentional (ΣZ) mit $T_{\mathcal{M}_1} \neq T_{\mathcal{M}_2}, \Sigma_{\mathcal{M}_1} = \Sigma_{\mathcal{M}_2}, Z_{\mathcal{M}_1} = Z_{\mathcal{M}_2}$;
3. temporal-intentional (TZ) mit $T_{\mathcal{M}_1} = T_{\mathcal{M}_2}, \Sigma_{\mathcal{M}_1} \neq \Sigma_{\mathcal{M}_2}, Z_{\mathcal{M}_1} = Z_{\mathcal{M}_2}$;
4. temporal-subjektiv ($T\Sigma$) mit $T_{\mathcal{M}_1} = T_{\mathcal{M}_2}, \Sigma_{\mathcal{M}_1} = \Sigma_{\mathcal{M}_2}, Z_{\mathcal{M}_1} \neq Z_{\mathcal{M}_2}$;

Während ΣT -, TZ - und $T\Sigma$ -Verwandtschaften die Ableitung neuer Modelle erlauben, ermöglichen es vollständig verwandte existierende Modelle zu bestätigen, zu verwerfen oder in Modelle mit geringerer temporaler Gültigkeit zu differenzieren.

Vererbung pragmatischer Eigenschaften. Für das Erlernen eines maschinellen Modells wird eine Menge vollständiger oder unvollständiger pragmatisch definierter vektorieller Modelle vorausgesetzt, die eine ΣT -, TZ - oder $T\Sigma$ -Verwandtschaft aufweisen. Die pragmatischen Eigenschaften des erlernten maschinellen Modells \mathcal{M}^* leiten sich dann aus denen der zugrunde liegenden vektoriellen Modelle $\mathcal{V}_0^*, \dots, \mathcal{V}_{n-1}^*$ ab:

$$\Sigma_{\mathcal{M}^*} = \bigcup_{i=0}^{n-1} \Sigma_{\mathcal{V}_i^*} \quad (5.10)$$

$$Z_{\mathcal{M}^*} = \bigcup_{i=0}^{n-1} Z_{\mathcal{V}_i^*} \quad (5.11)$$

$$T_{\mathcal{M}^*} = \left[\min(T_{\mathcal{V}_0^*}, \dots, T_{\mathcal{V}_{n-1}^*}), \max(T_{\mathcal{V}_0^*}, \dots, T_{\mathcal{V}_{n-1}^*}) \right] \quad (5.12)$$

Analog erfolgt die Vererbung für ein maschinelles Modell \mathcal{M}^* , das aus n Outputs anderer maschineller Modelle $\mathcal{M}_0, \dots, \mathcal{M}_{n-1}$ erlernt bzw. abstrahiert wird. Zusätzlich wird dann für \mathcal{M}^* eine höhere Wissens Ebene als für $\mathcal{M}_0, \dots, \mathcal{M}_{n-1}$ festgelegt.

Hierarchien. Analog zur Bloomschen Taxonomie wird hier eine zweidimensionale Ordnung maschineller Modelle genutzt. Für die erste Dimension werden die vier Bloomschen *Wissensdomänen* übernommen, d.h. es wird zwischen faktuellen, konzeptuellen, prozeduralen und metakognitiven Modellen unterschieden. In einer zweiten Dimension werden Ebenen kognitiver Prozesse bzw. *Wissensebenen* unterschieden, wobei angenommen wird, dass höhere Prozesse auf niedrigeren fußen. Ebene 0 repräsentiert vektorielle Modelle. Daraus erlernte maschinelle Modelle werden Ebene 1 zugeordnet, abstrahierte Modelle jeweils der nächsthöheren Ebene. Die Zahl aufeinander aufbauender Wissens Ebenen kann begrenzt werden (nach Bloom z.B. $n=7$).

5.2.3. Implementierungsaspekte

Bei etablierten maschinellen Lernverfahren ist das Speichern von Trainingsergebnissen oder strukturellen Eigenschaften vorgesehen. Eine Möglichkeit, pragmatische Eigenschaften eines maschinellen Modells \mathcal{M}^* bzw. vektoriellen Modells \mathcal{V}^* zu speichern und beim Lernen zu berücksichtigen, bieten diese jedoch nicht. Daher wird hier folgende Kodierung eingeführt:

- **Zeitintervall.** Für die Kodierung von $T_{\mathcal{M}}$ wird ein Timestamp verwendet. Der Unix-Timestamp etwa gibt die Anzahl der Sekunden seit Donnerstag, dem 1. Januar 1970 00:00 Uhr UTC an und kann auch für Zeitangaben, die auf Nicht-Unix-Geräten erhoben wurden, berechnet werden. Sofern diese zeitliche Auflösung ungenügend erscheint, lässt sich daraus mit minimalen Modifikationen ein Timestamp in Milli- oder Nanosekunden erzeugen. Obwohl hier nicht benötigt, gilt dies grundsätzlich auch für Zeiten vor dem 1. Januar 1970.
- **Subjektmenge.** Da kein Standard existiert, um ein Modellsubjekt $\sigma \in \Sigma_{\mathcal{M}}$ eindeutig zu kodieren, werden hier Abkürzungen als Unique Identifier (UID) verwendet. Für Modelle, die durch ein künstliches neuronales Netz erlernt wurden, wird die UID ANN verwendet, die UID RF analog für Random Forests. Synthetisierten Messdaten erhalten die UID 0.
- **Zweckmenge.** In der Zweckmenge $Z_{\mathcal{M}}$ werden vorherzusagende Zielvariablen mittels einer dreiteiligen UID aus Wissensdomäne, Wissensebene und Modellierungszweck definiert und gespeichert. So steht etwa C.1.k02 für die Unterscheidung zweier mittels k-Means bestimmter Cluster (k02) auf Ebene 1 der konzeptuellen Wissensdomäne (C). Vektorielle Modelle ohne gegebenen Zielwert erhalten die UID 0.0.0. Jedes neue pragmatisch definierte maschinelle Modell \mathcal{M} erhält automatisch eine neue UID.

Da der Modellzweck zwar als UID kodiert ist, aber offensichtlich mehrere Modelle für den gleichen Zweck existieren können, ist zusätzlich für jedes pragmatisch definierte maschinelle Modell ein eindeutiges Identifikationsmerkmal erforderlich. Im Rahmen dieser Arbeit wird dies durch eine dreiteilige Modell-ID realisiert, die sich aus Wissensdomäne, Wissensebene und einer fortlaufenden Nummerierung mittels Integer-Werten zusammensetzt. Die UID C.2.1 etwa identifiziert das erste Modell auf der zweiten Ebene der konzeptuellen Wissensdomäne.

Während pragmatische Eigenschaften von vektoriellen Modellen als Metadaten bzw. zusätzliche Features unmittelbar in Datensätze einfügbar sind, werden die Metadaten pragmatisch definierter maschineller Modelle in einer eigenen, unabhängigen Struktur verwaltet. Dies vermeidet nicht nur Abhängigkeit von den Implementierungen externer Algorithmen, sondern ist auch organisatorische Voraussetzung für die im nächsten Abschnitt beschriebenen Lernprozesse. Aus praktischen Erwägungen wird neben den drei pragmatischen Attributen $T_{\mathcal{M}}$, $\Sigma_{\mathcal{M}}$ und $Z_{\mathcal{M}}$ außerdem auch jeweils der Urbildraum, also die Eingabefeatures, als Metadatum kodiert (Tab. 5.1).

Attribut	Variablentyp	Beispielwerte
UID	zusammengesetzt	C.1.1
Urbildraum	Liste	0.0.2,0.0.5,0.0.7
Zeitintervall-Minimum $T_{\mathcal{M}min}$	Unix Timestamp	1765432345
Zeitintervall-Maximum $T_{\mathcal{M}max}$	Unix Timestamp	1765543456
Subjektteilmenge $\Sigma_{\mathcal{M}}$	Liste	ANN,RF
Zweckteilmenge $Z_{\mathcal{M}}$	Liste	C.1.k02

Tabelle 5.1.: Übersicht der Metadaten eines pragmatisch definierten maschinellen Modells.

5.3. Konstruktivistisches maschinelles Lernen

In der vorliegenden Arbeit wird Lernen als Konstruktion, Rekonstruktion oder Dekonstruktion von pragmatisch definierten maschinellen Modellen interpretiert. Diese Prozesse stellen die zentralen Lernprozesse konstruktivistischen maschinellen Lernens dar, dessen Ziel die Erzeugung einer hierarchisch geordneten Menge pragmatisch definierter maschineller Modelle ist. Als Ausgangspunkt wird eine Menge pragmatisch definierter vektorieller Modelle angenommen, im Folgenden Datenbasis genannt. Die vektoriellen Modelle können vollständig oder unvollständig sein.

5.3.1. Ablauf

Theoretisch soll jeder der drei Lernprozesse ausgehend von einem einzelnen vektoriellen Modell sowie durch schrittweise Hinzunahme weiterer vektorieller Modelle erfolgen. Da ein einzelnes vektorielles Modell in der Praxis jedoch nicht für die Anwendung maschineller Lernverfahren geeignet ist, erfolgt die Verarbeitung der Datenbasis in konsekutiv zu ziehenden Blöcke fester Größe (siehe auch Abb. 5.1). Innerhalb jedes gezogenen Blocks wird eine Teilmenge bzw. mehrere Teilmengen bestimmt, mit denen ein Lernprozess durchgeführt werden soll; die verbleibenden vektoriellen Modelle werden auf einer Halde abgelegt.

Identifikation und Selektion von Lernblöcken. TZ - oder $T\Sigma$ -Verwandtschaften innerhalb eines gezogenen Blocks werden mit einer einfachen Metadatenabfrage identifiziert. Da hier lediglich die unmittelbare Übereinstimmung beider Merkmale geprüft werden muss, ist dies mit minimalem algorithmischen Aufwand implementierbar. Enthält der gezogene Block hingegen eine Menge von ΣZ -verwandten vektoriellen Modellen, so muss es sich nicht zwingend um einen geeigneten Lernblock handeln; beispielsweise, weil größere zeitliche Lücken existieren. Identifizierte ΣZ -verwandte vektorielle Modelle werden daher anhand ihres zeitlichen Metadatum T geclustert (siehe auch Abschnitt 5.3.2). Übersteigt die Menge der identifizierten pragmatisch verwandten vektoriellen Modelle aus dem gezogenen Block einen benutzerdefinierten Schwellenwert, so wird diese Menge als geeignet für ein maschinelles Lernen betrachtet und als Lernblock bezeichnet. Enthält der gezogene Block mehrere Lernblöcke, wird jeweils der größte Lernblock zur Durchführung eines Lernprozesses verwendet.

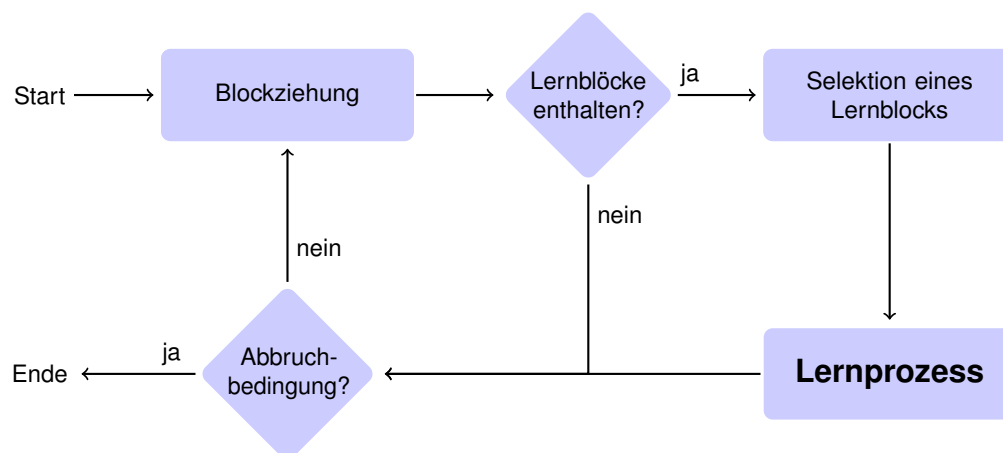


Abbildung 5.1.: Grundstruktur des konstruktivistischen maschinellen Lernens. Die Verarbeitung der Datenbasis erfolgt blockweise. Aus den vektoriellen Modellen eines Blocks werden mittels pragmatischer Verwandtschaft Lerndaten selektiert, anhand derer ein Lernprozess erfolgt.

5. Konstruktives maschinelles Lernen

Lernprozess. Anders als bei klassischen Ansätzen maschinellen Lernens bestimmt bei einem konstruktivistischen Ansatz nicht nur die Existenz vorgegebener Zielwerte, auf welche Weise mit einem Lernblock gelernt wird. Darüber hinaus ist es auch von dessen pragmatischen Eigenschaften abhängig, ob ein maschinelles Modell auf Basis eines identifizierten Lernblocks konstruiert oder rekonstruiert wird. Im Zuge eines Dekonstruktionsprozesses werden außerdem existierende maschinelle Modelle berücksichtigt. Grundlage dieses Vorgehens sind Faustregeln aus der konstruktivistischen Didaktik [204, S. 122ff], die hier wie folgt interpretiert werden:

1. *„So viel Konstruktion wie möglich!“* Jeder Lernblock wird zunächst auf eine Eignung zur Konstruktion neuer maschineller Modelle geprüft. Sind Zweck bzw. Zielparame-ter für die vektoriellen Modelle des Lernblocks unbekannt, so ist eine Konstruktion zwingend. Optional kann festgelegt werden, dass eine Konstruktion auch dann durchgeführt wird, wenn für alle vektoriellen Modelle des Lernblocks vorgegebene Zielwerte existieren.
2. *„Keine Rekonstruktionen um ihrer selbst willen!“* Das Ziel, einen bestimmten Parameter für einen gegebenen Datensatz im Sinne eines überwachten Lernens vorherzusagen, muss nicht in jedem Fall vollständig erreicht werden. Es ist folglich ein zulässiges Lernergebnis, wenn die angestrebte Rekonstruktion nur für einzelne der gegebenen Lernblöcke bzw. nur für einzelne Abschnitte der temporalen Ausdehnung des Datensatzes erfolgreich war.
3. *„Keine Konstruktionen ohne Ver-Störungen (sic!)“* Erfolgreich konstruierte maschinelle Modelle werden nicht nur mittels Rekonstruktion auf ihre Intersubjektivität geprüft, sondern im Zuge eines Dekonstruktionsprozesses auch auf Widersprüche zu bestehenden Modellen und gegebenenfalls auf Alternativen. Insbesondere können erfolgreich rekonstruierte maschinelle Modelle gegebenenfalls verändert oder verworfen werden.

Hieraus leitet sich ein Arbeitsablauf ab, um zu entscheiden, ob ein Lernblock in ein maschinelles Modell überführt und in die jeweilige Wissensdomäne integriert werden kann (Abb. 5.2).

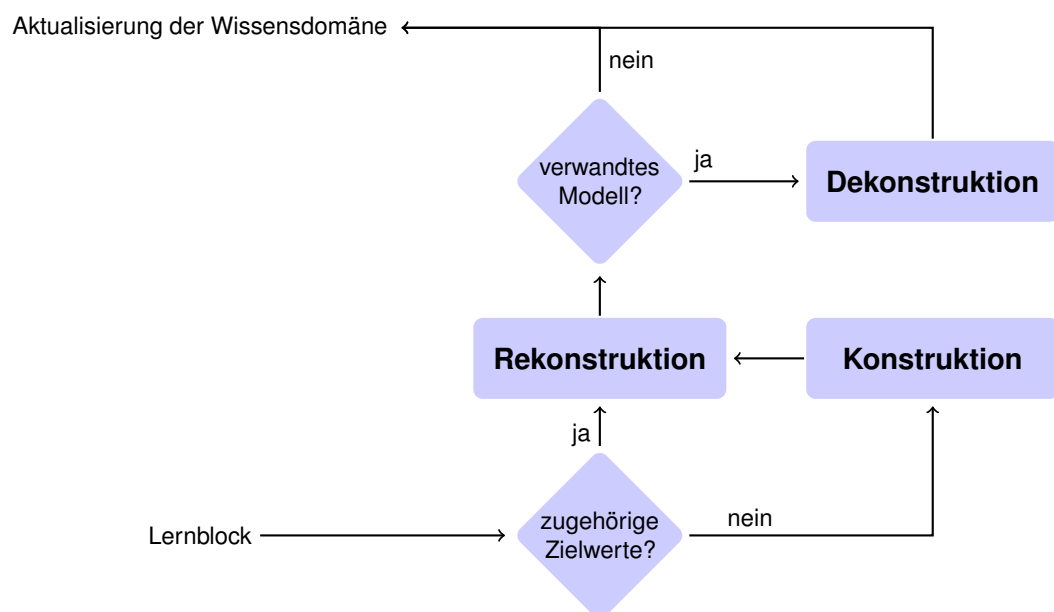


Abbildung 5.2.: Ablauf des Lernprozesses. Ein gegebener Lernblock wird entweder unmittelbar oder mittelbar durch den Rekonstruktionsprozess in eine Repräsentation als maschinelles Modell überführt. Dieses Modell wird dann auf mögliche Dekonstruktionen geprüft.

5.3.2. Implementierungsaspekte

Es wird vorausgesetzt, dass die zu analysierende Datenbasis in Form einer strukturierten Textdatei bzw. einer aus kommaseparierten Werten bestehenden Datei vorliegt. Der grundsätzliche Ablauf des konstruktivistischen maschinellen Lernens, soweit er nicht mittels *R*-Skripten realisiert ist, wird durch Terminalskripte gesteuert. Dies gilt insbesondere für die Ziehung von Blöcken sowie die Selektion von Lernblöcken und Lernprozessen. Um eventuelle Performance-Nachteile dieses Vorgehens zu minimieren, werden sämtliche Skripte im Shared Memory (`/dev/shm/`) ausgeführt und die Lernergebnisse erst am Ende in einen Zielordner geschrieben.

Lernblockidentifikation. Die Identifikation von Lernblöcken innerhalb einer Menge von ΣZ -verwandten vektoriellen Modellen erfolgt durch Clusteranalyse ihres zeitlichen Metadatum T , also eines eindimensionalen Vektors. Dazu wird zunächst die Dichteverteilung dieses Vektors mithilfe des *R*-Packages *KernSmooth*¹ abgeschätzt und mittels Schwellenwert (in %) abseitige T -Werte ausgeschlossen. Sofern sich die verbleibenden T -Werte in mehr als einen Cluster aufteilen lassen, wird mithilfe des *R*-Package *Ckmeans.1d.dp*² durchgeführt, das auf ein *k*-means-Clustering eindimensionaler Daten optimiert ist; anders als im allgemeinen Fall ist dabei im Eindimensionalen eine optimale Lösung möglich [261].

Abbruchbedingungen. Theoretisch kann konstruktivistisches maschinelles Lernen solange ausgeführt werden, bis alle vektoriellen Modelle der Datenbasis zur Erzeugung maschineller Modelle verwendet worden sind. In der Praxis stößt dies jedoch auf Zeit- und Ressourcengrenzen, weshalb ein benutzerdefinierter Mechanismus zum kontrollierten Abbruch implementiert wird. So kann einerseits die maximale Zahl von Blöcken definiert werden, die aus der Datenbasis gezogen und verarbeitet werden sollen. Zusätzlich kann definiert werden, ob bzw. mit wievielen Wiederholungen die Halde vektorieller Modelle auf Lernblöcke untersucht werden soll. Sobald eine der beiden Bedingungen erfüllt ist, wird das Lernen abgebrochen.

Wiederverwendbarkeit. Um erkannte Modelle auch für künftige Anwendungen nutzen zu können, muss die Konfiguration mindestens eines der beteiligten Lernverfahren nach der Rekonstruktion erhalten werden. Insbesondere müssen sämtliche gelernten Parameter gespeichert werden. Welches der verwendeten Verfahren (vgl. Abschnitt 5.4) hinterlegt wird, ist über einen Parameter definierbar. Die jeweiligen Architektur- und Gewichtsparameter werden in Unterverzeichnissen der Wissensdomäne gespeichert, die über die UID des Modells identifizierbar sind. Für das Modell mit der UID $C.1.3$ etwa werden die Information über das trainierte künstliche neuronale Netz im Unterordner des Ordners $C/1/3/ANN$ hinterlegt.

Eine Übersicht der entsprechenden Parameter für Ablauf und Blockziehung ist in Tab. 5.2 sowie in Anhang A zu finden.

	Parameter	Variablentyp	Default
Blockverarbeitung	Blockgröße	Integer	10.000
	Max. Anzahl Blocks	Integer	25
	Halden-Durchläufe	Integer	3
	Min. Lernblockgröße	Integer	2.000
	Cutoff für ΣZ -Verwandtschaft	Integer	20

Tabelle 5.2.: Übersicht der Blockverarbeitungsparameter im konstruktivistischen maschinellen Lernen.

¹Online abrufbar unter <https://cran.r-project.org/web/packages/KernSmooth/index.html>

²Online abrufbar unter <https://cran.r-project.org/web/packages/Ckmeans.1d.dp/index.html>

5.4. Konstruktion maschineller Modelle

Im didaktischen Kontext wird unter Konstruktion im Allgemeinen Kreativität, Innovation und Produktion, im Speziellen die Suche nach neuen Variationen, Kombinationen oder Transfers verstanden [205, S. 145]. Für maschinelle Modelle wird dies in der vorliegenden Arbeit als Identifikation bzw. Definition eines n -dimensionalen Outputs zu einer endlichen Menge unvollständiger vektorieller Modelle eines gegebenen Datensatzes interpretiert. Es wird also vorausgesetzt, dass Zielwerte entweder nicht bekannt sind oder nicht berücksichtigt werden.

Ziel einer Konstruktion ist somit die Identifikation eines gemeinsamen Zwecks für eine gegebene Menge von pragmatisch verwandten vektoriellen Modellen. Allerdings sind nicht alle Verwandtschaftsformen gleichermaßen für eine Modellkonstruktion geeignet. Insbesondere auf vollständig und TZ -verwandten vektoriellen Modellen basierende Konstruktionen bieten kaum Mehrwert. Vollständig verwandte vektorielle Modelle, die nicht redundant sondern divergent sind, stellen sogar eine schwerwiegende Fehlerquelle dar. TZ -verwandte Modelle erlauben grundsätzlich das Erzeugen neuer Modelle, was allerdings faktisch einen intersubjektiven Rekonstruktionsprozess darstellt. Für Konstruktionen werden daher entweder ΣZ -verwandte oder aus $T\Sigma$ -verwandten maschinellen Modellen gewonnene vektorielle Modelle verwendet.

A priori ist unklar, welches der aus einem Lernblock konstruierten Modelle den Anforderungen von Rekonstruktion und Dekonstruktion am gerechtsten wird. Während des Konstruktionsprozesses sollen daher möglichst viele möglichst unterschiedliche maschinelle Modelle identifiziert werden. Aus einem einzelnen Lernblock werden daher mehrere alternative Modelle konstruiert (Abb. 5.3). Der Konstruktionsprozess kann somit auch als Vervielfältigungs- oder Diversifizierungsmechanismus verstanden werden kann. Um aus dieser Menge das geeignetste Modell zu selektieren, werden diese konkurrierenden Modelle rekonstruiert und nach Vorhersagegenauigkeit gereiht. Das höchstgereichte Modell wird als Siegermodell entweder direkt in der Wissensdomäne gespeichert oder in den Dekonstruktionsprozess eingespeist (vgl. Abb. 5.2).

In der vorliegenden Arbeit werden Konstruktionsmechanismen für die konzeptuelle und prozedurale Wissensdomäne realisiert. Die hier eingesetzten Lernalgorithmen sind dabei als exemplarisches Minimalbeispiel zu sehen, nicht als abschließende Auswahl.

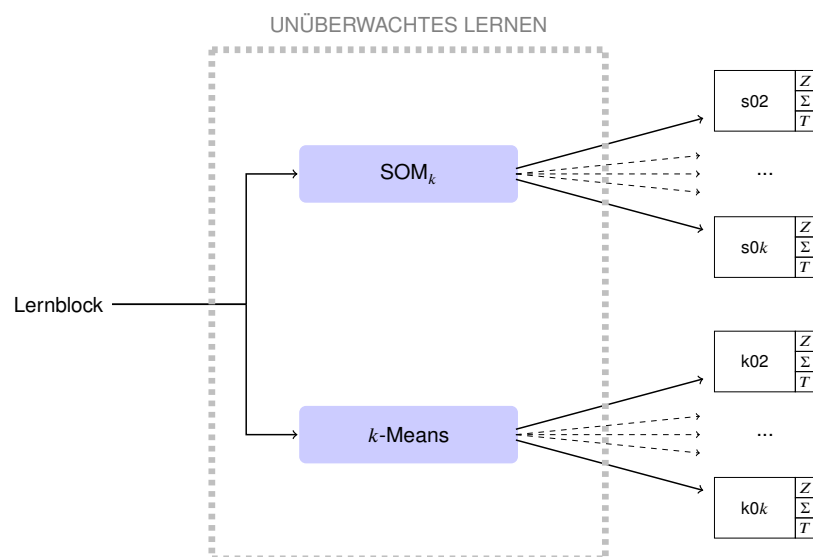


Abbildung 5.3.: Konstruktionsprozess für konzeptuelles Wissen. Unüberwachte Lernverfahren, hier: selbstorganisierende Karten (SOM) und k-Means, werden mit variierenden Parametern trainiert, um für einen Lernblock variierende Modelle $s_{02}, \dots, s_{0k}, k_{02}, \dots, k_{0k}$ zu konstruieren.

5.4.1. Konzeptuelles Wissen

Konzeptuelles Wissen im Sinne der Bloomschen Taxonomie kann beschrieben werden als Wissen um Klassifikationen, Kategorien und Strukturen. Maschinell lässt sich dies insbesondere durch Clustering- und Klassifizierungsalgorithmen umsetzen. Zweck des Clusterings ist in diesem Fall, innerhalb eines Lernblocks klar unterscheidbare Kategorien zu identifizieren.

Voraussetzung für den Einsatz von Clustering-Verfahren ist die vorherige Definition einer Zahl zu bestimmender Cluster. Üblicherweise werden daher mit dem gleichen Verfahren mehrere Durchläufe mit unterschiedlicher Cluster-Anzahl durchgeführt und die erhaltenen Clusterings mit einem externen Verfahren evaluiert [113]. Die maximale Cluster-Anzahl wird im Folgenden als maximale kategoriale Komplexität κ_k bezeichnet. Mit jedem Clustering-Verfahren werden somit $\kappa_k - 1$ maschinelle Modelle mit $k = \{2, \dots, \kappa_k\}$ Clustern oder Kategorien erzeugt.

Im Rahmen der vorliegenden Arbeit werden ein biologisch motiviertes sowie ein klassisches statistisches Clustering-Verfahren verwendet:

- **Selbstorganisierende Karte.** Kohonens Self-Organizing Map (SOM) bildet mit einer vorgegebenen Struktur künstlicher Neuronen Eingabevektoren ab. Dazu wird über die Euklidische Distanz dasjenige Neuron der SOM bestimmt, das am nächsten an einem gegebenen Eingangsvektor liegt. Der Index des Vektors wird mittels Minimumsuche bestimmt:

$$c = \arg \min_j (||X - W_j||) \quad (5.13)$$

wobei W_j der Gewichtsvektor ist.

Neben dem Gewinner-Neuron werden aber auch die im Sinne der Nachbarschaftsbeziehung benachbarten Neuronen angepasst; diese kann entweder konstant oder eine normalisierte Gaußfunktion sein. SOMs bewahren daher im Gegensatz zu vielen anderen unüberwachten Verfahren Nachbarschaftsbeziehungen und gelten insbesondere für Probleme mit mehreren Optima als gut geeignet [181].

Als Architektur werden hier eindimensionale SOMs verwendet. Dies hat einerseits den Vorteil, dass sich die konkrete Architektur unmittelbar aus der maximalen Modellkomplexität κ ableiten lässt. Andererseits wurde aber auch beobachtet, dass eindimensionale SOMs bessere Ergebnisse liefern als zweidimensionale [12].

Mit dem *R* package *kohonen*³ steht eine etablierte und freie Implementierung von Kohonens SOM zur Verfügung. Dieses erlaubt es nicht nur, SOMs schnell und effizient zu trainieren, sondern bietet auch umfangreiche Visualisierungsmöglichkeiten.

- **k-Means-Algorithmus.** Mit diesem Verfahren werden k Partitionen für eine Menge gegebener Datenpunkte x_j bestimmt, indem die Summe der quadratischen Abweichungen J zu den Clusterzentren μ_i minimiert werden:

$$J = \sum_{i=1}^k \sum_{x_j \in S_i} ||x_j - \mu_i||^2 \quad (5.14)$$

Aufgrund dieses vergleichsweise intuitiven und effizient implementierbaren Vorgehens konnte k-Means in den vergangenen Jahrzehnten in zahllosen Datenanalysen erfolgreich eingesetzt werden. Im Gegensatz zu anderen Verfahren erfordert die Durchführung beispielsweise außer der Anzahl der zu identifizierenden Cluster k keine weiteren Parameter. Relevanz und Popularität des k-Means-Verfahrens ist auch daran erkennbar, dass in der Programmierumgebung *R* für dessen Nutzung kein zusätzliches Paket geladen werden muss, sondern k-Means direkt im *R*-Kernel implementiert ist.

³Online abrufbar unter <http://cran.r-project.org/web/packages/kohonen/>

5.4.2. Prozedurales Wissen

Prozedurales Wissen im Sinne der Bloomschen Taxonomie kann beschrieben werden als Wissen darum, wie etwas gemacht wird – also als Wissen um subjektspezifische Fähigkeiten, Algorithmen oder Auswahlkriterien. Im Folgenden wird prozedurales Wissen daher als Gesamtmenge maschineller Modelle interpretiert, die Regressionsaufgaben realisieren. Die Konstruktion neuer prozeduraler Modelle für einen gegebenen Lernblock erfolgt daher durch Extraktion bislang unbekannter metrischer Features.

Im Rahmen der vorliegenden Arbeit werden ein teils biologisch motiviertes sowie ein statistisch orientiertes Verfahren verwendet:

- **Autoencoder.** Ein neuronales Netz mit symmetrischer Ein- und Ausgabeschicht wird verwendet, um die Abbildung der Eingangsdaten auf sich selbst zu lernen. Die ursprünglich vorgeschlagene Umsetzung von Autoencodern basierte auf einfachen Feedforward-Netzen und dem Standard-Backpropagation-Verfahren, später wurden jedoch auch komplexere mehrstufige Verfahren vorgeschlagen [98].

Ein Mittelweg stellt ein so genannter Sparsity Autoencoder, der nicht nur Feedforward-Netz und Backpropagation, sondern auch einen festen Sparsity-Parameter ρ nutzt:

$$E_{x \sim D} \left[a_i^{(2)} \right] \rightarrow \rho \quad (5.15)$$

wobei $a_i^{(2)}$ die Aktivierung des versteckten Neurons i repräsentiert. Die Gewichte des Neurons i werden nach jedem Durchgang so angepasst, dass sie diese Bedingung erfüllen.

Mit dem *R* package *autoencoder*⁴ steht eine Implementierung dieses Sparse Autoencoders zur Verfügung. Es erlaubt insbesondere es eine beliebige Wahl an versteckten Schichten und Neuronen sowie die direkte Ausgabe der Werte der versteckten Neuronenschicht.

- **Feature Extraction via Clustering.** Eine gegebene Menge von Features wird in k Feature-Cluster unterteilt, wobei stark korrelierende Features dem gleichen Cluster zugewiesen werden. Diese Unterteilung kann beispielsweise mittels hierarchischem Clustering oder mittels k-Means erfolgen [39]. Für jeden gefunden Cluster C_k wird dann mittels Centroid-Bildung ein synthetisches Feature c_k erzeugt:

$$c_k = \operatorname{argmax}_{u \in R^n} \left\{ \sum_{x_j \in C_k} r_{u,x_j}^2 + \sum_{y_j \in C_k} \eta_{u|y_j}^2 \right\} \quad (5.16)$$

wobei r^2 die quadrierte Pearson-Korrelation und für $\eta^2 \in [0, 1]$ gilt:

$$\eta_{u|y_j}^2 = \frac{\sum_{s \in \mathcal{M}_j} n_s (u_s - u)^2}{\sum_{i=1}^n (u_i - u)^2} \quad (5.17)$$

wobei \mathcal{M} die Kategorienmenge ist, n_s die Frequenz der Kategorie s , u der Mittelwert der Variable u und u_s der Mittelwert der Kategorie s von u .

Eine geeignete Implementierung eines solchen Feature-Clustering bietet das *R*-Package *ClustOfVar*⁵. Insbesondere erlaubt dies sowohl das Clustern von numerischen wie auch kategorialen Variablen. Als Clusterverfahren wird aufgrund der effizienteren Berechnung k-Means dem dort ebenfalls implementierten hierarchischen Clustering vorgezogen.

⁴Online abrufbar unter <http://cran.r-project.org/web/packages/autoencoder/>

⁵Online abrufbar unter <http://cran.r-project.org/web/packages/ClustOfVar/>

5.5. Rekonstruktion maschineller Modelle

Im didaktischen Kontext wird unter Rekonstruktion im Allgemeinen Anwendung, Wiederholung oder Nachahmung, im Speziellen die Suche nach Ordnung, Mustern oder Modellen verstanden [205, S. 145]. Analog dazu wird hier unter der Rekonstruktion eines maschinellen Modells das überwachte Lernen aus gegebenen Beispielen verstanden. Im Unterschied zu klassischen überwachten Lernverfahren werden jedoch mehrere konkurrierende maschinelle Modelle erzeugt und hinsichtlich ihrer intersubjektiven Gültigkeit bewertet.

Der Ablauf des Rekonstruktionsprozesses setzt sich dabei aus drei Teilschritten zusammen (Abb. 5.4). Im Rahmen einer Vorverarbeitung wird ein gegebener Lernblock bzw. eine gegebene Menge vektorieller Modelle zunächst auf die Zahl seiner Eingangsvariablen geprüft. Ist diese größer als die definierte maximale Modellkomplexität κ , findet eine algorithmische Feature-Auswahl statt, um die Komplexität des zu erlernenden Modells auf $\leq \kappa$ zu reduzieren.

Anschließend folgt das eigentliche überwachte Lernen, das hier mit zwei konkurrierenden Verfahren durchgeführt wird. Dabei werden zunächst nur die Ergebnisse jedes einzelnen Lernprozesses evaluiert. Im letzten Schritt wird die intersubjektive Übereinstimmung zwischen den Verfahren evaluiert und gegebenenfalls auf dieser Basis das geeignetste Modell selektiert.

5.5.1. Lernverfahren

Wie bereits zuvor beschrieben, sollen potentielle maschinelle Modelle nicht von parallel agierenden Agenten des gleichen überwachten Lernverfahrens erlernt werden, sondern von unterschiedlichen Verfahren. Um eine breite Anwendung zu ermöglichen, sollten diese Verfahren sowohl Klassifikations- als auch Regressionsaufgaben lösen können. Um die automatisierte Konfiguration zu erleichtern, sollten sie gleichzeitig nicht-parametrisch sein, also keine a-priori-Annahmen über die Dichteverteilung der Daten erfordern. Weiter ist eine grundsätzliche Diversität der Verfahren, im Sinne unterschiedlicher Verfahrensansätze, wünschenswert. Dazu werden hier ein biologisch inspiriertes und ein statistisch motiviertes Lernverfahren parallel angewendet.

Unter Berücksichtigung dieser Kriterien sowie der Verfügbarkeit geeigneter Implementierungen kommen für den Rekonstruktionsprozess folgende Verfahren zum Einsatz:

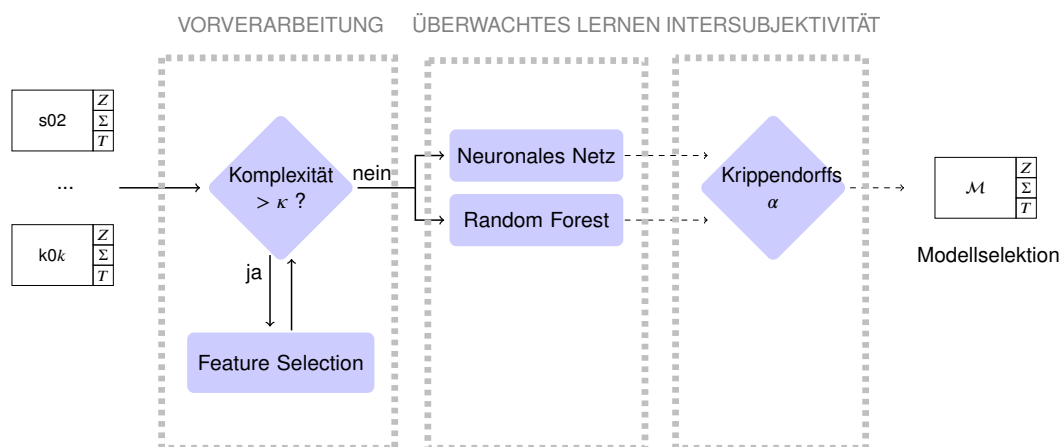


Abbildung 5.4.: Rekonstruktionsprozess. Nach Vorverarbeitung und überwachtem Lernen folgt eine Intersubjektivitätsbewertung mittels Reliabilitätskoeffizient. Wurde mehr als ein Modell erfolgreich rekonstruiert, wird mittels Modellselektion ein Sieger-Modell bestimmt.

5. Konstruktives maschinelles Lernen

- **Künstliche neuronale Netze.** Es wird ein klassisches Feedforward-Netz mit einer versteckten Schicht verwendet. Für die Neuronen der Ein- und Ausgabeschicht wird eine lineare Aktivierungsfunktion verwendet ($f(x) = x$), für diejenigen der versteckten Schicht eine sigmoide:

$$f(x) = \frac{1}{1 + e^{-x}} - \frac{1}{2} \quad (5.18)$$

wobei x jeweils die Summe der gewichteten Eingangssignale ist.

Als Trainingsalgorithmus wird Resilient Backpropagation (Rprop) verwendet, dessen Geschwindigkeit und Robustheit gegenüber klassischem Backpropagation optimiert wurden [211]; Rprop vereint damit die Vorteile von Manhattan-Training und Quickprop vereint. Die verwendete Architektur wird durch die Anzahl n der Eingabevariablen definiert, wobei pro Rekonstruktionsvorgang mehrere konkurrierende Netze mit den Architekturen $n-i-1$ mit $n \in \{2, \dots, 10\}$ und $i \in \{2, \dots, \lfloor n \cdot 1, 5 \rfloor\}$ trainiert werden; von diesen wird am Ende dasjenige mit der höchsten Vorhersagegenauigkeit für die weitere Analyse verwendet. Weitere Einzelheiten zu Parametern und Architekturen sind in Anhang A aufgeführt.

Mit der C-Bibliothek FANN⁶ steht ein leistungsfähiges Framework zur Verfügung, mit dem sich künstliche neuronale Netze schnell und effizient realisieren lassen. Insbesondere lassen sich so Rprop-Training und eine Automatisierung der Architekturauswahl realisieren.

- **Random Forests.** Als alternatives Lernverfahren wird ein Random Forest eingesetzt, also eine Menge von Entscheidungsbäumen, die aufgrund von Randomisierung nicht miteinander korreliert sind. Dazu verwenden diese Bäume einerseits statt des gesamten Datensatzes nur zufällige Stichproben, andererseits erhält jeder Baum statt aller Eingabevariablen nur eine zufällige Auswahl davon [30]. Aus diesem Grund können die Bäume unabhängig von einander trainiert werden, wobei das Random-Forest-Verfahren in hohem Maß parallelisierbar ist. Auch bei einer Evaluation von Testdaten erhalten die Entscheidungsbäume unabhängig voneinander die gleiche Eingabe, ihre potentiell unterschiedlichen Ergebnisse werden für die Gesamtvorhersage am Ende jedoch gemittelt.

Die Zahl der erforderlichen Entscheidungsbäume ist von der Komplexität des zu erlernenden Zusammenhangs abhängig, wird jedoch zumeist relativ hoch angesetzt. Soweit nicht anders angegeben, wird die Zahl der Bäume `ntree` in der vorliegenden Arbeit in einer moderaten Abhängigkeit von der maximalen Modellkomplexität mit `ntree = 50κ` gewählt.

Mit dem R package *randomForest*⁷ steht ein leistungsfähiges Framework zur Verfügung, das insbesondere eine parallelisierte Ausführung des Lernalgorithmus erlaubt.

5.5.2. Komplexitätsreduktion

Setzt sich der Urbildraum eines zu rekonstruierenden Modells aus mehr als κ Features zusammen, wird vor dem überwachten Lernen eine Komplexitätsreduktion mittels algorithmischer Feature-Auswahl durchgeführt. Prinzipiell kommen dafür sowohl Filter- als auch Wrapper- und Embedded-Verfahren in Frage. Ein wichtiges Auswahlkriterium stellt die Menge der zu verarbeitenden Daten bzw. der durchzuführenden Auswahlprozesse dar. Für große Datenmengen sind Filter-Methoden zwar aufgrund ihrer Effizienz und ihres vergleichsweise geringen Rechenaufwands zu bevorzugen, allerdings liefern diese nicht immer optimale Feature-Subsets. Eingebettete Feature-Auswahl-Algorithmen versprechen dagegen eine besonders sorgfältige Auswahl von Features, erfordern jedoch für große Datenmengen einen erheblich höheren Rechenaufwand als Filter.

⁶Online abrufbar unter <http://leenissen.dk/fann/>

⁷Online abrufbar unter <http://cran.r-project.org/web/packages/randomForest/>

Aus diesem Grund wird hier ein Hybrid-Ansatz verfolgt. Als Grundsatz gilt: Sofern ein Eingabedatensatz aus einer kleineren Menge von niederdimensionalen vektoriellen Modellen besteht, wird ein eingebettetes Verfahren bevorzugt. Im Umkehrschluss wird ein Filter verwendet, falls ein Eingabedatensatz aus einer großen Menge hochdimensionaler vektorieller Modelle besteht. Ob ein eingebettetes Verfahren zum Einsatz kommt, wird dabei mittels zweier benutzerdefinierter Hilfsparameter entschieden:

1. λ_x definiert die dafür maximal zulässige Zahl an Eingabedimensionen,
2. λ_y definiert die dafür maximal zulässige Zahl vektorieller Modelle.

Ist die Zahl der Eingabedimensionen $n > \lambda_x$ oder die Zahl der vektoriellen Modelle $m > \lambda_y$, wird ein Filter-Verfahren angewendet. Unter Berücksichtigung verfügbarer Implementierungen kommen hierfür folgende Methoden zum Einsatz:

1. **Correlation-based Feature Selection (CFS)**. Mit diesem multivariaten Filter-Verfahren werden nicht nur einzelne Features bewertet, sondern Feature-Subsets. Entsprechend ist das Ergebnis des Verfahrens keine Feature-Reihung, sondern die Angabe eines konkreten Feature-Subsets. Ist die Anzahl der ausgewählten Features größer als eine gegebene maximale Modellkomplexität κ , so findet unter Berücksichtigung von λ_x und λ_y eine erneute Komplexitätsreduktion statt.

Eine geeignete Implementierung dieses Verfahrens bietet das R-Package *FSelector*⁸, welches hier zum Einsatz kommt.

2. **Random Forests (RF)**. Dieses überwachte Lernverfahren kann als eingebettetes Verfahren verwendet werden, da bei jedem Lernvorgang implizit eine qualifizierte Reihung aller gegebenen Features erstellt wird. Um mit dieser Reihung eine Komplexitätsreduktion auf eine gegebene maximale Modellkomplexität κ zu erreichen, werden die ersten κ Features zur Erzeugung von κ Feature-Subsets mit 1 bis κ Dimensionen verwendet. Diese Subsets werden dann mittels überwachtem Lernen evaluiert, wobei mit dem Sieger-Subset anschließend ein neuer Lernblock erzeugt wird, der den bisherigen ersetzt.

Analog zur Verwendung im eigentlichen Rekonstruktionsprozess kommt hierfür das R-Package *randomForest*⁹ mit den oben genannten Parametern zum Einsatz.

5.5.3. Modellselektion

Während des Konstruktionsprozesses erzeugt ein unüberwachtes Lernverfahren mehrere maschinelle Modelle für einen gegebenen Lernblock (vgl. Abschnitt 5.3). Ziel des anschließenden Rekonstruktionsprozesses ist es einerseits, jedes dieser konstruierten Modelle auf seine Intersubjektivität hin zu evaluieren. Da die konstruierten Modelle in Konkurrenz zueinander stehen, ist außerdem zu entscheiden, welches der Modelle als Repräsentation des Lernblocks verwendet und in die Wissensdomäne integriert wird. Dazu werden die konstruierten Modelle entsprechend ihres Grades an Intersubjektivität gereiht und dasjenige Modell mit der besten Platzierung ausgewählt. Alle übrigen Modelle werden verworfen.

Bewertung der Intersubjektivität. Da innerhalb eines Rekonstruktionsprozesses zwei oder mehr Lernverfahren parallel eingesetzt werden, liegen für jeden Eingabevektor zwei oder mehr unabhängige, gegebenenfalls nicht-identische Zielwerte vor. Folglich stellt sich die Frage, inwieweit die Ergebnisse dieser konkurrierenden Verfahren insgesamt übereinstimmen. Um dies zu beantworten, wird analog zu empirischen Studien auf eine Interrater-Reliabilität zurückgegriffen.

⁸Online abrufbar unter <http://cran.r-project.org/web/packages/FSelector/>

⁹Online abrufbar unter <http://cran.r-project.org/web/packages/randomForest>

5. Konstruktives maschinelles Lernen

Unterschiede zwischen konkurrierenden Verfahren werden also nicht einfach durch Ersetzung durch einen Mittelwert nivelliert, sondern mithilfe des zufallskorrigierenden Reliabilitätskoeffizienten Krippendorffs α explizit quantifiziert und bewertet.

In Abhängigkeit von beobachteter Nicht-Übereinstimmung D_o und unter Zufallsbedingungen erwartbarer Nicht-Übereinstimmung D_e ist der Koeffizient auf dem Intervall $[-1, 1]$ definiert als:

$$\alpha = 1 - \frac{D_o}{D_e} \quad (5.19)$$

$$D_o = \frac{1}{n} \sum_c \sum_k o_{ck} \text{metric} \delta_{ck}^2 \quad (5.20)$$

$$D_e = \frac{1}{n(n-1)} \sum_c \sum_k n_c n_k \text{metric} \delta_{ck}^2 \quad (5.21)$$

wobei o_{ck} , n_c , n_k und n aus einer Koinzidenz-Matrix bestimmt werden (vgl. [138]). Diese Matrix wird aus den Ergebnissen aller konkurrierenden Bewerter erzeugt, wobei dafür je nach Definition entweder die Menge aller Samples oder eine definierte Teilmenge verwendet wird. In der Praxis wird aus Performanzgründen die Koinzidenzmatrix häufig lediglich für einen Teil der Eingabevektoren berechnet. In der vorliegenden Arbeit wird die Samplesbasis der Koinzidenzmatrix daher über einen benutzerdefinierten Parameter bestimmt.

Krippendorffs α lässt sich sowohl für nominale als auch metrische Skalen berechnen und kann daher sowohl für Klassifikations- als auch Regressionsaufgaben verwendet werden [138]. Im Gegensatz zu anderen Reliabilitätskoeffizienten kann er außerdem auf eine beliebige Anzahl konkurrierender Zielwerte bzw. Verfahren angewendet werden. Krippendorffs α wurde nicht nur wiederholt als Standardmaßzahl zur Quantifizierung einer Interrater-Reliabilität vorgeschlagen [139, 93], sondern auch in der Skriptsprache *R* als Teil des Packages *irr*¹⁰ implementiert.

Schwellenwertprüfung. Vor der Reihung aller rekonstruierten Modelle werden diejenigen verworfen, die zwar von den verwendeten überwachten Lernverfahren für sich betrachtet mit der benutzerdefinierten Präzision trainiert wurden, deren zugehöriger α -Wert jedoch einen benutzerdefinierten Schwellenwert nicht übersteigt. Bei der Wahl dieses Schwellenwerts ist zu beachten, dass $\alpha = 1$ eine optimale Reliabilität anzeigt, während $\alpha \leq 0$ impliziert, dass keinerlei Übereinstimmung der Bewertungen vorliegt. Übersteigt keines der rekonstruierten Modelle diesen α -Schwellenwert, wird der aktuelle Rekonstruktionsprozess insgesamt abgebrochen.

Modellreihung. Bleiben nach der Schwellenwertprüfung maschinelle Modelle erhalten, so bedeutet dies, dass sie von den eingesetzten überwachten Verfahren mit ausreichender Präzision erlernt wurden und eine ausreichende Interrater-Reliabilität aufweisen. Um die am wenigsten von spezifischen Verfahren abhängige Lernblock-Repräsentation zu identifizieren, werden diese Modelle anhand von Krippendorffs α absteigend gereiht. Da der maximale α -Wert das maximale Maß an Intersubjektivität impliziert, kann dasjenige Modell mit maximalem α -Wert im Sinne von Heinrich Hertz als deutlichstes Modell interpretiert werden¹¹. Sollten innerhalb der Reihung zwei oder mehr Modelle einen identischen α -Wert aufweisen, wird aus dieser neuen Teilmenge dasjenige Modell mit dem kleinsten Urbildraum ausgewählt. Im Sinne von Heinrich Hertz kann dies als Wahl des einfachsten Modells interpretiert werden¹².

Das auf diese Weise als optimale Lernblock-Repräsentation selektierte Modell wird in den anschließenden Dekonstruktionsprozess überführt. Alle übrigen Modelle werden verworfen.

¹⁰Online abrufbar unter <http://cran.r-project.org/web/packages/irr>

¹¹„Von zwei Bildern desselben Gegenstandes wird dasjenige das zweckmäßigere sein, welches mehr wesentliche Beziehungen des Gegenstandes widerspiegelt als das andere; welches, wie wir sagen wollen, das deutlichere ist.“ [97, S. 2f]

¹²„Bei gleicher Deutlichkeit wird von zwei Bildern dasjenige zweckmäßiger sein, welches neben den wesentlichen Zügen die geringere Zahl überflüssiger oder leerer Beziehungen enthält, welches also das einfachere ist.“ [97, S. 2f]

5.6. Dekonstruktion maschineller Modelle

In der Didaktik wird unter Dekonstruktion im Allgemeinen die Untersuchung einer bereits bestehenden Konstruktion auf Unvollständigkeit, auf Unvorgesehenes und Unbewusstes, im Speziellen die Suche nach möglichen Auslassungen, Vereinfachungen, Ergänzungen und Kritik verstanden [205, S. 145]. Im Folgenden wird daher die Dekonstruktion eines maschinellen Modells als gezieltes Infragestellen der Gültigkeit dieses Modells interpretiert.

Ablauf. Ausgehend von einem aus einem Lernblock rekonstruierten Modell \mathcal{M} wird jedes pragmatisch verwandte Modell \mathcal{M}_x aus der zugehörigen Wissensdomäne infrage gestellt. Für den Fall, dass für \mathcal{M} zwei oder mehr verwandte Modelle identifiziert werden, können diese je nach Benutzervorgabe entweder konsekutiv dekonstruiert werden oder der Dekonstruktionsprozess wird abgebrochen, sobald eine vollständige, ΣZ -, TZ - oder $T\Sigma$ -Dekonstruktion erfolgreich war. Abb. 5.5 zeigt den grundsätzlichen Ablauf eines solchen Dekonstruktionsprozesses.

Verwandtschaftsidentifikation. Aufgrund der besonderen Bedeutung einer vollständigen Verwandtschaft für bereits bestehende Modelle der Wissensdomäne, wird zunächst geprüft, ob eine solche vorliegt bzw. ob eine vollständige Dekonstruktion durchgeführt werden kann. Ist dies nicht der Fall, werden nacheinander ΣZ -, TZ - und $T\Sigma$ -Dekonstruktion geprüft. Sobald ein pragmatisch verwandtes maschinelles Modell \mathcal{M} identifiziert worden ist, wird weiter geprüft, ob sich beide zu einem neuen Modell $\mathcal{M}' = \mathcal{M} \cup \mathcal{M}_x$ vereinigen lässt.

Modellvereinigung. Ziel einer Dekonstruktion ist die Erweiterung bzw. Ersetzung bestehender Modelle \mathcal{M}_x . Es wird daher versucht, \mathcal{M} und \mathcal{M}_x zu einem neuen, hypothetischen Modell \mathcal{M}' zu vereinigen. Die Voraussetzungen für eine solche Modellvereinigung variieren je nach Art der pragmatischen Verwandtschaft zwischen \mathcal{M} und \mathcal{M}_x (siehe nachfolgende Abschnitte). Ist die Vereinigung erfolgreich, wird \mathcal{M}' mittels Rekonstruktion evaluiert; andernfalls wird die Dekonstruktion abgebrochen und \mathcal{M} als neues Modell in die Wissensdomäne integriert (Modellspeicherung). Bei erfolgreicher Rekonstruktion ersetzt \mathcal{M}' das bisherige Modell \mathcal{M}_x in der Wissensdomäne, andernfalls werden sowohl \mathcal{M} und \mathcal{M}_x verworfen (Modellentsorgung).

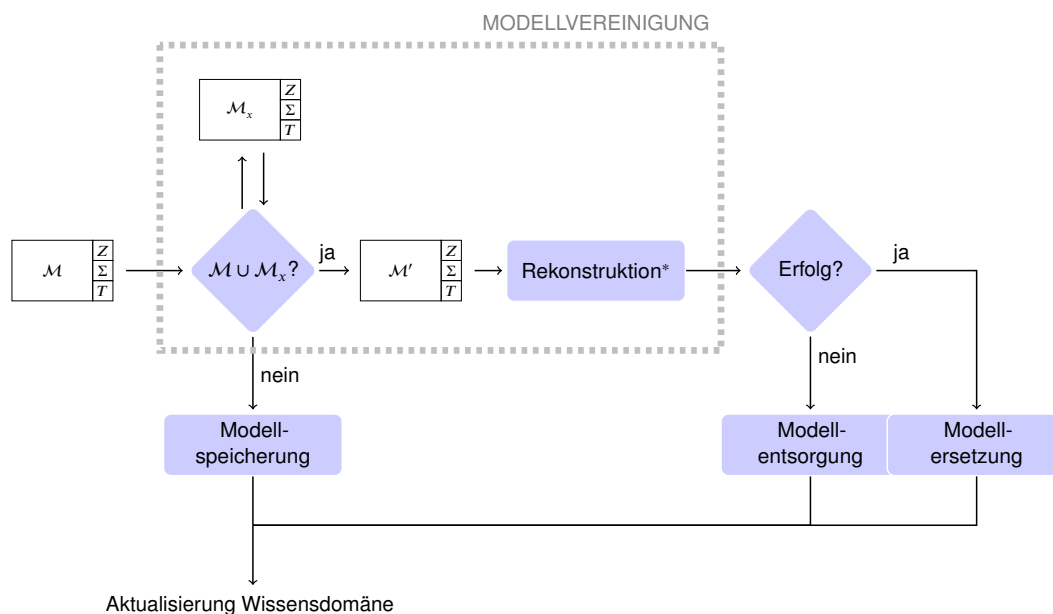


Abbildung 5.5.: Dekonstruktionsprozess. Für ein rekonstruiertes Modell \mathcal{M} werden Modellvereinigungen \mathcal{M}' mit Modellen \mathcal{M}_x aus der Wissensdomäne geprüft. Alle \mathcal{M}' werden anschließend erneut rekonstruiert (hier ohne Komplexitätsreduktion und Modellselektion).

5. Konstruktives maschinelles Lernen

Modellentsorgung. War die Vereinigung von \mathcal{M} und \mathcal{M}_x zu \mathcal{M}' erfolgreich, nicht jedoch die Rekonstruktion von \mathcal{M}' nicht erfolgreich, wird \mathcal{M}' verworfen. Der Umgang mit \mathcal{M} und \mathcal{M}_x erfolgt in diesem Fall benutzerdefiniert nach einer der folgenden Dekonstruktionsstrategien:

1. **Konservativ.** \mathcal{M}_x wird beibehalten, \mathcal{M} verworfen;
2. **Integrativ.** \mathcal{M}_x wird beibehalten, \mathcal{M} als neues Modell in die Domäne aufgenommen;
3. **Opportunistisch.** Falls \mathcal{M} auf einer größeren Menge vektorieller Modelle bzw. Trainingsdaten basiert als \mathcal{M}_x , wird \mathcal{M} als neues Modell in die Domäne aufgenommen und \mathcal{M}_x verworfen. Andernfalls wird \mathcal{M}_x beibehalten und \mathcal{M} verworfen.

5.6.1. ΣZ -Dekonstruktion

Mit einer ΣZ -Dekonstruktion wird geprüft, ob sich die zeitliche Gültigkeit von \mathcal{M}_x durch Vereinigung mit \mathcal{M} zu \mathcal{M}' erweitern lässt. Soweit nicht anders beschrieben, folgt dieser Prozess dem grundsätzlichen Ablaufschema eines Dekonstruktionsprozesses (Abb. 5.5).

Verwandtschaftsidentifikation. Voraussetzung für eine ΣZ -Dekonstruktion ist die vorherige Identifikation einer nicht-leeren Menge bereits in der Wissensdomäne gespeicherter maschineller Modelle \mathcal{M}_x mit ΣZ -Verwandtschaft zu \mathcal{M} . Da der zu erwartende Rekonstruktionserfolg umso geringer wird, je größer die zeitliche Distanz zwischen \mathcal{M}_x und \mathcal{M} ist, wird die Menge der zu berücksichtigenden verwandten Modelle mittels Schwellenwert eingeschränkt. Dies erfolgt durch Definition einer maximal erlaubten Distanz δT_{max} zwischen den Zeitintervallen ΣZ -verwandter Modelle \mathcal{M}_x und \mathcal{M} , die auf Basis eines Intervalls von 0 bis 1 drei unterschiedliche Voraussetzungen für die Durchführung einer ΣZ -Dekonstruktion ermöglicht:

1. $\delta T_{max} = 1$: Eine Expansion wird stets geprüft.
2. $\delta T_{max} = 0$: Eine Expansion wird nur geprüft, falls sich $T_{\mathcal{M}}$ und $T_{\mathcal{M}_x}$ temporal überlappen, d.h. falls entweder $\tau_{min}(\mathcal{M}_x) \geq \tau_{max}(\mathcal{M})$ oder $\tau_{min}(\mathcal{M}) \geq \tau_{max}(\mathcal{M}_x)$.
3. $0 < \delta T_{max} < 1$: Eine Expansion wird nur geprüft, falls die Distanz zwischen \mathcal{M}_x und \mathcal{M} das δT_{max} -fache der Spannweite $R_T(\mathcal{M}') = \tau_{max}(\mathcal{M}') - \tau_{min}(\mathcal{M}')$ nicht übersteigt, d.h. falls $\tau_{min}(\mathcal{M}_x) - \tau_{max}(\mathcal{M}) < \delta T_{max} \cdot R_T(\mathcal{M}')$ bzw. $\tau_{min}(\mathcal{M}) - \tau_{max}(\mathcal{M}_x) < \delta T_{max} \cdot R_T(\mathcal{M}')$

Modellvereinigung. Aus praktischen Erwägungen wird gefordert, dass die Urbildräume von \mathcal{M}_x und \mathcal{M} bei einer ΣZ -Dekonstruktion eine Schnittmenge mit zwei oder mehr Elementen besitzen müssen. Grund dafür ist, dass sonst keine kohärente Menge von Trainings- bzw. Testdaten für die Rekonstruktion von \mathcal{M}' mittels überwachter Lernverfahren erzeugt werden kann.

5.6.2. $T\Sigma$ -Dekonstruktion

Mit einer $T\Sigma$ -Dekonstruktion wird geprüft, ob sich aufbauend auf zwei $T\Sigma$ -verwandten Modellen \mathcal{M}_x und \mathcal{M} ein neues Modell auf der nächsthöheren Wissensebene konstruieren lässt. Voraussetzung dafür ist, dass das benutzerdefinierte Maximum an Wissensebenen nicht überschritten wird. Soweit nicht anders beschrieben, folgt eine $T\Sigma$ -Dekonstruktion dem grundsätzlichen Ablaufschema eines Dekonstruktionsprozesses (Abb. 5.5).

Modellvereinigung. Für diesen abstrahierenden Prozess wird ein erfolgreich rekonstruiertes \mathcal{M} zunächst unmittelbar als neues Modell in die Wissensdomäne integriert. Sofern \mathcal{M}_x und \mathcal{M} eine ausreichende Anzahl an Samples mit identischem zeitlichen Metadatum T aufweisen, werden die zugehörigen Zielparameter beider Modelle mittels Krippendorffs α auf Übereinstimmung getestet. Wird dabei der für eine Parameterübereinstimmung festgelegte Schwellenwert überschritten, so wird das vereinigte Modell verworfen. Andernfalls werden die Zielwerte der übereinstimmenden

Samples anschließend zum Urbildraum eines neuen Modells vereinigt und in einen neuen Zyklus aus Konstruktion, Rekonstruktion und Dekonstruktion eingespeist. Dieser neue Zyklus wird auf der nächsthöheren Wissensebene der jeweiligen Domäne ausgeführt.

5.6.3. TZ -Dekonstruktion

Mit einer TZ -Dekonstruktion wird geprüft, ob sich die Menge der gültigen Modellsubjekte von \mathcal{M}_x und \mathcal{M} erweitern lässt. Soweit nicht anders beschrieben, folgt eine TZ -Dekonstruktion dem grundsätzlichen Ablaufschema eines Dekonstruktionsprozesses (Abb. 5.5).

Modellvereinigung. Für diesen intersubjektivierenden Prozess wird ein erfolgreich rekonstruiertes \mathcal{M} zunächst unmittelbar als neues Modell in die Wissensdomäne integriert. Sofern \mathcal{M}_x und \mathcal{M} eine ausreichende die Anzahl an Samples mit identischem zeitlichen Metadatum T aufweisen, werden diese zu einem neuen Modell \mathcal{M}' vereinigt.

Rekonstruktion. Aufgrund des Verwandtschaftsgrade kann anstelle eines vollständigen Rekonstruktionsprozesses unmittelbar eine Überprüfung der zugehörigen Zielparameter beider Modelle mittels Krippendorffs α auf Übereinstimmung durchgeführt werden. Wird dabei der für eine Parameterübereinstimmung festgelegte Schwellenwert überschritten, so wird \mathcal{M}_x durch \mathcal{M}' ersetzt. Andernfalls bleiben alle Modelle der Wissensdomäne unverändert.

5.6.4. Vollständige Dekonstruktion

Im Gegensatz zu ΣZ -, $T\Sigma$ - und TZ -Dekonstruktionen ist eine vollständige Dekonstruktion nicht nur dazu geeignet, ein bestehendes Modell \mathcal{M}_x zu erweitern, sondern auch dieses zu falsifizieren. Daher realisiert dieser Dekonstruktionsmechanismus neben einer Modellerweiterung auch das Differenzieren bzw. gegebenenfalls das Verwerfen eines existierenden Modells. Eine vollständige Dekonstruktion folgt dem grundsätzlichen Ablaufschema eines Dekonstruktionsprozesses (Abb. 5.5), wobei jedoch einerseits die Modellvereinigung komplexer gestaltet ist sowie andererseits ein zusätzlicher Prozess zur Modelldifferenzierung vorgesehen ist (Abb. 5.6).

Modellvereinigung. Eine vollständige Verwandtschaft setzt nicht nur übereinstimmende Subjekt- und Zweckmengen voraus, sondern impliziert auch eine zeitliche Kongruenz von \mathcal{M}_x und \mathcal{M} . Um zwei vollständig pragmatisch verwandter Modelle zu vereinigen wird daher vorausgesetzt, dass für die Zeitintervalle $T_{\mathcal{M}_x}$ und $T_{\mathcal{M}}$ entweder gilt

$$\left(\tau_{min}(\mathcal{M}) \geq \tau_{min}(\mathcal{M}_x)\right) \wedge \left(\tau_{max}(\mathcal{M}) \leq \tau_{max}(\mathcal{M}_x)\right) \quad (5.22)$$

oder

$$\left(\tau_{min}(\mathcal{M}) \leq \tau_{min}(\mathcal{M}_x)\right) \wedge \left(\tau_{max}(\mathcal{M}) \geq \tau_{max}(\mathcal{M}_x)\right) \quad (5.23)$$

Ist eine dieser Bedingungen erfüllt, so wird überprüft, ob die Urbildräume von \mathcal{M}_x und \mathcal{M} eine Schnittmenge von zwei oder mehr Features besitzen. Ist dies nicht der Fall, wird weiter geprüft, ob aus einer Teilmenge $T\Sigma$ -verwandten Samples eine Art Submodell erzeugt werden kann. Ist dies der Fall wird mit Hilfe von Krippendorffs α überprüft, ob sich das $T\Sigma$ -verwandte Submodell auch tatsächlich auf den gleichen Modellzweck Z bezieht. In allen anderen Fällen wird die vollständige Dekonstruktion abgebrochen.

Modelldifferenzierung. Bleibt eine Rekonstruktion von \mathcal{M}' bzw. eines $T\Sigma$ -verwandten Submodell erfolglos, wird dies als Falsifizierung sowohl von \mathcal{M}' als auch von \mathcal{M}_x und \mathcal{M} gewertet. Dennoch werden diese Modelle in diesem Fall nicht unmittelbar verworfen, sondern zunächst geprüft, ob durch temporale Spaltung Teilmodelle aus \mathcal{M}' extrahiert werden können. Dieser als Modelldifferenzierung bezeichnete Prozess wird durch Clustering von T realisiert¹³. Werden

¹³Wiederum implementiert mittels des R -Packages *KernSmooth*.

5. Konstruktives maschinelles Lernen

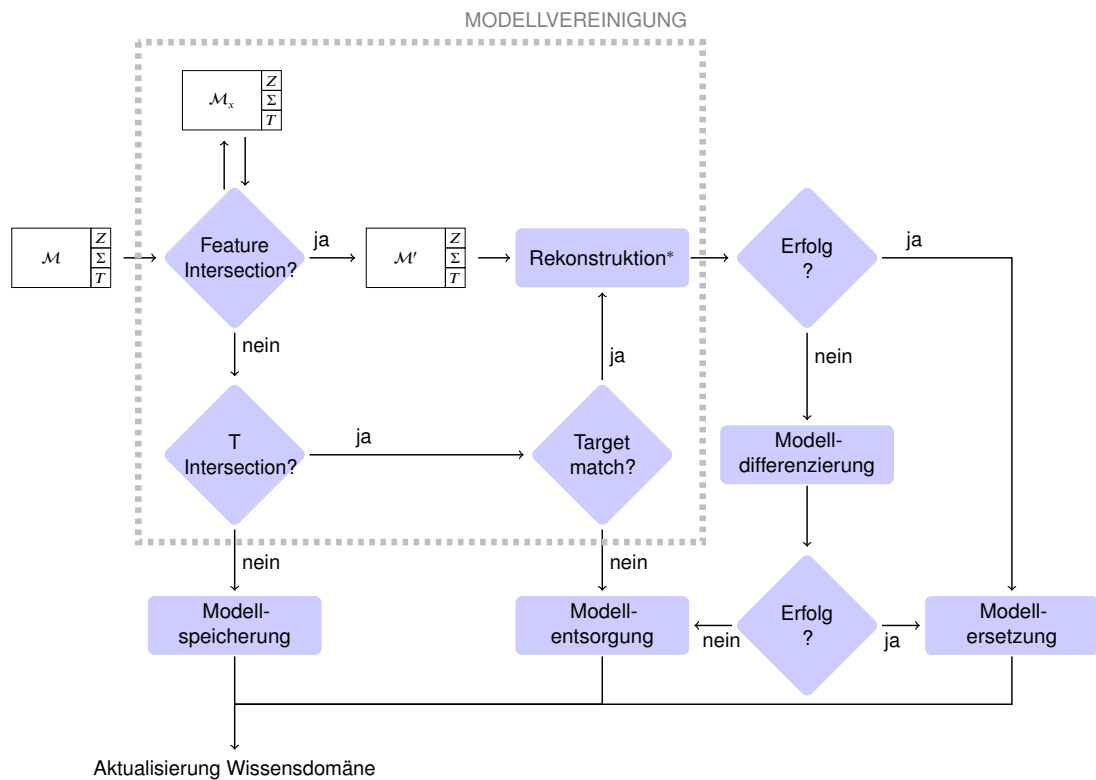


Abbildung 5.6.: Vollständige Dekonstruktion. Die Modellvereinigung erfolgt entweder via Feature-Schnittmenge oder zeitlich übereinstimmender Samples. Im Gegensatz zu den übrigen Dekonstruktionen ist eine Differenzierung oder Löschung von Domänenmodellen möglich.

unter Einhaltung eines Dichteverteilungsschwellenwerts sowie einer Mindestanzahl erforderlicher Samples für maschinelle Modelle zwei oder mehr temporale Cluster identifiziert, so werden mit diesen wiederum Rekonstruktionsprozesse durchgeführt.

Modellentsorgung. Wurden ein oder mehrere Modelle erfolgreich rekonstruiert, so werden diese als neue Modelle in die Wissensdomäne integriert. \mathcal{M}_x wird verworfen. Können weder \mathcal{M}' noch daraus abgeleitete Teilmodelle rekonstruiert werden, so werden sowohl diese als auch \mathcal{M}_x vollständig verworfen; sonstige Dekonstruktionsstrategien werden in diesem Fall nicht angewendet. Außerdem müssen weiter alle Modelle identifiziert und verworfen werden, die von \mathcal{M}_x abhängig sind. Sofern domänen-übergreifende Modellverknüpfungen vorhanden bzw. zugelassen sind, wird in weiteren Domänen analog verfahren.



6

Simulation von Impedanzspektroskopie-Messungen

Dieses Kapitel beschreibt eine repräsentative und skalierbare Nachbildung von Impedanzspektren, wie sie in Messungen an Epithelien oder epithelialen Zellkulturen gewonnen werden. Dies erfordert die Berücksichtigung zahlreicher Einflussfaktoren. Einerseits sind ein geeigneter Ersatzschaltkreis als Modell für die elektrischen Eigenschaften des Gewebes sowie geeignete Wertebereiche für dessen Komponenten zu bestimmen. Dies ist Voraussetzung, um theoretisch zu erwartende Impedanzen für das Gewebe errechnen zu können. Andererseits muss jedoch auch das Verhalten der verwendeten Messapparatur und damit verbundener Messfehler modelliert werden. Darüber hinaus werden zeitliche Zusammenhänge berücksichtigt, um einen pragmatisch definierten Datensatz generieren zu können. Abschließend muss eine ausreichende Übereinstimmung so generierter Datensätze mit gemessenen Daten sichergestellt werden.

6.1. Ersatzschaltkreise für epitheliales Gewebe

Wie im ersten Teil der Arbeit dargestellt, werden in elektrophysiologischen Messungen gemachte Beobachtungen häufig erklärbar, indem die untersuchten Zellen als elektrische Schaltkreise betrachtet werden. Anzahl und Auswahl der erforderlichen Schaltkreis-Komponenten variieren dabei je nach Fragestellung. Ebenso variieren die Wertebereiche dieser Komponenten sowohl von Zelltyp zu Zelltyp als auch unter Einfluss von Wirkstoffen. Die Plausibilität eines angenommenen Modells für eine Untersuchung muss somit stets neu geprüft werden.

Der einfachste überhaupt denkbare Ersatzschaltkreis für ein Epithelgewebe besteht aus einem einzelnen ohmschen Widerstand R und einem parallel dazu geschalteten Kondensator C , also einem so genannten RC-Glied, sowie einem weiteren, vorgeschalteten ohmschen Widerstand (Abb. 6.1a). Dieses im Folgenden als Ersatzschaltbild A bezeichnete Modell ist durch die Tatsache motiviert, dass epitheliales Gewebe neben Diffusionswegen für elektrisch geladene Teilchen grundsätzlich auch nicht-leitende kapazitive Regionen aufweist [84]. Darüber hinaus sind auch subepitheliales Gewebe und Basalmembran als zusätzliche Barriere für den Teilchenfluss repräsentiert. Ein solcher Ersatzschaltkreis erlaubt beispielsweise eine diagnostische Unterscheidung zwischen pathophysiologischen Veränderungen des Epithels und des subepithelialen Gewebes und ist Grundlage der so genannten Ein-Wege-Impedanzspektroskopie [77].

6. Simulation von Impedanzspektroskopie-Messungen

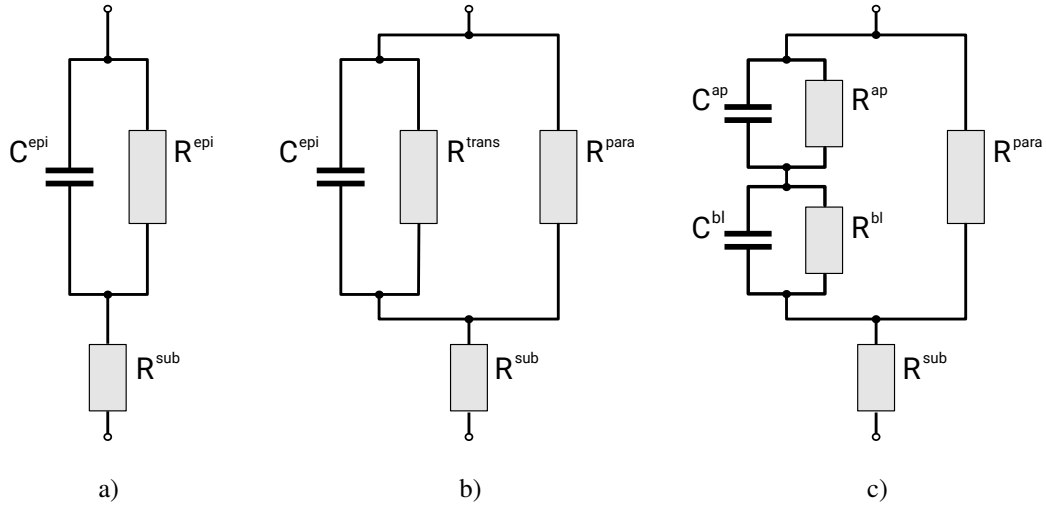


Abbildung 6.1.: Elektrische Modelle für epitheliales Gewebe. a) In Ersatzschaltung *A* setzt sich der Gewebewiderstand R^T aus epithelialelem Widerstand R^{epi} und dem Widerstand des subepithelialen Gewebes, R^{sub} , zusammen; kapazitive Membrananteile sind als epitheliale Kapazität C^{epi} repräsentiert. b) In Ersatzschaltung *B* wird R^{epi} in transzellulären und parazellulären Widerstand R^{trans} bzw. R^{para} unterteilt. c) Ersatzschaltung *C* unterscheidet zwischen apikalen und basolateralen Anteilen von R^{trans} und C^{epi} (R^{ap} und C^{ap} bzw. R^{bl} und C^{bl}).

Der einfachste Ersatzschaltkreis für ein Epithelgewebe, mit dem auch ein Ungleichgewicht zwischen der Diffusion elektrisch geladener Teilchen auf einem transzellulären Weg durch die Epithelzellen hindurch und der Diffusion auf einem parazellulären Weg durch die Zellzwischenräume abgebildet werden kann, erweitert den oben eingeführten Schaltkreis um einen zum RC-Glied parallelen ohmschen Widerstand (Abb. 6.1b). Die Unterscheidung zwischen transzellulärer und parazellulärer Leitfähigkeit in diesem als Ersatzschaltbild *B* bezeichneten Modell ist beispielsweise Grundlage der so genannten Zwei-Wege-Impedanzspektroskopie, mit der diese Größen für diagnostische Zwecke quantifiziert werden können [143].

Der einfachste Ersatzschaltkreis für ein Epithelgewebe, mit dem auch deutliche Differenzen in den Eigenschaften der nach innen und der nach außen gerichteten Gewebeseite abgebildet werden, ersetzt das RC-Glied im oben eingeführten Schaltkreis durch zwei in Reihe geschaltete RC-Glieder (Abb. 6.1c). Dieses als Ersatzschaltbild *C* bezeichnete Modell bildet somit die innere und die äußere Membran der Epithelzellen explizit ab und erlaubt insbesondere deren Eigenschaften zu unterscheiden. Obwohl mit klassischer Impedanzspektroskopie nur schwer bestimmbar, lassen sich die Parameter dieses Modells unter bestimmten Voraussetzungen durch Kombination von impedanzspektroskopischen Messungen und maschinellem Lernen bestimmen [230].

Die beiden zuerst eingeführten Ersatzschaltbilder *A* und *B* lassen sich unmittelbar aus dem Ersatzschaltbild *C* ableiten. Insbesondere gelten folgende Zusammenhänge:

$$R^{trans} = R^{ap} + R^{bl} \quad (6.1)$$

$$\frac{1}{R^{epi}} = \frac{1}{R^{trans}} + \frac{1}{R^{para}} \quad (6.2)$$

$$\frac{1}{C^{epi}} = \frac{1}{C^{ap}} + \frac{1}{C^{bl}} \quad (6.3)$$

Aufgrund dieser Überführbarkeit dient Ersatzschaltbild *C* im Folgenden durchgängig als Grundlage der Modellierung impedanzspektroskopischer Messkurven.

6.2. Eigenschaften epithelialer Zelllinien und ihrer funktionalen Zustände

Mit den Messdaten, die für diese Arbeit zugänglich waren¹, werden aufgrund ihrer prototypischen Eigenschaften die folgenden, vergleichsweise gut charakterisierten Zellkulturlinien modelliert:

- **HT-29/B6.** Beim Zelltyp HT-29 handelt es sich um eine mehrschichtige Drüsenkrebs-Zellkultur aus dem menschlichen Dickdarm [191], die Chlorid sezerniert [275]. Die Zelllinie HT-29/B6 ist ein davon abgeleiteter Monolayer-Klon [137], der statt durch Zugabe von Glukose durch Zugabe von Galaktose ausdifferenziert und ein “dichtes” Epithel bildet (engl. ‘tight’ epithelium). Unter physiologischen Bedingungen, also in Abwesenheit spezifischer Wirkstoffe, wird das Verhältnis von transzellulärer zu parazellulärer Leitfähigkeit, G^{trans}/G^{para} , auf 10:1 geschätzt [78]. Die Größenordnungen weiterer Modellparameter sind zum Teil aus publizierten Messergebnissen abschätzbar (Tab. 6.1).
- **IPEC-J2.** Beim Zelltyp IPEC-J2 handelt es sich eine aus dem Jejunum eines gesunden neugeborenen Schweins gewonnene einschichtige Zellkultur [21], die typische Eigenschaften von Enterozyten aus dem Dünndarm aufweist [207]. Sie zeigt meist einen transepithelialen Widerstand von mindestens 1 k Ω [21, 179]. Die Größenordnungen weiterer Modellparameter sind teils aus publizierten Messergebnissen abschätzbar (Tab. 6.2). Die Eigenschaften von IPEC-J2-Monolayern können stark variieren und hängen unter anderem von den verwendeten Zellkulturfiltern [228] und Wachstumsfaktoren [269] ab. Wegen ihrer Verbreitung als Modell in mikrobiologischen Studien [32] werden hier mit fötalem Rinderserum behandelte Kulturen betrachtet (auch als IPEC-J2/FBS bezeichnet [269]).
- **MDCK I.** Beim Zelltyp “Madin-Darby canine kidney” (MDCK) handelt es sich um eine Zellkulturlinie, die in den 1950er Jahren aus der Niere eines Cockerspaniels entnommen wurde [75]. Diese erwies sich jedoch als heterogen und erlaubte die Isolierung weiterer Subtypen [176]. MDCK I ist einer von zwei grundlegend verschiedenen Subtypen, in die sich die ursprüngliche Zelllinie unterteilen lässt [14, 209]. Unter physiologischen Bedingungen weisen MDCK-I-Zellen einen relativ hohen transepithelialen Widerstand mit Werten zwischen 1.500 und 14.000 Ohm auf [72]. Das Verhältnis von apikaler zu basolateraler Oberfläche wird auf 1:7,6 geschätzt [70]. Größenordnungen weiterer Modellparameter wurden mangels publizierter Messergebnisse aus unveröffentlichten Messungen abgeschätzt (Tab. 6.3).

Neben einem physiologischen Zellzustand ohne Wirkstoff-Einfluss wurden die Funktionsänderungen dieser drei Zelltypen unter dem Einfluss folgender Wirkstoffe modelliert:

- **EGTA.** Bei dem Wirkstoff Ethylenglycol-bis(aminoethylether)-N,N,N',N'-tetraessigsäure, kurz EGTA, handelt es sich um eine organische Säure, die als Chelator zur Bindung von Kalzium und Schwermetallen verwendet wird [56]. Dies ist beispielsweise zur Herstellung von Lösungen nützlich, mit denen physiologische Bedingungen in lebenden Zellen imitiert werden². In Versuchen an Epithelien führt die Reduktion der extrazellulären Konzentration freier Kalzium-Ionen mittels EGTA auf weniger als 1 Mikromol (μM) pro Liter dazu, dass sich die Tight Junctions der parazellulären Barriere öffnen. Dies erhöht die Leitfähigkeit der parazellulären Barriere bzw. reduziert deren Widerstand R^{para} . Gleichzeitig wird eine Reduktion von R^T sowie eine Änderung des Verhältnisses R^{trans}/R^{para} induziert. Hinweise auf Änderungen der kapazitiven Membraneigenschaften sind in der Literatur nicht zu finden.

¹Zur Verfügung gestellt vom Institut für Klinische Physiologie an der Charité Berlin.

²Bezogen auf eine intrazelluläre Konzentration von freien Kalzium-Ionen von etwa 100 Nanomol (nM).

6. Simulation von Impedanzspektroskopie-Messungen

A	Modell		Publizierte Messwerte		Modellierte Werte		Einheit
	B	C	Minimum	Maximum	Minimum	Maximum	
R^{sub}	R^{sub}	R^{sub}	—	$10,9 \pm 1,4$	1,0	30,0	$\Omega \cdot cm^2$
R^{epi}			—	$618,5 \pm 12,8$	152	1.498	$\Omega \cdot cm^2$
	R^{para}	R^{para}	—	2.694 ± 34	154	30.000	$\Omega \cdot cm^2$
	R^{trans}		—	≈ 1.300	153	20.000	$\Omega \cdot cm^2$
		R^{ap}	—	—	7,0	18.870	$\Omega \cdot cm^2$
		R^{bl}	—	—	4,0	19.375	$\Omega \cdot cm^2$
C^{epi}	C^{epi}		—	$4,63 \pm 0,03$	0,9	5,5	$\mu F/cm^2$
		C^{ap}	1,0	—	1,1	8,7	$\mu F/cm^2$
		C^{bl}	1,0	—	3,6	88,2	$\mu F/cm^2$

Tabelle 6.1.: Publierte [137, 143, 6] und modellierte Parameter für die Zelllinie HT-29/B6 unter physiologischen Bedingungen. Modellierte Werte sind aus publizierten abgeleitet.

A	Modell		Publizierte Messwerte		Modellierte Werte		Einheit
	B	C	Minimum	Maximum	Minimum	Maximum	
R^{sub}	R^{sub}	R^{sub}	—	—	0,3	50,0	$\Omega \cdot cm^2$
R^{epi}			1.840	8.500	913	8.567	$\Omega \cdot cm^2$
	R^{para}	R^{para}	2.200	12.900	956	15.000	$\Omega \cdot cm^2$
	R^{trans}		4.169 ± 967	24.700	3.000	20.000	$\Omega \cdot cm^2$
		R^{ap}	—	—	2.516	18.940	$\Omega \cdot cm^2$
		R^{bl}	—	—	2,0	16.978	$\Omega \cdot cm^2$
C^{epi}	C^{epi}		0,9	$1,12 \pm 0,06$	0,6	2,0	$\mu F/cm^2$
		C^{ap}	—	—	1,1	3,5	$\mu F/cm^2$
		C^{bl}	—	—	1,3	9,7	$\mu F/cm^2$

Tabelle 6.2.: Publierte [269, 84] und modellierte Parameter für die Zelllinie IPEC-J2 unter physiologischen Bedingungen. Modellierte Werte sind aus publizierten abgeleitet.

A	Modell		Publizierte Messwerte		Modellierte Werte		Einheit
	B	C	Minimum	Maximum	Minimum	Maximum	
R^{sub}	R^{sub}	R^{sub}	—	—	5,0	25,0	$\Omega \cdot cm^2$
R^{epi}			—	—	401	3.999	$\Omega \cdot cm^2$
	R^{para}	R^{para}	—	—	420	10.000	$\Omega \cdot cm^2$
	R^{trans}		—	—	420	15.000	$\Omega \cdot cm^2$
		R^{ap}	—	—	21,0	14.164	$\Omega \cdot cm^2$
		R^{bl}	—	—	2,0	14.654	$\Omega \cdot cm^2$
C^{epi}	C^{epi}		—	—	0,5	3,5	$\mu F/cm^2$
		C^{ap}	—	—	0,7	6,9	$\mu F/cm^2$
		C^{bl}	—	—	1,6	9,8	$\mu F/cm^2$

Tabelle 6.3.: Modellierte Parameter für die Zelllinie MDCK I unter physiologischen Bedingungen. Modellierte Werte sind aus unveröffentlichten Messungen abgeleitet.

6.3. Berechnung einer theoretischen Impedanz

- **Nystatin.** Bei diesem Wirkstoff handelt es sich um ein Antimykotikum, das ursprünglich aus dem Actinobacterium *Streptomyces noursei* isoliert wurde [102]. Durch Anlagerung von Nystatin an Zellmembranen entstehen darin Poren, die für Ionen durchlässig sind und somit die Gesamtdurchlässigkeit der Membran erhöhen. In Versuchen an Epithelien wird der Wirkstoff entweder auf der apikalen oder basolateralen Epithelseite appliziert, wodurch entsprechend R^{ap} oder R^{bl} reduziert wird. Hinweise auf Änderungen der kapazitiven Membraneigenschaften sind in der Literatur nicht zu finden. Unabhängig von der Applikationsseite ändern sich jedoch in beiden Fällen die Parameter R^{trans} und R^T ; darüber hinaus ändert sich ebenso das Verhältnis zwischen R^{trans} und R^{para} sowie das Verhältnis zwischen apikaler und basolateraler Zeitkonstante (τ_{ap} bzw. τ_{bl}).
- **EGTA+Nystatin.** Werden an einem Epithelgewebe beide zuvor beschriebenen Wirkstoffe gleichzeitig oder in kurzem zeitlichen Abstand voneinander appliziert, so hat dies zur Folge, dass sowohl R^{trans} als auch R^{para} reduziert werden. Außerdem beeinflusst dies auch hier alle von diesen beiden Modellparametern abhängigen Größen (siehe oben).

Insgesamt ergeben sich damit vier mögliche funktionale Zustände pro Zelltyp. Außer für MDCK I lagen jeweils für alle vier Zustände gemessene Daten vor, die einen Abgleich der Modellierung mit Messwerten an Epithelien erlaubten. Für eine weiterführende Parameterübersicht für alle Zelltypen und funktionalen Zustände siehe Anhang B.

6.3. Berechnung einer theoretischen Impedanz

Die unter idealisierten Bedingungen zu erwartende Impedanz eines Schaltkreises lässt sich aus den Impedanzen der einzelnen Schaltkreiskomponenten herleiten. Wie im vorherigen Abschnitt dargestellt, sind die hier relevanten Bauteile ohmsche Widerstände und Kondensatoren. Während die Impedanz Z_R eines ohmschen Widerstands R frequenzunabhängig ist und somit für jede Frequenz f bzw. Kreisfrequenz $\omega = 2\pi f$ stets

$$Z_R(\omega) = R \quad (6.4)$$

gilt, ist die Impedanz Z_C eines Kondensators C frequenzabhängig und gegeben durch

$$Z_C(\omega) = \frac{1}{i\omega C} \quad (6.5)$$

mit $i = \sqrt{-1}$.

Unter Anwendung der kirchhoffschen Sätze folgt für die Impedanz Z_{RC} eines RC-Glieds:

$$\frac{1}{Z_{RC}(\omega)} = \frac{1}{Z_R(\omega)} + \frac{1}{Z_C(\omega)} \quad (6.6)$$

und es ergibt sich weiter für die transepitheliale oder Gesamtimpedanz Z^T :

$$Z^T(\omega) = \frac{R^{epi}(1 - i \cdot \omega \cdot R^{epi} C^{epi})}{1 + (\omega R^{epi} C^{epi})^2} + R^{sub} \quad (6.7)$$

Da Z^T komplexwertig ist, lässt sich die Gesamtimpedanz auch als Summe als Kombination aus Realteil \Re und Imaginärteil \Im darstellen:

$$Z^T(\omega) = \Re(Z^T(\omega)) + i \cdot \Im(Z^T(\omega)) \quad (6.8)$$

6. Simulation von Impedanzspektroskopie-Messungen

wobei gilt:

$$\Re(Z^T(\omega)) = \frac{R^{epi}}{1 + (\omega R^{epi} C^{epi})^2} + R^{sub} \quad (6.9)$$

$$\Im(Z^T(\omega)) = \frac{-\omega(R^{epi})^2 C^{epi}}{1 + (\omega R^{epi} C^{epi})^2} \quad (6.10)$$

Durch Anwendung der Gleichungen 6.1 bis 6.3 lassen sich sowohl $\Re(Z^T)$ als auch $\Im(Z^T)$ für eine gegebene Frequenz ω direkt aus den sechs Modellparametern R^{sub} , R^{para} , R^{ap} , R^{bl} , C^{ap} und C^{bl} berechnen. Für eine vollständige Herleitung siehe Anhang B.

6.4. Modellierung von Messfehlern

Impedanzmessungen an Epithelien werden aus methodischen Gründen in Ussing-Kammern durchgeführt (vgl. Kapitel 4). Neben dem Verhalten des Gewebes wurde daher auch das Verhalten dieses besonderen Messaufbaus als gerätespezifisches Fehlermodell formalisiert, wodurch eine statistische Modellierung der zu erwartenden Abweichungen für $\Re(Z)$ - und $\Im(Z)$ -Werte in Abhängigkeit vom Gesamtwiderstand und der jeweils angewendeten Frequenz möglich wird.

Dieses Fehlermodell basiert auf Messungen, die an Ussing-Kammern des Instituts für Klinische Physiologie der Charité Berlin durchgeführt wurden. Dazu wurden verschiedene Polypropylen-Membranen mit einer Dicke von 0,2 Millimetern und einer Fläche von 0,6 cm² (kapazitiver Wert um 10 pF/cm²) genutzt. Diese wurden durch mechanische Bearbeitung soweit perforiert, dass sie einen Widerstand zwischen 100 und 2.000 Ωcm^2 aufwiesen. Pro Membran bzw. Widerstandswert wurden 30 bis 40 aufeinanderfolgende Spektren aufgenommen.

Um daraus ein formalisiertes Fehlermodell zu erstellen, wurden zunächst die für Kurven einer mittleren Größenordnung ($R^T \approx 525 \Omega\text{cm}^2$) ermittelten Abweichungen in Relation zum Widerstandswert gesetzt und anschließend pro Frequenz die relative Standardabweichungen für $\Re(Z)$ - und $\Im(Z)$ -Werte berechnet (Abb. 6.2). Die Frequenzabhängigkeit der Abweichungen wurde dann mittels frequenzweiser nicht-linearer Regression der Standardabweichungen modelliert.

Für $\Re(Z)$ bei einer Frequenz f wurde eine Fourier-Reihe zweiter Ordnung ($n=2$) als beste Näherung für die Standardabweichung σ_{\Re} bestimmt:

$$\sigma_{\Re}(f) = a_0 + \sum_{i=1}^n a_i \cdot \cos(nwf) + b_i \cdot \sin(nwf) \quad (6.11)$$

wobei mittels Fehlerminimierung $w = 5,353 \cdot 10^{-5}$, $a_0 = 4,848$, $a_1 = -4,11$, $b_1 = -0,8092$, $a_2 = -0,3583$ und $b_2 = 0,2014$ als günstigste Funktionsparameter bestimmt wurden.

Für $\Im(Z)$ bei einer Frequenz f wurde eine polynomiale Funktion vierten Grades ($n=4$) als beste Näherung für die Standardabweichung σ_{\Im} bestimmt:

$$\sigma_{\Im}(f) = a_0 + \sum_{i=1}^n a_i \cdot f^i \quad (6.12)$$

wobei mittels Fehlerminimierung $a_0 = 0,1889$, $a_1 = 0,0002737$, $a_2 = 1,863 \cdot 10^{-9}$, $a_3 = -1,906 \cdot 10^{-13}$ und $a_4 = 2,267 \cdot 10^{-18}$ als günstigste Funktionsparameter bestimmt wurden.

Um auch die Abhängigkeit des Messfehlers vom Gesamtwiderstand des Untersuchungsobjekts zu berücksichtigen, wurde eine vertikale Verschiebung von $\sigma_{\Re}(f)$ und $\sigma_{\Im}(f)$ mittels Exponentialfunktionen eingeführt. Diese simulieren die beobachteten Veränderungen bei 1,3 Hertz:

$$\sigma_{\Re}(1,3\text{Hz}) = 0,636 \cdot R^T^{-0,3278} \quad (6.13)$$

$$\sigma_{\Im}(1,3\text{Hz}) = 8,7008 \cdot R^T^{-0,8689} \quad (6.14)$$

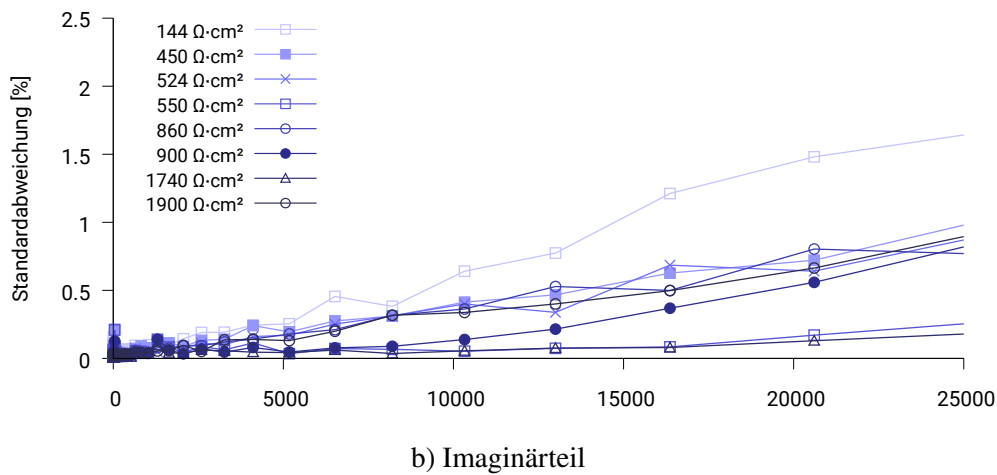
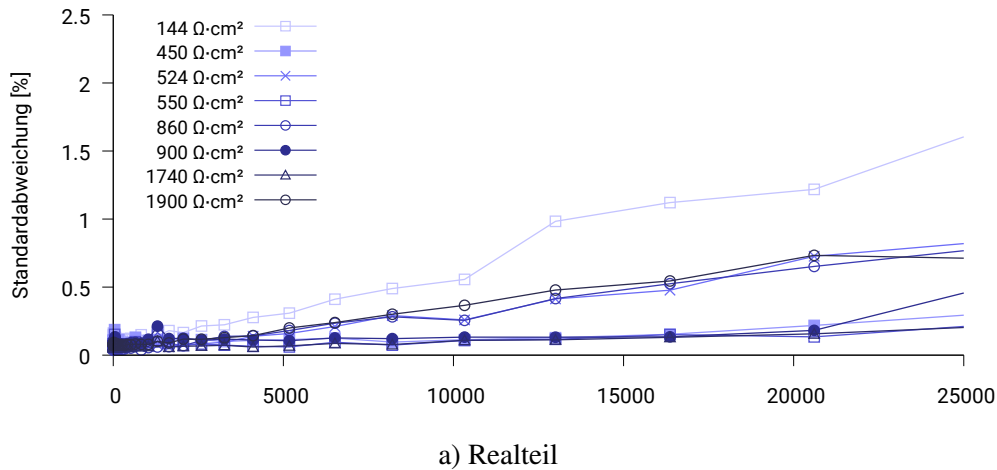


Abbildung 6.2.: Frequenz- und widerstandsabhängige Bias des Messaufbaus. Die Normalverteilung der Abweichung vom Erwartungswert ist für acht Referenz-Messungen mit ohmschen Widerständen zwischen 144 und 1.900 Ωcm^2 getrennt nach a) Real- und b) Imaginärteil komplexwertiger Impedanzen Z dargestellt. Die Messungen wurden vom Institut für Klinische Physiologie der Charité Berlin zur Verfügung gestellt.

Damit kann a_0 in $\sigma_{\Re}(f)$ als Funktion von R^T definiert werden:

$$a_0(R^T) = a_0 + \sigma_{\Re}(1, 3Hz) \quad (6.15)$$

und analog für a_0 in $\sigma_{\Im}(f)$:

$$a_0(R^T) = a_0 + \sigma_{\Im}(1, 3Hz) \quad (6.16)$$

Dieses Modell wurde im weiteren Verlauf dazu genutzt, Messfehler für theoretische Impedanzen nachzubilden, wobei R^T sowie die Frequenz f bzw. die Kreisfrequenz ω jeweils gegeben waren. Alle Abweichungen wurden getrennt nach Real- und Imaginärteil erzeugt und zu den zugehörigen theoretischen Werten addiert (für eine ausführliche Beschreibung siehe Anhang B).

Obwohl theoretisch denkbar, konnten in der Literatur keine Hinweise auf Messfehler als Folge kapazitiver Gewebe- bzw. Membrananteile gefunden werden. Daher werden diese in der vorliegenden Arbeit als vernachlässigbar betrachtet.

6.5. Systematische Synthetisierung von Datensätzen

Um für jede Zelllinie und jeden Zustand systematisch Impedanzspektren zu synthetisieren, wurden pseudorandomisierte Parameterwerte für Ersatzschaltkreis C (Abb. 6.1c) erzeugt. Unterschieden wurde dabei zwischen Parametern, die innerhalb einer Reihe aufeinanderfolgender Messungen veränderlich sind, und Parametern, die als konstant angenommen werden. R^{para} , R^{ap} und R^{bl} wurden in diesem Sinne als veränderlich betrachtet, da sie nicht nur natürlichen Schwankungen unterliegen, sondern sich auch aufgrund von Wirkstoffzugaben innerhalb einer zusammenhängenden Messreihe ändern können. Umgekehrt können R^{sub} , C^{ap} und C^{bl} zwar bei gleichem Zelltyp und gleichem funktionalem Zustand innerhalb ihrer Parametergrenzen variieren, für eine zusammenhängende Messreihe wurden sie jedoch als Konstanten modelliert.

Auf dieser Basis wurden jeweils Messreihen aus Impedanzspektren erzeugt, wobei die Parameterintervalle von R^{sub} , C^{ap} und C^{bl} mittels definierter Schrittweite variiert wurden; Intervallgrenzen für C^{epi} dienten dabei ausschließlich als statische Grenzen bzw. zur Überprüfung der Zulässigkeit der C^{ap} - und C^{bl} -Werte. Innerhalb einer zusammenhängenden Messreihe wurde außerdem der Parameterwert für R^{epi} pseudorandomisiert. Die Werte der von R^{epi} abhängigen, veränderlichen Parameter R^{para} , R^{ap} und R^{bl} wurden innerhalb der jeweiligen Parametergrenzen abgeleitet. Das algorithmische Vorgehen im Einzelnen ist in Anhang B dargestellt.

Motiviert ist dieses Vorgehen dadurch, dass die drei Parameterwerte für Modell A (Abb 6.1a) mittels Messung direkt bestimmt werden können und dadurch bestimmte Parametergrenzen als weniger fehlerbehaftet gelten als diejenigen für Modell B (Abb. 6.1b). Messwerte für Modell B wiederum wurden als zuverlässiger betrachtet als diejenigen für Modell C (Abb. 6.1c), da diese nur indirekt und unter Verwendung mehrstufiger Schätzverfahren ermittelt werden können.

Für die pragmatische Eigenschaft T einer erzeugten Kurve wurde darüber hinaus ein zeitlicher Zusammenhang entsprechend dieses Synthesemusters abgebildet. Dies wurde durch die inkrementellen Faktoren $\epsilon > 1$ innerhalb einer Messreihe sowie $\delta \gg \epsilon$ zwischen zwei Messreihen realisiert. Als Zweck Z kommt zwar grundsätzlich jeder der Modellparameter in Betracht, realistischer ist es jedoch, zunächst das Fehlen eines solchen gegebenen Zielparameters anzunehmen. Das Modellsjekt Σ wurde für alle synthetisierten Kurven als identisch angenommen (UID 0).

Aus jedem modelliertem Datensatz wurden 25.000 Samples gezogen, wodurch sich eine Datenbasis von insgesamt 275.000 Impedanzspektren ergab (Tab. 6.4). Diese wurde für die Analysen in Kapitel 7 und 8 in einen Trainingsdatensatz ($n=200.000$) und einen Testdatensatz ($n=75.000$) aufgeteilt. Eine weiterführende Charakterisierung der Datenbasis ist in Anhang C zu finden.

Zustand	Zellkulturlinie			gesamt
	HT-29/B6	IPEC-J2	MDCK I	
Kontrolle	25.000	25.000	25.000	75.000
Nystatin	25.000	25.000	25.000	75.000
EGTA	25.000	25.000	25.000	75.000
EGTA+Nystatin	25.000	25.000	—	50.000
gesamt	100.000	100.000	75.000	275.000

Tabelle 6.4.: Anzahl der modellierten Impedanzspektren pro Zellkulturlinie und funktionalem Zustand.

6.6. Abgleich modellierter und gemessener Daten

Um eine realitätsnahe Nachbildung impedanzspektroskopischer Messdaten sicherzustellen, wurde die Übereinstimmung zwischen modellierten und gemessenen Spektren für jede Zelllinie und

6.6. Abgleich modellierter und gemessener Daten

Zustand	Parameter	Zellkulturlinie			Minimum
		HT-29/B6	IPEC-J2	MDCK I	
Kontrolle	R^{sub}	100,0	91,7	97,1	91,7
	R^{epi}	94,3	96,7	88,5	88,5
	C^{epi}	100,0	91,7	91,1	91,1
Nystatin	R^{sub}	100,0	90,9	100,0	90,9
	R^{epi}	100,0	96,6	93,3	93,3
	C^{epi}	94,7	93,3	93,8	93,3
EGTA	R^{sub}	100,0	85,7	100,0	85,7
	R^{epi}	88,9	100,0	93,7	88,9
	C^{epi}	100,0	100,0	92,5	92,5
EGTA+	R^{sub}	100,0	100,0	—	100,0
Nystatin	R^{epi}	90,9	100,0	—	90,9
	C^{epi}	100,0	93,3	—	93,3
Minimum		88,9	85,7	88,5	85,7

Tabelle 6.5.: κ_B -Werte pro Zellkulturlinie und funktionalem Zustand [%].

jeden funktionalen Zustand überprüft. Eine besondere Herausforderung stellte in diesem Zusammenhang der gerätespezifische Messfehler dar (vgl. Abschnitt 6.4), da aufgrund dessen selbst bei identisch modellierten Parameterwerten des Ersatzschaltkreises kein geometrisch deckungsgleiches Spektrum vorliegt. Ein geometrisch motivierter Abgleich wird zusätzlich erschwert durch das asymmetrische Größenverhältnis zwischen einer Menge von jeweils 25.000 modellierten³ und einer Menge von maximal 200 gemessenen Impedanzspektren.

Stattdessen wurde die Übereinstimmung zwischen modellierten und gemessenen Daten hier indirekt anhand der drei Parameter des Ersatzschaltkreises A (R^{sub} , R^{epi} , C^{epi}) evaluiert. Diese sind zwar für gemessene Spektren naturgemäß nicht bekannt, können jedoch mit verschiedenen zuverlässigen Verfahren abgeschätzt werden (vgl. [230, 233]). Um einen eventuellen Verfahrensbias zu minimieren, wurden dazu je drei Verfahren parallel angewendet und die Differenzen zwischen den Zielwerten in speziellen zweidimensionalen Plots aufbereitet. Durch Binning dieser Plots und mithilfe eines relativen Fehlermaßes κ_B wurde dann die Kongruenz zwischen den Plots für modellierte und für gemessene Daten quantifiziert. Für eine detaillierte Beschreibung des Fehlermaßes und des Plot-Verfahrens siehe Anhang B (Abschnitt B.2).

Für die untersuchten 11 Teil-Datensätze der Zellkulturlinien bzw. -zustände (vgl. Tab. 6.4) wurde für die Parameter R^{sub} , R^{epi} und C^{epi} eine ausreichende Überlappung zwischen den zweidimensionalen Plots modellierter und gemessener Daten erreicht. In allen der 33 Abgleiche wurde ein κ_B -Wert von mindestens 85 Prozent erreicht (Tab. 6.5). Bei mehr als der Hälfte der Abgleiche betrug der κ_B -Wert mehr als 95 Prozent. Ein κ_B -Wert von 100 Prozent impliziert in diesem Zusammenhang, dass die jeweils gemessenen Impedanzspektren bezüglich R^{sub} , R^{epi} oder C^{epi} vollständig durch die abzugleichende Menge modellierter Spektren repräsentiert werden. Modellerte Impedanzspektren, die in dem Plot keine Überlappung mit den gemessenen aufweisen, sind somit überschüssig und gehen nicht in den κ_B -Wert ein.

Eine weiterführende Darstellung der Ergebnisse und Plots dieses Abgleichs für alle Zelllinien und funktionalen Zustände ist in Anhang B zu finden (Abschnitt B.3ff).

³Um ein für dieses Test-Sample geeignetes maschinelles Lernverfahren trainieren zu können, wurde zusätzlich ein unabhängiges Training-Sample aus 50.000 modellierten Spektren verwendet. Auswertung und Abgleich der Modellierungen mit Messdaten beziehen sich ausschließlich auf das Test-Sample.

6. Simulation von Impedanzspektroskopie-Messungen

7

Konzeptuelles Wissen für Impedanzspektren

In diesem Kapitel wird untersucht, inwieweit sich Impedanzspektren mittels konstruktivistischen maschinellen Lernens der zugehörigen epithelialen Zelllinie zuordnen lassen. Ziel ist die korrekte Unterscheidung von 75.000 zu testenden Spektren der Zelllinien HT-29/B6, IPEC-J2 und MDCK I. Dazu wird das konstruktivistische maschinelle Lernen einmal mit und einmal ohne vorgegebenem Zielwert auf die Trainingsdaten der Datenbasis (n=200.000) angewendet. Soweit nicht anders angegeben, werden dazu die durch Testläufe ermittelten Lernparameter in Tab. 7.1 angewendet. Für eine weiterführende Beschreibung der Datenbasis siehe Anhang C.

	Parameter	Wert
Konstruktion	Max. kategoriale Komplexität κ_k	5
	Min. Kategoriengröße	0,1
	Max. Fehler	0,25
Komplexitäts- reduktion	Max. Modellkomplexität κ	10
	Filter-Grenze λ_x	50
	Filter-Grenze λ_y	15.000
	Maximalreduktion erzwingen	Ja
Rekonstruktion	Min. Accuracy	0,8
	Samplegröße für Reliabilitätsprüfung	0,1
	Schwellenwert für Reliabilitätsprüfung α_{min}	0,8
	Redundanzvermeidung	Ja
Dekonstruktion	Umgang mit neuen Modellen	integrativ
	Abbruchbedingung	minimal
	Temporale Erweiterungstoleranz δT_{max}	1
	Toleranz bei vollständiger Dekonstruktion	0,1
	ΣZ -verwandte Modelle immer vereinigen	Ja
	Schwache Reliabilität erlaubt	Nein

Tabelle 7.1.: Übersicht der zur Exploration konzeptuellen Wissens verwendeten konstruktivistischen Lernparameter. Für eine weiterführende Beschreibung der Parameter siehe Anhang A.

7.1. Exploration einer konzeptuellen Wissensdomäne

Maschinelle Modelle, die Klassifikationsaufgaben repräsentieren, werden hier als konzeptuelles Wissen bzw. konzeptuelle maschinelle Modelle bezeichnet. Wie in Kapitel 5 beschrieben, wird durch Anwendung konstruktivistischen maschinellen Lernens auf eine gegebene Datenbasis eine Menge geordneter und gegebenenfalls miteinander verknüpfter Modelle erzeugt. Handelt es sich dabei um eine Menge konzeptueller Modelle wird diese analog zum menschlichen Wissen als konzeptuelle Wissensdomäne bezeichnet.

Wird das konstruktivistische maschinelle Lernen erstmalig auf Trainingsdaten angewendet, existiert zunächst noch keine konzeptuelle Wissensdomäne. Allerdings können Domänen, die auf anderen Daten bzw. in früheren Durchläufen erzeugt wurden, genutzt werden. Um etwa eine Konstruktion von maschinellen Modellen einer höheren Wissens Ebene aus maschinellen Modellen einer niedrigeren Wissens Ebene zu erreichen, kann eine bereits existierende Wissensdomäne einbezogen werden. Dadurch können insbesondere Bezüge zwischen den Modellen der gleichen Ebene identifiziert und zur Abstraktion von Modellen einer höheren Ebene genutzt werden.

Auf dieser Grundlage wird hier mittels eines Exploration genannten Prozesses eine konzeptuelle Wissensdomäne erzeugt, die als Vorwissen für eine spätere Adaption konzeptuellen Wissens genutzt werden kann. Hierfür werden Features bzw. Feature-Sets einer gegebenen Datenbasis getrennt voneinander und konsekutiv einem konstruktivistischen maschinellen Lernen unterworfen. Jedes der Feature-Sets wird dabei ohne vorgegebenen Zielparameter bzw. Modellzweck verwendet, was dazu führt, dass für jeden untersuchten Lernblock automatisch ein Konstruktionsprozess ausgeführt wird. Dieses sequentielle, quasi-unüberwachte Lernen führt im Ergebnis zu einer selbstständigen Exploration der Trainingsdaten und erhöht sowohl die Wahrscheinlichkeit, pragmatisch verwandte Modelle zu identifizieren, als auch deren wahrscheinliche Anzahl.

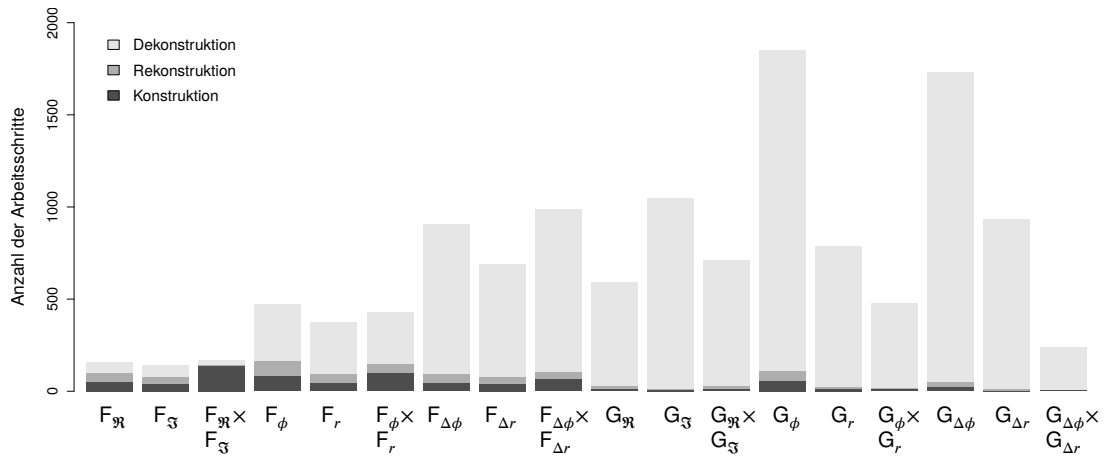
7.1.1. Explorationsverlauf

Die Exploration der Trainingsdaten erfolgt in 18 aufeinander aufbauenden Teilschritten, in denen jeweils eine Teilmenge der Eingabefeatures verwendet wird. Die für die Datenbasis definierten 18 Feature-Sets reflektieren dabei sowohl die kartesische Darstellung der Impedanzen (Feature-Sets $F_{\mathfrak{R}}$, $F_{\mathfrak{I}}$ sowie $F_{\mathfrak{R}} \times F_{\mathfrak{I}}$) als auch die polare (Feature-Sets F_{ϕ} , F_r sowie $F_{\phi} \times F_r$). Weiter werden aus der Polarkoordinaten-Darstellung Feature-Sets abgeleitet, die Differenzen zwischen den Features eines einzelnen Impedanzspektrums beschreiben ($F_{\Delta\phi}$, $F_{\Delta r}$, $F_{\Delta\phi} \times F_{\Delta r}$). Zusätzlich werden für diese 6 Feature-Sets statistische Eigenschaften berechnet und als 6 Sets globaler Features eines Impedanzspektrums zusammengefasst ($G_{\mathfrak{R}}$, $G_{\mathfrak{I}}$, $G_{\mathfrak{R}} \times G_{\mathfrak{I}}$, G_{ϕ} , G_r , $G_{\phi} \times G_r$, $G_{\Delta\phi}$, $G_{\Delta r}$, $G_{\Delta\phi} \times G_{\Delta r}$). Für eine detaillierte Beschreibung der Feature-Sets siehe Anhang C.

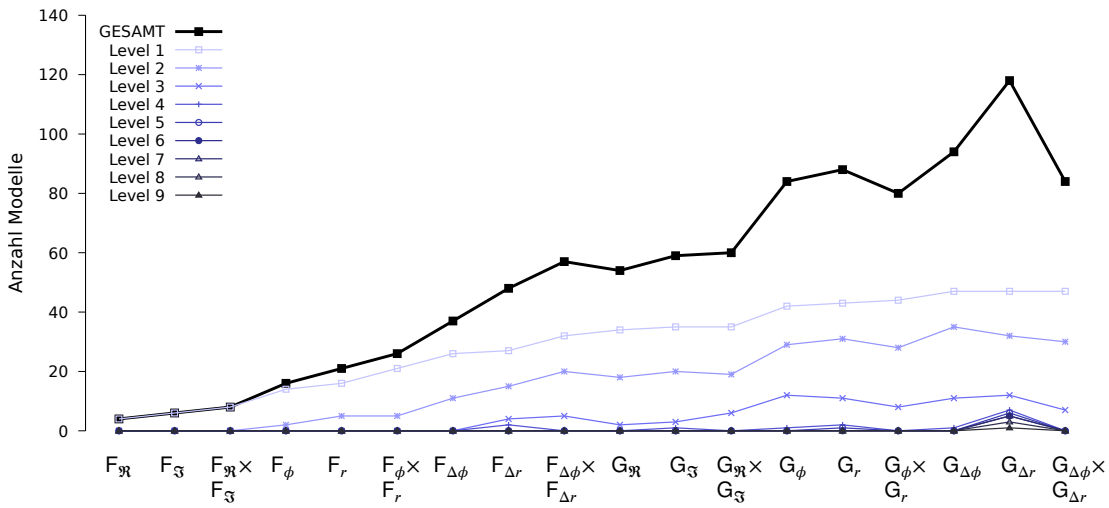
Pro Teilschritt wird zwar nur eine Teilmenge der Eingabefeatures verwendet, die Anzahl der Samples entspricht jedoch der gesamten Trainingsdaten ($n=200.000$). Pro Block wird zunächst ein Konstruktionsprozess durchgeführt, auf den im Erfolgsfall ein Rekonstruktionsprozess und gegebenenfalls ein Dekonstruktionsprozess folgt. Die Menge der im Rahmen dieser Prozesse ausgeführten Operationen variiert dabei von Feature-Set zu Feature-Set (Abb. 7.1a). Die Anzahl der Modelle der Wissensdomäne nimmt im Verlauf der Exploration von Feature-Set zu Feature-Set stetig zu (Abb. 7.1b). Veränderungen innerhalb eines Feature-Set-Durchlaufs sind hier zugunsten der Übersichtlichkeit nicht im Explorationsverlauf abgebildet.

Die Hierarchie der Domäne verändert sich während der 18-teiligen Exploration kontinuierlich. So enthält diese etwa während der Exploration der ersten drei Feature-Sets, die eine Darstellung der Impedanzen in kartesischen Koordinaten repräsentieren, ausschliesslich maschinelle Modelle der Wissens Ebene 1. Danach werden auch Modelle der Ebenen 2, 3 und 4 erzeugt. Im weiteren Verlauf werden die Modelle dieser höheren Ebenen teils wieder gelöscht. Die durchschnittliche Komplexität der Modelle der Domäne steigt anfangs und fällt danach wieder ab (Abb. 7.1c).

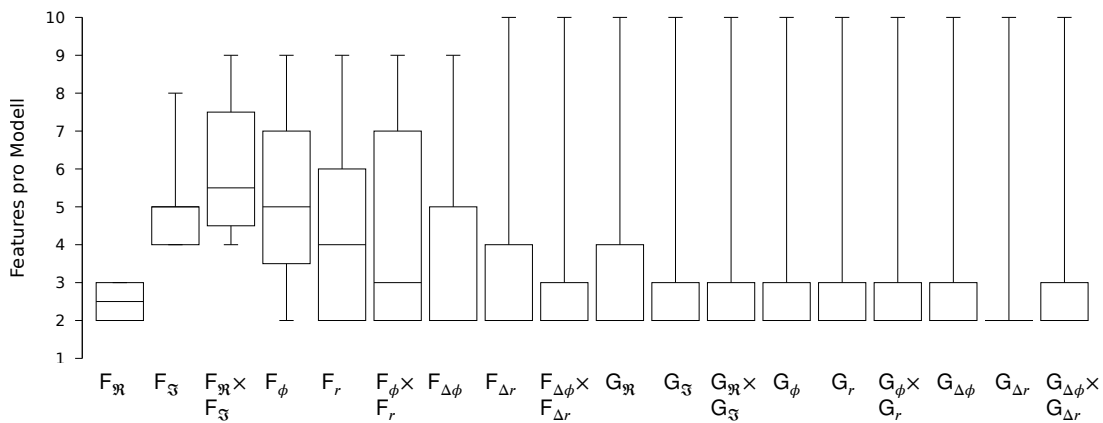
7.1. Exploration einer konzeptuellen Wissensdomäne



(a) Ausgeführte Arbeitsschritte



(b) Größe der Wissensdomäne



(c) Komplexität κ der Domänenmodelle

Abbildung 7.1.: Verlauf der Exploration der konzeptuellen Wissensdomäne. Die Exploration erfolgt in 18 Teilschritten, dargestellt ist jeweils der Zustand am Ende eines Feature-Set-Durchlaufs.

7.1.2. Charakterisierung der explorierten Wissensdomäne

Die nach den 18 Teilschritten der Exploration vorliegende konzeptuelle Wissensdomäne setzt sich aus insgesamt 84 maschinellen Modellen zusammen. Diese verteilen sich in einer hierarchischen Anordnung über die Wissensebenen 1 bis 3, wobei Modelle der Ebenen 2 und 3 auf Modellen der Ebenen 1 und 2 aufbauen. Die Anzahl der Modelle pro Ebene nimmt dabei ab, je höher sich die Ebene innerhalb der Hierarchie befindet (Abb. 7.2a). 47 der 84 Domänenmodelle befinden sich auf Ebene 1, 30 Modelle auf Ebene 2 und sieben auf Ebene 3 (Tab. 7.2).

70 Modelle der Domäne sind 2-fach-Klassifikatoren, d.h. sie realisieren eine Zuordnung zu einer von zwei Klassen. Weitere 11 maschinelle Modelle sind 3-fach-Klassifikatoren, die die jeweilige Eingabe einer von drei Klassen zuordnen. Drei Modelle realisieren eine Unterscheidung von vier Klassen. 33 der 84 Modelle der Domäne liegen k-Means-basierte Cluster zugrunde, d.h. die dadurch realisierte Unterscheidung von zwei, drei oder vier Klassen wurde durch Identifikation von zwei, drei oder vier Clustern mittels k-Means konstruiert. Im Gegensatz dazu wurden 34 Modelle von einer selbstorganisierenden Karte konstruiert. Die übrigen 17 Modelle wurden durch Vereinigung je eines k-Means- und eines SOM-basierten Modells erzeugt.

Als Eingabefeatures eines Modells dienen jeweils die Ausgaben der Modelle der darunterliegenden Ebene, wobei die Ebene 0 die Features der Datenbasis repräsentiert. Im Rahmen der Exploration wurde die Modellkomplexität auf $\kappa \leq 10$ begrenzt, so dass die Modelle der konzeptuellen Wissensdomäne maximal zehn Eingabefeatures besitzen können. Innerhalb der explorierten Domäne nutzt nur ein Modell tatsächlich zehn Eingabefeatures, während 61 der 84 Modelle lediglich zwei Eingabefeatures verwenden. Somit liegt der Median von κ bei 2. Die sieben maschinellen Modelle der Ebene 3 nutzen alle je zwei Eingabefeatures.

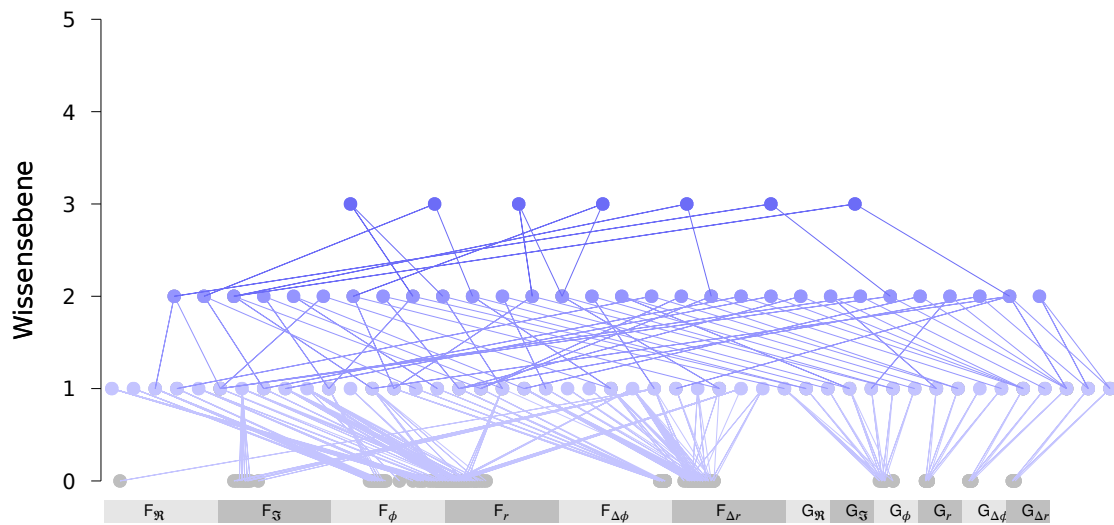
Abb. 7.2a zeigt in einer Frontalaufsicht die Menge aller Verknüpfungen zwischen den maschinellen Modellen der Wissensdomäne. Verbindungen zur Ebene 0 visualisieren dabei die Nutzung der Features der zur Exploration genutzten Trainingsdaten. Durch diese Konnektom-Darstellung wird ersichtlich, dass die Features der Trainingsdaten von den Modellen der explorierten Domäne nicht gleichmäßig genutzt werden. So werden Features der Sets F_ϕ , F_r und $F_{\Delta r}$ häufig genutzt. Umgekehrt nutzt kein einziges Modell ein Feature der Sets $G_{\mathfrak{R}}$ und $G_{\mathfrak{S}}$.

Eine Analyse der temporalen Gültigkeit der Domänenmodelle ergibt, dass diese im Mittel 13,1 Prozent der Gesamtzeitspanne der Datenbasis abdecken. Wie jedoch Abb. 7.2b sowie Abb. 7.3a-c zeigen, variiert diese Abdeckung und somit die Gültigkeit des durch die Modelle repräsentierten Wissens teils stark. Auf Ebene 1 decken sechs der 47 Domänenmodelle mehr als 50 Prozent der Zeitspanne ab und repräsentieren damit im Verhältnis zur Datenbasis zeitlich weit gefasstes Wissen. Während jedoch auf Ebene 1 die zeitliche Abdeckung im Durchschnitt 18,8 Prozent und im Maximum 86,4 Prozent beträgt, ist die Abdeckung der darüberliegenden Ebenen geringer. So decken die 30 Modelle der Ebene 2 im Durchschnitt 5,6 Prozent und die sieben Modelle der Ebene 3 im Durchschnitt 6,6 Prozent der Gesamtzeitspanne der Datenbasis ab. Die Ebenen 2 und 3 bilden somit zwar abstrahiertes, aber zeitlich eng begrenztes Wissen ab.

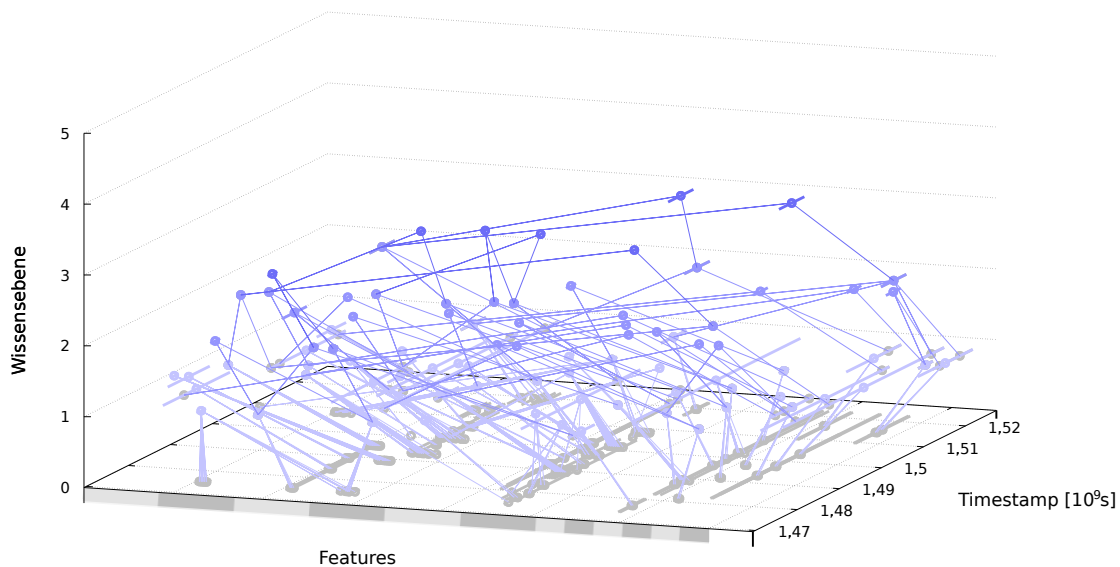
	2-fach-Klassifikatoren	3-fach-Klassifikatoren	4-fach-Klassifikatoren	Gesamt
Ebene 1	45	1	1	47
Ebene 2	19	10	1	30
Ebene 3	6	0	1	7
Domäne	70	11	3	84

Tabelle 7.2.: Übersicht über die Modellzwecke der explorierten konzeptuellen Wissensdomäne.

7.1. Exploration einer konzeptuellen Wissensdomäne



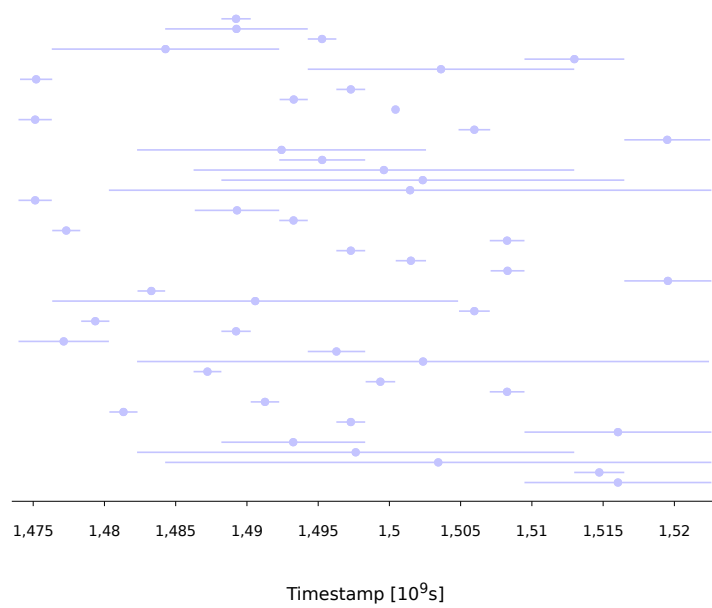
a) Modellhierarchie



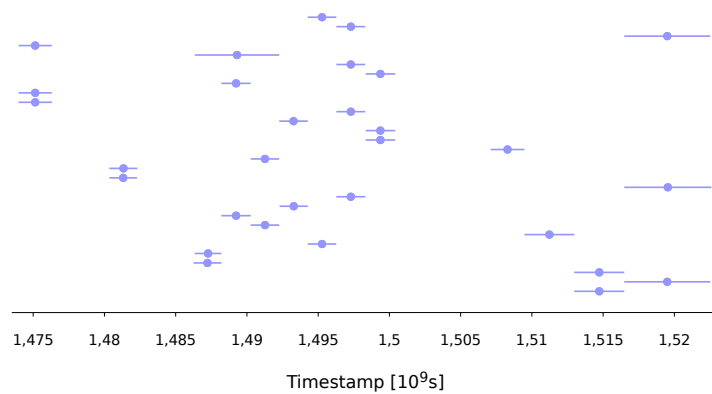
b) Dreidimensionale Aufsicht

Abbildung 7.2.: Überblick über die explorierte konzeptuelle Wissensdomäne. Modelle der gleichen Wissensebene sind jeweils im gleichen Farbton dargestellt. Die Feature-Sets der Trainingsdaten, die sich auf Ebene 0 befinden, sind entlang der x-Achse gekennzeichnet. a) Modellhierarchie. In dieser Frontal-Aufsicht werden Verknüpfungen zwischen den Domänenmodellen durch vertikale Geraden repräsentiert; die Zeitachse wird in dieser Darstellung vernachlässigt. b) Dreidimensionale Aufsicht. Zusätzlich zur Darstellung als Modellhierarchie ist die zeitliche Ausdehnung jedes Modells auf der z-Achse repräsentiert.

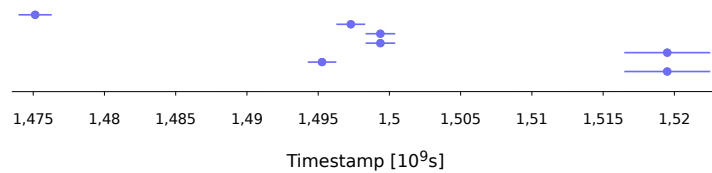
7. Konzeptuelles Wissen für Impedanzspektren



a) Zeitausdehnung der Modelle der Ebene 1 (n=47)



b) Zeitausdehnung der Modelle der Ebene 2 (n=30)



c) Zeitausdehnung der Modelle der Ebene 3 (n=7)

Abbildung 7.3.: Temporale Gültigkeit der Modelle der konzeptuellen Wissensdomäne. Dargestellt ist die zeitliche Ausdehnung der 47 Modelle der Ebene 1 (a), der 30 Modelle der Ebene 2 (b) sowie der sieben Modelle der Ebene 3 (c). Die Modelle sind jeweils anhand ihrer ID von oben nach unten aufsteigend in 1er-Schritten auf einer einheitslosen y-Achse angeordnet.

7.2. Adaption konzeptuellen Wissens

Im vorherigen Abschnitt wurden Impedanzspektren ohne einen gegebenen Zielparameter analog zu einem unüberwachten Lernen exploriert. Hier wird dagegen im Rahmen einer Adaption ein Trainingsdatensatz ($n=200.000$) genutzt, um einen vorgegebenen Zielparameter analog zu einem überwachten Lernen zu erlernen. Der Lernerfolg wird anschließend mit einem Testdatensatz ($n=75.000$) überprüft. Anders als bei klassischem überwachten Lernen erfolgt dies jedoch blockweise und beinhaltet zusätzlich einen Dekonstruktionsprozess.

Darüber hinaus wird die konzeptuelle Wissensdomäne zur Abstraktion mittels Feature-Transformation genutzt. Dazu wird ein Lernblock bei erfolgloser Rekonstruktion in Features der nächsthöheren Wissensebene transformiert und zusammen mit dem definierten Zielparameter einem erneuten Rekonstruktionsprozess unterworfen. Die zur Transformation verwendeten Modelle der Zielebene werden dabei anhand der pragmatischen Eigenschaften identifiziert. Schlägt die Rekonstruktion fehl, wird der Lernblock solange auf die nächsthöhere Ebene transformiert, bis die höchste Ebene der Domäne erreicht ist.

Als Beispiel für die Adaption konzeptuellen Wissens wird die Unterscheidung impedanzspektroskopischer Messungen hinsichtlich der zugehörigen epithelialen Zelllinie untersucht. Dies entspricht einer Klassifikationsaufgabe, bei der Samples in eine von drei Klassen eingeordnet werden sollen. Der zugehörige Zielwert wird als kategoriale Variable L kodiert:

$$L = \begin{cases} 1 & \text{falls "HT-29/B6"} \\ 2 & \text{falls "IPEC J-2"} \\ 3 & \text{falls "MDCK I"} \end{cases} \quad (7.1)$$

Ziel dieses Adaptionsprozesses ist die Abbildung der Klassifikationsaufgabe in Form eines oder mehrerer maschineller Modelle innerhalb einer konzeptuellen Wissensdomäne. Im Gegensatz zu klassischem überwachten Lernen bildet eine solche Repräsentation sowohl eine temporale als auch eine intersubjektive Gültigkeit explizit ab. Dabei soll eine Accuracy von mindestens 0.75 sowie ein Intercoder-Reliabilitätskoeffizient von mindestens 0.75 erreicht sein. Weiter soll die Modellkomplexität κ die Anzahl der verwendeten Frequenzen ($n=42$) nicht übersteigen. Tab. 7.3 gibt einen Überblick über die verwendeten Lernparameter.

	Parameter	Wert
Komplexitätsreduktion	Max. Modellkomplexität κ	42
	Filter-Grenze λ_x	350
	Filter-Grenze λ_y	15.000
	Maximalreduktion erzwingen	Ja
Rekonstruktion	Min. Accuracy	0,75
	Samplegröße für Reliabilitätsprüfung	0,1
	Schwellwert für Reliabilitätsprüfung α_{min}	0,75
	Redundanzvermeidung	Ja
Dekonstruktion	Umgang mit neuen Modellen	integrativ
	Abbruchbedingung	minimal
	Temporale Erweiterungstoleranz δT_{max}	1
	Toleranz bei vollständiger Dekonstruktion	0,1
	ΣZ -verwandte Modelle immer vereinigen	Ja
	Schwache Reliabilität erlaubt	Nein

Tabelle 7.3.: Übersicht der zur Adaption konzeptuellen Wissens verwendeten konstruktivistischen Lernparameter. Für eine ausführliche Beschreibung der Parameter siehe Anhang A.

7. Konzeptuelles Wissen für Impedanzspektren

7.2.1. Adaptionsverlauf

Im Folgenden wird dargestellt, wie die gegebene Klassifikationsaufgabe durch das konstruktivistische maschinelle Lernen adaptiert wird. Die Adaption wird blockweise sowie unter Verwendung der explorierten konzeptuellen Wissensdomäne durchgeführt.

Tab. 7.4 zeigt, dass die identifizierten Lernblöcke des jeweiligen Blocks (Wissensebene 0) durchgängig sowohl von Random Forests als auch neuronalen Netzen mit ausreichender Accuracy gelernt werden. Allerdings liegt die Inter-Coder-Reliabilität nur für die Lernblöcke der Blöcke 10 und 19 über dem gesetzten Schwellenwert von 0,75. Beide maschinellen Modelle werden unverändert in die Wissensdomäne integriert, da weder eine vollständige noch eine ΣZ -Dekonstruktion durchgeführt wird.

Für diejenigen Blöcke, für die keine Rekonstruktion möglich war, werden die Samples anschließend von der Ebene 0 auf die Ebene 1 transformiert und ein erneuter Rekonstruktionsprozess angestoßen. Eine solche Rekonstruktion ist für Block 18 erfolgreich. Dieser wurde in einen Block der Ebene 1 transformiert, der sich aus zehn Features zusammensetzt und 10.000 Samples umfasst. Der transformierte Block konnte sowohl von einem Random Forest als auch von einem neuronalen Netz erfolgreich und mit einer Inter-Coder-Reliabilität von 0,78 rekonstruiert werden.

Für die übrigen auf Ebene 1 transformierten Blöcke war keine Rekonstruktion möglich. Die anschließenden Transformationen auf Ebene 2 und 3 waren durchgängig nicht erfolgreich und sind zugunsten der Übersichtlichkeit nicht in Tab. 7.4 dargestellt.

Block	Rekonstruktion (Ebene 0)				Vollständige/ ΣZ -Dekonstruktion				Rekonstruktion (Ebene 1)			
	κ	Größe	RF/NN	α	κ	Größe	RF/NN	α	κ	Größe	RF/NN	α
1	32	10.000	+/+	0,53	—	—	—	—	3	9.956	-/-	—
2	41	10.000	+/+	0,50	—	—	—	—	4	9.833	+/-	—
									2	10.000	-/-	—
3	35	10.000	+/+	0,53	—	—	—	—	4	9.666	+/-	—
									3	10.000	-/-	—
4	41	10.000	+/+	0,56	—	—	—	—	4	9.833	+/-	—
									3	10.000	-/-	—
5	34	10.000	+/+	0,47	—	—	—	—	7	9.833	+/-	—
									6	10.000	-/-	—
6	32	10.000	+/+	0,53	—	—	—	—	8	10.000	-/+	—
7	36	10.000	+/+	0,46	—	—	—	—	11	9.666	-/+	—
									10	10.000	-/-	—
8	22	10.000	+/+	0,53	—	—	—	—	14	10.000	-/+	—
9	42	10.000	+/+	0,62	—	—	—	—	13	10.000	-/+	—
10	27	10.000	+/+	0,83	—	—	—	—	—	—	—	—
11	42	10.000	+/+	0,53	—	—	—	—	13	10.000	-/+	—
12	40	10.000	+/+	0,47	—	—	—	—	15	10.000	-/+	—
13	28	10.000	+/+	0,47	—	—	—	—	10	9.833	-/+	—
									6	10.000	-/-	—
14	36	10.000	+/+	0,57	—	—	—	—	10	9.833	-/+	—
									9	10.000	-/-	—
15	36	10.000	+/+	0,61	—	—	—	—	8	10.000	-/+	—
16	28	10.000	+/+	0,47	—	—	—	—	9	9.666	-/+	—
									7	10.000	-/-	—
17	42	10.000	+/+	0,48	—	—	—	—	10	9.666	-/+	—
									9	10.000	-/-	—
18	25	10.000	+/+	0,62	—	—	—	—	10	10.000	+/+	0,78
19	2	10.000	+/+	1,00	—	—	—	—	—	—	—	—

Tabelle 7.4.: Verlauf der Adaption konzeptuellen Wissens. Jede Zeile entspricht einem Block bzw. Lernblock; die Größe beschreibt die Zahl der enthaltenen Samples. Ob Random Forests (RF) bzw. neuronale Netzen (NN) erfolgreich trainiert wurden, ist durch + bzw. — gekennzeichnet.

7.2.2. Anwendung auf Testdaten

Im Rahmen der Adaption wurden drei konzeptuelle Modelle in die Wissensdomäne integriert, die eine intersubjektiv nachvollziehbare Unterscheidung von Impedanzspektren hinsichtlich der zugrundeliegenden Zelllinie realisieren. Deren pragmatischen Eigenschaften Σ und Z sind somit definiert als $\Sigma = \{NN, RF\}$ und $Z = \{L\}$. Temporal decken diese Modelle zusammen 11.002.274 Sekunden bzw. 22,5 Prozent der Zeitspanne der gesamten Datenbasis ab. Zwei der Modelle befinden sich auf der Wissensebene 1. Das dritte Modell greift auf Modelle der Ebene 1 der explorierten Wissensdomäne zurück und wurde daher auf Ebene 2 integriert. Für eine weiterführende Darstellung der adaptierten Modelle siehe Anhang D.

Werden diese konzeptuellen Modelle auf Testdaten angewendet, so werden zunächst deren pragmatische Eigenschaften mit denen der Testdaten abgeglichen. Insbesondere werden die erlernten konzeptuellen Modelle ausschließlich auf Testdaten angewendet, die innerhalb deren temporalen Geltungsbereichs liegen. Dies trifft hier auf insgesamt 11.138 Samples bzw. 14,9 Prozent der Testdaten der Datenbasis ($n=75.000$) zu. Umgekehrt bleiben 85,1 Prozent der verfügbaren Testdaten im Folgenden unberücksichtigt.

Die abgedeckten Samples wurden für jedes der adaptierten Modelle getrennt mit dem zugehörigen zuvor trainierten künstlichen neuronalen Netz und dem zugehörigen zuvor trainierten Random Forest evaluiert. Für Modelle der Wissensebene 1 werden die relevanten Features der Testdaten genutzt, für Modelle einer Ebene $n \geq 2$ werden die Testdaten unter Anwendung der dazwischen liegenden Modelle in die entsprechenden Features der Ebene $n - 1$ transformiert. Für die nachfolgende Auswertung werden die Ergebnisse der drei konzeptuellen Modelle entsprechend des verwendeten Lernverfahrens aggregiert dargestellt.

Zusammengenommen klassifizierten die neuronalen Netze 9.836 der abgedeckten 11.138 Test-Samples korrekt, während 1.301 der abgedeckten Test-Samples falsch zugeordnet wurden. Dies entspricht einer Accuracy von 0,88. Die Random Forests klassifizierten zusammen 9.325 der 11.138 Test-Samples richtig, während 1.813 Samples falsch klassifiziert wurden. Dies entspricht einer Accuracy von 0,84. Tab. 7.5 zeigt die Konfusionsmatrizen für die Klassifizierung durch die künstlichen neuronalen Netze und die Random Forests.

Zielwert	NN-Klassifizierung			RF-Klassifizierung		
	HT-29/B6	IPEC-J2	MDCK I	HT-29/B6	IPEC-J2	MDCK I
HT-29/B6	5.230	151	1	3.614	1.766	2
IPEC-J2	916	3.595	62	18	4.552	3
MDCK I	6	142	1.035	1	23	1.159

Tabelle 7.5.: Konfusionsmatrizen für die Klassifizierung nach zugehöriger Zelllinie. Dargestellt sind die mittels künstlicher neuronaler Netze (links) und mittels Random Forests (rechts) erzielten Zuordnungen für die abgedeckten Testdaten ($n=11.138$). In beiden Fällen sind die angegebenen Werte für die Evaluierungen der drei verwendeten konzeptuellen Modelle aggregiert.

7. Konzeptuelles Wissen für Impedanzspektren

8

Prozedurales Wissen für Impedanzspektren

In diesem Kapitel wird untersucht, in wieweit sich konstruktivistisches maschinelles Lernen nutzen lässt, um metrische Eigenschaften von Epithelien anhand impedanzspektroskopischer Messungen zu quantifizieren. Anlog zum Erlernen konzeptuellen Wissens wird zunächst eine Exploration der Trainingsdaten ohne Zielparameter durchgeführt. Anschließend wird die so erzeugte Wissensdomäne zur Adaption eines klinisch relevanten Parameters genutzt. Soweit nicht anders angegeben, werden die durch Testläufe ermittelten konstruktivistischen Lernparameter in Tab. 8.1 angewendet. Für eine weiterführende Charakterisierung der Datenbasis siehe Anhang C.

	Parameter	Wert
Konstruktion	Max. prozedurale Komplexität κ_p	2
	Max. Fehler	0,25
Komplexitäts- reduktion	Max. Modellkomplexität κ	10
	Filter-Grenze λ_x	50
	Filter-Grenze λ_y	15.000
	Maximalreduktion erzwingen	Ja
Rekonstruktion	Max. durchschnittliche Abweichungen [%]	10
	Max. Maximalabweichung [%]	100
	Samplegröße für Reliabilitätsprüfung	0,1
	Schwellwert für Reliabilitätsprüfung α_{min}	0,8
	Redundanzvermeidung	Ja
Dekonstruktion	Umgang mit neuen Modellen	integrativ
	Abbruchbedingung	minimal
	Temporale Erweiterungstoleranz δT_{max}	1
	Toleranz bei vollständiger Dekonstruktion	0,1
	ΣZ -verwandte Modelle immer vereinigen	Ja
	Schwache Reliabilität erlaubt	Nein

Tabelle 8.1.: Übersicht der zur Exploration prozeduralen Wissens verwendeten konstruktivistischen Lernparameter. Für eine weiterführende Beschreibung der Parameter siehe Anhang A.

8.1. Exploration einer prozeduralen Wissensdomäne

Wie in Abschnitt 5.4.2 beschrieben, werden maschinelle Modelle, die auf Regression basieren, hier als prozedurales Wissen bzw. prozedurale maschinelle Modelle bezeichnet. Durch Anwendung des konstruktivistischen maschinellen Lernens zur Erzeugung prozeduraler maschineller Modelle auf eine gegebene Menge Trainingsdaten entsteht eine hierarchische Ordnung dieser Modelle. Insbesondere werden mit diesem Prozess Bezüge zwischen den Modellen identifiziert und gespeichert. Eine solche Menge strukturierter und miteinander verknüpfter prozeduraler Modelle kann als Äquivalent zum prozeduralen Wissen nach Bloom betrachtet werden und wird daher im Folgenden analog als prozedurale Wissensdomäne bezeichnet.

Analog zum Erlernen konzeptuellen Wissens (Kapitel 7) kann das Erlernen prozeduralen Wissens sowohl unmittelbar auf den Trainingsdaten als auch unter Einbeziehung einer existierenden Wissensdomäne erfolgen. Dies schafft die Voraussetzung dafür, Modelle, die auf anderen Daten bzw. in früheren Durchläufen erzeugt wurden, zur Erzeugung maschineller Modelle einer höheren Wissensebene zu nutzen. Wie bereits mit der Exploration einer konzeptuellen Wissensdomäne gezeigt, erhöht ein sequentielles Lernen auf unterschiedlichen Feature-Sets sowohl die Wahrscheinlichkeit, innerhalb der Trainingsdaten Modelle höherer Wissensebenen zu identifizieren, als auch deren wahrscheinliche Anzahl.

8.1.1. Explorationsverlauf

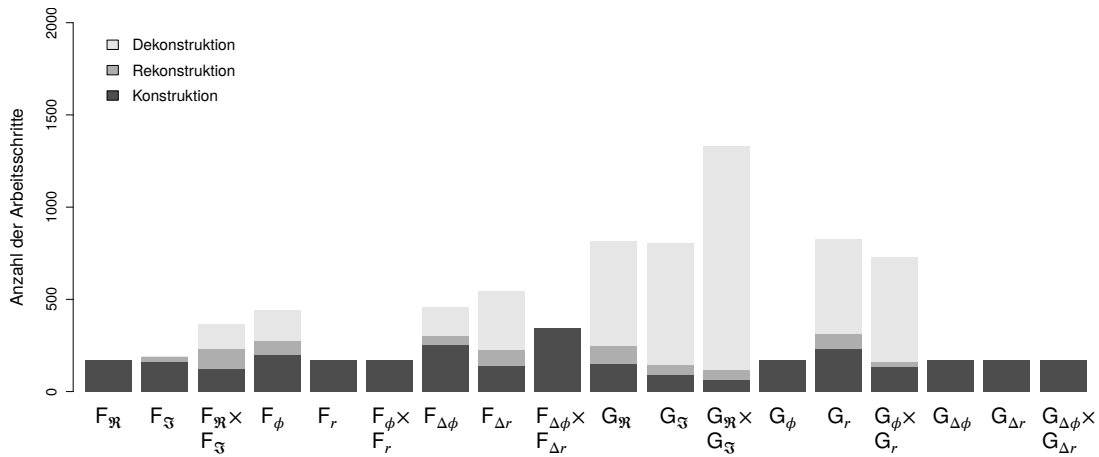
Die prozedurale Exploration der Datenbasis erfolgt analog zur konzeptuellen Exploration in 18 Teilschritten. In diesen wird jeweils eines der Feature-Sets der kartesischen Darstellung der Impedanzen ($F_{\mathfrak{R}}, F_{\mathfrak{S}}, F_{\mathfrak{R}} \times F_{\mathfrak{S}}$), der polaren Darstellung ($F_{\phi}, F_r, F_{\phi} \times F_r$) sowie der davon abgeleiteten Feature-Sets ($F_{\Delta\phi}, F_{\Delta r}, F_{\Delta\phi} \times F_{\Delta r}, G_{\mathfrak{R}}, G_{\mathfrak{S}}, G_{\mathfrak{R}} \times G_{\mathfrak{S}}, G_{\phi}, G_r, G_{\phi} \times G_r, G_{\Delta\phi}, G_{\Delta r}, G_{\Delta\phi} \times G_{\Delta r}$) als Eingabe verwendet. Die Feature-Sets werden dabei konsekutiv zum Training genutzt, wobei die jeweils zuletzt erzeugte Wissensdomäne wiederum als Ausgangspunkt für den darauf folgenden Teilschritt verwendet wird. Die Anzahl der Samples entspricht jeweils dem gesamten Trainingsdatensatzes ($n=200.000$). Analog zur Exploration konzeptuellen Wissen werden die Testdaten nicht zur Exploration verwendet.

Aus den Trainingssamples werden jeweils Blöcke von 10.000 Samples gezogen und dem Lernprozess unterworfen. Die Anzahl der erfolgreichen Konstruktionen variiert dabei ebenso von Feature-Set zu Feature-Set wie die Anzahl der ausgeführten Lernoperationen für Rekonstruktions- und Dekonstruktionsprozesse (Abb. 8.1a). Für die Lernblöcke der Feature-Sets $F_{\mathfrak{R}}, F_r, F_{\phi} \times F_r, F_{\Delta\phi} \times F_{\Delta r}, G_{\phi}, G_{\Delta\phi}, G_{\Delta r}$ und $G_{\Delta r} \times G_{\Delta r}$ wurde ausschließlich der Konstruktionsprozess durchlaufen. Infolge der erfolglosen Konstruktionen konnte aus diesen Lernblöcken keine Modelle rekonstruiert bzw. in die Wissensdomäne integriert werden.

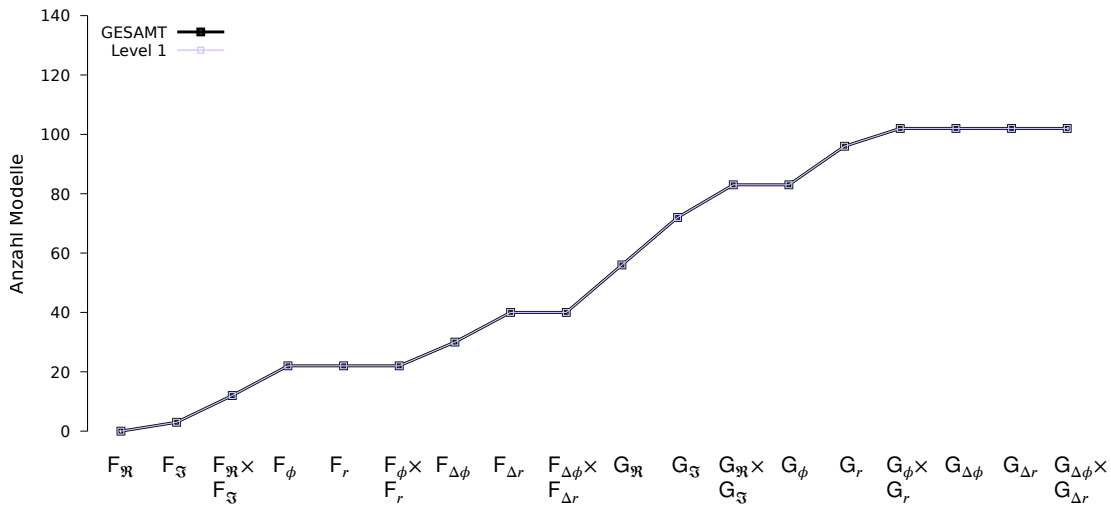
Im Verlauf der Exploration nimmt die Anzahl der Domänenmodelle von Feature-Set zu Feature-Set kontinuierlich zu (Abb. 8.1b). Die Hierarchie der prozeduralen Domäne bleibt im Verlauf der Exploration unverändert. Während sie nach Durchlauf des Feature-Sets $F_{\mathfrak{R}}$ noch keine Modelle enthält, weist die Wissensdomäne nach Durchlauf des Feature-Sets $F_{\mathfrak{S}}$ Modelle auf Ebene 1 auf. Diese flache Hierarchie wird – anders als während der Exploration der konzeptuellen Domäne – im weiteren Verlauf beibehalten, so dass die prozedurale Domäne auch nach Abschluss der Exploration ausschließlich maschinelle Modelle der Ebene 1 aufweist (Abb. 8.1b).

Die Anzahl der Eingabefeatures pro maschinell Modell liegt im Verlauf der Exploration überwiegend zwischen 2 und 5. Die Komplexität der Modelle nimmt damit in der ersten Hälfte der Exploration leicht zu (Abb. 8.1c). Für die letzten sieben Feature-Sets, beginnend ab $G_{\mathfrak{R}} \times G_{\mathfrak{S}}$, sind trotz einer weiteren Zunahme der Zahl der Domänenmodelle keine statistischen Veränderungen in der Modellkomplexität mehr erkennbar.

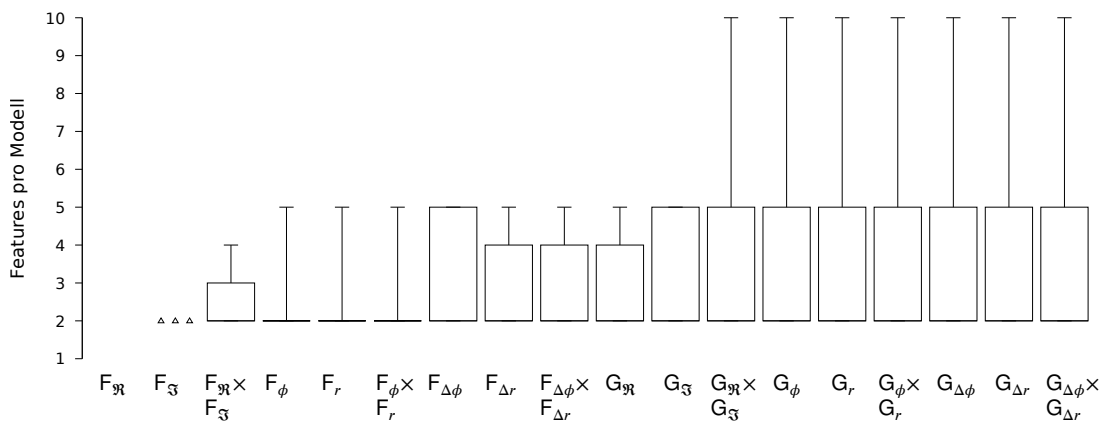
8.1. Exploration einer prozeduralen Wissensdomäne



(a) Ausgeführte Arbeitsschritte



(b) Größe der Wissensdomäne



(c) Komplexität κ der Domänenmodelle

Abbildung 8.1.: Verlauf der Exploration der prozeduralen Wissensdomäne. Die Exploration erfolgt in 18 Teilschritten, dargestellt ist jeweils der Zustand am Ende eines Feature-Set-Durchlaufs.

8.1.2. Charakterisierung der explorierten Wissensdomäne

Nach Ende der aus 18 Teilschritten bestehenden Exploration setzt sich die prozedurale Wissensdomäne aus 102 maschinellen Modellen zusammen. Diese befinden sich alle auf der ersten Wissensebene. Eine Abstraktion von Modellen der Ebene 2 oder höher liegt damit nicht vor. Der Verlauf der Exploration lässt erkennen, dass dies auch zu keinem früheren Zeitpunkt, d.h. nach Abschluss früherer Teilschritte, der Fall war (Abb. 8.1b).

Alle 102 Modelle, die in die Wissensdomäne aufgenommen wurden, wurden mittels Autoencoder konstruiert. Sie stellen somit eindimensionale Repräsentationen des jeweiligen Feature-Sets dar, die einen Lernblock oder mehrere Lernblöcke abdecken. Die Definitionsbereiche der konstruierten Zielparameter der Modelle sind in Abb. 8.2 dargestellt. Urheber der konstruierten Modelle sind Autoencoder mit entweder einem oder zwei Neuronen pro versteckter Schicht, wobei im Falle zweier Neuronen jeweils die Zielwerte des mit besseren Ergebnissen rekonstruierten Neurons verwendet wurden. Mittels Feature-Clustering konstruierte Modelle wurden nicht in der prozeduralen Wissensdomäne gespeichert.

Da sich alle Modelle der Domäne auf Wissensbene 1 befinden, dienen ausschließlich die Features der Trainingsdaten (Ebene 0) als Eingabefeatures. Pro Modell werden maximal zehn Eingabefeatures verwendet, im Durchschnitt 3,3 Features. Minimum, erstes Quartil und Median liegen bei 2 Features, das dritte Quartil bei 5 Features. Wie aus Abb. 8.3 hervorgeht, werden die Feature-Sets der Datenbasis nicht gleichmäßig genutzt. Insbesondere wurden die Features der Sets F_r , G_ϕ , $G_{\Delta\phi}$ und $G_{\Delta r}$ nicht zur Konstruktion von Modellen der Ebene 1 genutzt.

In der zeitlichen Ausdehnung deckt ein Modell der explorierten Domäne im Mittel 8,5 Prozent der Zeitspanne der Datenbasis ab. Allerdings schwankt diese temporale Abdeckung zwischen einem Minimum von 3,9 und einem Maximum von 53,8 Prozent. Die Domäne bzw. die Modelle der Ebene 1 bilden somit sowohl Wissen mit enger als auch weiterer zeitlicher Begrenzung ab. Abb. 8.4 zeigt die Ausdehnung der prozeduralen Modelle der Ebene 1 unter Berücksichtigung der tatsächlichen zeitlichen Start- und Endpunkte.

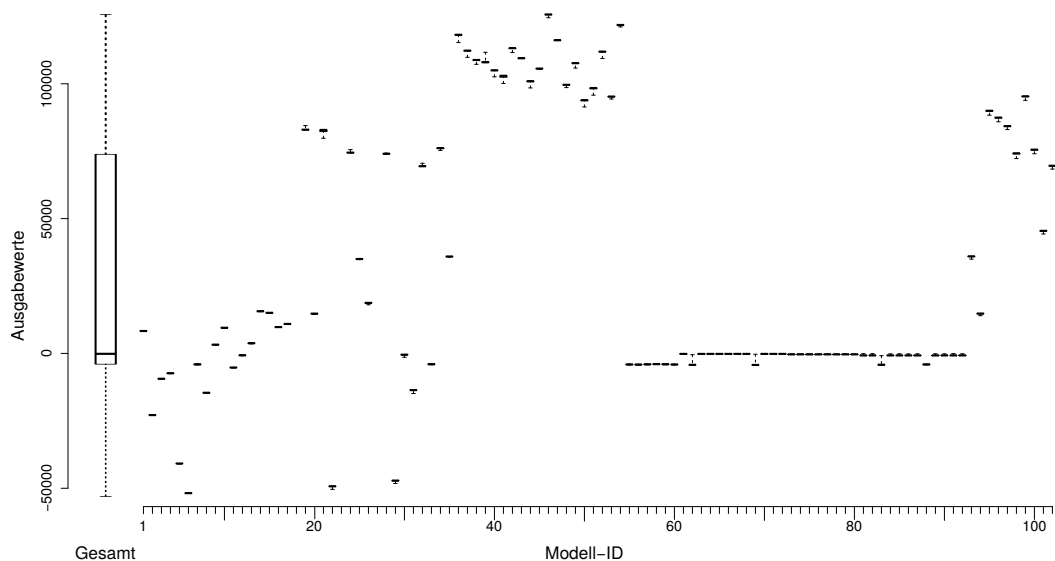
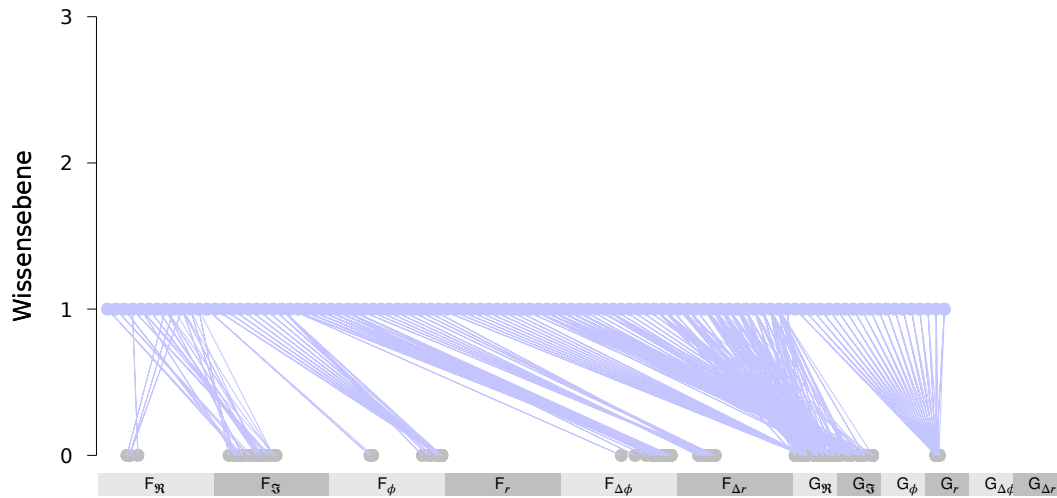
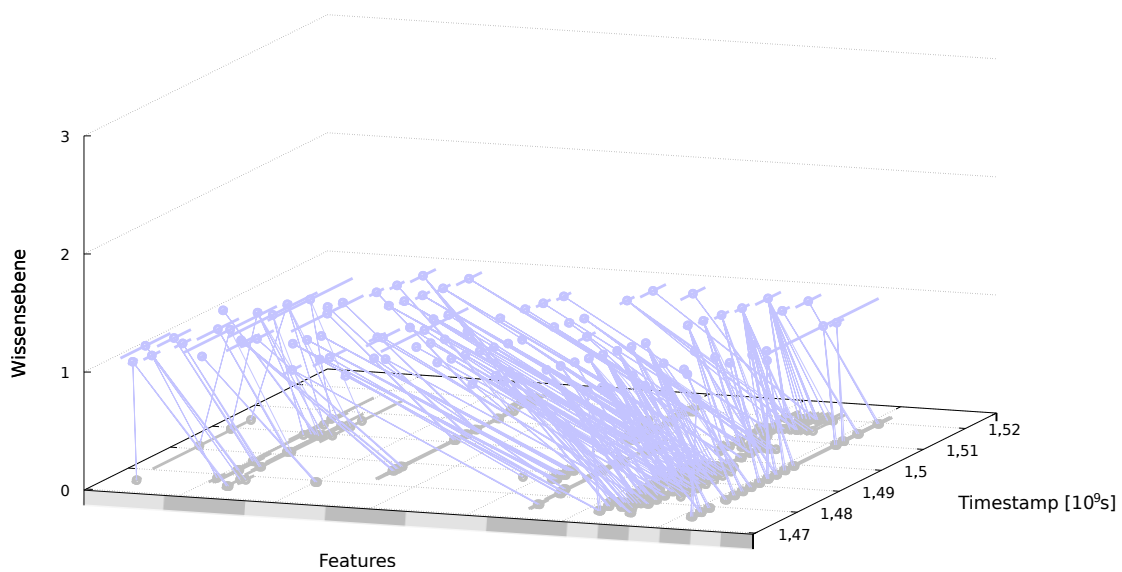


Abbildung 8.2.: Wertebereiche der Modelle der explorierten prozeduralen Wissensdomäne. Dargestellt sind Boxplots der Zielwerte der gesamten Domäne (links) sowie der einzelnen prozeduralen Modelle (rechts), die anhand ihrer ID auf einer einheitslosen x-Achse angeordnet sind.

8.1. Exploration einer prozeduralen Wissensdomäne



a) Modellhierarchie



b) Dreidimensionale Aufsicht

Abbildung 8.3.: Überblick über die explorierte prozedurale Wissensdomäne. Modelle der gleichen Wissensebene sind jeweils im gleichen Farbton dargestellt. Die Feature-Sets der Trainingsdaten, die sich auf Ebene 0 befinden, sind entlang der x-Achse gekennzeichnet. a) Modellhierarchie. In dieser Frontal-Aufsicht werden Verknüpfungen zwischen den Domänenmodellen durch vertikale Geraden repräsentiert; die Zeitachse wird in dieser Darstellung vernachlässigt. b) Dreidimensionale Aufsicht. Zusätzlich zur Darstellung als Modellhierarchie ist die zeitliche Ausdehnung jedes Modells auf der z-Achse repräsentiert.

8. Prozedurales Wissen für Impedanzspektren

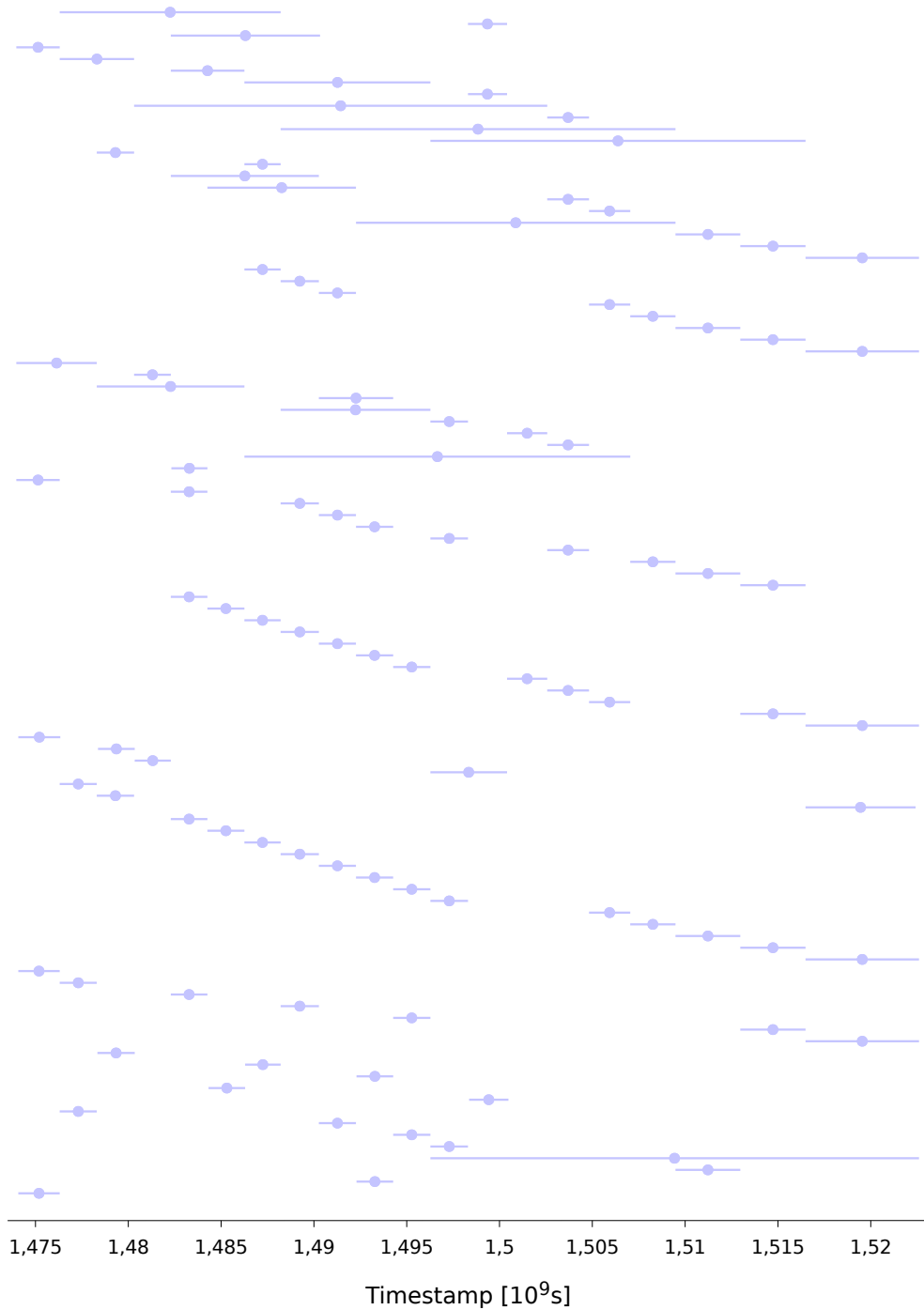


Abbildung 8.4.: Temporale Gültigkeit der Modelle der prozeduralen Wissensdomäne. Dargestellt ist die zeitliche Ausdehnung der 102 Modelle der Ebene 1. Die Modelle sind anhand ihrer ID von oben nach unten aufsteigend in 1er-Schritten auf einer einheitslosen y-Achse angeordnet.

8.2. Adaption prozeduralen Wissens

Das Erlernen von Regressionsaufgaben mittels konstruktivistischen maschinellen Lernens erfolgt analog zur Adaption konzeptuellen Wissens. Der Trainingsdatensatz ($n=200.000$) und die explorierte prozedurale Wissensdomäne werden dazu verwendet, um einen gegebenen Zielparameter abzubilden. Anschließend wird der Lernerfolg mit einem Testdatensatz ($n=75.000$) evaluiert. Im Gegensatz zur Adaption konzeptuellen Wissens wird hier jedoch statt einem kategorialen ein metrischer Zielparameter erlernt.

Als Beispiel für das Erlernen prozeduralen Wissens wird die Quantifizierung der epithelialen Kapazität C^{epi} untersucht. Dabei handelt es sich um einen globalen Parameter des Gewebes, der die kapazitiven Eigenschaften der apikalen und basolateralen Membran zusammenfasst (vgl. Abschnitt 6.1). Innerhalb der verwendeten Datenbasis variiert dieser Parameter insgesamt zwischen 0,5 und 5,5 μF , wobei dessen Durchschnittswerte für HT-29/B6-Zellen bei rund 3,0 μF , bei IPEC-J2-Zellen bei rund 1,5 μF und bei MDCK I bei rund 2,0 μF liegen (vgl. Anhang C).

Ziel des Adaptionprozesses ist hier die Abbildung der Regressionsaufgabe in Form eines oder mehrerer maschineller Modelle innerhalb der prozeduralen Wissensdomäne. Insbesondere soll so ermöglicht werden, die temporale und die intersubjektive Gültigkeit zu bewerten. Bezüglich der Regressionsgenauigkeit wird eine durchschnittliche relative Abweichung vom Zielwert von weniger als 10 Prozent und eine maximale relative Abweichung von weniger als 100 Prozent angestrebt. Die Modellkomplexität κ soll die Anzahl der verwendeten Frequenzen ($n=42$) nicht übersteigen. Tab. 8.2 gibt einen vollständigen Überblick über die verwendeten Parameter.

8.2.1. Adaptionsverlauf

Im Folgenden wird dokumentiert, wie die gegebene Regressionsaufgabe durch das konstruktivistische maschinelle Lernen unter Berücksichtigung der explorierten prozeduralen Wissensdomäne adaptiert wird. Tab. 8.3 zeigt den blockweisen Ablauf, wobei $T\Sigma$ - und TZ -Dekonstruktionen zugunsten der Übersichtlichkeit und mangels Relevanz für die Adaption nicht dargestellt sind.

Die Lernblöcke der Wissensebene 0 wurden für 17 der 20 Blöcke der Datenbasis erfolgreich rekonstruiert. Random Forests und neuronale Netze lieferten dabei relative Abweichungen von den

	Parameter	Wert
Komplexitätsreduktion	Max. Modellkomplexität κ	42
	Filter-Grenze λ_x	350
	Filter-Grenze λ_y	10.000
	Maximalreduktion erzwingen	Ja
Rekonstruktion	Max. durchschnittliche Abweichungen [%]	20
	Max. Maximalabweichung [%]	100
	Samplegröße für Reliabilitätsprüfung	5000
	Schwellenwert für Reliabilitätsprüfung α_{min}	0,75
	Redundanzvermeidung	Ja
	Dekonstruktion	Umgang mit neuen Modellen
Abbruchbedingung		minimal
Temporale Erweiterungstoleranz δT_{max}		1
Toleranz bei vollständiger Dekonstruktion		0,1
ΣZ -verwandte Modelle immer vereinigen		Ja
Schwache Reliabilität erlaubt		Nein

Tabelle 8.2.: Übersicht der zur Adaption prozeduralen Wissens verwendeten konstruktivistischen Lernparameter. Für eine weiterführende Beschreibung der Parameter siehe Anhang A.

8. Prozedurales Wissen für Impedanzspektren

Zielwerten, die die gesetzten Schwellenwerte nicht überschritten. Die Inter-Coder-Reliabilität lag zwischen 0,96 und 0,98. Die für die Blöcke 2, 5, 6, 8, 9, 13 und 14 rekonstruierten maschinellen Modelle wurden im Rahmen von ΣZ -Dekonstruktionen jeweils mit einem existierenden Modell der Wissensdomäne vereinigt, das infolgedessen durch das vereinigte Modell ersetzt wurde. Für Block 14 wurde eine erfolglose vollständige Dekonstruktion und anschließend eine erfolgreiche Modelldifferenzierung durchgeführt. Die Lernblöcke der Blöcke 10, 11 und 12 konnten auf Ebene 0 nicht erfolgreich rekonstruiert werden und wurden daher auf die Ebene 1 transformiert. Die Rekonstruktion der transformierten Lernblöcke verlief hier jedoch ebenfalls erfolglos.

Block	Rekonstruktion (Ebene 0)				Vollständige/ ΣZ -Dekonstruktion				Rekonstruktion (Ebene 1)			
	κ	Größe	RF/NN	α	κ	Größe	RF/NN	α	κ	Größe	RF/NN	α
1	40	9.857	+/+	0,97	—	—	—	—	—	—	—	—
2	42	9.900	+/+	0,97	22	9.900	+/+	0,98	—	—	—	—
3	42	9.900	+/+	0,97	16	9.900	-/-	—	—	—	—	—
4	39	9.900	+/+	0,97	17	9.900	-/-	—	—	—	—	—
					18	9.900	+/-	—	—	—	—	—
5	42	9.900	+/+	0,97	16	9.900	-/-	—	—	—	—	—
					18	9.900	+/-	—	—	—	—	—
					18	9.900	+/+	0,95	—	—	—	—
6	42	9.900	+/+	0,97	15	9.900	-/-	—	—	—	—	—
					22	9.900	+/-	—	—	—	—	—
					16	9.900	+/+	0,96	—	—	—	—
7	42	9.900	+/+	0,96	{10,16,15}	9.900	-/-	—	—	—	—	—
8	41	9.900	+/+	0,97	{10,14}	9.900	-/-	—	—	—	—	—
					13	9.900	+/-	—	—	—	—	—
					23	9.900	+/+	0,97	—	—	—	—
9	36	9.900	+/+	0,97	{13,14,14}	9.900	-/-	—	—	—	—	—
					17	9.900	+/+	0,97	—	—	—	—
10	30	9.900	+/-	—	—	—	—	—	{12,2}	8.833	-/+	—
									4	9.000	+/-	—
									6	9.900	-/-	—
11	39	9.900	+/-	—	—	—	—	—	{10,4}	8.900	-/-	—
									6	9.900	-/-	—
12	40	9.900	+/-	—	—	—	—	—	{11,7}	8.800	+/-	—
									4	9.900	-/-	—
13	40	9.900	+/+	0,97	17	9.900	-/+	—	—	—	—	—
					23	9.900	+/+	0,98	—	—	—	—
14	42	9.900	+/+	0,96	18	FULL	-/+	—	—	—	—	—
						DIFF	+/+	0,99	—	—	—	—
						DIFF	+/+	0,98	—	—	—	—
15	40	9.900	+/+	0,96	{17,13,11,16,16}	9.900	-/-	—	—	—	—	—
16	40	9.900	+/+	0,96	{14,14,13,13}	9.900	-/-	—	—	—	—	—
					15	9.900	+/-	—	—	—	—	—
17	41	9.900	+/+	0,97	{13,14,12,12,19}	9.900	-/-	—	—	—	—	—
					14	9.900	+/-	—	—	—	—	—
18	42	9.900	+/+	0,97	{12,14,11,11,15}	9.900	-/-	—	—	—	—	—
					16	9.900	+/-	—	—	—	—	—
19	41	9.900	+/+	0,97	{12,14,12,12,15}	9.900	-/-	—	—	—	—	—
					15	9.900	+/-	—	—	—	—	—
20	41	9.900	+/+	0,98	{11,14,14,11,11,17}	9.900	-/-	—	—	—	—	—
					17	9.900	-/+	—	—	—	—	—

Tabelle 8.3.: Verlauf der Adaption prozeduralen Wissens. Lernblöcke eines Blocks sind entweder als eigene Zeile oder aggregiert in einer Mengendarstellung der Modellkomplexität κ repräsentiert; die Größe beschreibt die Zahl der enthaltenen Samples. Ob Random Forests (RF) bzw. neuronale Netze (NN) erfolgreich trainiert wurden, ist durch + bzw. – gekennzeichnet.

8.2.2. Anwendung auf Testdaten

Im Rahmen der Adaption wurden acht prozedurale Modelle in die Wissensdomäne integriert, die eine intersubjektiv nachvollziehbare Quantifizierung der epithelialen Kapazität realisieren. Die pragmatischen Eigenschaften Σ und Z sind somit definiert als $\Sigma = \{NN, RF\}$ und $Z = \{C^{epi}\}$. Temporal decken diese Modelle zusammen 40.518.505 Sekunden bzw. 82,8 Prozent der Zeitspanne der gesamten Datenbasis ab. Alle acht Modelle nutzen unmittelbar die Features der Datenbasis und wurden daher auf Ebene 1 der Wissensebene integriert. Für eine ausführliche Darstellung der adaptierten Modelle siehe Anhang D.

Analog zur Anwendung von konzeptuellen Modellen auf Testdaten (Abschnitt 7.2), setzt die Anwendung erlernter prozeduraler Modelle eine Überlappung der pragmatischen Eigenschaften der Modelle und der Testdaten voraus. Nur Testdaten, deren temporale Gültigkeit innerhalb des temporalen Geltungsbereichs der erlernten prozeduralen Modelle liegen, werden evaluiert. Dies trifft hier auf insgesamt 63.246 Samples oder 84,4 Prozent der Testdaten ($n=75.000$) zu, die im Folgenden ausgewertet werden. 15,6 Prozent der Testdaten weisen dagegen keine Überlappung mit den adaptierten Modellen auf und gehen daher nicht in die Auswertung ein.

Für die nachfolgende Auswertung werden die Ergebnisse der acht prozeduralen Modelle entsprechend des angewendeten Lernverfahrens aggregiert betrachtet. Abb. 8.5 zeigt die Genauigkeit von künstlichen neuronalen Netzen und Random Forests als Boxplot. Im Mittel bestimmen die neuronalen Netze die epitheliale Kapazität mit einem relativen Fehler von $\pm 6,6$ Prozent, wobei dieser das Verhältnis der absoluten Abweichung zum tatsächlichen Wert ausdrückt. Maximal beträgt der relative Fehler $\pm 180,3$ Prozent. Die Random Forests bestimmen die epitheliale Kapazität im Mittel mit einem relativen Fehler von $\pm 2,4$ Prozent. Der maximale relative Fehler beträgt hier $\pm 90,9$ Prozent. Weitere statistische Kenngrößen können Tab. 8.4 entnommen werden.

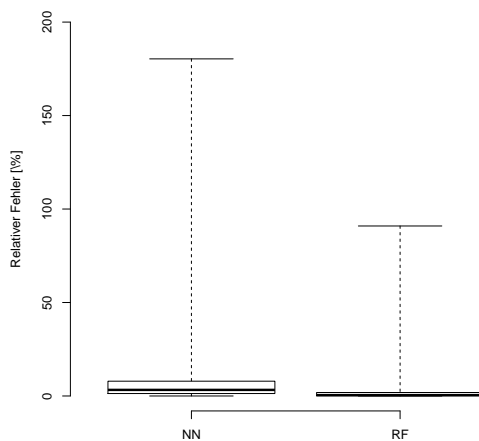


Abbildung 8.5.: Boxplot des relativen Fehlers der bestimmten epithelialen Kapazität [%].

	Rel. Fehler [$\pm\%$]	
	NN	RF
Maximum	180,3	90,9
3. Quartil	7,9	1,9
Arth. Mittel	6,6	2,4
Median	3,3	0,5
1. Quartil	1,3	0,1
Minimum	0,0	0,0

Tabelle 8.4.: Relativer Fehler der bestimmten epithelialen Kapazität [%].

8. Prozedurales Wissen für Impedanzspektren

Teil III.

DISKUSSION UND SCHLUSSFOLGERUNGEN

*Unser Wissen ist ein kritisches Raten,
ein Netz von Hypothesen,
ein Gewebe von Vermutungen.*

— Karl Popper (1902-1994)



9

Epitheliale Impedanzanalyse

Das vorgeschlagene Verfahren eines konstruktivistischen maschinellen Lernens wird in der vorliegenden Arbeit zur Unterscheidung (Kapitel 7) bzw. Quantifizierung (Kapitel 8) epithelialer Eigenschaften eingesetzt. Die dabei erzielten Ergebnisse werden in diesem Kapitel in Beziehung zu anderen Methoden gesetzt. Darüber hinaus wird die Übertragbarkeit des Analyseverfahrens auf weitere Fragestellungen der epithelialen Impedanzanalyse sowie die Bedeutung und Qualität der Simulation von Impedanzspektren diskutiert.

9.1. Modellierung impedanzspektroskopischer Messungen

In der Physiologie ist es seit langem übliche Praxis, das Verhalten von Epithelien in Form ihrer elektrischen Eigenschaften zu beschreiben. Obwohl dies insbesondere für die Modellierung von Ersatzschaltkreisen gilt, geht der hier beschriebene Ansatz sowohl qualitativ als auch quantitativ deutlich über traditionelle Verfahren hinaus. Er erlaubt nicht nur die Entwicklung differenzierter physiologischer Hypothesen, sondern aufgrund der Skalierbarkeit auch einen effektiven Einsatz von Methoden des maschinellen Lernens. Bei ausschließlicher Verwendung von manuell erhobenen Messdaten wird die hierfür anzustrebende Menge von mehreren zehntausend bis hunderttausend Samples in der Regel um mindestens ein bis zwei Größenordnungen verfehlt.

Anders als die Ein-Wege-Impedanzspektroskopie [77] oder die Zwei-Wege-Impedanzspektroskopie [143] unterscheidet der hier angenommene Ersatzschaltkreis zwischen apikaler und basolateraler Seite des Epithels. Wie bereits mit früheren Arbeiten gezeigt wurde [230], erlaubt dies nicht nur eine detailliertere Modellierung epithelialer Funktionen, sondern auch pathologisch relevanter Dysfunktionen. Gleichzeitig verwendet diese Modell Parameter, die – anders als bei Ein- und Zwei-Wege-Impedanzspektroskopie – nicht unmittelbar mittels Messungen bestimmbar bzw. überprüfbar sind. Die Herausforderung dieses Ansatzes besteht also insbesondere darin, die verwendeten Modellparameter auf ihre Plausibilität zu überprüfen.

Kernelement des hier angewendeten Modellierungsansatzes stellt die Definition spezifischer Parametergrenzen für Zelllinien und -zustände dar. Die hier verwendeten Parametergrenzen bauen sowohl auf publizierten Messwerten als auch auf zuvor definierten Parametergrenzen auf. Sie wurden darüber hinaus im Rahmen dieser Arbeit intensiv mit Messdaten verglichen,

9. Epitheliale Impedanzanalyse

die vom Institut für Klinische Physiologie der Charité Berlin zur Verfügung gestellt werden. Im Ergebnis kann insbesondere nicht nur die Annahme unterschiedlicher Parametergrenzen für die drei Zellkulturlinien HT-29/B6, IPEC-J2 sowie MDCK I als gesichert gelten, sondern ebenso die Annahme unterschiedlicher Parametergrenzen für die untersuchten funktionalen Zustände.

Ferner berücksichtigt der Modellierungsansatz explizit den gerätespezifischen Messfehler. Die gewählte mathematische Beschreibung stellt dabei eine Balance zwischen einer detaillierten und einer performanten Modellierung dar. Einerseits, weil die Beschreibung lediglich auf den tatsächlich in den Messungen verwendeten Frequenzumfang (bis 16,35 kHz) beschränkt ist, andererseits, weil für die Wahrscheinlichkeit jedes einzelnen Fehlers eine Normalverteilung angenommen wird. Dennoch ist das Modell, das zweistufig sowie getrennt für Real- und Imaginärteil berechnet wird, ausreichend komplex, um tatsächlich gemessene Daten zu erklären.

Belegt wird die Erklärkraft des Modellierungsansatzes durch eine hohe Übereinstimmung von modellierten und gemessenen Impedanzspektren. Der indirekte Abgleich auf Basis von Parametern, die aus den Spektren bestimmt werden, wurde bereits in einer früheren Arbeit als zweckmäßige Methodik etabliert [230]. Hier wurde dieser auf insgesamt drei voneinander unabhängige Parameter erweitert sowie durch Anwendung etablierter Binning-Algorithmen soweit präzisiert, dass sich für jeden der drei Parameter sowohl ein grafischer Abgleich als auch ein rechnerisches Fehlermaß κ_B ableiten lässt. Insgesamt ergibt sich so eine umfassende automatisierte Quantifizierung der indirekten Übereinstimmung, wodurch erstmals ein schneller und zuverlässiger Abgleich einer beliebigen Menge modellierter und gemessener Impedanzspektren möglich ist.

Ab einem κ_B -Wert von 0,85 wird die Nachbildung empirisch beobachteter Eigenschaften epithelialer Zelllinien unter einem gegebenen funktionalen Zustand als erfolgreich gewertet. Insbesondere die in der Literatur angenommenen Effekte der Wirkstoffe EGTA und Nystatin konnten mittels geeigneter Parametergrenzen für die Komponenten des Ersatzschaltkreises realistisch nachgebildet werden. Für die modellierte Anwendung von EGTA lag die niedrigste Übereinstimmung bei einem κ_B -Wert von 85,7 über alle Vergleichsparameter und Zelllinien hinweg, für die modellierte Anwendung von Nystatin bzw. die Kombinationsanwendung lag der niedrigste κ_B -Wert bei 90,9. Ebenso konnten die charakteristischen Unterschiede der untersuchten Zelllinien realistisch nachgebildet werden, was sich in einem κ_B -Wert von 88,5 oder mehr für alle Vergleichsparameter der modellierten Kontrollbedingungen von HT-29/B6, IPEC-J2 und MDCK I widerspiegelt.

Über bestehende Verfahren hinaus geht die Modellierung pragmatischer Metadaten, die hier erstmals für Impedanzspektren durchgeführt wurde. Der Fokus lag dabei auf dem temporalen Metadatum in Form eines UNIX-Timestamps (vgl. Abschnitt 6.5 bzw. 5.2.3). Die zugrundeliegenden Annahmen über den zeitlichen Ablauf von Impedanzmessungen basieren auf Erfahrungswerten des Instituts für Klinische Physiologie der Charité Berlin. Insgesamt handelt es sich bei der temporalen Modellierung zunächst um einen Proof of Concept. In künftigen Messungen sollen die tatsächlichen Erhebungszeitpunkte dokumentiert und genutzt werden. Für das Subjekt-Metadatum wurden durchweg ein einheitlicher Wert verwendet, der die hier angewendete computergestützte Synthetisierung repräsentiert. Da in diesem Kontext kein eindeutiger Zweck für jede Messung definierbar ist, wurde dieser einheitlich mit dem Default-Wert 0 als fehlend markiert.

Aufgrund der systematischen Synthetisierung der Impedanzspektren war es möglich, eine ebenso realistische wie repräsentative Datenbasis zu synthetisieren. Insbesondere minimiert die gleichmäßige Zusammensetzung aus jeweils 25.000 synthetisierten Messungen pro Zellkulturlinie und funktionalem Zustand (Tab. 6.4) einen eventuellen Bias in nachgelagerten Analysen. Einzige Ausnahme stellt die Zelllinie MDCK I dar, für die der Einfluss der Wirkstoffkombination EGTA und Nystatin mangels Messdaten nicht modelliert werden konnte. Insgesamt ergibt sich so ein realistischer, repräsentativer und umfangreicher Referenzdatensatz, der sich nicht nur für die computergestützte epitheliale Impedanzanalyse, sondern aufgrund der Modellierung pragmatischer Metadaten ebenso zur Ableitung pragmatisch definierter maschineller Modelle eignet.

9.2. Unterscheidung von Zelllinien

Als Beispiel für die Adaption konzeptuellen Wissens wurde in Abschnitt 7.2 untersucht, inwieweit sich Impedanzspektren mittels konstruktivistischen maschinellen Lernens hinsichtlich der zugehörigen Zelllinie unterscheiden lassen. Wie dargestellt, weisen die Zelllinien HT-29/B6, IPEC-J2 und MDCK I Unterschiede in Aufbau und Funktionsweise auf. Herausfordernd ist deren Unterscheidung anhand von Impedanzspektren insbesondere dann, wenn neben Messungen unter Kontrollbedingungen auch Messungen anderer funktionaler Zustände als Zelltyp-Repräsentation betrachtet werden. Ein Impedanzspektrum mit einem epithelialen Widerstand von rund $250 \Omega\text{cm}^2$ etwa kann einerseits sowohl Kontrollbedingungen, Nystatin-Zugabe und EGTA-Zugabe bei HT-29/B6- als auch MDCK-I-Zellen entstammen, andererseits aber auch einer kombinierten Zugabe von Nystatin und EGTA bei HT-29/B6 oder IPEC-J2 (Anhang B). Ein menschlicher Betrachter kann die Zuordnung eines solchen Spektrums zu einer Zelllinie in der Regel ohne Kenntnis von Zelllinie und Messbedingungen nicht zuverlässig aus der Geometrie des Spektrums ableiten. Ein computergestütztes Verfahren könnte daher sowohl die Qualitätssicherung in klinischen Screenings verbessern als auch eine schnelle Zuordnung unbekannter Proben ermöglichen. Entsprechende Ansätze sind in der Literatur jedoch bislang nicht zu finden.

In der vorliegenden Arbeit ist ein zentrales Ergebnis der Adaption konzeptuellen Wissens die implizite Differenzierung zwischen Impedanzspektren, die eine Unterscheidung nach Zelllinien erlauben, und solchen, für die dies nicht der Fall ist. Im Rahmen der Adaption konnten drei konzeptuelle Modelle in die Wissensdomäne integriert werden, mit denen eine Zuordnung zu einer Zelllinie für insgesamt 14,9 Prozent der zu testenden Impedanzspektren realisiert wird. Von den mit diesen Modellen abgedeckten Spektren konnten mittels künstlicher neuronaler Netze 88,3 Prozent und mittels Random Forests 83,7 Prozent korrekt zugeordnet werden. Die übrigen 85,1 Prozent der zu testenden Impedanzspektren konnten zwar blockweise jeweils von einem künstlichen neuronalen Netz und einem Random Forest mit der geforderten Accuracy rekonstruiert werden, allerdings nicht mit der zuvor definierten Intersubjektivität von $\alpha \geq 0,75$ (Tab. 7.4).

Die zu erlernende Klassifikationsaufgabe kann daher nur für bestimmte Spektren als zuverlässig lösbar betrachtet werden. Vergleicht man die abgedeckten mit den nicht abgedeckten Spektren, so fällt etwa auf, dass die Zelllinie MDCK I unter den Impedanzspektren, die von den adaptierten konzeptuellen Modellen abgedeckt werden, seltener repräsentiert ist als in der Gruppe der nicht abgedeckten; die funktionalen Zustände sind im Vergleich dazu gleichmäßiger repräsentiert. Dies lässt den Schluss zu, dass zwar HT-29/B6-Zellen überwiegend mit der vorgegebenen Zuverlässigkeit gegen IPEC-J2-Zellen abgegrenzt werden können, nicht jedoch gegen MDCK-I-Zellen. Tatsächlich sind für Impedanzspektren der Zelllinien HT-29/B6 und IPEC-J2 physiologisch

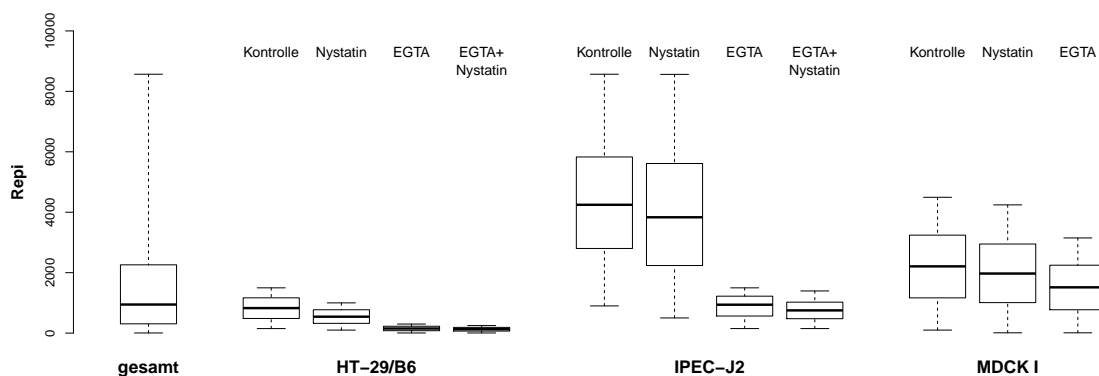


Abbildung 9.1.: Charakteristische Parameter der Datenbasis. R^{epi}

9. Epitheliale Impedanzanalyse

bedingte Geometrie-Unterschiede bekannt. Während etwa für IPEC-J2-Zellen unter Kontrollbedingungen Werte zwischen 1.800 und 8.500 Ωcm^2 für den epithelialen Widerstand publiziert wurden, finden sich in der Literatur für HT-29/B6-Zellen Werte bis maximal 650 Ωcm^2 . Die angenommenen Werte für MDCK I überschneiden sich mit diesen beiden Wertebereichen.

Ein nicht unerheblicher Teil der zu analysierenden Impedanzspektren ist mit gleicher Wahrscheinlichkeit der Zelllinie HT-29/B6, IPEC-J2 oder MDCK I zuzuordnen. Grund für diese Mehrdeutigkeit ist einerseits, dass biologische Zellen jenseits von prototypischen Exemplaren grundsätzlich eine Varianz in ihren Eigenschaften aufweisen, wodurch fließende Übergänge entstehen. Andererseits induziert die Anwendung von Wirkstoffen eine zusätzliche Varianz sowohl der elektrischen Eigenschaften als auch der Geometrie davon abhängiger Impedanzspektren. Je fließender die Übergänge, desto weniger eindeutig sind jedoch Zuordnungen zu einer Zelllinie möglich. Die durch das konstruktivistische maschinelle Lernen erfolgte Begrenzung der adaptierten Modelle auf einen Teil der zu erlernenden Impedanzspektren ergibt sich daher unmittelbar aus der Nebenbedingung, dass die Klassifikationsaufgabe möglichst zuverlässig erfolgen soll.

9.3. Quantifizierung der epithelialen Kapazität

Als Beispiel für adaptiertes prozedurales Wissen wurde in Abschnitt 8.2 untersucht, inwieweit sich die Kapazität C^{epi} eines Epithelgewebes mittels konstruktivistischen maschinellen Lernens quantifizieren lässt. In dem minimalen Ersatzschaltbild *A* (Abb. 6.1a) repräsentiert dieser Parameter die aggregierten kapazitiven Anteile der inneren und äußeren Lipiddoppelschicht der zusammenhängenden Barriere, die Epithelien mit ihren Zell-Zell-Kontakten realisieren. Die dadurch etablierte Kapazität ist innerhalb einer Zellkultur relativ konstant. Für die untersuchten Zelllinien wurden Werte zwischen 0,5 und 5,5 $\mu\text{F}/\text{cm}^2$ modelliert (Tab. 6.1, 6.2 und 6.3).

Obwohl die Kapazität eines Epithels typischerweise unempfindlich gegenüber Wirkstoffen und physiologischen Gradienten ist, kann sich diese Eigenschaft unter pathophysiologischen Bedingungen wie Norovirus-Infekten oder einer chronischen Darmentzündung infolge von Zöliakie dennoch ändern. Wird etwa infolge von Zellschäden sowie kompensatorischer Zellteilung die Fläche der Mucosa größer [270] oder kleiner [251, 250], beeinflusst dies auch C^{epi} . Darüber hinaus kann die Kenntnis von C^{epi} in der klinischen Praxis aufwendige morphometrische Analysen zur Abschätzung der Oberfläche der Mucosa ersetzen. Zur Abschätzung von C^{epi} aus Impedanzspektren wurden computergestützte Auswertungen vorgeschlagen [22, 233], die jedoch durch gerätespezifische Messfehler im Hochfrequenz-Bereich erheblich erschwert werden.

In der vorliegenden Arbeit wurde C^{epi} mit einer durchschnittlichen Abweichung von weniger als zehn Prozent vom Zielwert bestimmt (Abschnitt 8.2.2). Ähnlich wie in vorangehenden Arbeiten [232], wiesen bei dieser Regressionaufgabe Random Forests sowohl im Mittel (2,4 Prozent) als auch im Maximum (90,9 Prozent) eine niedrigere Fehlerate auf als künstliche neuronale Netze (Mittel: 6,6 Prozent; Maximum: 180,3 Prozent). Die Random-Forests-Auswertung zeigt damit eine Genauigkeit, die mit früheren Arbeiten hierzu vergleichbar ist. So berichteten Bertrand et al. eine Bestimmung von C^{epi} mittels computergestütztem Fitting mit Abweichungen von ± 1 Prozent [22], Schmid et al. bestimmten C^{epi} durch Kombination von Modellierung und maschinellem Lernen und erzielten Abweichungen von bis zu 37,1 Prozent [232].

Im Unterschied zu diesen beiden Vergleichsarbeiten wird hier jedoch aufgrund des konstruktivistischen maschinellen Lernens die Gültigkeit der C^{epi} -Quantifizierung explizit und automatisiert ermittelt. Im Ergebnis ist erkennbar, dass dieser Parameter für einen Teil der Trainings- bzw. Testdaten nicht unabhängig voneinander mit der geforderten Genauigkeit bestimmt werden kann (Abschnitt 8.2). Umgekehrt ist das Verfahren aufgrund der automatisierten Modellierung nicht auf eine einzelne Zelllinie mit zwei [232] bzw. drei [22] funktionalen Zuständen beschränkt. Einerseits wurde die Quantifizierung von C^{epi} hier für drei Zelllinien mit je drei bzw. vier funktionalen

Zuständen adaptiert, was einen erheblich allgemeineren Ansatz darstellt. Andererseits ermöglichen die Prinzipien eines konstruktivistischen maschinellen Lernens jederzeit die Ergänzung des prozeduralen Wissens durch Adaption weiterer Zelllinien und funktionaler Zustände.

9.4. Generalisierbarkeit des Analyseverfahrens

Die hier vorgeschlagene Methode zur Unterscheidung von epithelialen Zelllinien und zur Quantifizierung ihrer Kapazität beruht neben einem konstruktivistischen maschinellen Lernen auch auf einer detaillierten mathematischen Modellierung. Während das Lernverfahren bei entsprechender Konfiguration auch auf andere Fragestellungen angewendet werden kann, ist eine Übertragung der hier beschriebenen Modellierung von Impedanzspektren und der Ableitung von epithelialen Eigenschaften auf andere Fragestellungen der epithelialen Impedanzanalyse sowie die konkrete Anwendung auf reale Messdaten nur unter bestimmten Voraussetzungen möglich.

Derzeit ist das Analyseverfahren spezifisch für den am Institut für Klinische Physiologie der Charité Berlin verwendeten Messaufbau definiert. Grund dafür ist der für jedes Messgerät typische gerätespezifische Messfehler. Für eine Erweiterung auf andere Messaufbauten ist daher ein individuelles gerätespezifisches Fehlermodell zu erstellen, mit dem analog zur hier verwendeten Modellierung real- und imaginärteilspezifische Fehlerwahrscheinlichkeiten errechnet werden können. Der in der vorliegenden Arbeit verfolgte Modellierungsansatz erfordert dabei lediglich eine zweistellige Anzahl an Leerkammer-Messungen bzw. Messungen an Referenz-Widerständen, auf deren Basis dann Regressionsanalysen durchgeführt werden. Ist für ein gegebenes Mess-Setup eine ausreichend detaillierte Fehlermodellierung gefunden worden, ist die Annahme der Gültigkeit für alle zugehörigen gemessenen und synthetisierten Impedanzspektren plausibel.

Unter Berücksichtigung des hier genutzten Ansatzes, insbesondere der Modellierung des gerätespezifischen Messfehlers, kann das Analyseverfahren grundsätzlich auf weitere Zelllinien ausgedehnt werden. Voraussetzung dafür ist die Identifikation belastbarer Parametergrenzen für jede neue Zelllinie. Dies erfordert analog zum hier beschriebenen Vorgehen eine ausreichende Anzahl an Messdaten sowie einen Abgleich zwischen gemessenen und modellierten Impedanzspektren. Mithilfe der im Rahmen dieser Arbeit etablierten Methodik kann dieser Abgleich weitgehend automatisiert erfolgen. Auch der Einfluss der untersuchten Wirkstoffe EGTA und Nystatin wurde bislang nur für HT-29/B6-, IPEC-J2- und MDCK-I-Zellen modelliert und bestätigt. Eine Übertragung auf weitere Zelllinien würde hier ebenfalls eine belastbare Identifikation der Modellparameter erfordern. Insgesamt ist das Vorgehen bei der Modellierung der Parameter übertragbar, nicht jedoch die hier konkret modellierten Parametergrenzen.

Darüber hinaus ist eine Erweiterung des Analyseverfahrens auf weitere Fragestellungen vorstellbar. So könnten beispielsweise analog zur Bestimmung der epithelialen Kapazität mittels Adaption prozeduralen Wissens Modelle identifiziert werden, mit denen das Verhältnis zwischen apikaler und basolateraler Zeitkonstante intersubjektiv nachvollziehbar quantifiziert werden kann. Dieses Verhältnis ist ein diagnostischer Parameter von pathophysiologischer Relevanz, der sich nicht unmittelbar aus impedanzspektroskopischen Messungen ablesen lässt [233]. Eine andere klinische relevante Erweiterung wäre eine automatisierte Unterscheidung zwischen dichten und lecken Epithelien, die analog zur Adaption konzeptuellen Wissens durchgeführt werden könnte.

9. *Epitheliale Impedanzanalyse*



10 Implementierungs- und Verfahrensaspekte

Der Einsatz des vorgeschlagenen konstruktivistischen maschinellen Lernens in den untersuchten Anwendungen (Kapitel 7 und 8) erlaubt neben einer Bewertung der Anwendungsergebnisse (Kapitel 9) auch eine Diskussion der Umsetzung. In diesem Kapitel werden daher die zentralen Designentscheidungen zur Automatisierung und Parallelisierung von Lernverfahren sowie zum Lernen in und zur Organisation von Wissensdomänen diskutiert.

10.1. Automatisierung von unüberwachtem und überwachtem Lernen

Die Realisierung des hier vorgeschlagenen Verfahrens setzt voraus, dass alle eingebundenen Lernalgorithmen automatisiert ausgeführt werden. Dies erfordert jedoch implizit eine automatisierte Bestimmung geeigneter Lernparameter und stellt damit eine besondere Herausforderung dar. Dieser Herausforderung wird beim Training überwachter Lernverfahren mit anderen Strategien begegnet als beim Training unüberwachter Verfahren.

Unüberwachtes Lernen. Obwohl die zielführende Nutzung unüberwachter Lernverfahren in der Praxis häufig eine Herausforderung darstellt, ist deren Automatisierung hier aufgrund des Gesamtdesigns vergleichsweise einfach realisierbar. Im Konstruktionsprozess werden unüberwachte Verfahren eingesetzt, um eine Menge konkurrierender maschineller Modelle zu erzeugen. Dazu wird jeder eingesetzte Algorithmus mit einer benutzerdefinierten Anzahl unterschiedlicher Konfigurationen ausgeführt (vgl. Abschnitt 5.4), deren Ergebnisse dann im Rekonstruktionsprozess weiterverarbeitet werden. Die Erzeugung dieser Konfigurationen ist dabei leicht realisierbar, insbesondere bei der Erzeugung konzeptueller Modelle mittels k-Means und prozeduraler Modelle mittels k-Means-basiertem Feature-Clustering. In beiden Fällen ist als einziger Konfigurationsparameter die Anzahl der zu identifizierenden Cluster k anzugeben, die entsprechend einer Nutzervorgabe variiert wird. Da selbstorganisierende Karten und der hier eingesetzte Sparse Autoencoder mehrere Parameter benötigen, wurde jeweils eine Standardkonfiguration definiert, die nur in der Anzahl der zu erzeugenden Outputs variiert wurde. Für selbstorganisierende Karten wurde eine eindimensionale Topologie vorgegeben, für welche die Anzahl der Knoten variiert wurde. Für den Autoencoder wurden die Standardparameter des *R*-Package *autoencoder* übernommen und

10. Implementierungs- und Verfahrensaspekte

die Anzahl der versteckten Knoten entsprechend der Benutzervorgabe variiert. Motiviert ist dies durch gute Praxiserfahrungen mit beiden Konfigurationen sowie einen angestrebten Ausgleich zwischen optimalem Ergebnis und Minimierung des Rechenaufwands.

Überwachtes Lernen. Ziel des Rekonstruktionsprozesses ist es zunächst, für ein gegebenes Modell mit jedem der eingesetzten überwachten Lernverfahren ein optimales Trainingsergebnis zu erzielen. Da gleichzeitig aber auch eine automatisierte Bewertung des Trainings erfolgen soll, werden hierfür bei Regressionen die mittlere prozentuale Abweichung vom Zielwert und bei Klassifikationen die so genannte Accuracy verwendet (siehe auch Abschnitt 3.5). Mithilfe eines nutzerspezifischen Schwellenwerts wird so zunächst entschieden, ob das Training eines Modells erfolgreich bzw. welche Konfiguration des jeweiligen Algorithmus am erfolgreichsten war. Anschließend wird auf Basis der Intercoder-Reliabilität das intersubjektivste beziehungsweise am erfolgreichsten rekonstruierte Modell ausgewählt und die übrigen Modelle verworfen.

Die Auswahl geeigneter bzw. optimaler Parameter ist für Random Forests leichter umsetzbar als für künstliche neuronale Netze. Denn die Fehlerrate von Random Forests hängt neben der Auswahl geeigneter Features bzw. Samples für die einzelnen Bäume insbesondere von deren Anzahl ab. Mit einer zunehmenden Zahl von Bäumen sinkt die Fehlerrate von Random Forests, wobei diese gleichzeitig gegen einen festen Wert konvergiert [30]. Somit kann durch Wahl einer ausreichend hohen Anzahl von Bäumen sichergestellt werden, dass ein Random Forest eine minimale Fehlerrate aufweist. In der vorliegenden Arbeit wurde dafür auf Basis praktischer Erfahrungen eine Menge von 500 Bäumen pro Random Forest festgelegt.

Für Feedforward-Netze mit Backpropagation hingegen sind nicht nur eine Netzarchitektur, sondern auch mehrere Lernparameter wie die Weite der Lernschritte oder die Anzahl der Lerniterationen zu definieren. Für die Nutzung des Rprop-Verfahrens spricht daher, dass dieser eine automatisierte und gleichzeitig hoch effektive Optimierung der Lernschrittweite realisiert. Gleichzeitig handelt es sich dabei um einen der am schnellsten konvergierenden Trainingsalgorithmen für künstliche neuronale Netze [210, 107]. Aufgrund dieser Eigenschaften wurde hier das Verfolgen einer Early-Stopping-Strategie als plausibel betrachtet, bei der das Rprop-Training nach jeweils 2.500 Schritten automatisch beendet wird (vgl. z.B. [197, 37, 271]).

Zentrales zu optimierendes Merkmal der verwendeten Feedforward-Netze bleibt somit deren Architektur. Diese wird sowohl durch die Modellkomplexität κ als auch durch die eindimensionale Ausgabe bereits eng begrenzt. Anders als im allgemeinen Fall begrenzt dies gemeinsam mit der Beschränkung auf eine einzelne versteckte Schicht mit $\lfloor n \cdot 1,5 \rfloor$ Neuronen (vgl. Abschnitt 5.5) die Anzahl der möglichen Architekturen auf eine endliche Menge niedriger Mächtigkeit. Folglich ist eine Evaluation aller in Frage kommenden Architekturen im Sinne einer Brute-Force-Strategie in endlicher Zeit möglich. Gegenüber klassischen Netzwerkarchitekturalgorithmien wie FlexNet [171] hat dies nicht nur den Vorteil, dass eine geringere Wahrscheinlichkeit für Overfitting besteht, sondern insbesondere, dass eine effiziente Parallelisierung möglich ist.

Netzwerk-Konstruktion. Bereits in den 1980er und 1990er Jahren wurden Lernverfahren vorgeschlagen, die nicht nur die Bedeutung von Verknüpfungen mittels Gewichtsanzpassung erlernen, sondern auch deren Struktur selbst verändern können. Bekannte Vertreter solcher Algorithmen zur Konstruktion von Netzwerkarchitekturen waren etwa Cascade Correlation [61] oder FlexNet [171], die innerhalb eines gegebenen künstlichen neuronalen Netzes Knoten hinzufügen oder entfernen konnten. Im Gegensatz zum vorliegenden Ansatz handelt es sich dabei jedoch um Verfahren, die ausschließlich nach den Regeln eines überwachten Lernens funktionieren. So beruht insbesondere bei FlexNet die Bewertung von Strukturveränderungen auf dem Nutzen, den diese zur Vorhersage der gegebenen Zielwerte bringen. Konstruktivistisches maschinelles Lernen hingegen erlaubt auch die Exploration unbekannter Datensätze im Sinne eines Lernens ohne Zielparameter. Darüber hinaus erlauben klassische Netzwerkkonstruktionsalgorithmen die Berücksichtigung von Metadaten ebenso wenig wie klassische überwachte Lernverfahren.

10.2. Abstraktion und Wissensebenen

Eine Wissensdomäne ist hier in Form einer hierarchischen Ordnung von Modellen realisiert. Dies ist motiviert durch die Annahme unterschiedlicher Abstraktionsebenen menschlicher Kognition (vgl. Abschnitt 2.4), weist aber auch Analogien zur Schichten-Anordnung von Neuronen wie etwa im visuellen Kortex auf. In dieser Hinsicht stehen die didaktisch motivierten Wissens-ebenen nach Bloom und der hier verfolgte konstruktivistische Ansatz nicht im Widerspruch zu neurophysiologischen Strukturen. Gleichzeitig ermöglicht die Einführung einer solchen Modell-hierarchie unter dem gegebenen Lernparadigma überhaupt erst die Abbildung komplexerer Zusammenhänge. Denn Konstruktion und Rekonstruktion sind in der vorliegenden Implementierung darauf ausgelegt, Modelle möglichst geringer Komplexität zu erzeugen.

Abstraktionskaskaden. Zu Beginn der Exploration konzeptuellen Wissens enthält die Wissensdomäne ausschließlich maschinelle Modelle der Wissensebene 1, im weiteren Verlauf werden auch Modelle der darüberliegenden Ebenen 2 bis 9 abstrahiert und hinzugefügt (Abb. 7.1b). Dies erfolgt durch $T\Sigma$ -Dekonstruktion und kann ausgehend von einem gegebenen Lernblock zu einer Kaskade von Abstraktionen führen, in deren Folge Modelle mehrerer Ebenen erzeugt werden können (Abb. 10.1). Dies ist bedingt durch den grundsätzlichen Ablauf des konstruktivistischen Lernens. Denn mittels $T\Sigma$ -Dekonstruktion werden die Ausgabewerte eines Modells der Ebene n zu Eingabewerten eines neuen Lernblocks der Ebene $n+1$ mit dann unbekanntem Zielwert (vgl. Abschnitt 5.6). Dieser Lernblock der Ebene $n+1$ wiederum wird – wie jeder Lernblock ohne bekannte Zielwerte – einem Konstruktionsprozess und das konstruierte maschinelle Modell anschließend einem Rekonstruktionsprozess unterworfen (vgl. 5.2). Bestehen für ein solches rekonstruiertes Modell wiederum $T\Sigma$ -Verwandtschaften innerhalb der zugehörigen Wissensdomäne, werden diese unmittelbar geprüft. Existieren also bereits weitere $T\Sigma$ -verwandte Modelle auf Ebene $n+1$, so kann ein durch Abstraktion neu erzeugtes Modell der Ebene $n+1$, das erfolg-

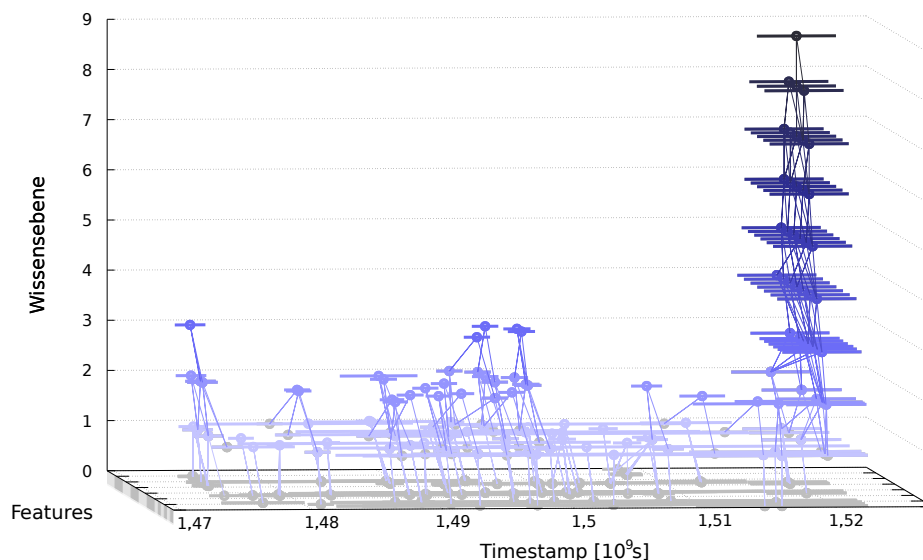


Abbildung 10.1.: Ergebnis einer Abstraktionskaskade. Dargestellt ist der temporäre Zustand der konzeptuellen Wissensdomäne nach Abschluss der Exploration von Feature-Set $G_{\Delta r}$. Die Modelle der Ebenen 4 bis 9 wurden aus dem letzten Block abgeleitet und bauen aufeinander auf. Sie umfassen je 9.725 Samples und weisen eine identische temporale Gültigkeit auf.

10. Implementierungs- und Verfahrensaspekte

reich rekonstruiert wurde, zu einer erneuten $T\Sigma$ -Dekonstruktion und gegebenenfalls zu einem neuen Modell auf Ebene $n+2$ führen. Begrenzt werden können solche aufeinander aufbauenden Abstraktionen durch Definition einer maximal zulässigen Wissensebene.

Dekonstruktionskaskaden. Während die Zahl der Modelle einer Wissensdomäne im Verlauf einer Exploration tendenziell zunimmt, kann die Zahl der damit abgedeckten Wissensebenen Schwankungen unterliegen. So werden etwa bei Exploration der konzeptuellen Domäne nach Training mit dem Feature-Set $G_{\Delta r}$ neun Ebenen abgedeckt, während es nach dem anschließenden Training mit Feature-Set $G_{\Delta \phi}$ nur noch drei Ebenen sind (Abb. 7.1b). Bei diesem Effekt handelt es sich um eine implizite Folge des hier eingeführten Lernens durch Dekonstruktion, mit dem Modelle der Wissensdomäne verändert oder gegebenenfalls gelöscht werden (vgl. Abschnitt 5.6). Denn diese Modifikationen betreffen implizit auch Domänenmodelle, die die Ausgabe eines veränderten oder gelöschten Modells als Eingabe verwenden. Bei Löschung eines Modells können solche Abhängigkeiten nur durch Löschung aller abhängigen Modelle sinnvoll aufgelöst werden. In der aktuellen Implementierung ist diese Strategie jedoch auch für Änderungen umgesetzt, mit denen existierende Modelle der Wissensdomäne erweitert werden. Dadurch wird ein Modell \mathcal{M}_i nach einer temporal erweiternden Modellvereinigung gelöscht und das vereinigte Modell als neues Modell \mathcal{M}_j gespeichert. Diese Strategie stellt sicher, dass die Wissensdomäne zu keinem Zeitpunkt Modelle ohne Verbindung zur Datenbasis enthält. Andererseits kann dies jedoch zu einer Dekonstruktionskaskade führen, in deren Verlauf zahlreiche bereits identifizierte Modelle höherer Ebenen wieder gelöscht werden. Um diesen tendenziell destruktiven Vorgang abzufedern, wird daher als unmittelbare Folge von erfolgreichen vollständigen oder ΣZ -Dekonstruktionen eine Prüfung möglicher $T\Sigma$ -Dekonstruktionen forciert, die im Idealfall zu einer Abstraktionskaskade führt und die Dekonstruktionskaskade in umgekehrter Richtung wieder aufbaut.

Prozedurale Abstraktion. Anders als für die konzeptuelle Wissensdomäne konnten für die prozedurale im Verlauf der Exploration keine Modelle oberhalb der ersten Wissensebene rekonstruiert werden (Abb. 8.1b). Nach Abschluss der Exploration befinden sich somit sämtliche Domänenmodelle auf Ebene 1 (Abb. 8.3). Die Konstruktion bzw. Rekonstruktion von prozeduralen Modellen der Ebene 2 aus den Zielwerten von Modellen der Ebene 1 schlägt dagegen fehl. Wie Abb. 8.2 zeigt, wird eine solche Modellkonstruktion auf Basis der Modelle der Ebene 1 durch die zum Teil sehr unterschiedlichen Wertebereiche der jeweiligen Zielwerte erschwert. Diese Werte werden im Fall der hier eingesetzten Autoencoder deren versteckten Neuronen entnommen, ihre Streuung und Wertebereichsgrenzen stellen also Repräsentationen der Eingabedaten dar. Eine Nivellierung der dadurch ausgedrückten Modellunterschiede mittels Normierung auf einen festen Wertebereich für alle Modelle der Ebene 1 ist daher trotz möglicherweise positiver Effekte auf die Abstraktionsfähigkeit als problematisch zu betrachten. Eine weitere Schwierigkeit der prozeduralen Abstraktion stellt eine Designentscheidung der aktuellen Implementierung dar. Für die Konstruktion von Modellen ab Ebene 2 werden derzeit stets nur zwei Eingabefeatures der darunterliegenden Ebene verwendet, wodurch die Anzahl und Komplexität extrahierbarer Zusammenhänge ebenfalls limitiert ist. Darüber hinaus könnten Autoencoder und Feature-Clustering möglicherweise keine universell erfolgreich einsetzbaren Konstruktionsmechanismen für prozedurales Wissen darstellen. Um die Abstraktionsfähigkeit für prozedurale Modelle zu erhöhen, sollte daher in künftigen Implementierungen bzw. Erweiterungen die Nutzung weiterer unüberwachter Regressionsverfahren zur Modellkonstruktion geprüft werden.

10.3. Modellverwaltung

Mittels Dekonstruktion wird ein rekonstruiertes Modell mit Modellen der Wissensdomäne abgeglichen und dort gegebenenfalls gespeichert. Die hier gewählte Umsetzung dieser Domäne innerhalb des Dateisystems des Betriebssystems ermöglicht eine intuitive und effiziente Einbindung der

beteiligten Lernalgorithmen (siehe Abschnitte 5.4 und 5.5). Sie erfordert jedoch auch eine explizite Verwaltung der Datenstrukturen und weist einen vergleichsweise hohen Speicherbedarf auf.

Fragmentierung. Der Dekonstruktionsprozess erlaubt entweder eine Überprüfung aller verwandten Modelle oder einen Abbruch nach der ersten erfolgreich durchgeführten Dekonstruktion (vgl. Anhang A). Ein solcher Abbruch wird hier als minimale Dekonstruktion bezeichnet und ist die Standardeinstellung in der aktuellen Implementierung. Motiviert ist dies durch den Zeit- und Rechenaufwand, den eine vollständige Überprüfung erfordert. Eine minimale Dekonstruktion stellt somit einen Kompromiss zwischen minimalem Rechenaufwand und optimalem Ergebnis dar. Sie kann allerdings zu einer Art Fragmentierung der Wissensdomäne führen, etwa weil ΣZ -verwandte Modelle nicht überprüft werden, deren Vereinigung zu einem neuen Modell mit größerer temporaler Ausdehnung führen würde. Zu dieser Modellvereinigung kommt es nur dann, wenn sie in einem nachfolgenden Dekonstruktionsprozess überprüft werden. Um das Abstrahieren von Modellen hoher temporaler Ausdehnung zu begünstigen, wird die minimale Dekonstruktion künftig entweder eliminiert oder um eine Art heuristische Garbage Collection erweitert, mit der am Ende jeder Iteration pragmatische Verwandtschaften gesammelt überprüft werden.

Speicherbedarf. Mit jedem Modell werden standardmäßig sowohl Konfigurations- als auch Trainingsdaten aller beteiligten Lernverfahren in der Wissensdomäne gespeichert. Mit der Zahl der Domänenmodelle steigt daher auch der Speicherbedarf. So benötigt etwa die konzeptuelle Wissensdomäne nach Exploration rund 400 Megabyte, was rund 75 Prozent der Trainings-Datenbasis entspricht (Abb. 10.2). Um den Speicherbedarf zu reduzieren, können statt der Konfigurationsdaten aller nur die eines Lernverfahrens gespeichert werden, hier also entweder eines neuronale Netzes oder eines Random Forests. Größeren Anteil am Speicherbedarf haben jedoch die Trainingsdaten. Während deren Speicherung in den Domänenmodellen der Ebene 1 noch eine redundante Speicherung der Datenbasis darstellt, handelt es sich auf den höheren Ebenen um abgeleitete, neue Daten. Ein Verzicht auf deren Speicherung wäre daher nur dann möglich, wenn die Trainingsdaten sämtlicher Domänenmodelle bei Bedarf aus der Datenbasis abgeleitet werden, was zwar realisierbar, aber wenig effizient wäre. Um künftige Implementierungen hinsichtlich des Speicherbedarfs zu optimieren, bietet sich daher am ehesten die Nutzung eines Datenbank-Managementsystems an.

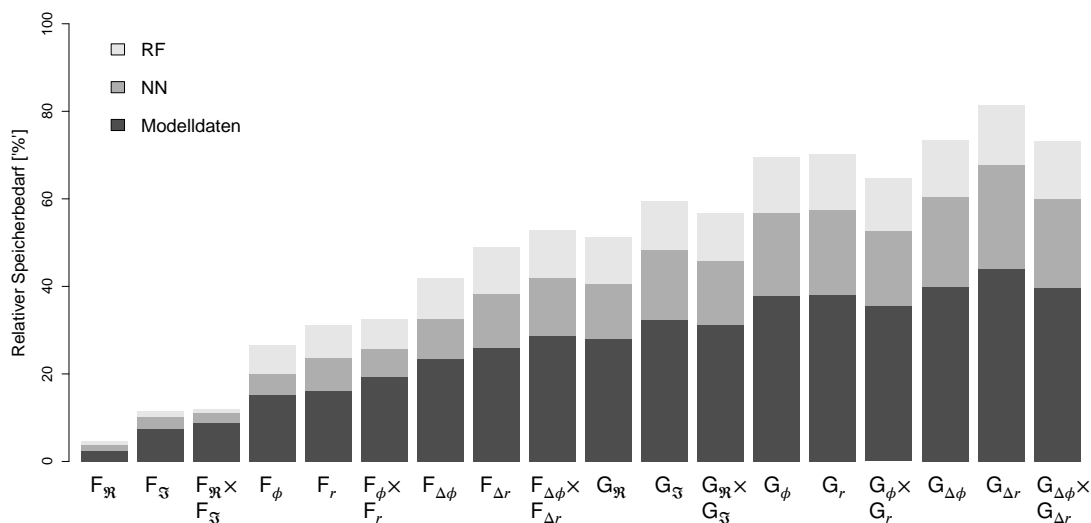


Abbildung 10.2.: Entwicklung des Speicherbedarfs der konzeptuellen Wissensdomäne im Verlauf der Exploration. Dargestellt ist der nach Training mit dem jeweiligen Feature-Set benötigte Speicherplatz im Verhältnis zur Größe der verwendeten Datenbasis.

10.4. Parallelisierung

Das hier vorgeschlagene Verfahren eines konstruktivistischen maschinellen Lernens setzt als einen zentralen Mechanismus auf einen Vergleich konkurrierender Lernverfahren und ist darüber hinaus an einigen Stellen mit einer Art Exhaustions- oder Brute-Force-Suche vergleichbar. Während ein solches Konzept auf den ersten Blick Redundanz und erhöhten Rechenaufwand erwarten lässt, erlaubt es jedoch gleichzeitig, eine Reihe intuitiver Parallelisierungen zu implementieren. Zwar ist die inherent konsekutive Abfolge von Konstruktions-, Rekonstruktions- und Dekonstruktionsprozess nicht sinnvoll auflösbar. Jedoch bieten diese Teilprozesse selbst eine Reihe von Ansatzpunkten für Single-Instruction-Multiple-Data- (SIMD) und Multiple-Instruction-Single-Data-Parallelisierungen (MISD) nach Flynn [65]. Um einen Kompromiss zwischen effizienter Ausführung und effizienter Implementierung zu erreichen, wurden hierfür Prozess-Parallelisierungen genutzt, die eine Parallel-Ausführung auch auf Standard-Desktop-Rechnern erlauben. Aufgrund des damit bereits erreichten Speedups wurde auf weitergehende Parallelisierungen, etwa mittels Rechnerclustern oder unter Ausnutzung von GPUs, verzichtet.

Die Konstruktion neuer Modelle aus einem gegebenen Lernblock erfolgt grundsätzlich unabhängig voneinander (vgl. Abschnitt 5.4). Im Sinne einer MISD-Parallelisierung können daher beispielsweise SOM und k-Means parallel ausgeführt werden, um konzeptuelle Modelle zu konstruieren. Da aber in beiden Fällen k Modelle gesucht werden, die aus bis zu k Clustern besteht, kann eine MISD-Parallelisierung auch auf k-Means selbst bzw. auf k verschiedene SOMs angewendet werden. In der hier genutzten R -Umgebung konnte dieser Ansatz mittels des *dopar*-Befehls sowie des Packages *doMC* schnell und effizient umgesetzt werden. Analog zur Konstruktion konzeptueller Modelle lässt sich auf diesen beiden Ebenen auch für die zur Konstruktion prozeduraler Modelle eingesetzten Autoencoder- und Feature-Clustering-Verfahren MISD-Parallelisierung realisieren. Bei einer Erweiterung des Konstruktionsprozesses um ein zusätzliches unüberwachtes Lernverfahren wäre dieses ebenfalls unabhängig von den bestehenden und könnte daher parallel zu k-Means und SOM ausgeführt werden.

Der Rekonstruktionsprozess bietet aufgrund des Ablaufs ebenfalls eine Reihe von Ansatzpunkten für eine parallelisierte Implementierung. So legt die Ausführung voneinander unabhängiger überwachter Lernverfahren auf dem gleichen Datensatz (vgl. Abschnitt 5.5) erneut eine MISD-Parallelausführung nahe. Darüber hinaus sind auch die Algorithmen selbst wiederum parallelisiert implementiert. Im Falle der in R implementierten Random Forests etwa ermöglicht der *dopar*-Befehl eine SIMD-Parallelisierung jeder Instanz des Algorithmus. In Verbindung mit dem Package *doMC* lassen sich dadurch die strukturell voneinander unabhängigen Entscheidungsbäume effizient über alle verfügbaren Prozessorkerne verteilen, wodurch das Random-Forest-Verfahren unmittelbar von einer höheren Anzahl an Kernen profitiert. Für die als C-Binaries genutzten künstlichen neuronalen Netze kann mittels einfacher Skript-Forks eine MISD-Parallelisierung realisiert werden, da diese eigene Threads darstellen und somit über die verfügbaren Prozessorkerne verteilt werden können. Die zu Beginn eines Rekonstruktionsprozesses ausgeführte hybride Komplexitätsreduktion wird nur parallel ausgeführt, wenn Random Forests zur Feature-Auswahl verwendet werden; das alternativ verwendete CFS bietet keine effizienten Ansatzpunkte für eine Parallelisierung. Die abschließende Modellreihung ist ein im wesentlichen sequentieller Vorgang und bietet ebenfalls keine Ansatzpunkte für eine Thread-Parallelisierung.

Mit den implementierten Parallelisierungen lassen sich bei entsprechender Rechnerkapazität sowohl bei überwachtem wie unüberwachtem Lernen deutliche Performanzgewinne erzielen. Vergleicht man etwa die Exploration einer konzeptuellen Wissensdomäne auf einem 4-Kern-Prozessor mit 8 parallelen Threads mit der Ausführung auf einem 16-Kern-Prozessor mit 32 parallelen Threads, so ergibt sich ein Gesamt-Speedup von rund 2,5. Wie Abb. 10.3 zeigt, variiert der Speedup jedoch pro Feature-Set. Eine Analyse der tatsächlich ausgeführten Lernoperationen

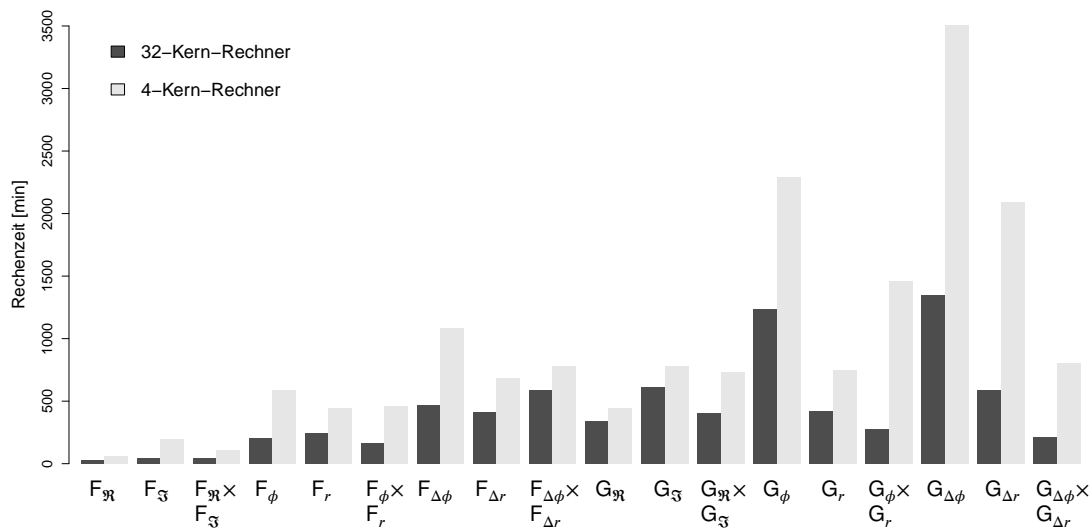


Abbildung 10.3.: Rechenzeit für die Exploration der konzeptuellen Wissensdomäne. Für jedes der 18 Feature-Sets ist die tatsächlich verstrichene Zeit bei Ausführung auf einem 4-Kern-Rechner einerseits und einem 16-Kern-Rechner andererseits dargestellt.

zeigt, dass hier insbesondere die Zahl der Dekonstruktionsoperationen variiert (Abb. 7.1a). Dies wiederum erklärt sich durch die stetig wachsende Zahl von Modellen in der Wissensdomäne (Abb. 7.1b), mit der auch die Zahl der zu prüfenden pragmatischen Verwandtschaften zunimmt. Diese Verwandtschaftsprüfungen wurden für den hier erstmals umgesetzte Dekonstruktionsprozess jedoch noch nicht parallelisiert. In der aktuellen Implementierung wird der Dekonstruktionsprozess nur dort parallelisiert durchgeführt, wo er auf Konstruktions- oder Rekonstruktionsprozesse, also unüberwachte oder überwachte Lernverfahren, zurückgreift (vgl. Abschnitt 5.6).

Da der Dekonstruktionsprozess somit im Sinne des Amdahlschen Gesetzes den seriellen Flaschenhals des Verfahrens darstellt [7], sind von dessen Optimierung in künftigen Implementierungen weitere deutliche Performanzgewinne zu erwarten. Potenzial für Parallelisierungen bieten hier zunächst Lernblöcke, aus deren Rekonstruktionen mehrere voneinander unabhängige Modelle in den Dekonstruktionsprozess überführt werden. Da dies einen jedoch eher einen Sonderfall darstellt, ist der davon zu erwartende Speedup vergleichsweise begrenzt. Deutlichere Performanzgewinne sind hingegen von einer Parallelisierung der bislang sequentiellen Ausführung von ΣZ -, $T\Sigma$ -, TZ - und vollständiger Dekonstruktion für ein gegebenes rekonstruiertes Modell zu erwarten. So könnten beispielsweise $T\Sigma$ -Verwandtschaften parallel geprüft werden, da diese voneinander unabhängig sind und sich nicht gegenseitig beeinflussen.

10. Implementierungs- und Verfahrensaspekte



11

Chancen und Risiken konstruktivistischen Lernens

Das hier vorgeschlagene Lernverfahren folgt einem anderen Paradigma als bestehende Verfahren. Im Gegensatz zu biologisch motivierten Algorithmen dienen hier Denk- statt neurophysiologische Strukturen als Vorbild, entscheiden Metadaten mit über den Lernfortschritt und spielen Zweifel eine zentrale Rolle. Darüber hinaus definiert dieser neue Ansatz Grenzen der Generalisierbarkeit und macht diese für jedes erzeugte Modell transparent. In diesem Kapitel werden Vor- und Nachteile dieses neuen Ansatzes erörtert sowie noch offene Fragen skizziert.

11.1. Umgang mit Mehrdeutigkeiten

In empirischen Daten ist Mehrdeutigkeit ebenso allgegenwärtig wie in der menschlichen Wahrnehmung. Das hier vorgeschlagene Verfahren ermöglicht erstmals einen expliziten Umgang mit Datensätzen, die mehrdeutig interpretiert werden können. Denn die Einführung des Konzepts maschineller Modelle auf Basis der Allgemeinen Modelltheorie nach Stachowiak eröffnet neben der bloßen Feststellung, ob ein Algorithmus einen Zusammenhang zwischen Eingabe und Ausgabe erlernt hat oder nicht, weitere Bewertungskriterien. Für ein Kippbild etwa würden sich damit anhand der pragmatischen Eigenschaften T , Σ und Z mehrere Modelle mit unterschiedlicher Gültigkeit konstruieren lassen (vgl. Kapitel 1). Rekonstruktions- und Dekonstruktionsprozess dagegen identifizieren und eliminieren gezielt solche Modelle, die im Widerspruch zu anderen stehen und als Mehrdeutigkeiten interpretiert werden könnten. Während der Rekonstruktionsprozess dies über den Weg der Intersubjektivität erreicht, identifiziert und korrigiert der Dekonstruktionsprozess Widersprüche zu bereits in einer Wissensdomäne gespeicherten Modellen.

Voraussetzung dafür, Mehrdeutigkeit innerhalb einer gegebenen Datenbasis mittels konstruktivistischen maschinellen Lernens auflösen zu können, ist eine initiale Erhebung der pragmatischen Eigenschaften T , Σ und Z für jedes darin enthaltene vektorielle Modell, d.h. für jeden Trainings- oder Testdatenvektor. Dies erfordert, für jedes gemessene Sample Erhebungszeitpunkt, Kennzeichnung des Erhebers sowie gegebenenfalls Erhebungszweck zu dokumentieren. In der vorliegenden Arbeit ist dies aufgrund der durch Simulation erzeugten Datenbasis leicht umsetzbar. Aber auch für künftige Anwendungen stellt diese Voraussetzung keine Hürde dar. So ist etwa in modernen Steuerungs- oder Messgeräten, die auf einem eingebetteten System basieren, die Speicherung des

11. Chancen und Risiken konstruktivistischen Lernens

zugehörigen Timestamps in einigen Anwendungen bereits üblich. Eine Speicherung der Geräte-Seriennummer und einer Messgrößen-Kennung ist ebenfalls in vielen industriellen Anwendungen die Regel, so dass in der Praxis letztlich nur noch ein gezieltes Mapping dieser pragmatischen Eigenschaften auf ein dreidimensionales Metadatenset $\{T, \Sigma, Z\}$ erforderlich wäre, das jedes Sample bzw. vektorielle Modell eindeutig beschreibt.

Die Ablaufsteuerung des konstruktivistischen Lernens erfolgt an ihren zentralen Punkten in Abhängigkeit von den pragmatischen Eigenschaften T , Σ und Z . Diese Eigenschaften bestimmen den Verwandtschaftsgrad zwischen Modellen (vgl. Abschnitt 5.2) und beeinflussen das Lernen bereits vor Ausführung von Konstruktions-, Rekonstruktions- oder Dekonstruktionsprozessen. Grund ist die Nutzung dieser Verwandtschaften zur Identifikation von Lernblöcken aus einer gegebenen Menge vektorieller Modelle, wodurch jedes aus einem Lernblock abgeleitete maschinelle Modell über eindeutig definierte pragmatische Eigenschaften verfügt. Im Kontext des anschließenden Konstruktionsprozesses werden aus einem gegebenen Lernblock gezielt konkurrierende Modelle mit identischer temporaler Gültigkeit T , aber unterschiedlichem Zweck Z konstruiert. Ein Betrachter Σ ist für ein solches Modell zu diesem Zeitpunkt formal noch nicht definiert.

Ziel des Rekonstruktionsprozesses ist sowohl die Subjektivierung eines konstruierten Modells mittels überwachten Lernens als auch die Eliminierung dadurch erzeugter Mehrdeutigkeiten. Wurde ein konstruiertes Modell \mathcal{M} mit definierten pragmatischen Eigenschaften $T(\mathcal{M})$ und $Z(\mathcal{M})$ von einem überwachten Lernverfahren L_i entsprechend des jeweiligen Fehler-Schwellenwerts erfolgreich erlernt, so gilt $\Sigma(\mathcal{M}) = \{L_i\}$ mit $i \in \{0, \dots, n - 1\}$; in der vorliegenden Untersuchung gilt aufgrund der parallelen Verwendung eines künstlichen neuronalen Netzes und eines Random Forests $n=2$. Um subjektbezogene Mehrdeutigkeiten zu eliminieren, werden die Ergebnisse des überwachten Lernens miteinander verglichen. Die Nutzung von Krippendorffs α als Inter-Kodierer-Reliabilitätskoeffizient erlaubt dabei eine Bewertung der Intersubjektivität. Als mehrdeutig – und damit verwerfbar – werden solche Modelle interpretiert, die ein benutzerdefiniertes Minimum des Reliabilitätskoeffizienten nicht überschreiten. Von den verbleibenden Modellen wird dasjenige mit der höchsten Inter-Kodierer-Reliabilität als optimale Repräsentation des Lernblocks selektiert. Der Rekonstruktionsprozess stellt somit sicher, dass aus einem gegebenen Lernblock eine eindeutig definierte und intersubjektiv nachvollziehbare Repräsentation ausgewählt wird.

Einmal rekonstruierte Modelle speichert das hier vorgeschlagene Verfahren in einer Wissensdomäne. Dabei stellt der Dekonstruktionsprozess sicher, dass infolge dieser Speicherung keine Mehrdeutigkeiten innerhalb der Domäne entstehen (vgl. Abschnitt 5.6). Dies wäre etwa der Fall, wenn die Wissensdomäne zwei oder mehr vollständig verwandte maschinelle Modelle enthalten würde, die für ein enthaltenes vektorielles Modell unterschiedliche Zielwerte vorsehen. Ein Widerspruch ergäbe sich aber auch dann, wenn vollständig verwandte Modelle unterschiedliche Wertebereiche aufweisen würden. Für jede vollständige Verwandtschaft eines neuen Modells mit einem Modell der Wissensdomäne wird daher nicht nur geprüft, ob eine Modellvereinigung möglich ist, sondern bei deren Scheitern auch, ob eine Modelldifferenzierung eine Rekonstruktion erlaubt (vgl. Abschnitt 5.6.4). Ist beides unmöglich, wird dies als Mehrdeutigkeit interpretiert und sowohl das neue als auch das in der Wissensdomäne gespeicherte Modell gelöscht.

Insgesamt reduzieren diese Mechanismen uneindeutige und fehlerhafte Interpretationen, die durch mehrdeutige Datensätze entstehen können. Während die Berücksichtigung pragmatischer Eigenschaften dazu dient, Modelle formal eindeutig zu unterscheiden, wird mittels Rekonstruktion eine ausreichende Intersubjektivität und mittels Dekonstruktion eine widerspruchsfreie Speicherung der Modelle sicher gestellt. Soweit eine Datenbasis Mehrdeutigkeiten enthält, werden diese somit identifiziert und exkludiert. Bei der Evaluation von Testdaten unter Verwendung einer bestehenden Wissensdomäne stellt die Berücksichtigung pragmatischer Verwandtschaften darüber hinaus sicher, dass nur eindeutige Vorhersagen gemacht werden. Daten, für welche die Wissensdomäne keine pragmatisch verwandten Modelle enthält, werden daher exkludiert.

11.2. Nachvollziehbarkeit

In modernen Anwendungen maschinellen Lernens stellen transparente Strukturen und nachvollziehbare Ergebnisse zentrale Anforderungen dar. Das hier vorgeschlagene Verfahren wird diesen Anforderungen dank konstruktivistischem Ansatz und pragmatischem Modellbegriff in mehrfacher Hinsicht gerecht. Diese Nachvollziehbarkeit wird dadurch ermöglicht, dass sich nicht nur die Lernabläufe selbst an menschlichen Denkstrukturen orientieren, sondern auch die Repräsentation des erlernten Wissens. Das vorgeschlagene Verfahren unterscheidet sich hierin von vielen etablierten überwachten und unüberwachten Verfahren, die Anwendern zwar Lösungen liefern, jedoch keinen intuitiv interpretierbaren Lösungsweg.

Die hierarchische Organisation erlernter Modelle in einer Wissensdomäne erlaubt es, deren Funktionen detailliert nachzuvollziehen. Denn durch Abbildung des Modellzwecks als pragmatische Eigenschaft Z ist beispielsweise für jedes konzeptuelle Modell unmittelbar bestimmbar, ob es sich um einen Klassifikator handelt, der zwischen zwei, drei oder vier Klassen unterscheidet. Die konzeptuelle Wissensdomäne als Ganzes wird so als Konnektom von Klassifikatoren interpretierbar, die unidirektional miteinander verknüpft sind. Eine solche Interpretation erlaubt einerseits eine intuitive Rückverfolgung des individuellen Lösungsweges eines Domänenmodells, mit der eine von diesem Modell getroffene Klassifizierung auf Basis vorgelagerter Klassifikatoren auch einzelfallbezogen erklärt werden kann. Ferner lässt sich eine konzeptuelle Wissensdomäne dadurch als leicht interpretierbares Klassifikatoren-Konnektom visualisieren (Abb. 11.1).

Die hier eingeführte Definition einer temporalem Gültigkeit erlernten Wissens macht transparent, dass zu testende Daten dahingehend unterschieden werden müssen, ob sie durch ein bereits erlerntes Modell abgedeckt sind oder nicht. Dies kann anhand des Metadatum T abgeglichen werden und für die Interpretation der Ergebnisse hilfreich sein. Während etwa Vorhersagen für zeitlich abgedeckte Testdaten im Wesentlichen eine Interpolation bekannter Trainingsdaten darstellt,

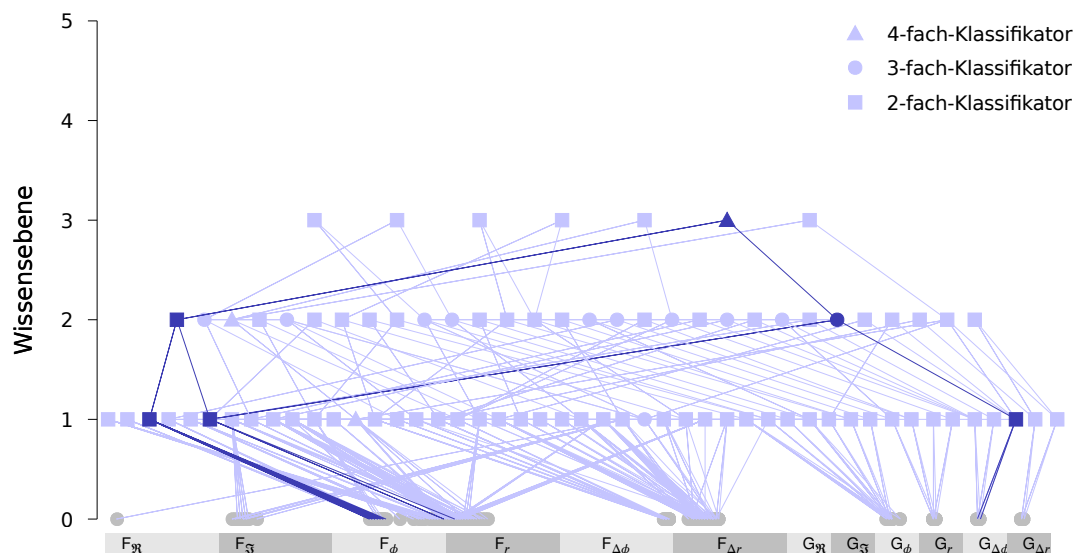


Abbildung 11.1.: Klassifikatoren-Konnektom der explorierten konzeptuellen Wissensdomäne. Der Modellzweck Z wird in einer Frontal-Ansicht auf die Domäne für jedes Modell grafisch repräsentiert. Exemplarisch farblich hervorgehoben ist der 4-fach-Klassifikator C.3.6 sowie alle Modelle der darunter liegenden Ebenen, auf denen dieser aufbaut. Auf der x-Achse sind die Feature-Sets der Datenbasis gekennzeichnet (Ebene 0).

11. Chancen und Risiken konstruktivistischen Lernens

werden Vorhersagen zeitlich nicht abgedeckter Testdaten vor diesem Hintergrund als Extrapolationen mit stärker spekulativem Charakter erkennbar. Analog dazu machen die pragmatischen Eigenschaften Σ und Z transparent, ob Vorhersagen auf Basis von subjekt- und zweckbezogener Erfahrung oder auf Basis evidenzloser Generalisierungen erfolgen. Diese neue Qualität des erlernten Wissens erlaubt insbesondere, automatisiert zu entscheiden, ob Vorhersagen für neue, dem System bislang unbekannte Testdaten durchgeführt werden sollen oder nicht.

Obwohl Modelle und Domäne beide einem kontinuierlichen Wandel durch Lernen unterliegen, ist die Entstehung der jeweils aktuellen Ausprägung vollständig nachvollziehbar und begründbar. Ermöglicht wird dies durch eine umfassende und regelorientierte Verfahrensprotokollierung. Bereits die vor Konstruktion, Rekonstruktion und Dekonstruktion stattfindende Identifikation von Lernblöcken dokumentiert nicht nur den gewählten Lernblock, sondern alle potentiellen Kandidaten. Ein weiterer Grund ist die individuelle und unveränderliche Kennzeichnung jeder Modellausprägung mittels UID. Im Rahmen einer aus mehreren Teilschritten bestehenden Exploration lässt sich somit beispielsweise Zeitpunkt und Ursache einer ΣZ -Dekonstruktion zweier Modelle \mathcal{M}_i und \mathcal{M}_j zu einem neuen Modell \mathcal{M}_k nachträglich identifizieren. Auch die Ursachen von Abstraktions- und Dekonstruktionskaskaden lassen sich bereits auf Basis der Verfahrensprotokollierung nachvollziehen, also ohne detaillierte Analyse der Daten selbst.

Nicht zuletzt ist auch die in empirischen Untersuchungen grundsätzlich geforderte intersubjektive Nachvollziehbarkeit in dem hier vorgeschlagenen konstruktivistischen maschinellen Lernen verankert. Insbesondere für kategoriale Bewertungen, also Klassifikationsaufgaben, ist sicherzustellen, dass diese nicht nur rein subjektiv erfolgen. Grundlage für die algorithmische Umsetzung dieses Prinzips ist die parallele Ausführung konkurrierender überwachter Lernverfahren während des Rekonstruktionsprozesses und die anschließende Bewertung mittels Inter-Kodierer-Reliabilitätskoeffizient. Der benutzerdefinierte Schwellenwert für den Koeffizienten Krippendorffs α stellt sicher, dass Modelle verworfen werden, die nicht in einem ausreichenden Maß intersubjektiv nachvollziehbar sind. Dieser Mechanismus führt dazu, dass die damit erzeugte Wissensdomäne ausschließlich Modelle mit hoher intersubjektiver Nachvollziehbarkeit enthält und verringert die Wahrscheinlichkeit, dass Modelle aus nicht-kausalen Korrelationen darin integriert werden.

11.3. Modellierungslücken

Die in dieser Arbeit operationalisierten pragmatischen Eigenschaften von Modellen stellen Definitionsgrenzen dar. Dies impliziert, dass ein derart definiertes maschinelles Modell nicht unbeschränkt gültig ist. Im Umkehrschluss existieren mit jeder Wissensdomäne gleichzeitig Definitionsbereiche, welche die Modelle der Domäne explizit nicht abbilden. Die prinzipiell unendliche Menge aller denkbaren Modellsubjekte Σ ist beispielsweise in der vorliegenden Arbeit in jedem Domänenmodell auf künstliche neuronale Netze und Random Forests beschränkt ($n=2$).

In der Praxis werden werden diese Definitionsgrenzen insbesondere bei Betrachtung der temporalen Gültigkeit T deutlich. Wie die Ergebnisse für die konzeptuellen und prozeduralen Wissensdomänen zeigen, decken die Modelle dieser Domänen nie die gesamte Zeitspanne ab, die durch die angewendete Datenbasis gegeben ist. So decken die adaptierten konzeptuellen Modelle gemeinsam nur 22,5 Prozent und die adaptierten prozeduralen Modelle gemeinsam nur 82,8 Prozent der Zeitspanne der Datenbasis ab. Daraus folgt, dass für einen gegebenen Modellzweck Z nicht zu jedem durch die Datenbasis abgebildeten Zeitpunkt eine modellbasierte Evaluation möglich ist. Insbesondere ist durch die Gesamtheit aller Modelle eines gegebenen Z ein Maximum der temporalen Gültigkeit T von Z gegeben. Über dieses Maximum hinaus erlaubt die jeweilige Wissensdomäne keine unmittelbar gültigen Vorhersagen. Das kohärente und zielgerichtete Extrapolieren erworbenen Wissens auf neue, nicht erlernte Situationen und Aufgaben stellt im Sinne der Bloomschen Taxonomie eine eigene Form von Wissen dar, die gesondert zu erlernen und in

einer metakognitiven Domäne zu repräsentieren wäre.

Während Modellierungslücken also einerseits in einem konstruktivistischem maschinellen Lernen per Definition unvermeidbar sind, stellt sich andererseits dennoch die Frage, inwieweit Exploration und Adaption deren Entstehung begünstigen oder minimieren können. Eine temporale Modellierungslücke kann beispielsweise entstehen, wenn aus einem Lernblock ein neues Modell erzeugt wird, das zwar eine ΣZ -Verwandtschaft mit einem bereits in der Wissensdomäne enthaltenen Modell aufweist, aber keine temporale Überlappung. Falls dann die ΣZ -Dekonstruktion fehlschlägt, also keine Modellvereinigung möglich ist, wird die Domäne anschließend zwei Modelle mit übereinstimmenden Eigenschaften Σ und Z , aber unterschiedlichen temporalen Grenzen T enthalten. Eine Modellvereinigung kann insbesondere durch die tatsächliche Struktur der durch die Modelle abgebildeten Daten verhindert werden, etwa durch unterschiedliche Wertebereiche der Zielparameter prozeduraler Modelle. Neben einer fehlgeschlagenen ΣZ -Dekonstruktion kann aber auch eine minimale ΣZ -Dekonstruktion Modellierungslücken verursachen, sofern in der Wissensdomäne mehr als ein ΣZ -verwandtes Modell enthalten ist.

Indirekt haben auch die vor der eigentlichen ΣZ -Dekonstruktion durchgeführten Prozesse Einfluss auf deren Ausgang. So entscheiden etwa der Rekonstruktionsprozess und die damit verbundene Modellselektion darüber, welches Modell in den Dekonstruktionsprozess eingespeist wird. In der aktuellen Implementierung setzt das Verfahren dafür sowohl voraus, dass beide überwachten Lernverfahren Vorhersagen mit einer gegebenen Mindestgenauigkeit treffen, als auch dass ein gegebener Schwellenwert für den Reliabilitätskoeffizienten α erreicht wird. Insbesondere die Mindestgenauigkeit der Vorhersagen stellt, sofern sie nicht auf einen trivial niedrigen Wert gesetzt wird, eine Herausforderung für das Verfahren insgesamt dar. Denn optimierungsbedürftige Lernverfahren wie künstliche neuronale Netze erzielen oft mit automatisierten Abläufen nicht das bestmögliche Ergebnis. Der Lösungsansatz der aktuellen Implementierung sieht dafür vor, mehrere künstliche neuronale Netze mit unterschiedlichen Konfigurationen parallel zu trainieren und anschließend dasjenige mit der höchsten Vorhersagegenauigkeit zu verwenden. Wird dabei der gegebene Schwellenwert nicht erreicht, bricht der Rekonstruktionsprozess auch dann ab, wenn ein parallel dazu trainierter Random Forest eine ausreichende Genauigkeit aufweist. Eine künftige Weiterentwicklung könnte daher weitere überwachte Lernverfahren wie etwa Support-Vektor-Maschinen integrieren und die Bewertung der Intersubjektivität eines Modells auf Basis einer paarweise statt für drei Verfahren berechneten Interrater-Reliabilität vornehmen.

Nicht zu vernachlässigen für das Ergebnis der Lernprozesse ist der Einfluss von Zufallsinitialisierungen der eingesetzten überwachten und unüberwachten Lernverfahren. Ob der k-Means-Algorithmus ein Clustering mit dem globalen oder lediglich einem lokalen Minimum identifiziert, ist beispielsweise von dessen Initialisierung abhängig [113]. Auch bei Autoencodern bzw. künstlichen neuronalen Netzen beeinflusst die Zufallsinitialisierung der Verbindungsgewichte den Trainingsfortschritt und kann beispielsweise dazu führen, dass das Gradientenabstiegsverfahren statt dem globalen nur ein lokales Minimum identifiziert. Das Lernergebnis von Random Forests ist ebenfalls von zufallsinitialisierten Parametern abhängig [30]. In der Folge kann das hier vorgeschlagene automatisierte Vorgehen dazu führen, dass potentielle Domänenmodelle nicht konstruiert bzw. rekonstruiert werden können. Innerhalb des Konstruktionsprozesses könnte dies etwa durch Hinzunahme weiterer parallel ausgeführter unüberwachter Lernverfahren ausgeglichen werden. Ein anderer Ansatz zur Minimierung des Einflusses von Zufallsinitialisierungen wäre die mehrmalige Ausführung jedes Algorithmus, was entsprechend die Zahl der zu rekonstruierenden Modelle erhöhen würde. Ein ähnlicher Effekt lässt sich in der aktuellen Implementierung bereits über die Inkrementierung der Iterationsdurchläufe erreichen. Analog könnten auch innerhalb des Rekonstruktionsprozesses die jeweiligen überwachten Lernverfahren mehrfach ausgeführt und nur der am besten trainierte Durchgang gewertet werden. Ein anderer Ansatz wäre die Hinzunahme weiterer überwachter Lernverfahren wie etwa Support-Vektor-Maschinen.

11.4. Fazit und Ausblick

Die vorliegende Arbeit führt ein maschinelles Lernen unter Berücksichtigung konstruktivistischer Prinzipien sowie unter Verwendung des pragmatischen Modellbegriffs nach Stachowiak ein. Auf dieser Basis wurde eine konkrete Implementierung nicht nur von Konstruktions- und Rekonstruktionsprozessen, sondern auch eines Dekonstruktionsprozesses realisiert. Dieser neue, ergänzende Prozess erfordert nicht nur eine Formalisierung dessen, was beim Mensch Zweifel genannt wird, sondern insbesondere ein übergeordnetes Konzept, das diese drei Lernprozesse in Beziehung zueinander setzt. Stachowiaks Modellbegriff hat sich hierfür als ideale Grundlage erwiesen, anhand deren sich nicht nur der Dekonstruktionsprozess, sondern auch das Konzept eines konstruktivistischen maschinellen Lernens insgesamt effektiv operationalisieren lässt. Die Organisation von maschinellen Modellen in einer hierarchisch geordneten Wissensdomäne stellt dabei eine zentrale Voraussetzung für den Dekonstruktionsprozess dar.

Gleichzeitig bildet eine Wissensdomäne explizit auch Beziehungen zwischen erlernten Modellen ab, wodurch diese unmittelbar nachvollziehbar werden. Diese Beziehungen werden nicht nur durch Verknüpfungen von Ein- und Ausgaben der Modelle definiert, sondern insbesondere durch die pragmatischen Eigenschaften T , Σ und Z . Deren konsequente Nutzung in Form von Metadaten verbessert jedoch nicht nur die Nachvollziehbarkeit des Lernens und der Ergebnisse, sondern erlaubt auch einen effektiven Umgang mit Mehrdeutigkeiten. Einerseits ermöglicht sie eine formale Definition und Überprüfbarkeit dessen, was unter Mehrdeutigkeit zu verstehen ist. Andererseits stellt sie auch einen leicht operationalisierbaren Weg dar, diese zu erkennen und zu differenzieren. Im Umkehrschluss ist aufgrund dieser Überprüfbarkeit auch festlegbar, dass ausschließlich eindeutige Modelle in eine Wissensdomäne aufgenommen werden sollen.

Gemeinsam mit den Prinzipien eines konstruktivistischen maschinellen Lernens realisiert der vorgeschlagene Ansatz auch Ideen eines kritischen Rationalismus. Dazu zählt neben der durch Dekonstruktionen implizierten Vorläufigkeit von Erkenntnis insbesondere eine Abkehr vom Prinzip der Verifikation, also der Suche nach Belegbarkeit, hin zum Prinzip der Falsifikation, also der Suche nach Widerlegbarkeit. So werden im Rahmen des Rekonstruktionsprozesses grundsätzlich alle Modelle verworfen, welche die benutzerdefinierten Schwellenwerte für subjektiven Lernerfolg oder für intersubjektive Übereinstimmung nicht erreichen. Ein Modell, das nach Durchlaufen von Rekonstruktions- und Dekonstruktionsprozess in eine Wissensdomäne integriert wird, ist somit aus Sicht des konstruktivistischen maschinellen Lernens zum gegenwärtigen Zeitpunkt nicht widerlegbar bzw. nicht durch ein alternatives Modell ersetzbar.

Für die Anwendung auf empirische Daten wie impedanzspektroskopische Messungen bietet das hier vorgeschlagene konstruktivistische maschinelle Lernen somit klare methodische Vorteile gegenüber klassischen überwachten bzw. unüberwachten Lernverfahren. Sowohl bei Exploration von Daten ohne definierten Zweck als auch bei Adaption eines gegebenen Zwecks wahrt das Verfahren die Prinzipien von Nachvollziehbarkeit und Falsifikation. Im Ergebnis führt die Exploration der hier genutzten Datenbasis zu ausdifferenzierten konzeptuellen und prozeduralen Wissensdomänen. Analog führt auch die Adaption von konzeptuellem und prozeduralem Wissen nicht zu einheitlichen, sondern temporal differenzierten Modellen für die Datenbasis. In früheren Arbeiten zur Analyse von Impedanzspektren musste vergleichbare Fallunterscheidungen bzw. Differenzierungen von Hypothesen noch manuell getroffen werden [230, 231, 232, 233].

Über den konkreten Anwendungsfall hinaus liefert diese Arbeit auch eine Reihe von Ansatzpunkten für die künftige Erforschung und Optimierung eines konstruktivistischen maschinellen Lernens. So sieht die überarbeitete Bloomsche Taxonomie neben dem hier untersuchten konzeptuellen und prozeduralen Wissen auch Faktenwissen und metakognitives Wissen als relevante Dimensionen menschlicher Kognition an. In künftigen Arbeiten wäre somit noch zu klären, wie und mit welchem Nutzen sich eine faktuelle und eine metakognitive Wissensdomäne realisieren

lassen. Eine Aufgabe metakognitiver Modelle könnte etwa das spekulative Extrapolieren von Modellen der konzeptuellen oder prozeduralen Domäne über deren pragmatisch definierten Grenzen hinaus sein. Dies könnte eine Übertragung des unmittelbar erlernten Wissens auf neue Situationen erlauben, wie es dem Menschen möglich ist. Neben diesen vier Wissensarten aus der kognitiven Domäne der Bloomschen Taxonomie könnte für künftige Anwendungen möglicherweise auch die Untersuchung bzw. Realisierung einer affektiven oder psychomotorischen Domäne Relevanz besitzen. Inwieweit diese sich im Rahmen eines konstruktivistischen maschinellen Lernens realisieren lassen, ist zum gegenwärtigen Zeitpunkt ebenfalls noch unklar.

Konstruktivistisches maschinelles Lernen erlaubt sowohl den aktiven Erwerb von unbekanntem und bekanntem Wissen mittels Exploration und Adaption als auch dessen Anwendung als Reaktion auf bestimmte Eingaben. Offen ist hingegen, ob ein solches Lernen auch ein aktives, selbstständiges Handeln auf Basis von erworbenem Wissen ermöglichen kann. Für das menschliche Gehirn wird ein solches Handeln häufig als mentaler Prozess beschrieben, für aktiv handelnde Computersysteme wurde in der Vergangenheit häufig der Begriff Agent verwendet. Im Kontext des hier vorgeschlagenen Verfahrens wäre beispielsweise denkbar, dass ein selbstständiges Handeln in Form von Modellen einer metakognitiven Wissensdomäne realisiert wird, die auf Basis von konzeptuellen und prozeduralen Modelle zur selbstständigen Konstruktion von Handlungszielen und Absichten fähig ist. Sollte dies gelingen, so würde damit eine Art konstruktivistischer Agent entstehen, der eine Ich-Konzept über sich selbst konstruiert und somit dem, was man heute über Denken und Lernen des Menschen weiß, im Kern sehr nahe kommen könnte.

11. Chancen und Risiken konstruktivistischen Lernens

ANHANG



A

Implementierung und Konfiguration

Die Teilschritte konstruktivistischen maschinellen Lernens werden mittels unüberwachten und überwachten Lernverfahren realisiert. Diese werden weitgehend mithilfe existierender Implementierungen der Programmiersprache *R* angewendet, die im Folgenden gemeinsam mit den verwendeten Trainingsparametern dargestellt werden. Daneben sind außerdem die Parameter des konstruktivistischen maschinellen Lernens selbst sowie die damit verbundenen Konfigurationsmöglichkeiten zusammenfassend dargestellt.

A.1. Unüberwachte Lernverfahren

Alle unüberwachten Lernverfahren werden über *R*-Skripte gesteuert, die im Skript-Unterordner `construction` der Implementierung liegen.

Selbstorganisierende Karten. Es kommen eindimensionale Self-Organizing Maps (SOMs) zum Einsatz. Im Rahmen eines Konstruktionsprozesses werden dabei mehrere SOMs mit jeweils mindestens 2 und maximal κ_k Neuronen trainiert. Der Parameter κ_k ist dabei die maximale kategoriale Komplexität, die als grundlegender Parameter konstruktivistischen maschinellen Lernens bereits gegeben ist. Innerhalb des Konstruktionsprozesses ist κ_k der einzig variable Parameter. Die übrigen Trainingsparameter bleiben fix (siehe Tab. A.1).

Alle SOMs sind mithilfe des *R*-Packages *kohonen*¹ realisiert. Für das zu erzeugende `som`-Objekt ist zunächst über ein `somgrid` die Anzahl von Gitterzeilen und -reihen sowie die Art des Gitters zu definieren. Nach dem Training wird für ein Sample der Index desjenigen Neurons als Clusterzuordnung gespeichert, dem es von der SOM zugeordnet worden ist. Die so konstruierten Zielwerte sind dabei unmittelbar im zugehörigen *R*-Objekt hinterlegt.

k-Means-Algorithmus. Es kommt der Standard-k-Means-Algorithmus zum Einsatz. Im Rahmen eines Konstruktionsprozesses werden dabei jeweils mindestens 2 und maximal κ_k Clusterzentren angenommen. Es werden also κ_k verschiedene Clustermöglichkeiten identifiziert.

In *R* ist k-Means ein fester Sprachbestandteil, der nicht über ein Package geladen werden muss. Der zugehörige Befehl `kmeans` ist im Sprachkern implementiert und direkt aufrufbar. Nach dem Training wird der Index desjenigen Clusters, dem das Sample durch den k-Means-Algorithmus

¹Online abrufbar unter <http://cran.r-project.org/web/packages/kohonen/>

A. Implementierung und Konfiguration

	Parameter	Wert
SOM	<code>somgrid\$xdim</code>	1
	<code>somgrid\$ydim</code>	κ_k
	<code>somgrid\$topo</code>	hexagonal
k-Means	<code>k</code>	κ_k
Autoencoder	<code>nl</code>	3
	<code>N.hidden</code>	κ_p
	<code>beta</code>	6
	<code>epsilon</code>	0.001
	<code>lambda</code>	0.0002
	<code>rho</code>	0.01
	<code>max.iterations</code>	1000
	<code>rescale.flag</code>	TRUE
	<code>rescaling.offset</code>	0.001
ClustOfVar	<code>init</code>	κ_p

Tabelle A.1.: Übersicht der Trainingsparameter der unüberwachten Lernverfahren. κ_k bezeichnet die maximal zulässige kategoriale, κ_p die maximal zulässige prozedurale Komplexität.

zugeordnet wurde, gespeichert. Diese zu konstruierenden Zielwerte sind dabei bereits unmittelbar im zugehörigen *R*-Objekt hinterlegt.

Autoencoder. Es kommen Autoencoder mit einer versteckten Schicht und einer variierenden Zahl versteckter Neuronen zum Einsatz. Die Zahl der versteckten Neuronen ist durch den übergeordneten Parameter `MaxModelTargets` definiert (vgl. Tab. A.3). Für die übrigen Parameter werden Defaultwerte des verwendeten *R*-Packages *autoencoder*² übernommen (vgl. Tab. A.1).

Da es sich bei einem Autoencoder um ein neuronales Netz handelt, müssen sämtliche Ein- und Ausgabewerte dem Wertebereich der verwendeten Aktivierungsfunktion entsprechen. Für eine entsprechende Anpassung stellt das Package *autoencoder* eine Rescaling-Funktion zur Verfügung, die bei jeder Durchführung angewendet wird.

Um die durch den Autoencoder erzeugte Repräsentation niedrigerer Dimension extrahieren zu können, ist es erforderlich, statt der Werte der Ausgabeschicht die Werte der versteckten Schicht auszugeben. Das Package *autoencoder* stellt dafür einen Parameter `hidden.output` zur Verfügung, mit dem dies ermöglicht wird. Während bei selbstorganisierenden Karten und dem k-Means-Algorithmus die zu konstruierenden Zielwerte unmittelbar im zugehörigen *R*-Objekt hinterlegt sind, müssen diese bei *autoencoder* wie bei einem überwachten Lernverfahren mittels des Befehls `predict` explizit erzeugt werden.

Clustering of Variables (ClustOfVar). Es kommt ein Feature-Clustering-Verfahren zum Einsatz, das auf dem k-Means-Algorithmus basiert. Die Zahl der zu unterscheidenden Feature-Cluster ist dabei durch die maximale prozedurale Komplexität κ_p gegeben und wird über den Parameter `init` festgelegt.

Alle Feature-Clusterings sind mithilfe des *R*-Packages *ClustOfVar*³ implementiert. Dieses erzeugt über den Befehl `kmeansvar` zunächst ein Objekt, das Informationen darüber enthält, welche Features welchem der κ_p Feature-Cluster zugeordnet sind (vgl. Tab. A.3). Analog zum Autoencoder müssen die zu konstruierenden Zielwerte anschließend wie bei einem überwachten Lernverfahren mittels des Befehls `predict` erzeugt werden. Zusätzlich werden die zugehörigen Features exportiert und in den resultierenden Modellen berücksichtigt.

²Online abrufbar unter <http://cran.r-project.org/web/packages/autoencoder/>

³Online abrufbar unter <http://cran.r-project.org/web/packages/ClustOfVar/>

A.2. Überwachte Lernverfahren

Die überwachten Lernverfahren werden mittels *C*- oder *R*-Implementierungen realisiert, die im Skript-Unterordner `reconstruction` der Implementierung liegen. Alle Eingabedaten werden randomisiert im Verhältnis 2:1 geteilt, wodurch zwei Drittel der Daten als Trainingsdaten und ein Drittel als Testdaten verwendet werden.

Künstliche neuronale Netze. Es werden ausschließlich Feedforward-Netze mit je einer versteckten Schicht und einer variierenden Zahl versteckter und Eingabeneuronen eingesetzt. Für das Training wird Resilient Backpropagation (Rprop) verwendet [211]. Dieser Algorithmus vereint die Vorteile von Manhattan-Training und Quickprop, wodurch es sowohl in Geschwindigkeit als auch Robustheit dem klassischen Backpropagation überlegen ist. Die Rprop-Parameter Δ_{min} , Δ_{max} und η^+ wurden auf Basis etablierter Standards gewählt [211], die Parameter Δ_0 und η^- auf Basis publizierter Optimierungsvorschläge [115, 107] (siehe Tab. A.2).

Die Implementierung erfolgt als Binaries auf Basis der *C*-Bibliothek FANN⁴, die ein leistungsfähiges und variables Framework darstellt. Als Kompromiss zwischen einer rein statischen und einer rein dynamischen Realisierung wird eine Menge fest definierter Netzarchitekturen mit je einer versteckten Neuronenschicht und einer variierenden Anzahl n an Eingabeneuronen vorgehalten. In Abhängigkeit von n werden pro Rekonstruktionsvorgang drei konkurrierende neuronale Netze mit den Architekturen $n-i-1$ mit $n \in \{2, \dots, 10\}$ trainiert, von denen nur dasjenige mit der höchsten Vorhersagegenauigkeit weiterverwendet wird. In der vorliegenden Implementierung gilt $n \leq \text{MaxFeatures} = 10$ (Default-Wert, vgl. Abschnitt A.4) sowie $i \in \{2, \dots, [n \cdot 1, 5]\}$,

Da eine sigmoide Aktivierungsfunktion verwendet wird (FANN_SIGMOID_SYMMETRIC), müssen sämtliche Ein- und Ausgabewerte dem Wertebereich $(-1, 1)$ entsprechen. Für eine entsprechende Anpassung stellt die Bibliothek FANN eine Rescaling-Funktion zur Verfügung. Mit dieser werden die Eingabedaten jedes Netzes jeweils auf einen Bereich zwischen $-0,8$ und $0,8$ und wieder zurück in das Ursprungsintervall skaliert.

⁴Online abrufbar unter <http://leenissen.dk/fann/>

	Parameter	Wert
FANN	<code>desired_error</code>	0.0001
	<code>max_epochs</code>	2500
	<code>activation_steepness_hidden</code>	0.9
	<code>activation_steepness_output</code>	0.9
	<code>activation_function_hidden</code>	FANN_SIGMOID_SYMMETRIC
	<code>activation_function_output</code>	FANN_LINEAR
	<code>train_stop_function</code>	FANN_STOPFUNC_MSE
	<code>training_algorithm</code>	FANN_TRAIN_RPROP
	<code>rprop_increase_factor</code>	1.2
	<code>rprop_decrease_factor</code>	0.7
	<code>rprop_delta_max</code>	50.0
	<code>rprop_delta_min</code>	0.000001
	<code>rprop_delta_zero</code>	0.07
	<code>scaling_params</code>	$(-0.8, 0.8, -0.8, 0.8)$
Random Forest	<code>ntrees</code>	50
	<code>mtry</code>	Klass.: \sqrt{n} ; Regr.: $n/3$
	<code>nodesize</code>	Klass.: 1; Regr.: 5
	<code>norm.votes</code>	FALSE

Tabelle A.2.: Übersicht der Trainingsparameter der überwachten Lernverfahren. Soweit nicht anders angegeben, gelten für Klassifikation und Regression die gleichen Werte.

A. Implementierung und Konfiguration

Random Forests. Es kommen Random Forests zum Einsatz, deren Größe in einem moderaten Verhältnis zur maximalen Modellkomplexität κ steht, Die Zahl der verwendeten Entscheidungsbäume ist mit `ntree=500` festgelegt, d.h. $\kappa_k < ntree \leq 50 \cdot \kappa \forall \kappa_k \in \{2, \dots, 10\}$.

Die Implementierung erfolgt mithilfe des R-Packages *randomForest*⁵, das insbesondere eine Option zur parallelisierten Ausführung des Lernalgorithmus bietet. Mit diesem Package kann sowohl Klassifikation als auch Regression realisiert werden, wobei die Vorhersage von Zielwerten für Testdaten die explizite Definition der Zielwerte als Klassifikationszielwerte (`as.factor`) oder als Regressionszielwerte (`as.vector`) erfordert.

Darüber hinaus wird in der Implementierung in einigen weiteren Punkten zwischen Klassifikation und Regression unterschieden. Dies betrifft einerseits die Zahl der zufällig auszuwählenden Variablen für jeden Split, die entsprechend der Default-Einstellung auf \sqrt{n} für Klassifikation und auf $n/3$ für Regression gesetzt ist. Auch die Mindestgröße der Endknoten beträgt entsprechend der Default-Werte für Klassifikation 1 und für Regression 5 (siehe auch Tab. A.2).

A.3. Ablaufsteuerung

Der in Kapitel 5.3 beschriebene Ablauf eines konstruktivistischen maschinellen Lernens sowie der damit verbundenen Lernprozesse wird über ein Basisskript `learn.sh` und eine Konfigurationsdatei `learn.cfg` gesteuert. Die weitere Ablaufsteuerung ist in einer Ordnerstruktur organisiert, wobei für Konstruktion, Rekonstruktion und Dekonstruktion jeweils Unterordner `construction`, `reconstruction` und `deconstruction` existieren.

Mit dem Basisskript `learn.sh` werden insbesondere folgende zentrale Aufgaben realisiert:

1. Auslesen der Konfigurationsdatei. Gemäß Benutzervorgabe wird ein temporäres Verzeichnis angelegt, in dem der Lernvorgang und damit alle folgenden Lernprozesse ausgeführt werden. Um eine performante Ausführung zu erreichen, wird empfohlen, für dieses Lernverzeichnis einen Ordner im Shared Memory (`/dev/shm`) zu wählen.
2. Iterationsweises Abarbeiten der Datenhalde. Die Halde verwaltet alle noch nicht in ein Modell integrierten vektoriellen Modelle und wird pro Iteration mit einem nachgelagerten Skript (`learn_iteration.sh`) verarbeitet.

Mit dem Iterationsskript `learn_iteration.sh` werden u.a. folgende Aufgaben realisiert:

1. Blockweise Verarbeitung der Eingabedaten. Gemäß Benutzervorgabe werden von der Halde gleichgroße Blöcke gezogen und auf geeignete Lernblöcke untersucht. Die Identifikation von Lernblöcken wird in ein eigenes Skript (`select_learn_block`) ausgelagert.
2. Auswahl eines geeigneten Lernprozesses. Für jeden identifizierten Lernblock wird entschieden, ob Konstruktion, Rekonstruktion oder Dekonstruktion durchgeführt werden. Diese sind in Skripten in den Unterordnern `construction`, `reconstruction` und `deconstruction` ausgelagert.
3. Verwaltung der Wissensdomäne. Jedes konstruierte Modell erhält hier vor dem Übergang in den Rekonstruktionsprozess eine vorläufige ID. Wird das Modell durch Re- oder Dekonstruktion verworfen, kann diese ID später erneut vergeben werden. Andernfalls wird auf dieses Basis die Integration in die Wissensdomäne durchgeführt.

Tab. A.3 gibt einen Überblick über die Parameter, die über die Konfigurationsdatei gesetzt werden können. Für eine ausführliche Beschreibung der Parameter siehe Abschnitt A.4.

⁵Online abrufbar unter <http://cran.r-project.org/web/packages/randomForest/>

	Parameter	Bedeutung
Allgemein	InputFile	Eingabedaten
	LearnDirectory	Temporäres Lernverzeichnis
	MaxLearnDir	Maximale Auslastung des Lernverzeichnisses
	UseExistingModels	Übernehmen einer existierenden Wissensdomäne
Preprocessing	SetFeatures	Vorgelagerte manuelle Feature-Selektion
	SetTargets	Überschreiben gesetzter Zielparameter
	SortTimestamp	Chronologische Ordnung der Samples
	CutTimestamp	Prüfen der Timestamps auf Kürzung
Blockverarbeitung	BlockSize	Blockgröße
	MaxBlocks	Max. Anzahl Blocks
	StackIterations	Halden-Durchläufe
	LearnBlockMinimum	Min. Lernblockgröße
	SigmaZetaCutoff	Cutoff für ΣZ -Verwandtschaft
Konstruktion	MaxCategories	Max. kategoriale Komplexität κ_k
	MinCategorySize	Mindestgröße der Kategorien
	MaxModelTargets	Max. prozedurale Komplexität κ_p
	MaxTargetError	Max. Fehler
Komplexitätsreduktion	MaxFeatures	Max. Modellkomplexität $\kappa_{\mathcal{M}}$
	MaxFilterX	Filter-Grenze λ_x
	MaxFilterY	Filter-Grenze λ_y
	MaxModelReduction	Maximalreduktion erzwingen
Rekonstruktion	MaxTestErrorAvg	Max. durchschnittliche Abweichungen [%]
	MaxTestErrorMax	Max. Maximalabweichung [%]
	MinTestAccuracy	Min. Accuracy
	ReliabilitySample	Samplegröße für Reliabilitätsprüfung
	MinReliability	Schwellenwert für Reliabilitätsprüfung α_{min}
	ReduceModelRedundancy	Redundanzvermeidung
Dekonstruktion	DeconstStrategy	Umgang mit neuen Modellen
	DeconstMode	Abbruchbedingung
	DeconstMaxDistanceT	Temporale Erweiterungstoleranz δT_{max}
	DeconstFullTolerance	Toleranz bei vollständiger Dekonstruktion
	ForceTimeExpansion	ΣZ -verwandte Modelle immer vereinigen

Tabelle A.3.: Übersicht der konstruktivistischen Lernparameter. Die Werte der Parameter werden mithilfe einer Konfigurationsdatei für das konstruktivistische maschinelle Lernen definiert. Für eine ausführliche Beschreibung der einzelnen Parametern siehe Abschnitt A.4.

A.4. Konfigurationsreferenz

Im Folgenden sind alle Parameter erläutert, die in der Konfigurationsdatei gesetzt werden können (vgl. Tab. A.3). Diese wird als Textdatei mit tabulator-separated values (TSV) erwartet. Die Standard-Konfigurationsdatei heißt `learn.cfg` und liegt im Basisordner.

Kommentarzeilen können in die Konfigurationsdatei eckigen Klammern eingefügt werden: `[Kommentar]`. Inline-Kommentare können hinter Parameter und zugehörigen Wert eingefügt werden und werden durch das `#`-Symbol gekennzeichnet.

Allgemein

<code>InputFile <string></code>	Voller Dateipfad zu den zu erlernenden Daten. Diese müssen in einer einzigen TSV-Textdatei vorliegen, wobei die pragmatischen Metadaten am Ende jeder Zeile erwartet werden. Die <i>T</i> -Spalte muss die drittletzte, die Σ -Spalte die vorletzte und die <i>T</i> -Spalte die letzte sein. Default: <code>./learn.dat</code>
<code>LearnDirectory <string></code>	Voller Dateipfad zu einem Ordner, in den die zu lernenden Daten sowie alle Skripte kopiert werden. Auf dem Datenträger, auf dem der Ordner liegt, sollte mindestens noch soviel Speicher frei sein wie im Startverzeichnis. Default: <code>/dev/shm</code>
<code>MaxLearnDir <int></code>	Maximal erlaubte Auslastung des Datenträgers, auf dem das <code>LearnDirectory</code> liegt. Dieser Wert wird vor Beginn jeder Iteration überprüft. Default: <code>80 [%]</code> .
<code>UseExistingModels <boolean></code>	Optionale Übernehmen einer existierenden Wissensdomäne. Falls <code>true</code> , wird eine im Startverzeichnis enthaltene Wissensdomäne wiederverwendet. Ansonsten wird diese überschrieben. Default: <code>false</code> .

Preprocessing

<code>SetFeatures <string></code>	Auswahl derjenigen Spalten des TSV-formatierten <code>InputFile</code> , die als Eingabevariablen verwendet werden. Zulässige Formatierungen sind <code>1, 3, 5</code> oder <code>1-10</code> . Diese Option muss immer gesetzt werden, ein Default-Wert ist nicht implementiert.
<code>SetTargets <string></code>	Optionales Überschreiben der <i>Z</i> -Spalte des TSV-formatierten <code>InputFile</code> mit einem festen String. Dies kann zum Beispiel nützlich sein, wenn die ursprünglichen Werte fehlerhaft sind. Wird die <i>Z</i> -Spalte mit dieser Option auf <code>0.0.0</code> gesetzt, erzwingt dies für alle Lernblöcke Konstruktionsprozesse.
<code>SortTimestamp <boolean></code>	Optionales Sortieren der Samples entsprechend der <i>T</i> -Spalte des TSV-formatierten <code>InputFile</code> . Die Sortierung wird einmalig bei Programmstart ausgeführt und erfolgt in aufsteigender Ordnung.
<code>CutTimestamp <boolean></code>	Optionales Kürzen des Formats der <i>T</i> -Spalte des TSV-formatierten <code>InputFile</code> . In den entsprechenden Timestamps ein redundantes Prefix identifiziert wird, das allen Samples gemeinsam ist, wird dieses entfernt. Ansonsten bleibt die <i>T</i> -Spalte unverändert.

Blockverarbeitung

BlockSize <int>	Größe der aus der Halde zu ziehenden Blöcke. Der Wert für <code>BlockSize</code> beeinflusst die Anzahl der auszuführenden Lernprozesse und sollte größer gewählt sein als <code>LearnBlockMinimum</code> . Dieser Wert muss immer gesetzt werden, ein Default-Wert ist nicht implementiert.
MaxBlocks <int>	Maximale Anzahl von Blöcken, die aus der Halde gezogen werden. Spätestens bei Überschreiten dieses Schwellenwerts wird das Verfahren insgesamt abgebrochen. Dies Anzahl der identifizierten Lernblöcke ist dabei unerheblich. Default: 10.
StackIterations <int>	Maximale Anzahl von Halden-Durchläufen. Spätestens bei Überschreiten dieses Schwellenwerts wird das Verfahren insgesamt abgebrochen. Default: 1.
LearnBlockMinimum <int>	Minimal erforderliche Zahl an Samples pro Lernblock. Werden innerhalb eines gezogenen Blocks Lernblöcke identifiziert, die weniger Samples enthalten, werden diese Lernblöcke verworfen. Dieser Wert muss immer gesetzt werden, ein Default-Wert ist nicht implementiert.
SigmaZetaCutoff <double>	Schwellenwert für die Identifizierung ΣZ -verwandter Lernblöcke. Bestimmt welche Dichteverteilung des Metadatum T in einem gezogenen Block mindestens erreicht werden muss, um damit ein eindimensionales Clustering durchzuführen (vgl. Abschnitt 5.3.1). Bereiche, die diesen Schwellenwert unterschreiten, werden nicht betrachtet. Default: 20 [%].

Konstruktion

MaxCategories <int>	Maximale kategoriale Komplexität κ_k . Bestimmt innerhalb der konzeptuellen Wissensdomäne die Anzahl der Cluster, die mittels Clustering-Verfahren bestimmt werden. Dieser Wert muss immer gesetzt werden, ein Default-Wert ist nicht implementiert.
MinCategorySize <int>	Mindestgröße einer Kategorie. Legt innerhalb der konzeptuellen Wissensdomäne die erforderliche Zahl an Samples pro Cluster fest. Wird diese nicht von allen Clustern eines Clusterings erreicht, wird das Clustering insgesamt verworfen. Statt eines Integers kann auch ein Dezimalwert zwischen 0 und 1 angegeben werden, mit dem die Sampleszahl aus der Gesamtzahl der Samples des Lernblocks berechnet wird. Dieser Wert muss immer gesetzt werden, ein Default-Wert ist nicht implementiert.
MaxModelTargets <int>	Maximale prozedurale Komplexität κ_p . Bestimmt innerhalb der prozeduralen Wissensdomäne die Zahl der Zielparameter, die konstruiert werden. Falls $\kappa_p > 1$, wird während des Rekonstruktionsprozesses die Redundanz der Zielparameter überprüft (vgl. <code>ReduceModelRedundancy</code>). Default: 1.

A. Implementierung und Konfiguration

MaxTargetError <double>

Schwellenwert für überwachtes Konstruieren prozeduralen Wissens. Bestimmt auf dem Intervall $[0,1]$, wie hoch der durchschnittliche Trainingsfehler eines Autoencoders höchstens sein darf. Default: 0.25 .

Komplexitätsreduktion

MaxFeatures <int>

Maximal zulässige Modellkomplexität κ . Mit diesem Wert wird festgelegt, wieviele Eingabewerte ein Modell höchstens haben darf. Ab $\text{MaxFeatures}+1$ Eingabewerten wird eine Komplexitätsreduktion durchgeführt. Default: 10 .

MaxFilterX <int>

Filter-Grenze λ_x . Dieser Wert regelt, wann ein eingebettetes Feature-Selection-Verfahren und wann ein Filterverfahren zum Einsatz kommt. Enthält ein Modell mehr als MaxFilterX Eingabewerte, wird ein Filterverfahren eingesetzt; ansonsten ein eingebettetes Verfahren. Dieser Wert muss immer gesetzt werden, ein Default-Wert ist nicht implementiert.

MaxFilterY <int>

Filter-Grenze λ_y . Dieser Wert regelt, wann ein eingebettetes Feature-Selection-Verfahren und wann ein Filterverfahren zum Einsatz kommt. Enthält ein Modell mehr als MaxFilterY Samples, wird ein Filterverfahren eingesetzt; ansonsten ein eingebettetes Verfahren. Dieser Wert muss immer gesetzt werden, ein Default-Wert ist nicht implementiert.

MaxModelReduction <boolean>

Erzwingen einer maximalen Komplexitätsreduktion. Ist dieser Wert auf `true` gesetzt, werden solange Komplexitätsreduktionsschritte durchgeführt, bis keine Features mehr entfernt werden können. Andernfalls wird die Komplexitätsreduktion abgebrochen, sobald erstmals weniger als MaxFeatures Eingabevariablen erreicht wurden. Default: `false`.

Rekonstruktion

MinTestAccuracy <double>

Schwellenwert für die Accuracy-Prüfung. Legt innerhalb der konzeptuellen Wissensdomäne fest, welchen Accuracy-Wert Modelle für Vorhersagen von Klassifikationen mindestens erreicht werden müssen. Wird MinTestAccuracy nicht erreicht, wird das Modell verworfen. Default: 0.8 .

MaxTestErrorAvg <double>

Schwellenwert für den durchschnittlichen Vorhersagefehler. Legt innerhalb der prozeduralen Wissensdomäne fest, welche durchschnittliche prozentuale Abweichung vom Zielwert für Regressionsvorhersagen maximal erreicht werden darf. Wird MaxTestErrorAvg überschritten, wird das Modell verworfen. Default: 10 [%].

MaxTestErrorMax <double>

Schwellenwert für den maximalen Vorhersagefehler. Legt innerhalb der prozeduralen Wissensdomäne fest, welche maximale prozentuale Abweichung vom Zielwert für Regressionsvorhersagen maximal erreicht werden darf. Wird MaxTestErrorAvg überschritten, wird das Modell verworfen. Default: 100 [%].

A.4. Konfigurationsreferenz

ReliabilitySample <double>	Samplegröße für Reliabilitätsprüfung. Bestimmt den Anteil der Samples, anhand dessen der Reliabilitätskoeffizient eines Modells bestimmt wird. Der angegebene Wert muss zwischen 0 und 1 liegen. Default: 0.1.
MinReliability <double>	Schwellenwert für Reliabilitätsprüfung. Bestimmt ab welchem Wert des Interreliabilitätskoeffizienten eine Übereinstimmung angenommen wird. Default: 0.8.
ReduceModelRedundancy <boolean>	Redundanzvermeidung. Falls <code>true</code> und <code>MaxModelTargets > 1</code> , wird stets nur eines der κ_p konstruierten Modell weiterverwendet. Dazu werden diese anhand der Interreliabilität und der Anzahl der Features gereiht. Default: <code>false</code> .
Dekonstruktion	
DeconstStrategy <string>	Umgang mit neuen Modellen. Bestimmt, ob neue erfolgreich dekonstruierte Modelle in die Wissensdomäne aufgenommen werden (vgl. Abschnitt 5.6. Mögliche Werte: <code>conservative</code> , <code>integrative</code> , <code>opportunistic</code> . Default: <code>conservative</code> .
DeconstMode	Abbruchbedingung. Der Dekonstruktionsprozess wird entweder abgebrochen, sobald eine vollständige, eine ΣZ -, eine TZ - oder eine $T\Sigma$ -Dekonstruktion erfolgreich war (<code>minimal</code>) oder erst dann, wenn keine weiteren pragmatisch verwandten Modelle mehr identifiziert werden können (<code>full</code>); siehe auch Abschnitt 5.6. Default: <code>minimal</code> .
DeconstMaxDistanceT	Temporale Erweiterungstoleranz δT_{max} . Definiert, wie groß die temporale Lücke zwischen ΣZ -verwandten Modelle maximal sein darf, damit eine Expansion geprüft wird (vgl. Abschnitt 5.6.1). Ist dieser Wert auf 1 gesetzt, wird stets geprüft. Ist er auf 0 gesetzt, wird nur geprüft, wenn eine Überlappung vorliegt. Ansonsten darf die temporale Lücke die Größe des durch δT_{max} gesetzten Anteil an der Gesamtspannweite nicht überschreiten. Default: 1.
DeconstFullTolerance	Temporale Toleranz bei vollständiger Verwandtschaft. Definiert, um welchen Anteil das Zeitintervall T erweitert werden kann, um eine Verwandtschaft zu konstituieren. Default: 0.1.
ForceTimeExpansion <boolean>	Vereinigung von ΣZ -verwandten Modellen. Falls <code>true</code> , werden ΣZ -verwandte Modelle bevorzugt behandelt bzw. vereinigt. Default: <code>true</code> .

A. Implementierung und Konfiguration

B Modellierung epithelialer Zellkulturen

Für die Nachbildung von Messdaten für die epithelialen Zellkulturlinien HT-29/B6, IPEC-J2 und MDCK I wird der elektrische Ersatzschaltkreis C (siehe Abschnitt 6.1) sowie je nach funktionellem Zustand fest definierte Wertebereiche für die Schaltelemente angenommen. Auf dieser Grundlage wird dann die theoretisch zu erwartende Impedanz für eine gegebene Messfrequenz berechnet und anschließend ein gerätespezifischer Messfehler aufaddiert.

Pro nachgebildetem Impedanzspektrum werden 42 invariante Frequenzen f_0, \dots, f_{41} zwischen 1,3 und 16,35 kHz angenommen (Tab. B.1), d.h. es werden pro Spektrum auf Basis von sechs Schaltkreisparametern je 42 frequenzabhängige Impedanzen Z berechnet:

$$(Z_0, \dots, Z_{41}) \sim (R_s, R_p, R_a, R_b, C_a, C_b) \quad (\text{B.1})$$

Real- und Imaginärteil der komplexwertigen Impedanzen $Z = \Re(Z) + i \cdot \Im(Z)$ werden dabei getrennt verarbeitet und berechnet (siehe Abschnitt B.1):

$$(Z_0, \dots, Z_{41}) = ((\Re(Z_0), \Im(Z_0)), \dots, (\Re(Z_{41}), \Im(Z_{41}))) \quad (\text{B.2})$$

Gleiches gilt für das Aufaddieren des modellierten gerätespezifischen Messfehlers σ_{\Re} bzw. σ_{\Im} (siehe Abschnitt 6.4):

$$\Re(Z_i) = \Re(Z^T(\omega_i)) + \sigma_{\Re}(f_i) \quad (\text{B.3})$$

$$\Im(Z_i) = \Im(Z^T(\omega_i)) + \sigma_{\Im}(f_i) \quad (\text{B.4})$$

wobei $\omega_i = 2\pi f_i$ sowie $i \in \{0, \dots, 41\}$ gilt.

Die an diese Berechnungen anschließende systematische Synthetisierung von Messungen an Epithelien erfolgt individuell für jede Zellkulturlinie und jeden funktionalen Zustand, wobei jeweils individuelle Grenzen für die Modellparameter definiert werden (siehe Abschnitt B.3 ff). Die Berechnung entsprechender Impedanzspektren sowie deren pragmatischer Eigenschaften T , Σ und Z folgt dann Alg. B.1 und Alg. B.2 (siehe auch Abschnitt 6.5).

Alle modellierten Datensätze, d.h. alle berechneten Impedanzspektren, wurden mit tatsächlich gemessenen Spektren abgeglichen. Das Abgleichsverfahren wird im Einzelnen im Abschnitt B.2 beschrieben. Die Ergebnisse des Abgleichs sind getrennt nach Zellkulturlinien und funktionalen Zellzustände in den nachfolgenden Abschnitten dargestellt.

B.1. Berechnung der transepithelialen Impedanz Z^T

Der angenommene Ersatzschaltkreis C besteht aus zwei seriellen RC-Gliedern a und b mit den Zeitkonstanten $\tau_a = R_a \cdot C_a$ und $\tau_b = R_b \cdot C_b$, einem dazu parallelen Widerstand R_p sowie einem weiteren, diesen Schaltkreiselementen vorgeschalteten seriellen Widerstand R_s (vgl. Kapitel 6.1)¹.

Für die von der Kreisfrequenz ω abhängige Gesamtimpedanz bzw. transepitheliale Impedanz Z^T gilt aufgrund der kirchhoffschen Gesetze:

$$Z^T(\omega) = \frac{(R_p R_a + R_p R_b) + i\omega(R_p R_a \tau_b + R_p R_b \tau_a)}{(R_p + R_a + R_b - \omega^2 R_p \tau_b \tau_a) + i\omega(R_p \tau_a + R_p \tau_b + R_a \tau_b + R_b \tau_a)} + R_s \quad (\text{B.5})$$

Durch Anwendung der dritten binomischen Formel $(x + y)(x - y) = x^2 - y^2$ und durch Multiplikation mit dem Term

$$\frac{(R_p + R_a + R_b - \omega^2 R_p \tau_b \tau_a) - i\omega(R_p \tau_a + R_p \tau_b + R_a \tau_b + R_b \tau_a)}{(R_p + R_a + R_b - \omega^2 R_p \tau_b \tau_a) - i\omega(R_p \tau_a + R_p \tau_b + R_a \tau_b + R_b \tau_a)} = 1$$

ergibt sich folgende Darstellung für Z^T :

$$Z^T(\omega) = \frac{((R_p R_a + R_p R_b) + i\omega(R_p R_a \tau_b + R_p R_b \tau_a))(R_p + R_a + R_b - \omega^2 R_p \tau_b \tau_a) - i\omega(R_p \tau_a + R_p \tau_b + R_a \tau_b + R_b \tau_a)}{(R_p + R_a + R_b - \omega^2 R_p \tau_b \tau_a)^2 - i^2 \omega^2 (R_p \tau_a + R_p \tau_b + R_a \tau_b + R_b \tau_a)^2} + R_s \quad (\text{B.6})$$

$$= \frac{((R_p R_a + R_p R_b)(R_p + R_a + R_b - \omega^2 R_p \tau_b \tau_a) + (R_p R_a + R_p R_b)(-i\omega(R_p \tau_a + R_p \tau_b + R_a \tau_b + R_b \tau_a)))}{(R_p + R_a + R_b - \omega^2 R_p \tau_b \tau_a)^2 + \omega^2 (R_p \tau_a + R_p \tau_b + R_a \tau_b + R_b \tau_a)^2} + \frac{(i\omega(R_p R_a \tau_b + R_p R_b \tau_a))(R_p + R_a + R_b - \omega^2 R_p \tau_b \tau_a) + (i\omega(R_p R_a \tau_b + R_p R_b \tau_a))(-i\omega(R_p \tau_a + R_p \tau_b + R_a \tau_b + R_b \tau_a))}{(R_p + R_a + R_b - \omega^2 R_p \tau_b \tau_a)^2 + \omega^2 (R_p \tau_a + R_p \tau_b + R_a \tau_b + R_b \tau_a)^2} + R_s \quad (\text{B.7})$$

Für Real- und Imaginärteil der komplexwertigen Impedanz $Z^T(\omega) = \Re(Z^T(\omega)) + i \cdot \Im(Z^T(\omega))$ ergibt sich dann:

$$\Re(Z^T(\omega)) = \frac{R_p(R_a + R_b)(R_p + R_a + R_b - \omega^2 R_p \tau_b \tau_a) + \omega^2 (R_p R_a \tau_b + R_p R_b \tau_a)(R_p \tau_a + R_p \tau_b + R_a \tau_b + R_b \tau_a)}{(R_p + R_a + R_b - \omega^2 R_p \tau_b \tau_a)^2 + \omega^2 (R_p \tau_a + R_p \tau_b + R_a \tau_b + R_b \tau_a)^2} + R_s \quad (\text{B.8})$$

$$\Im(Z^T(\omega)) = \omega \frac{R_p(R_a \tau_b + R_b \tau_a)(R_p + R_a + R_b - \omega^2 R_p \tau_b \tau_a) - R_p(R_a + R_b)(R_p \tau_a + R_p \tau_b + R_a \tau_b + R_b \tau_a)}{(R_p + R_a + R_b - \omega^2 R_p \tau_b \tau_a)^2 + \omega^2 (R_p \tau_a + R_p \tau_b + R_a \tau_b + R_b \tau_a)^2} \quad (\text{B.9})$$

¹Um die Darstellung übersichtlicher zu halten, werden hier verkürzte Bezeichnungen für die Schaltkreis-Parameter verwendet (R_a für R^{wp} , R_b für R^{bl} , R_p für R^{para} , etc.).

B.1. Berechnung der transepithelialen Impedanz Z^T

Z	f [Hz]	ω [1/s]
Z_0	1.3	8.16814
Z_1	1.6366	10.2831
Z_2	2.0604	12.9459
Z_3	2.5938	16.2973
Z_4	3.2654	20.5171
Z_5	4.111	25.8302
Z_6	5.1754	32.518
Z_7	6.5154	40.9375
Z_8	8.202	51.5347
Z_9	10.326	64.8802
Z_{10}	13	81.6814
Z_{11}	16.366	102.831
Z_{12}	20.603	129.452
Z_{13}	25.938	162.973
Z_{14}	32.654	205.171
Z_{15}	41.109	258.295
Z_{16}	51.753	325.174
Z_{17}	65.154	409.375
Z_{18}	82.02	515.347
Z_{19}	103.26	648.802
Z_{20}	130	816.814
Z_{21}	163.66	1028.31
Z_{22}	206.03	1294.52
Z_{23}	259.38	1629.73
Z_{24}	326.54	2051.71
Z_{25}	411.09	2582.95
Z_{26}	517.53	3251.74
Z_{27}	651.53	4093.68
Z_{28}	820.2	5153.47
Z_{29}	1032.6	6488.02
Z_{30}	1300	8168.14
Z_{31}	1636.6	10283.1
Z_{32}	2060.3	12945.2
Z_{33}	2593.8	16297.3
Z_{34}	3265.4	20517.1
Z_{35}	4110.9	25829.5
Z_{36}	5175.3	32517.4
Z_{37}	6515.3	40936.8
Z_{38}	8202.5	51537.8
Z_{39}	10326	64880.2
Z_{40}	12999	81675.1
Z_{41}	16350	102730

Tabelle B.1.: Frequenzen der modellierten und gemessenen Impedanzspektren.

B. Modellierung epithelialer Zellkulturen

```

T ← TIMESTAMP
Σ ← 0
Z ← 0
for Cap ∈ [Cminap, Cmaxap] do
  for Cbl ∈ [Cminbl, Cmaxbl] do
    Cepi ← (Cap · Cbl) / (Cap + Cbl)
    if Cepi ∈ [Cminepi, Cmaxepi] then
      for Repi ∈ [Rminepi, Rmaxepi] do
        if Repi ∈ [Rminepi, Rmaxtrans] then
          for Rtrans ∈ [Rmintrans, Rmaxtrans] do
            for Rap ∈ [Rminap, Rmaxap] do
              Rbl ← Rtrans - Cap
              Rpara ← (Rtrans · Repi) / (Rtrans - Repi)
              if Rpara ∈ [Rminpara, Rmaxpara] then
                print_error_measurement()
                T++
              T = T + ε
            T = T + δ
          T = T + δ
        T = T + δ
      T = T + δ
    T = T + δ
  T = T + δ

```

Algorithmus B.1: Variante A einer systematische Synthetisierung von Impedanzspektren (nehme gegebene R^{epi} , R^{ap} und R^{bl} und berechne R^{para} entsprechend (mit $\epsilon > 1$ und $\delta \gg \epsilon$).

```

T ← TIMESTAMP
Σ ← 0
Z ← 0
for Cap ∈ [Cminap, Cmaxap] do
  for Cbl ∈ [Cminbl, Cmaxbl] do
    Cepi ← (Cap · Cbl) / (Cap + Cbl)
    if Cepi ∈ [Cminepi, Cmaxepi] then
      for Repi ∈ [Rminepi, Rmaxepi] do
        if Repi ∈ [Rminepi, Rmaxpara] then
          for Rpara ∈ [Rminpara, Rmaxpara] do
            Rtrans ← (Rpara · Repi) / (Rpara - Repi)
            if Rtrans ∈ [Rmintrans, Rmaxtrans] then
              for Rap ∈ [Rminap, Rmaxap] do
                Rbl ← Rtrans - Cap
                print_error_measurement()
                T++
              T = T + ε
            T = T + δ
          T = T + δ
        T = T + δ
      T = T + δ
    T = T + δ
  T = T + δ

```

Algorithmus B.2: Variante B einer systematische Synthetisierung von Impedanzspektren (nehme gegebenes R^{epi} und R^{para} und berechne R^{ap} und R^{bl} entsprechend (mit $\epsilon > 1$ und $\delta \gg \epsilon$).

B.2. Abgleich modellierter und gemessener Daten

Alle modellierten Datensätze wurden mit Impedanzspektren abgeglichen, die an der zugehörigen Zellkultur unter dem jeweils modellierten funktionalen Zustand gemessen wurden. Das hier beschriebene Verfahren ist eine vollautomatisierte Weiterentwicklung eines früheren Ansatzes zur Quantifizierung der Ähnlichkeit zwischen modellierten und gemessenen Spektren [230].

Zu den besonderen Herausforderungen eines solchen Vergleichs zählt insbesondere, dass a) tatsächliche Messungen gerätebedingt fehlerbehaftet sind (insbesondere bei hohen Frequenzen), und b) vergleichsweise geringe Fallzahlen pro Zelllinie und funktionellem Zustand verfügbar sind ($n < 200$). Da dies eine rein statistisch-geometrische Bewertung der Ähnlichkeit zwischen den diskreten Kurven erschwert, wurden stattdessen die Übereinstimmungen folgender voneinander unabhängiger Messgrößen als Bewertungsmaßstab herangezogen:

- subepithelialer Widerstand mit $R^{sub} = \lim_{\omega \rightarrow \infty} \Re(Z(\omega))$
- epithelialer Widerstand mit $R^{epi} = \lim_{\omega \rightarrow 0} \Re(Z(\omega))$
- epitheliale Kapazität mit $C^{epi} = \frac{1}{\omega_{min} R^{epi}}$

Im Gegensatz zu modellierten, ist jedoch für gemessene Spektren eine exakte Bestimmung des wahren Parameterwerts unmöglich. Die hier betrachteten Verfahren zur Parameterabschätzung weisen insbesondere je nach Kurvenform eine unterschiedliche Zuverlässigkeit auf. Um damit verbundene Ungenauigkeiten als Fehlerquelle zu minimieren, wurden daher je drei Verfahren (siehe Abschnitte B.2.1 bis B.2.3) parallel angewendet und die normierten Differenzen zwischen den ermittelten Werten zur Quantifizierung herangezogen (siehe Abschnitt B.2.4).

B.2.1. Methode M1: Diskrete Approximation

In der klinischen Praxis werden die Parameter R^{sub} , R^{epi} und C^{epi} häufig unmittelbar auf Basis der Real- und Imaginärteile eines Impedanzspektrums angenähert. Dann folgt:

- $R^{sub} = \Re(Z_{41})$
- $R^{epi} = \Re(Z_0)$
- $C^{epi} = \omega_i R^{epi}$ mit $\Im(Z_i) = \min(\{\Im(Z_0), \dots, \Im(Z_{41})\})$

B.2.2. Methode M2: Cole-Cole-Fit

Mittels Cole-Cole-Fit [45] können die Parameter R^{sub} , R^{epi} und C^{epi} bestimmt werden, indem eine Kreisfunktion auf dem komplexwertigen Plot eines Impedanzspektrums gefittet wird.

Aus mehreren konkurrierenden Fit-Verfahren wurde hier die Methode von Kasa [117] gewählt und angewendet. Diese hat neben guten statistischen Eigenschaften [255] auch den Vorteil, dass sie als frei verfügbare Implementierung im R-Package *conicfit* zur Verfügung steht².

R^{sub} und R^{epi} werden dann über die Schnittpunkte x_1 und x_2 (mit $x_1 < x_2$) des Kreises mit der x-Achse berechnet, C^{epi} über die Frequenz der am nächsten am Kreismittelpunkt $c = (x_c, y_c)$ liegenden Impedanz:

- $R^{sub} = x_1$
- $R^{epi} = x_2 - x_1$
- $C^{epi} = \omega_i R^{epi}$ mit $|\Re(Z_i) - x_c| = \min(\{|\Re(Z_0) - x_c|, \dots, |\Re(Z_{41}) - x_c|\})$

²Online abrufbar unter <https://cran.r-project.org/web/packages/conicfit/>

B.2.3. Methode M3: Maschinelle Lernverfahren

Die Bestimmung von R^{sub} , als auch von R^{epi} und C^{epi} aus Impedanzspektren mithilfe maschineller Lernverfahren hat sich als zuverlässig und robust erwiesen [231, 230, 233]. Hier wird dazu jeweils ein künstliches neuronales Netz mit Feedforward-Architektur verwendet. Als Trainingsdaten werden je 50.000 Samples verwendet, als Testdaten 25.000 davon verschiedene Samples.

Als Input-Features zur Bestimmung von R^{sub} dienen analog zu publizierten Verfahren [230] Real- und Imaginärteil der Impedanzen der zehn höchsten Frequenzen, also $\{Z_{32}, \dots, Z_{41}\}$; umgekehrt zur Bestimmung von R^{epi} der Impedanzen der zehn niedrigsten Frequenzen, $\{Z_0, \dots, Z_9\}$, wobei deren Realteile zuvor um den R^{sub} -Wert korrigiert werden. Als Input-Features zur Bestimmung von C^{epi} dienen analog die dazwischen liegenden 20 Impedanzen $\{Z_{12}, \dots, Z_{31}\}$.

In Abhängigkeit von der Zahl der verwendeten Input-Features wird für R^{sub} und R^{epi} eine 20-5-1-Architektur mit einer versteckten Schicht verwendet, für C^{epi} eine 40-20-5-1-Architektur mit zwei versteckten Schichten. Für die Neuronen der Ein- und Ausgabeschicht wird eine lineare Aktivierungsfunktion verwendet, für diejenigen der versteckten Schichten eine sigmoide.

Als Trainingsalgorithmus dient Resilient Backpropagation (Rprop), das dem klassischen Backpropagation sowohl in Geschwindigkeit als auch Robustheit überlegen ist. Alle Neuronalen Netze wurden mit der C-Bibliothek FANN³ realisiert, mit dem sich Netzarchitekturen effizient realisieren und Lernprozesse leicht automatisieren lassen. Die verwendeten Parameter für Rprop bzw. FANN stimmen mit denen aus Anhang A (Tab. A.2) überein.

B.2.4. Ähnlichkeitsdiagramme

Um die Ähnlichkeit zwischen modellierten und gemessenen Spektren trotz ungleicher Fallzahlen plausibel abschätzen zu können, wird ein Binningverfahren verwendet. Die Kongruenz der Mess-Bins mit Modell-Bins dient als relatives Ähnlichkeitsmaß κ_B :

$$\kappa_B = \frac{n_\kappa}{n} \quad (\text{B.10})$$

wobei n_κ die Anzahl der kongruenten Bins und n die Gesamtanzahl der Mess-Bins ist.

Um die Bins zu erhalten, werden drei unabhängige Parameterschätzungen für R^{sub} , R^{epi} oder C^{epi} jeweils in eine zweidimensionale Darstellung überführt. Dazu werden zunächst die Schätzwerte m_1 , m_2 und m_3 für den jeweiligen Parameter mit dem arithmetischen Mittel normiert:

$$m_i = \frac{m_i - \bar{m}_{arithm}}{\bar{m}_{arithm}} \quad (\text{B.11})$$

mit $\bar{m}_{arithm} = (\sum m_i)/3$ und $i \in \{1, 2, 3\}$.

Im nächsten Schritt werden die Differenzen $\delta_1 = m_2 - m_3$, $\delta_2 = m_1 - m_2$ und $\delta_3 = m_1 - m_3$ für alle Mess- und alle Modellspektren gebildet. Sowohl für Mess- als auch Modellspektren spannt dann die Menge aller Tupel $\Delta_\delta = (\delta_1, \delta_2, \delta_3)$ eine Ebene im dreidimensionalen Raum auf, durch deren Drehen auf die Y-Z-Achse man eine zweidimensionale Darstellung erhält. In den nachfolgenden Abschnitten ist die Menge der jeweils 25.000 Differenz-Tupel der Modell-Ebene stets blau dargestellt, diejenigen der Mess-Ebene dagegen schwarz.

Die beiden so erzeugten Punktebenen werden danach auf die maximale Ausdehnung der Mess-Ebene zugeschnitten und mithilfe des R-Packages *hexbin*⁴ jeweils in hexagonale Bins transformiert. Mittels der Package-Funktion *hdiffplot* können dann beide Bin-Plots überlagernd dargestellt und abgeglichen werden. Dabei werden kongruente Bins, in denen also in beiden Ebenen Punkte liegen, grün dargestellt (●), Mess-Bins ohne entsprechenden Modell-Bin dagegen rot (●). Modell-Bins ohne entsprechenden Mess-Bin sind türkis eingefärbt (●).

³Online abrufbar unter <http://leenissen.dk/fann/>

⁴Online abrufbar unter <https://cran.r-project.org/web/packages/hexbin/>

B.3. HT-29/B6

B.3.1. Kontrollbedingungen

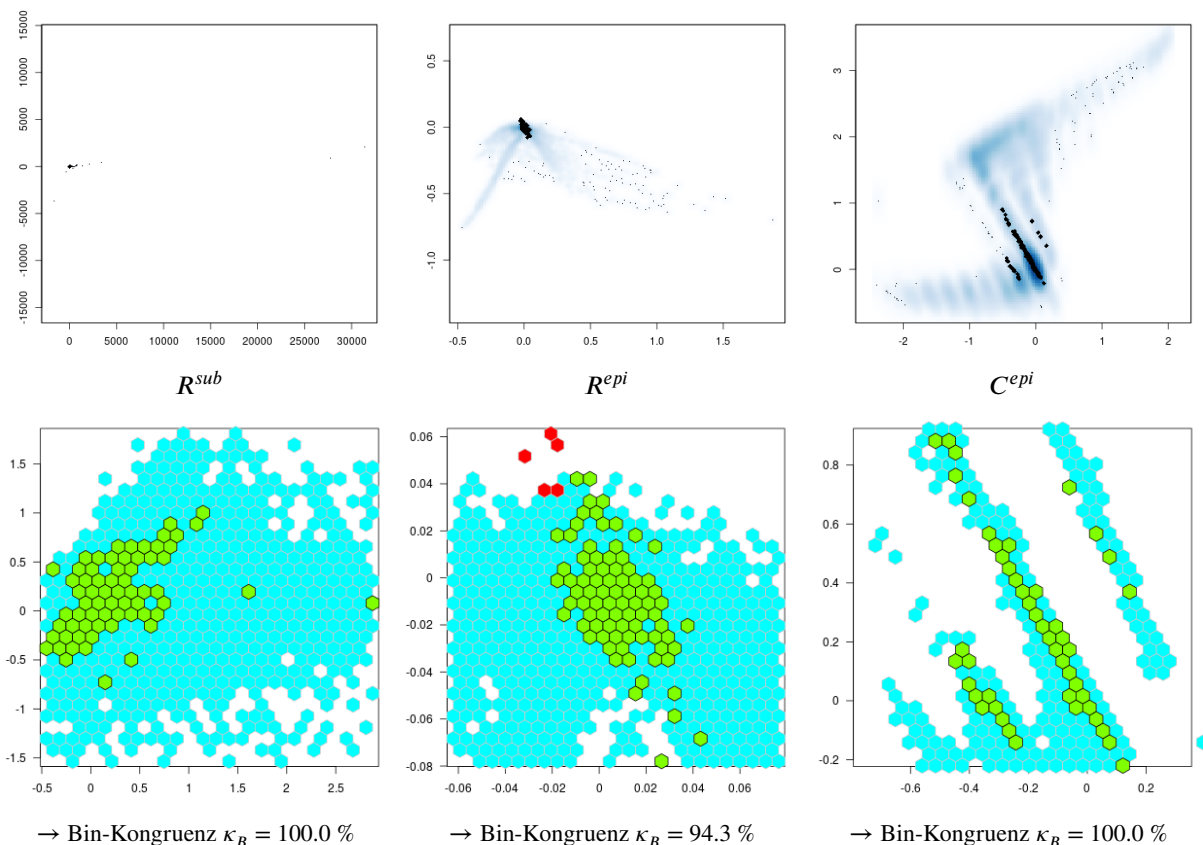
Wertebereiche der Modellparameter wurden auf Basis publizierter Daten abgeleitet (vgl. Tab. 6.1).

Ersatzschaltbild			Modellierter Wertebereich		Schrittweite		Einheit
A	B	C	Minimum	Maximum	R^{sub}	R^{epi}	
R^{sub}	R^{sub}	R^{sub}	1.0	30.0	0.0002	—	$\Omega \cdot cm^2$
R^{epi}			150	1498	—	2.7	$\Omega \cdot cm^2$
	R^{para}	R^{para}	152	30000	—	—	$\Omega \cdot cm^2$
	R^{trans}		151	20000	—	—	$\Omega \cdot cm^2$
		R^{ap}	1	19500	—	—	$\Omega \cdot cm^2$
		R^{bl}	1	19999	—	—	$\Omega \cdot cm^2$
C^{epi}	C^{epi}		1.0	5.5	—	—	$\mu F/cm^2$
		C^{ap}	1.3	8.1	—	—	$\mu F/cm^2$
		C^{bl}	5.0	86.2	—	—	$\mu F/cm^2$

Tabelle B.2.: Modellierte Parameter für die Zelllinie HT-29/B6 unter physiologischen Bedingungen.

Abgleich Modellierung / Messungen

Differenzen der Schätzungen für ungelernete modellierte Impedanzspektren (Sample 001, n=25000) wurden mit Differenzen für gemessene Spektren ohne Wirkstoffzugabe (teils mit Farbstoff Fluoreszin) verglichen (n=171).



B. Modellierung epithelialer Zellkulturen

B.3.2. Manipulation mit Nystatin

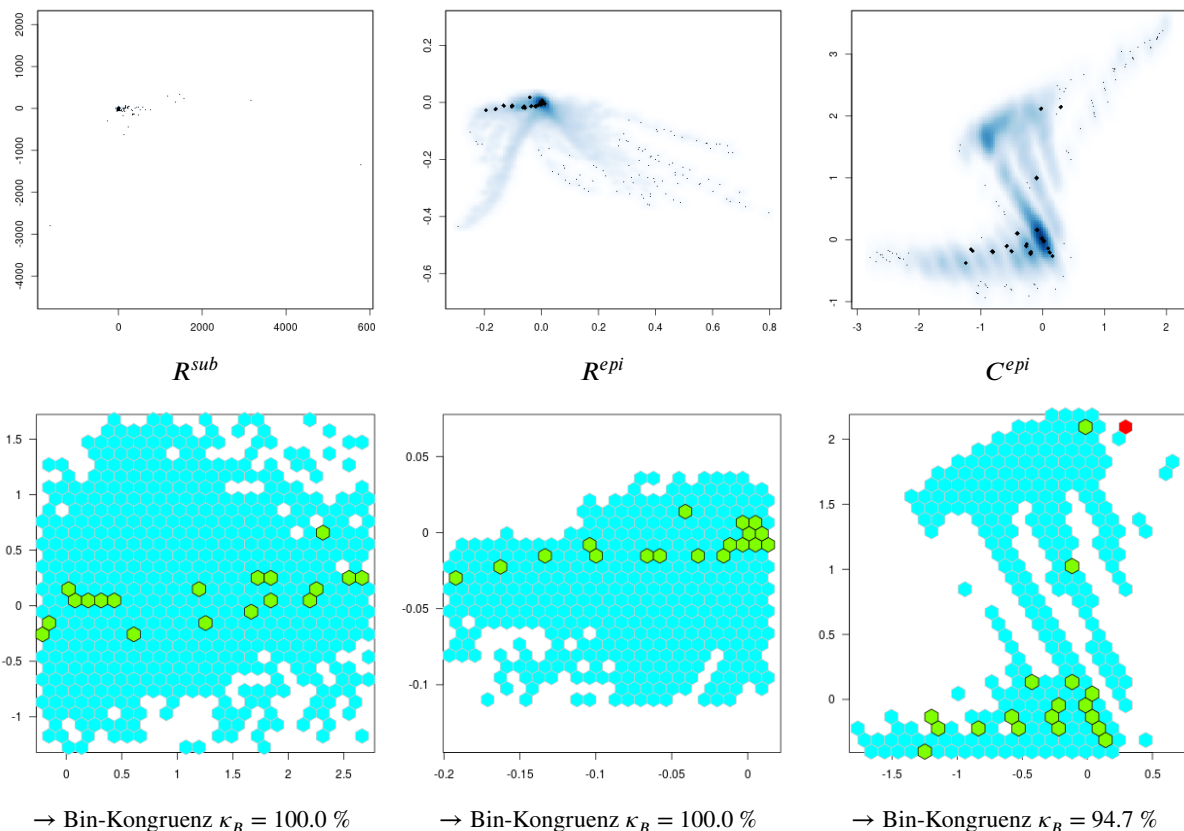
Aus publizierten und gemessenen Werten ist bekannt, dass der Wertebereich von R^{epi} unter Nystatin-Einfluss im Vergleich zu Kontrollbedingungen verkleinert ist. Als Ursache wird modelliert, dass durch den R^{trans} -Modulator Nystatin der mögliche Wertebereich des Parameters R^{trans} verkleinert wird (vgl. Abschnitt 6.2).

Ersatzschaltbild			Modellierter Wertebereich		Schrittweite		Einheit
A	B	C	Minimum	Maximum	R^{sub}	R^{epi}	
R^{sub}	R^{sub}	R^{sub}	1.0	30.0	0.0002	—	$\Omega \cdot cm^2$
	R^{epi}		100	1000	—	1.8	$\Omega \cdot cm^2$
	R^{para}	R^{para}	101	30000	—	—	$\Omega \cdot cm^2$
	R^{trans}		101	20000	—	—	$\Omega \cdot cm^2$
		R^{ap}	1	19500	—	—	$\Omega \cdot cm^2$
		R^{bl}	1	19999	—	—	$\Omega \cdot cm^2$
C^{epi}	C^{epi}		1.0	5.3	—	—	$\mu F/cm^2$
		C^{ap}	1.2	7.2	—	—	$\mu F/cm^2$
		C^{bl}	5.0	87.1	—	—	$\mu F/cm^2$

Tabelle B.3.: Modellierte Parameter für die Zelllinie HT-29/B6 unter Einfluss von Nystatin.

Abgleich Modellierung / Messungen

Differenzen der Schätzungen für ungelernete modellierte Impedanzspektren (Sample 003, n=25000) wurden mit Differenzen für gemessene Spektren unter apikaler oder basolateraler Nystatin-Zugabe verglichen (n=23).



B.3.3. Manipulation mit EGTA

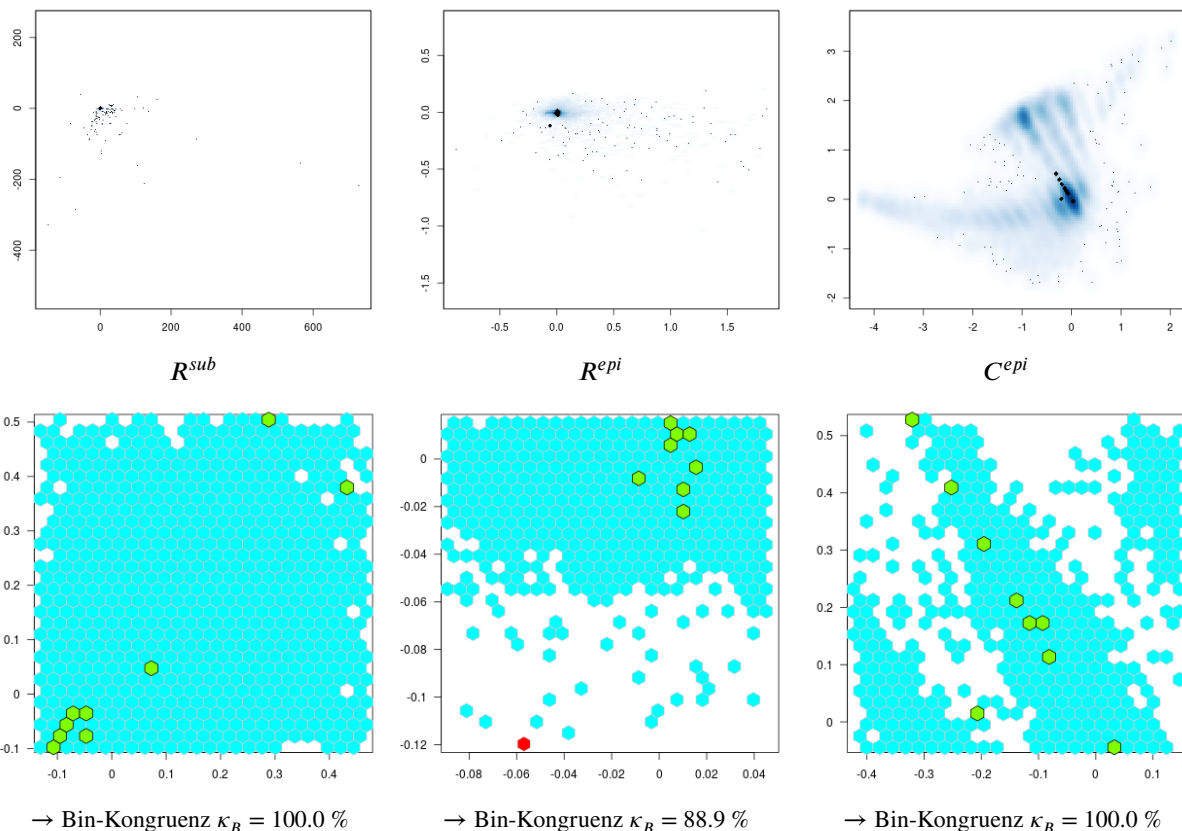
Aus publizierten und gemessenen Werten ist bekannt, dass der Wertebereich von R^{epi} unter EGTA-Einfluss im Vergleich zu Kontrollbedingungen verkleinert ist. Als Ursache wird angenommen, dass durch Wirkung des R^{para} -Modulators EGTA der mögliche Wertebereich des Parameters R^{para} verkleinert ist (vgl. Abschnitt 6.2).

Ersatzschaltbild			Modellierter Wertebereich		Schrittweite		Einheit
A	B	C	Minimum	Maximum	R^{sub}	R^{epi}	
R^{sub}	R^{sub}	R^{sub}	1.0	30.0	0.0002	—	$\Omega \cdot cm^2$
	R^{epi}		5	300	—	0.6	$\Omega \cdot cm^2$
	R^{para}	R^{para}	5	15000	—	—	$\Omega \cdot cm^2$
	R^{trans}		10	20000	—	—	$\Omega \cdot cm^2$
		R^{ap}	1	19500	—	—	$\Omega \cdot cm^2$
		R^{bl}	1	19999	—	—	$\Omega \cdot cm^2$
C^{epi}	C^{epi}		1.2	5.4	—	—	$\mu F/cm^2$
		C^{ap}	1.3	6.0	—	—	$\mu F/cm^2$
		C^{bl}	6.6	86.7	—	—	$\mu F/cm^2$

Tabelle B.4.: Modellierter Parameter für die Zelllinie HT-29/B6 unter Einfluss von EGTA.

Abgleich Modellierung / Messungen

Differenzen der Schätzungen für ungelernete modellierte Impedanzspektren (Sample 004, n=25000) wurden mit Differenzen für gemessene Spektren unter Zugabe von EGTA verglichen (n=10).



B. Modellierung epithelialer Zellkulturen

B.3.4. Manipulation mit EGTA und Nystatin

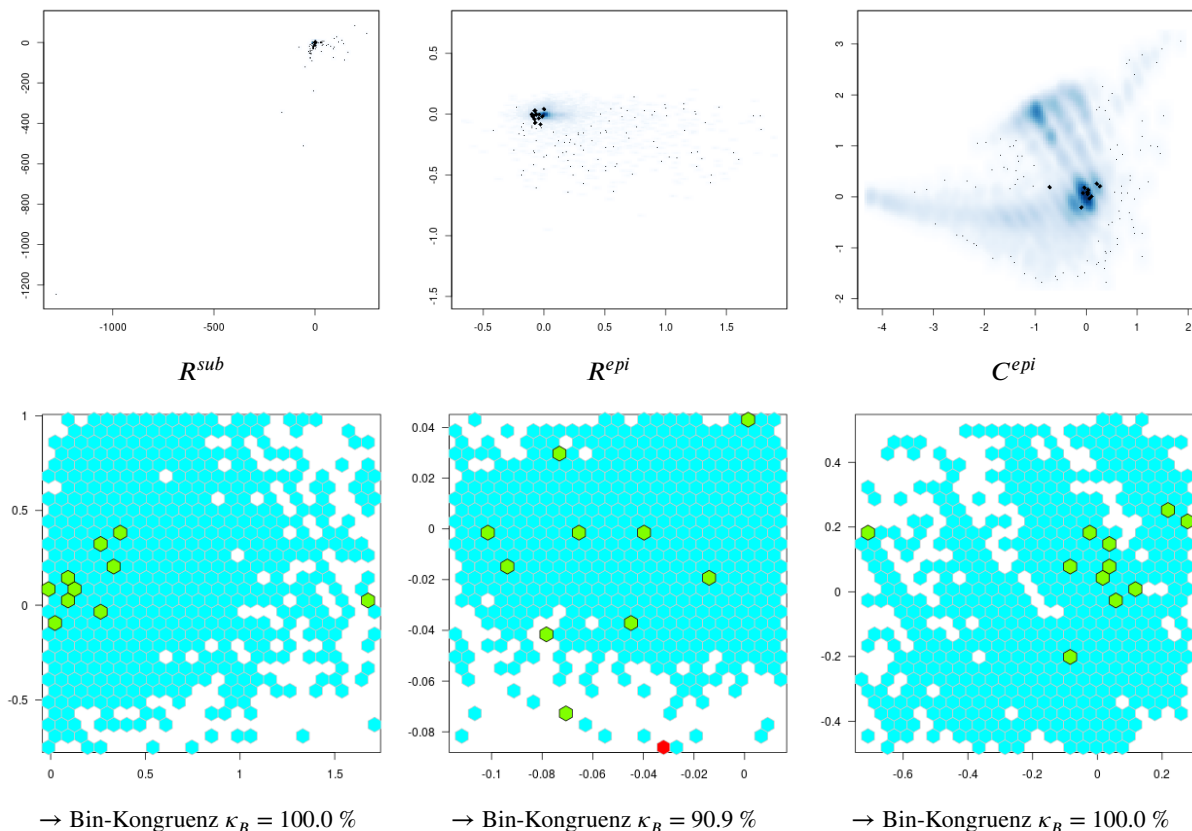
Unter Einfluss von EGTA und Nystatin ist der Wertebereich von R^{epi} im Vergleich zu Kontrollbedingungen verkleinert. Als Ursache wird modelliert, dass durch EGTA der mögliche Wertebereich des Parameters R^{para} und durch Nystatin der mögliche Wertebereich des Parameters R^{trans} verkleinert ist (vgl. Abschnitt 6.2).

Ersatzschaltbild			Modellierter Wertebereich		Schrittweite		Einheit
A	B	C	Minimum	Maximum	R^{sub}	R^{epi}	
R^{sub}	R^{sub}	R^{sub}	1.0	30.0	0.0002	—	$\Omega \cdot cm^2$
		R^{epi}	5	250	—	0.5	$\Omega \cdot cm^2$
	R^{para}	R^{para}	5	15000	—	—	$\Omega \cdot cm^2$
	R^{trans}		5	20000	—	—	$\Omega \cdot cm^2$
		R^{ap}	1	19500	—	—	$\Omega \cdot cm^2$
		R^{bl}	1	19999	—	—	$\Omega \cdot cm^2$
C^{epi}	C^{epi}		1.3	5.3	—	—	$\mu F/cm^2$
		C^{ap}	1.4	7.4	—	—	$\mu F/cm^2$
		C^{bl}	8.0	83.1	—	—	$\mu F/cm^2$

Tabelle B.5.: Modellierte Parameter für die Zelllinie HT-29/B6 unter Einfluss von EGTA und Nystatin.

Abgleich Modellierung / Messungen

Differenzen der Schätzungen für ungelernete modellierte Impedanzspektren (Sample 003, n=25000) wurden mit Differenzen für gemessene Spektren unter Zugabe von EGTA und Nystatin verglichen (n=11).



B.4. IPEC-J2

B.4.1. Kontrollbedingungen

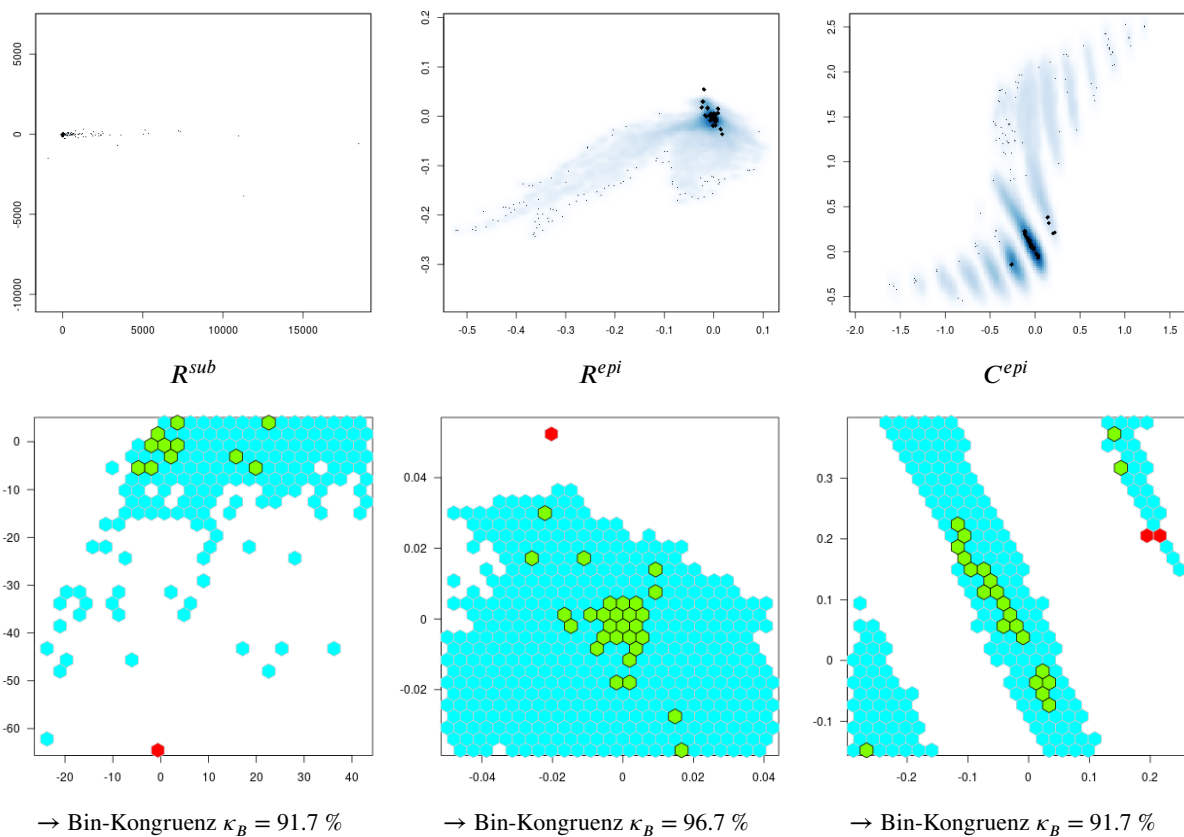
Wertebereiche der Modellparameter wurden auf Basis publizierter Daten abgeleitet (vgl. Tab. 6.2).

Ersatzschaltbild			Modellierter Wertebereich		Schrittweite		Einheit
A	B	C	Minimum	Maximum	R^{sub}	R^{epi}	
R^{sub}	R^{sub}	R^{sub}	0.3	50.0	0.0004	—	$\Omega \cdot cm^2$
	R^{epi}		900	8567	—	16.0	$\Omega \cdot cm^2$
	R^{para}	R^{para}	944	15000	—	—	$\Omega \cdot cm^2$
	R^{trans}		3000	20000	—	—	$\Omega \cdot cm^2$
		R^{ap}	2500	19000	—	—	$\Omega \cdot cm^2$
		R^{bl}	1	17500	—	—	$\Omega \cdot cm^2$
C^{epi}	C^{epi}		0.7	2.0	—	—	$\mu F/cm^2$
		C^{ap}	1.1	3.3	—	—	$\mu F/cm^2$
		C^{bl}	1.5	9.3	—	—	$\mu F/cm^2$

Tabelle B.6.: Modellierte Parameter für die Zelllinie IPEC-J2 unter physiologischen Bedingungen.

Abgleich Modellierung / Messungen

Differenzen der Schätzungen für ungelernete modellierte Impedanzspektren (Sample 001, n=25000) wurden mit Differenzen für gemessene Spektren ohne Wirkstoffzugabe (teils mit Farbstoff Fluoreszin) verglichen (n=42).



B. Modellierung epithelialer Zellkulturen

B.4.2. Manipulation mit Nystatin

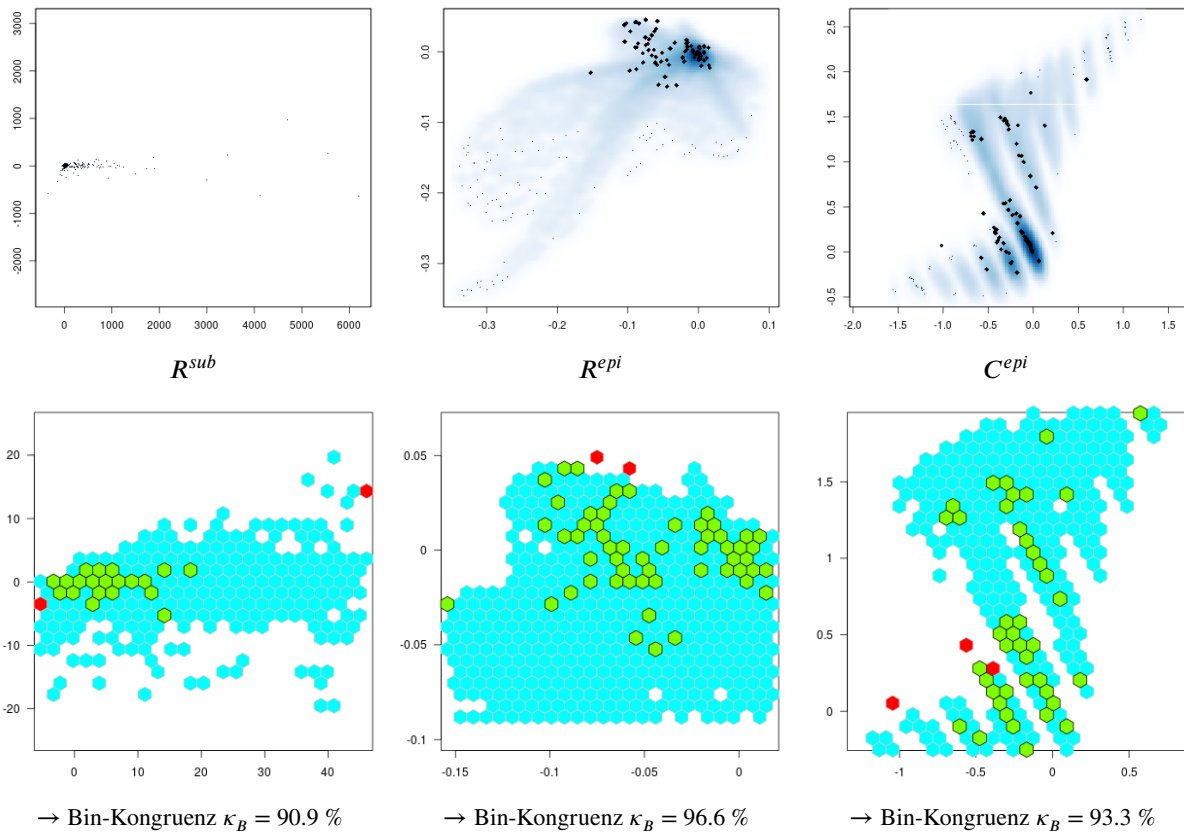
Aus publizierten und gemessenen Werten ist bekannt, dass der Wertebereich von R^{epi} unter Nystatin-Einfluss verkleinert ist. Als Ursache wird angenommen, dass durch den R^{trans} -Modulator Nystatin der mögliche Wertebereich des Parameters R^{trans} verkleinert wird (vgl. Abschnitt 6.2).

Ersatzschaltbild			Modellierter Wertebereich		Schrittweite		Einheit
A	B	C	Minimum	Maximum	R^{sub}	R^{epi}	
R^{sub}	R^{sub}	R^{sub}	0.3	50.0	0.0004	—	$\Omega \cdot cm^2$
	R^{epi}		502	8562	—	17.0	$\Omega \cdot cm^2$
	R^{para}	R^{para}	900	15000	—	—	$\Omega \cdot cm^2$
	R^{trans}		520	20000	—	—	$\Omega \cdot cm^2$
		R^{ap}	350	19000	—	—	$\Omega \cdot cm^2$
		R^{bl}	1	19650	—	—	$\Omega \cdot cm^2$
C^{epi}	C^{epi}		0.7	2.0	—	—	$\mu F/cm^2$
		C^{ap}	1.1	3.3	—	—	$\mu F/cm^2$
		C^{bl}	1.4	9.6	—	—	$\mu F/cm^2$

Tabelle B.7.: Modellierte Parameter für die Zelllinie IPEC-J2 unter Einfluss von Nystatin.

Abgleich Modellierung / Messungen

Differenzen der Schätzungen für ungelernete modellierte Impedanzspektren (Sample 001, n=25000) wurden mit Differenzen für gemessene Spektren unter apikaler oder basolateraler Nystatin-Zugabe verglichen (n=72).



B.4.3. Manipulation mit EGTA

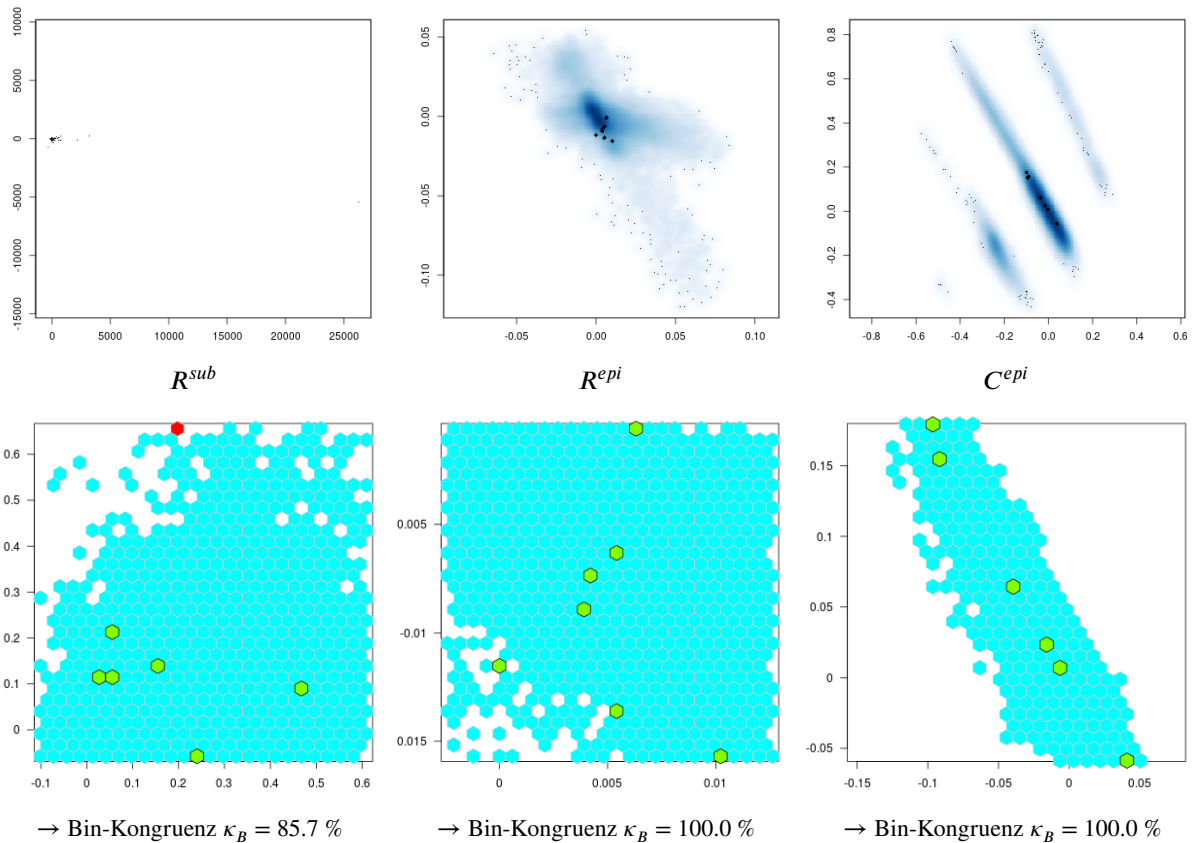
Aus publizierten und gemessenen Werten ist bekannt, dass der Wertebereich von R^{epi} unter EGTA-Einfluss im Vergleich zu Kontrollbedingungen verkleinert ist. Als Ursache wird modelliert, dass durch Wirkung des R^{para} -Modulators EGTA der mögliche Wertebereich des Parameters R^{para} verkleinert ist (vgl. Abschnitt 6.2).

Ersatzschaltbild			Modellierter Wertebereich		Schrittweite		Einheit
A	B	C	Minimum	Maximum	R^{sub}	R^{epi}	
R^{sub}	R^{sub}	R^{sub}	0.3	50.0	0.0004	—	$\Omega \cdot cm^2$
	R^{epi}		150	1498	—	2.7	$\Omega \cdot cm^2$
	R^{para}	R^{para}	151	2000	—	—	$\Omega \cdot cm^2$
	R^{trans}		3000	20000	—	—	$\Omega \cdot cm^2$
		R^{ap}	2500	19000	—	—	$\Omega \cdot cm^2$
		R^{bl}	1	17500	—	—	$\Omega \cdot cm^2$
C^{epi}	C^{epi}		0.7	2.0	—	—	$\mu F/cm^2$
		C^{ap}	1.1	3.3	—	—	$\mu F/cm^2$
		C^{bl}	1.6	9.6	—	—	$\mu F/cm^2$

Tabelle B.8.: Modellierte Parameter für die Zelllinie IPEC-J2 unter Einfluss von EGTA.

Abgleich Modellierung / Messungen

Differenzen der Schätzungen für ungelernete modellierte Impedanzspektren (Sample 005, n=25000) wurden mit Differenzen für gemessene Spektren unter Zugabe von EGTA verglichen (n=7).



B. Modellierung epithelialer Zellkulturen

B.4.4. Manipulation mit EGTA und Nystatin

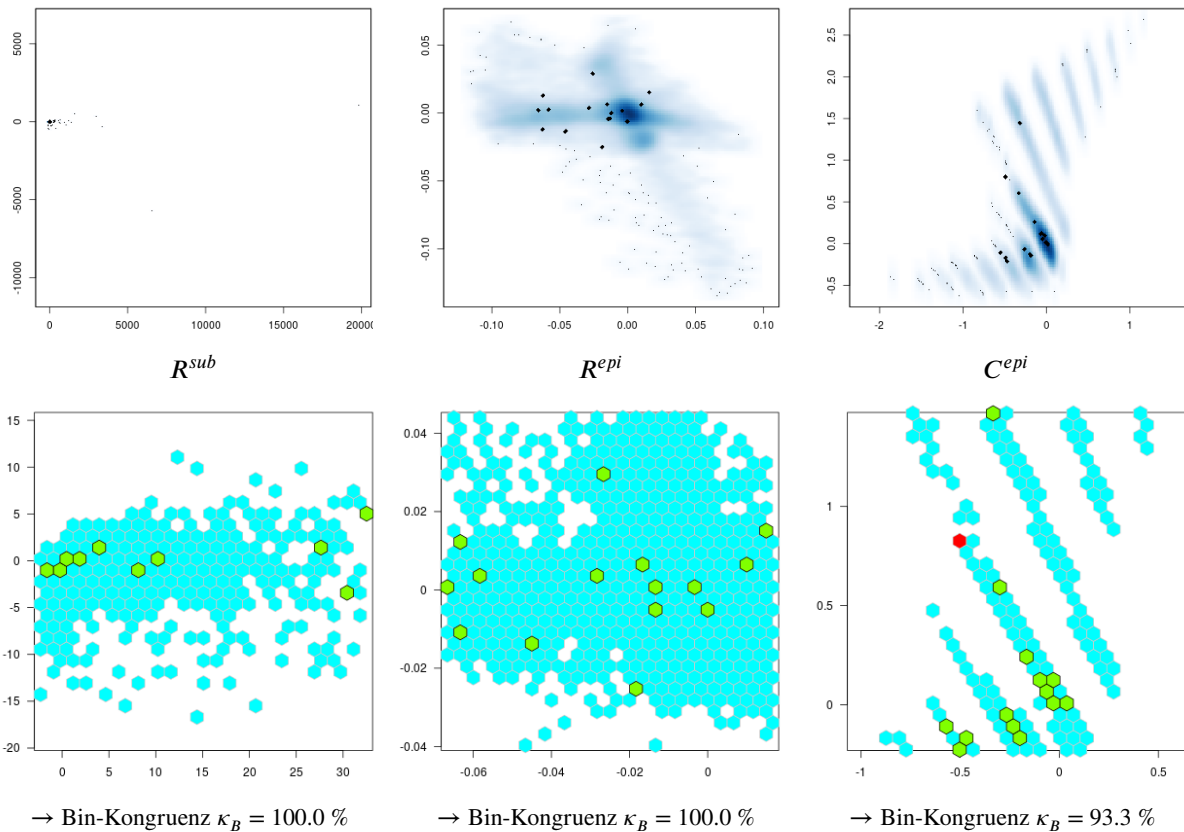
Unter Einfluss von EGTA und Nystatin ist der Wertebereich von R^{epi} im Vergleich zu Kontrollbedingungen verkleinert. Als Ursache wird modelliert, dass durch EGTA der mögliche Wertebereich des Parameters R^{para} und durch Nystatin der mögliche Wertebereich des Parameters R^{trans} verkleinert ist (vgl. Abschnitt 6.2).

Ersatzschaltbild			Modellierter Wertebereich		Schrittweite		Einheit
A	B	C	Minimum	Maximum	R^{sub}	R^{epi}	
R^{sub}	R^{sub}	R^{sub}	0.3	50.0	0.0004	—	$\Omega \cdot cm^2$
	R^{epi}		150	1395	—	2.7	$\Omega \cdot cm^2$
	R^{para}	R^{para}	152	1500	—	—	$\Omega \cdot cm^2$
	R^{trans}		500	20000	—	—	$\Omega \cdot cm^2$
		R^{ap}	350	19000	—	—	$\Omega \cdot cm^2$
		R^{bl}	1	19650	—	—	$\Omega \cdot cm^2$
C^{epi}	C^{epi}		0.7	1.9	—	—	$\mu F/cm^2$
		C^{ap}	1.1	3.2	—	—	$\mu F/cm^2$
		C^{bl}	1.5	9.6	—	—	$\mu F/cm^2$

Tabelle B.9.: Modellierter Parameter für die Zelllinie IPEC-J2 unter Einfluss von EGTA und Nystatin.

Abgleich Modellierung / Messungen

Differenzen der Schätzungen für ungelernete modellierte Impedanzspektren (Sample 001, n=25000) wurden mit Differenzen für gemessene Spektren unter Zugabe von EGTA und Nystatin verglichen (n=17).



B.5. MDCK I

B.5.1. Kontrollbedingungen

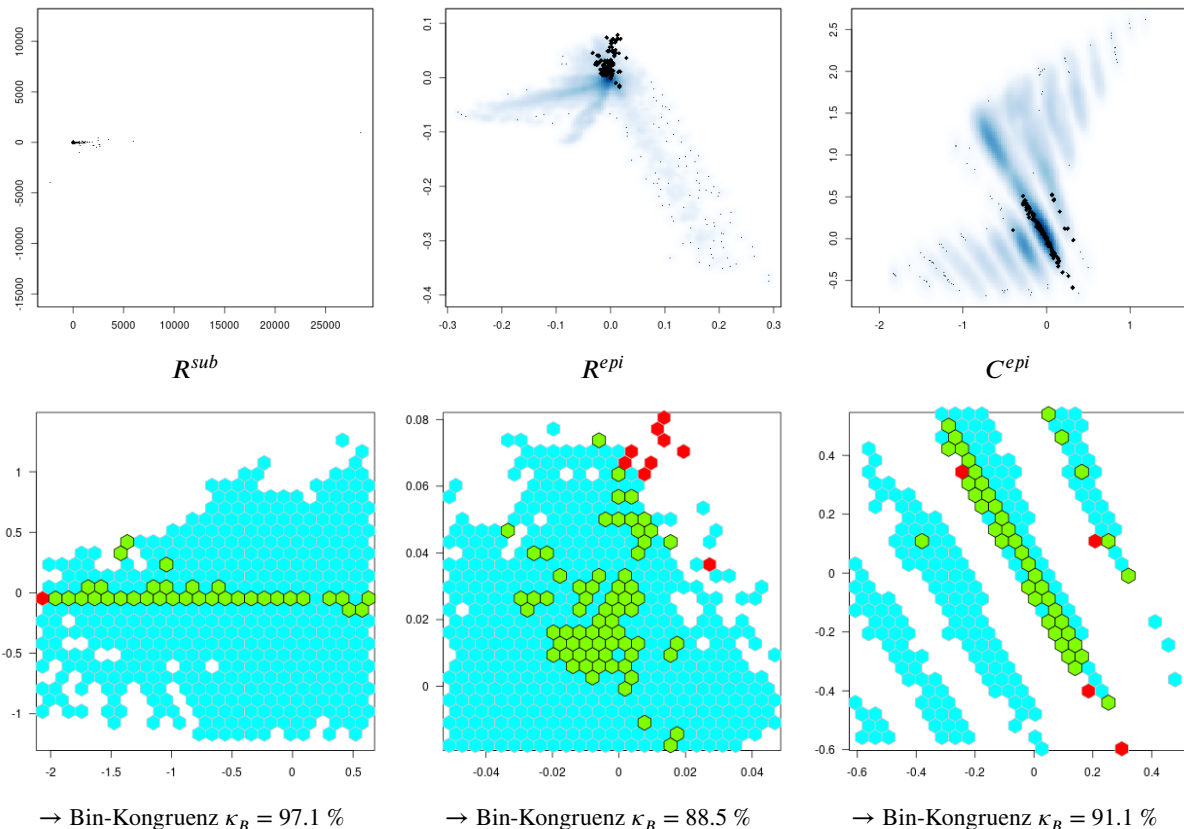
Wertebereiche der Modellparameter wurden auf Basis publizierter Daten abgeleitet (vgl. Tab. 6.3).

Ersatzschaltbild			Modellierter Wertebereich		Schrittweite		Einheit
A	B	C	Minimum	Maximum	R^{sub}	R^{epi}	
R^{sub}	R^{sub}	R^{sub}	5.0	25.0	0.00015	—	$\Omega \cdot cm^2$
	R^{epi}		100	4495	—	8.75	$\Omega \cdot cm^2$
	R^{para}	R^{para}	102	10000	—	—	$\Omega \cdot cm^2$
	R^{trans}		101	15000	—	—	$\Omega \cdot cm^2$
		R^{ap}	10	14500	—	—	$\Omega \cdot cm^2$
		R^{bl}	1	14990	—	—	$\Omega \cdot cm^2$
C^{epi}	C^{epi}		0.8	3.3	—	—	$\mu F/cm^2$
		C^{ap}	0.9	5.7	—	—	$\mu F/cm^2$
		C^{bl}	2.9	9.5	—	—	$\mu F/cm^2$

Tabelle B.10.: Modellierte Parameter für die Zelllinie MDCK I unter physiologischen Bedingungen.

Abgleich Modellierung / Messungen

Differenzen der Schätzungen für ungelernete modellierte Impedanzspektren (Sample 005, n=25000) wurden mit Differenzen für gemessene Spektren ohne Wirkstoffzugabe (teils mit Farbstoff Fluoreszin) verglichen (n=115).



B. Modellierung epithelialer Zellkulturen

B.5.2. Manipulation mit Nystatin

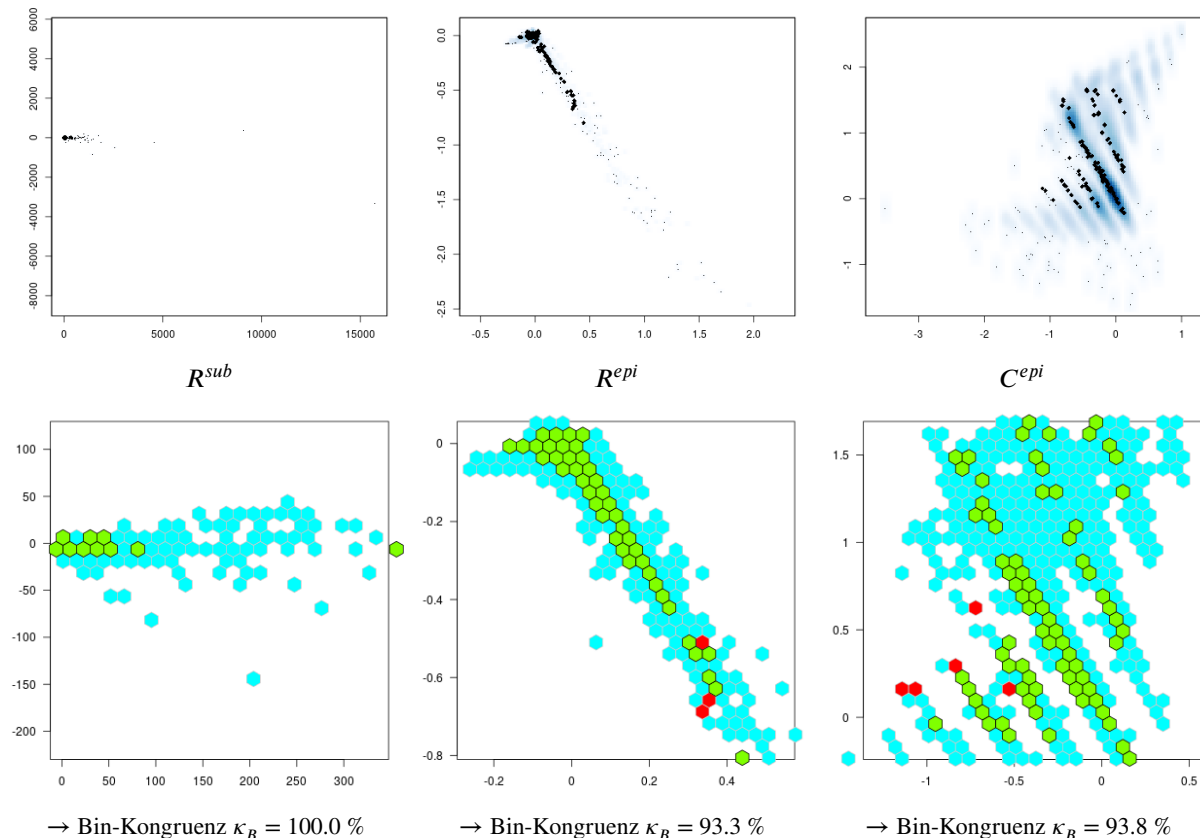
Aus publizierten und gemessenen Werten ist bekannt, dass der Wertebereich von R^{epi} unter Nystatin-Einfluss verkleinert ist. Als Ursache wird angenommen, dass durch den R^{trans} -Modulator Nystatin der mögliche Wertebereich des Parameters R^{trans} verkleinert wird (vgl. Abschnitt 6.2).

Ersatzschaltbild			Modellierter Wertebereich		Schrittweite		Einheit
A	B	C	Minimum	Maximum	R^{sub}	R^{epi}	
R^{sub}	R^{sub}	R^{sub}	5.0	25.0	0.00015	—	$\Omega \cdot cm^2$
	R^{epi}		10	4243	—	8.25	$\Omega \cdot cm^2$
	R^{para}	R^{para}	100	10000	—	—	$\Omega \cdot cm^2$
	R^{trans}		10	10000	—	—	$\Omega \cdot cm^2$
		R^{ap}	5	9500	—	—	$\Omega \cdot cm^2$
		R^{bl}	1	9995	—	—	$\Omega \cdot cm^2$
C^{epi}	C^{epi}		0.8	3.5	—	—	$\mu F/cm^2$
		C^{ap}	0.9	5.9	—	—	$\mu F/cm^2$
		C^{bl}	2.7	9.6	—	—	$\mu F/cm^2$

Tabelle B.11.: Modellierter Parameter für die Zelllinie MDCK I unter Einfluss von Nystatin.

Abgleich Modellierung / Messungen

Differenzen der Schätzungen für ungelernete modellierte Impedanzspektren (Sample 001, n=25000) wurden mit Differenzen für gemessene Spektren unter apikaler oder basolateraler Nystatin-Zugabe verglichen (n=135).



B.5.3. Manipulation mit EGTA

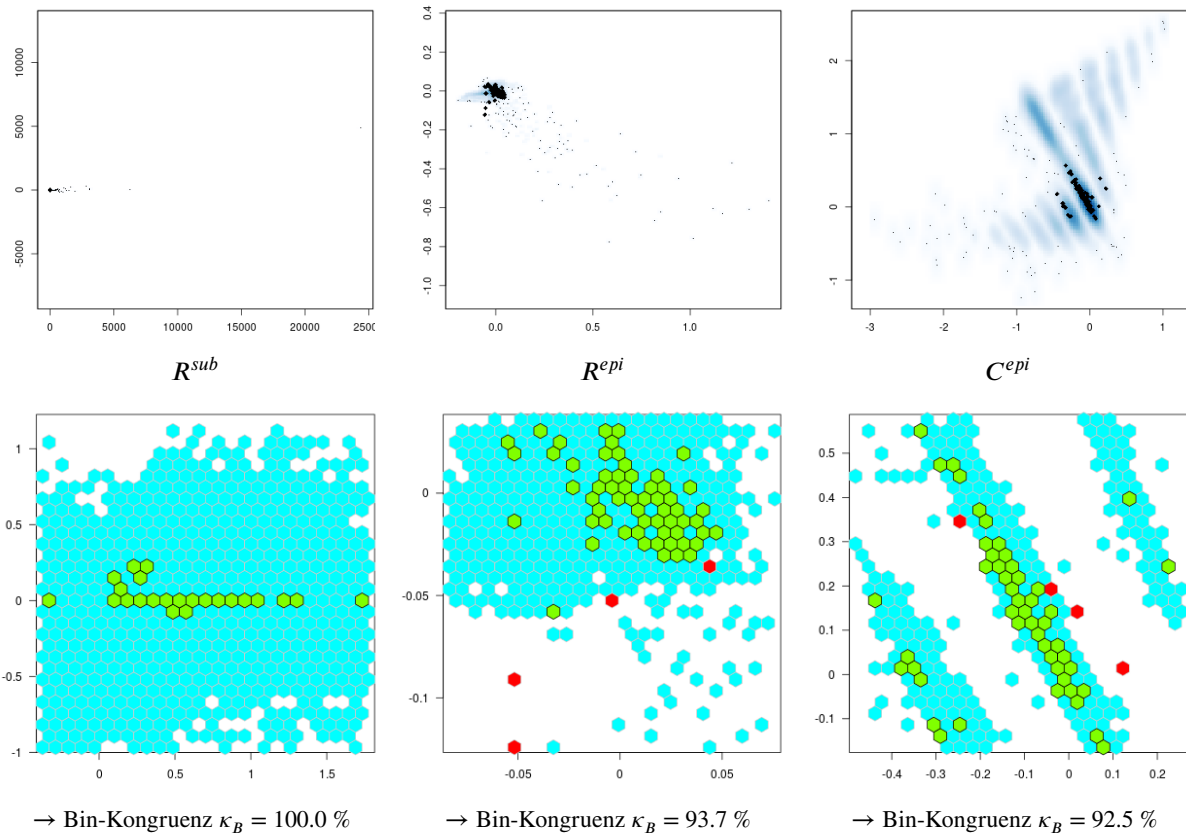
Aus publizierten und gemessenen Werten ist bekannt, dass der Wertebereich von R^{epi} unter EGTA-Einfluss im Vergleich zu Kontrollbedingungen verkleinert ist. Als Ursache wird modelliert, dass durch Wirkung des R^{para} -Modulators EGTA der mögliche Wertebereich des Parameters R^{para} verkleinert ist (vgl. Abschnitt 6.2).

Ersatzschaltbild			Modellierter Wertebereich		Schrittweite		Einheit
A	B	C	Minimum	Maximum	R^{sub}	R^{epi}	
R^{sub}	R^{sub}	R^{sub}	5.0	25.0	0.00015	—	$\Omega \cdot cm^2$
	R^{epi}		10	3149	—	6.0	$\Omega \cdot cm^2$
	R^{para}	R^{para}	10	5000	—	—	$\Omega \cdot cm^2$
	R^{trans}		100	15000	—	—	$\Omega \cdot cm^2$
		R^{ap}	10	14500	—	—	$\Omega \cdot cm^2$
		R^{bl}	1	14990	—	—	$\Omega \cdot cm^2$
C^{epi}	C^{epi}		0.7	3.4	—	—	$\mu F/cm^2$
		C^{ap}	0.9	6.3	—	—	$\mu F/cm^2$
		C^{bl}	2.6	9.6	—	—	$\mu F/cm^2$

Tabelle B.12.: Modellierte Parameter für die Zelllinie MDCK I unter Einfluss von EGTA.

Abgleich Modellierung / Messungen

Differenzen der Schätzungen für ungelernete modellierte Impedanzspektren (Sample 001, n=25000) wurden mit Differenzen für gemessene Spektren unter Zugabe von EGTA verglichen (n=121).



B. Modellierung epithelialer Zellkulturen



C

Charakterisierung der Datenbasis

Die zu analysierende Datenbasis besteht aus Stichproben aller modellierten Impedanz-Datensätze (siehe Kapitel 6). Die Stichprobengröße liegt dabei jeweils bei $n = 25.000$, so dass sich zusammengenommen eine zu analysierende Datenbasis aus 275.000 Impedanzspektren ergibt (vgl. Tab. 6.4). Diese werden unterteilt in einen Trainingsdatensatz von 200.000 Spektren, der für Exploration und Adaption genutzt wird, und einen Testdatensatz von 75.000 Spektren, der zur Evaluation der Adaption genutzt wird.

Neben den komplexwertigen Impedanzen dieser Spektren werden mittels rechnerischer Transformation alternative Darstellungen bzw. zusätzliche Input-Features erzeugt. So sind etwa die Darstellung in kartesischen und Polarkoordinaten äquivalent und zur Analyse von Impedanzspektren grundsätzlich gleichermaßen geeignet [84]. Sowohl aus kartesischen als auch Polarkoordinaten lassen sich weiter statistische Kennzahlen zu Real- und Imaginärteil bzw. Magnitude und Phasenwinkel berechnen, die als Input-Features für maschinelle Lernverfahren verwendet werden können [232, 233].

Als potentiell zu erlernende Zielparameter oder Output-Features stehen neben den zugrundeliegenden Ersatzschaltkreis-Parametern auch implizite kategoriale Zuordnungen zur Verfügung, etwa zu einer Zellkultur oder einem Zellzustand. Darüber hinaus lassen sich physiologisch relevante Parameter wie das Verhältnis zwischen trans- und parazellulärem Widerstand oder das Verhältnis zwischen den Zeitkonstanten der beiden RC-Glieder des Ersatzschaltkreises berechnen und als Zielparameter verwenden [232].

C.1. Input-Features

Die verwendeten Input-Features sind entweder explizit messbar oder implizit aus messbaren Features berechenbar. Sowohl explizite als auch implizite Features werden im Folgenden als separate Sets betrachtet.

Explizite Feature-Sets stellen dabei insbesondere Sets aus Real- und Imaginärteil der komplexwertigen Impedanzen Z_i dar, die für die $n = 42$ verwendeten Kreisfrequenzen ω_i berechnet wurden (Tab. B.1):

$$S_{\Re} = \{\Re(Z_0), \dots, \Re(Z_{n-1})\} \quad (\text{C.1})$$

$$S_{\Im} = \{\Im(Z_0), \dots, \Im(Z_{n-1})\} \quad (\text{C.2})$$

Der in der Modellierung enthaltene Parameter R^{sub} , der eine Verschiebung des Gesamtspektrums entlang der Achse des Realteils bewirkt, wird in der Analyse nicht betrachtet. Daher wird dieser vor der weiteren Verarbeitung aus dem Feature-Set S_{\Re} herausgerechnet, und dieses durch das Set $S_{\Re'}$ ersetzt:

C. Charakterisierung der Datenbasis

Feature	Berechnung	Bezeichnung
$\min(\mathcal{S}_{\mathfrak{R}})$	$\{x : x \leq x_i \ \forall x_i \in \mathcal{S}_{\mathfrak{R}}\}$	Minimum
$\max(\mathcal{S}_{\mathfrak{R}})$	$\{x : x_i \leq x \ \forall x_i \in \mathcal{S}_{\mathfrak{R}}\}$	Maximum
$R(\mathcal{S}_{\mathfrak{R}})$	$\max(\mathcal{S}_{\mathfrak{R}}) - \min(\mathcal{S}_{\mathfrak{R}})$	Spannweite
$P_{0,1}(\mathcal{S}_{\mathfrak{R}}^*)$	$x_{[0,1 \cdot n+1]}$	1. Perzentil (für das sortierte Set $\mathcal{S}_{\mathfrak{R}}^*$)
$Q_{0,25}(\mathcal{S}_{\mathfrak{R}}^*)$	$x_{[0,25 \cdot n+1]}$	1. Quartil (für das sortierte Set $\mathcal{S}_{\mathfrak{R}}^*$)
$Q_{0,75}(\mathcal{S}_{\mathfrak{R}}^*)$	$x_{[0,75 \cdot n+1]}$	3. Quartil (für das sortierte Set $\mathcal{S}_{\mathfrak{R}}^*$)
$P_{0,9}(\mathcal{S}_{\mathfrak{R}}^*)$	$x_{[0,9 \cdot n+1]}$	9. Perzentil (für das sortierte Set $\mathcal{S}_{\mathfrak{R}}^*$)
$R_{IQ}(\mathcal{S}_{\mathfrak{R}})$	$Q_{0,75}(\mathcal{S}_{\mathfrak{R}}^*) - Q_{0,25}(\mathcal{S}_{\mathfrak{R}}^*)$	Interquartilsabstand
$R_{IP}(\mathcal{S}_{\mathfrak{R}})$	$Q_{0,9}(\mathcal{S}_{\mathfrak{R}}^*) - Q_{0,1}(\mathcal{S}_{\mathfrak{R}}^*)$	Interperzentilsabstand
$\bar{x}_{med}(\mathcal{S}_{\mathfrak{R}}^*)$	$\begin{cases} x_{(\frac{n+1}{2})}, & \text{falls } n \text{ ungerade} \\ \frac{1}{2}(x_{(\frac{n}{2})} + x_{(\frac{n}{2}+1)}), & \text{falls } n \text{ gerade} \end{cases}$	Median (für das sortierte Set $\mathcal{S}_{\mathfrak{R}}^*$)
$\bar{x}_{arithm}(\mathcal{S}_{\mathfrak{R}})$	$\frac{1}{n} \sum_{i=1}^n x_i$	Arithmetisches Mittel
$\bar{x}_{geom}(\mathcal{S}_{\mathfrak{R}})$	$\sqrt[n]{\prod_{i=1}^n x_i}$	Geometrisches Mittel
$\bar{x}_{harm}(\mathcal{S}_{\mathfrak{R}})$	$\frac{n}{\sum_{i=1}^n \frac{1}{x_i}}$	Harmonisches Mittel
$s^2(\mathcal{S}_{\mathfrak{R}})$	$\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x}_{arithm}(\mathcal{S}_{\mathfrak{R}}))^2$	Varianz
$s(\mathcal{S}_{\mathfrak{R}})$	$\sqrt{s^2(\mathcal{S}_{\mathfrak{R}})}$	Standardabweichung
$R_{MM}(\mathcal{S}_{\mathfrak{R}})$	$\bar{x}_{med}(\mathcal{S}_{\mathfrak{R}}) - \bar{x}_{arithm}(\mathcal{S}_{\mathfrak{R}})$	Abstand zwischen Median und arithmetischem Mittel

Tabelle C.1.: Globale Features für das Feature-Set $\mathcal{S}_{\mathfrak{R}}$ mit $x_i \in \mathcal{S}_{\mathfrak{R}}$ und $n = \#(\mathcal{S}_{\mathfrak{R}})$.

$$\mathcal{S}_{\mathfrak{R}'} = \{\Re(Z_0) - R^{sub}, \dots, \Re(Z_{n-1}) - R^{sub}\} \quad (C.3)$$

Die Polarkoordinaten-Darstellung der Impedanzen wird hier aus der komplexwertigen Darstellung, bzw. aus $\mathcal{S}_{\mathfrak{R}'}$ und \mathcal{S}_{\Im} , errechnet. Da Magnitude $r(Z)$ und Phasenwinkel $\phi(Z)$ jedoch auch direkte Messgrößen sind, werden die zugehörigen Sets hier ebenfalls als explizite Feature-Sets bezeichnet:

$$\mathcal{S}_{\phi} = \{\phi(Z_0), \dots, \phi(Z_{n-1})\} \quad (C.4)$$

$$\mathcal{S}_r = \{r(Z_0), \dots, r(Z_{n-1})\} \quad (C.5)$$

Aus diesen Polarkoordinaten-Features werden zusätzlich zwei Sets aus quasi-expliziten Features abgeleitet, die Veränderungen innerhalb eines zugehörigen expliziten Feature-Sets beschreiben. Dazu werden jeweils $n - 1$ Features berechnet, welche die Differenzen *zwischen* den einzelnen Features eines gegebenen expliziten Sets aus n Features repräsentieren:

$$\mathcal{S}_{\Delta\phi} = \{\Delta\phi \mid \Delta\phi_i = \phi(Z_{i+1}) - \phi(Z_i), 0 \leq i < n - 1\} \quad (C.6)$$

$$\mathcal{S}_{\Delta r} = \{\Delta r \mid \Delta r_i = r(Z_{i+1}) - r(Z_i), 0 \leq i < n - 1\} \quad (C.7)$$

Aus den expliziten Feature-Sets $\mathcal{S}_{\mathfrak{R}'}$, \mathcal{S}_{\Im} , \mathcal{S}_{ϕ} , und \mathcal{S}_r sowie den quasi-expliziten Feature-Sets $\mathcal{S}_{\Delta\phi}$ und $\mathcal{S}_{\Delta r}$ wurden zusätzlich Sets impliziter globaler Features extrahiert. Dazu wurden für jedes explizite Feature-Set \mathcal{S} jeweils 16 univariate statistische Lagemaße bzw. Kennzahlen berechnet und als neues implizites

Feature-Set $G(S)$ zusammengefasst:

$$G(S) = \{min, max, R, P_{0,1}, Q_{0,25}, Q_{0,75}, P_{0,9}, R_{IQ}, R_{IP}, \bar{x}_{med}, \bar{x}_{arithm}, \bar{x}_{geom}, \bar{x}_{harm}, s^2, s, R_{MM}\} \quad (C.8)$$

Die genauen Berechnungsgrundlagen für $G(S)$ sind in Tab. C.1 exemplarisch für das Set $G(S_{\mathfrak{R}'})$ aufgeschlüsselt. Die Berechnung der weiteren impliziten Feature-Sets $G(S_{\mathfrak{S}})$, $G(S_{\phi})$, $G(S_r)$, $G(S_{\Delta\mathfrak{R}'})$, $G(S_{\Delta\mathfrak{S}})$, $G(S_{\Delta\phi})$, $G(S_{\Delta r})$ erfolgt analog dazu.

Daraus ergibt sich eine Menge I_e aus 250 expliziten und eine Menge I_i aus 96 impliziten Features:

$$I_e = S_{\mathfrak{R}'} \times S_{\mathfrak{S}} \times S_{\phi} \times S_r \times S_{\Delta\phi} \times S_{\Delta r} \quad (C.9)$$

$$I_i = G(S_{\mathfrak{R}'}) \times G(S_{\mathfrak{S}}) \times G(S_{\phi}) \times G(S_r) \times G(S_{\Delta\phi}) \times G(S_{\Delta r}) \quad (C.10)$$

Die Gesamtmenge I aller Input-Features wird schließlich als Kombination $I = I_e \times I_i$ aus allen expliziten und impliziten Feature-Sets definiert. Pro Impedanzspektrum steht dadurch eine Menge von insgesamt 346 Input-Features zur Verfügung.

C.2. Zielparameter

Als potentielle Zielparameter werden hier ein kategoriales und ein metrisches Feature der Datenbasis genutzt (vgl. Kapitel 7 und 8). Aufgrund der Modellierung ist dabei sichergestellt, dass ein tatsächlicher Zusammenhang zwischen den Impedanzspektren und den Output-Features besteht.

Als kategorialer Zielparameter wird die Kennung der zugehörigen Zelllinie L definiert:

$$L = \begin{cases} 0 & \text{falls "HT-29/B6"} \\ 1 & \text{falls "IPEC-J2"} \\ 2 & \text{falls "MDCK I"} \end{cases} \quad (C.11)$$

Als prozeduraler Zielparameter wird die epitheliale Kapazität C^{epi} verwendet:

$$C^{epi} = \frac{C^{ap} \cdot C^{bl}}{C^{ap} + C^{bl}} \quad (C.12)$$

Aufgrund der Zusammensetzung der Datenbasis aus Teildatensätzen (vgl. Abschnitt 6.5) treten die Klassen des kategorialen Zielparameters L nicht mit gleicher Häufigkeit auf. Die Klassen 0 und 1 haben eine Häufigkeit von $\frac{4}{11}$, die Klasse 2 von $\frac{3}{11}$. Die Verteilung für den metrischen Zielwert C^{epi} sind in Fig. C.1 dargestellt.

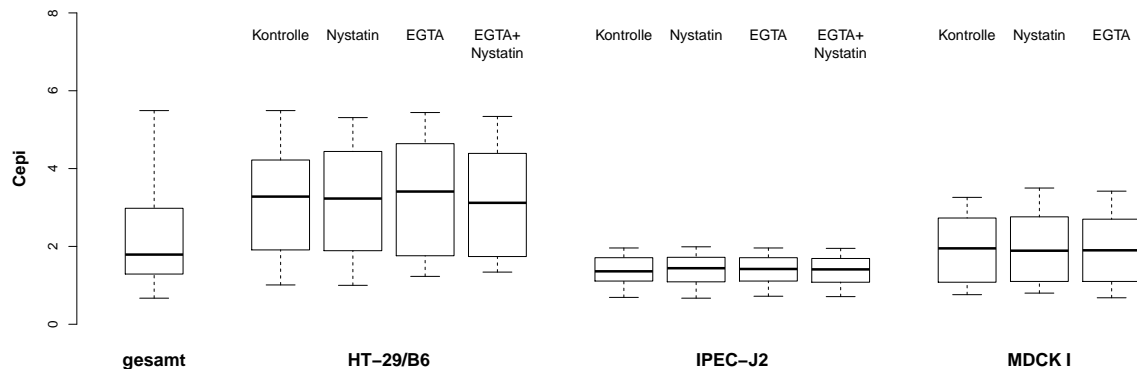


Abbildung C.1.: Verteilung der epithelialen Kapazität C^{epi} über die Datenbasis.

C.3. Metadaten

Um ein konstruktivistisches maschinelles Lernen mit einem Datensatz durchführen zu können, ist für jedes der Samples die Kenntnis der pragmatischen Eigenschaften T , Σ und Z erforderlich (vgl. Kapitel 5). Entsprechend sind diese als Metadaten Teil der Modellierung bzw. der zu analysierenden Datenbasis.

Die temporale Eigenschaft T wird dabei pro erzeugtem Impedanzspektrum nach festen Regeln hochgezählt. Initialwert ist stets der aktuelle Timestamp zum Beginn der jeweiligen Modellierung, also insbesondere der Modellierung einer spezifischen Zellkulturlinie und eines spezifischen funktionalen Zustands. Nach jedem erzeugten Impedanzspektrum wird T um 12 Sekunden erhöht, zusätzlich nach jeder potentiell zusammenhängenden Messreihe jeweils um $\epsilon = 120s$ bzw. $\delta = 1200s$ (siehe auch Alg. B.1 und B.2).

Über die gesamte Datenbasis bewegt sich T zwischen den Timestamps 1473865512 und 1522789715, was einem Modell-Zeitraum zwischen September 2016 und März 2018 entspricht (Abb. C.2). Die T -Werte der einzelnen Zellkulturlinien und -zustände erstrecken sich weitgehend ebenfalls über diesen Zeitraum, was einer Modellierung paralleler Messungen über einen längeren Zeitraum entspricht.

Da sämtliche Kurven synthetisiert sind, wird allen Spektren die gleiche Subjekteigenschaft Σ zugewiesen; die dafür verwendete ID lautet 0. Der zugehörige Zweck Z wird als unbekannt angenommen und ebenfalls einheitlich mit der ID 0 versehen.

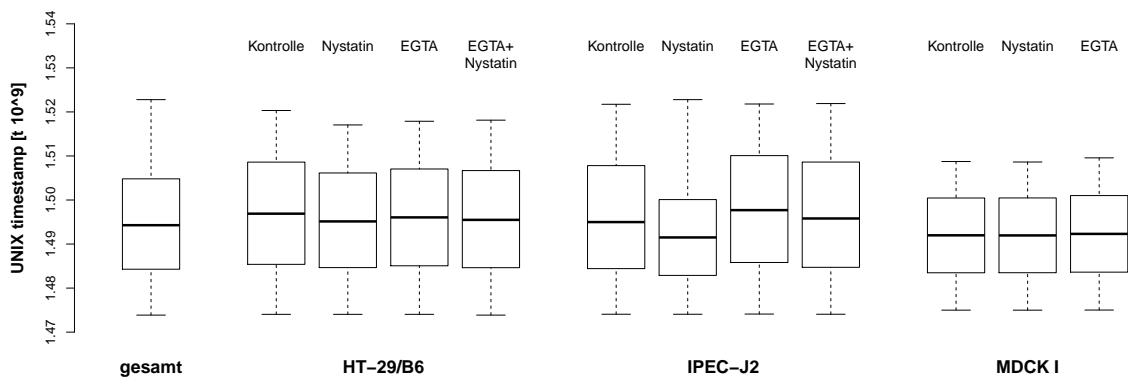


Abbildung C.2.: Verteilung des zeitlichen Metadatum T über die Datenbasis



D Charakterisierung des adaptierten Wissens

Im Folgenden werden die maschinellen Modelle beschrieben, die durch Adaption von konzeptuellem und prozeduralem Wissen erzeugt wurden. Ziel und Verlauf der Adaptionen sind in Kapitel 7 und 8 dargestellt.

D.1. Konzeptuelles Wissen

Zur Unterscheidung von Impedanzspektren hinsichtlich der zugrundeliegenden epithelialen Zelllinie L wurden insgesamt drei konzeptuelle Modelle L_1 , L_2 und L_3 erzeugt. Gemeinsam decken sie 30.000 Samples des Trainingsdatensatzes ($n=200.000$) und 11.138 Samples des Testdatensatzes ($n=75.000$) ab (Tab. D.1). Die einzelnen Modelle decken zwischen 4,1 und 7,2 Prozent der Zeitspanne der gesamten Datenbasis ab (Abb. D.2). 81,6 Prozent der Zeitspanne der gesamten Datenbasis werden nicht abgedeckt.

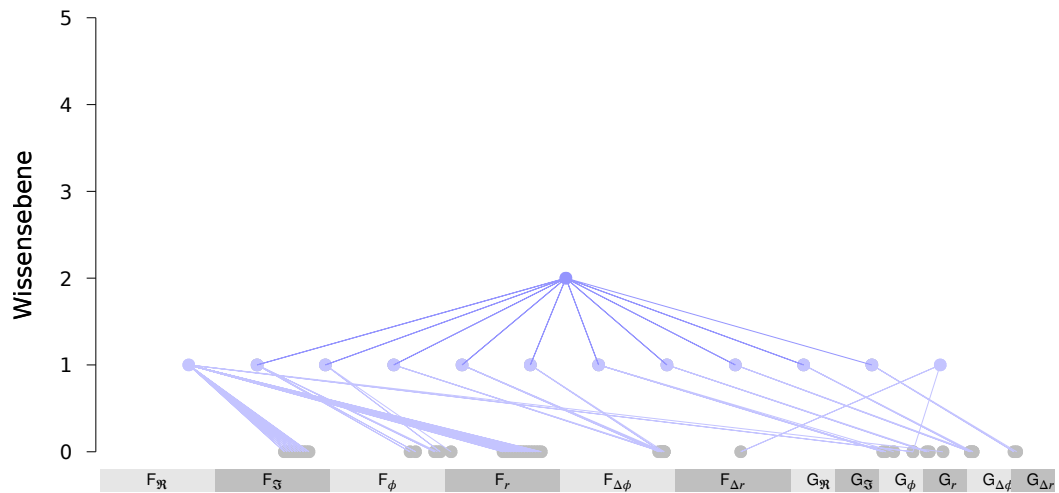
L_1 und L_3 sind Modelle der Ebene 1 und greifen somit unmittelbar auf Features der Datenbasis zurück. Modell L_2 befindet sich auf Ebene 2 und nutzt als Eingabe-Features die Ausgaben von zehn Modellen, die im Rahmen der Exploration auf der Ebene 1 integriert wurden. Abb. D.1a zeigt das Konnektom, Abb. D.1b eine dreidimensionale Aufsicht für L_1 , L_2 und L_3 . Modelle, die in der explorierten Domäne enthalten sind, aber nicht zur Adaption verwendet wurden, sind hier nicht dargestellt.

Tab. D.2 gibt einen Überblick über die verwendeten Features der Datenbasis. Zu unterscheiden ist dabei zwischen expliziten Features, die sich unmittelbar aus einer bzw. zwei Impedanzen ableiten, und impliziten Features, die sich jeweils auf eine Gruppe von Impedanzen beziehen (vgl. Anhang C). Die expliziten Features ergeben sich unmittelbar aus den ineinander überführbaren Darstellungen als kartesische oder Polarkoordinaten. Die Berechnungsgrundlage der impliziten Features kann Tab. C.1 entnommen werden. Die Zuordnung der Impedanz Z_i zur Frequenz f_i bzw. ω_i ist in Tab. B.1 abgebildet.

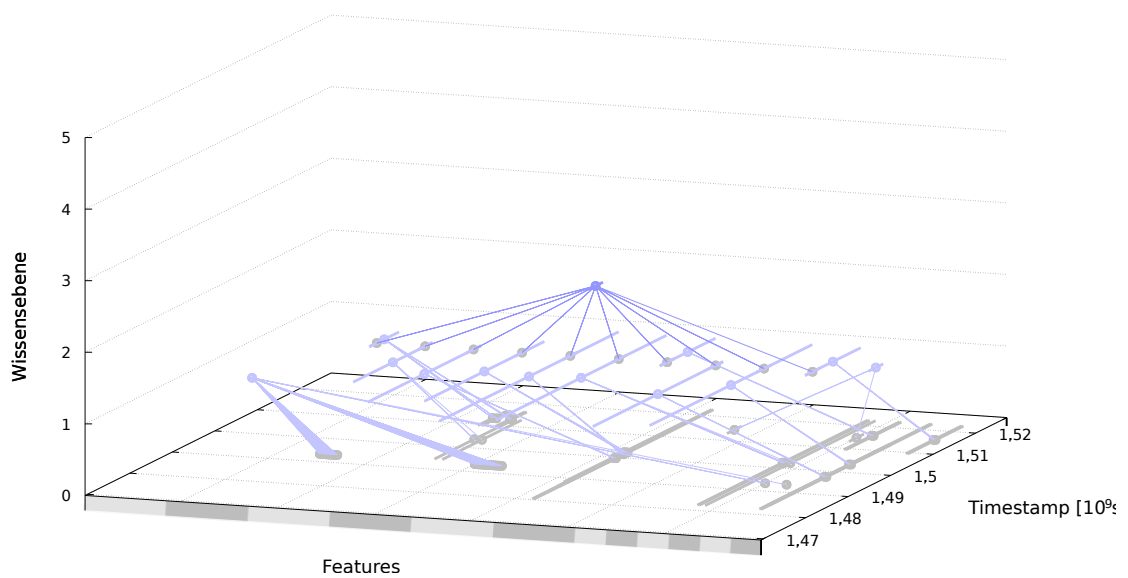
Modell	T	Σ	Z	Anzahl Features	Trainings-samples	Test-samples
L_1	[1492265315,1494268633]	{ANN, RF}	L	27	10.000	3.683
L_2	[1509480257,1512971023]	{ANN, RF}	L	18	10.000	3.745
L_3	[1512971514,1516483013]	{ANN, RF}	L	2	10.000	3.710

Tabelle D.1.: Überblick über die Eigenschaften der adaptierten konzeptuellen Modelle

D. Charakterisierung des adaptierten Wissens



a) Modellhierarchie



b) Dreidimensionale Aufsicht

Abbildung D.1.: Überblick über die adaptierte konzeptuelle Wissensdomäne. Modelle der gleichen Wissensebene sind jeweils im gleichen Farbton dargestellt. Die Feature-Sets der Trainingsdaten, die sich auf Ebene 0 befinden, sind entlang der x-Achse gekennzeichnet. a) Modellhierarchie. In dieser Frontal-Aufsicht werden Verknüpfungen zwischen den Domänenmodellen durch vertikale Geraden repräsentiert; die Zeitachse wird in dieser Darstellung vernachlässigt. b) Dreidimensionale Aufsicht. Zusätzlich zur Darstellung als Modellhierarchie ist die zeitliche Ausdehnung jedes Modells als Parallele zur z-Achse repräsentiert.

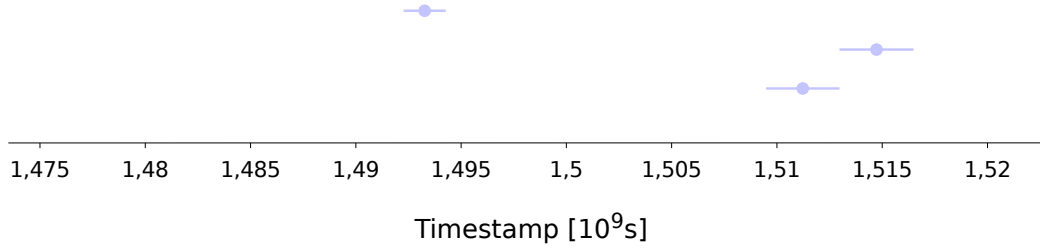


Abbildung D.2.: Temporale Gültigkeit des adaptierten konzeptuellen Wissens. Dargestellt ist die zeitliche Ausdehnung der beiden Modelle der Ebene 1 sowie des Modells der Ebene 2. Die Modelle sind dabei entsprechend ihrer Ebene bzw. ihrer ID in 1er-Schritten auf einer einheitslosen y-Achse in absteigender Reihenfolge angeordnet.

Feature	L_1	L_2	L_3
1	$\mathfrak{S}(Z_{27})$	$\phi(Z_1)$	$r(Z_{28}) - r(Z_{27})$
2	$\mathfrak{S}(Z_{28})$	$\phi(Z_3)$	$\max(S_r)$
3	$\mathfrak{S}(Z_{29})$	$\phi(Z_{40})$	
4	$\mathfrak{S}(Z_{30})$	$\phi(Z_{41})$	
5	$\mathfrak{S}(Z_{31})$	$r(Z_0)$	
6	$\mathfrak{S}(Z_{32})$	$r(Z_4)$	
7	$\mathfrak{S}(Z_{33})$	$\phi(Z_{39}) - \phi(Z_{38})$	
8	$\mathfrak{S}(Z_{34})$	$\phi(Z_{40}) - \phi(Z_{39})$	
9	$\mathfrak{S}(Z_{35})$	$\phi(Z_{41}) - \phi(Z_{40})$	
10	$\mathfrak{S}(Z_{36})$	$s(S_\phi)$	
11	$r(Z_{23})$	$s^2(S_\phi)$	
12	$r(Z_{24})$	$R(S_\phi)$	
13	$r(Z_{25})$	$s(S_r)$	
14	$r(Z_{26})$	$s^2(S_r)$	
15	$r(Z_{27})$	$s(S_{\Delta\phi})$	
16	$r(Z_{28})$	$s^2(S_{\Delta\phi})$	
17	$r(Z_{29})$	$s(S_{\Delta r})$	
18	$r(Z_{30})$	$s^2(S_{\Delta r})$	
19	$r(Z_{31})$		
20	$r(Z_{32})$		
21	$r(Z_{33})$		
22	$r(Z_{34})$		
23	$r(Z_{35})$		
24	$r(Z_{36})$		
25	$r(Z_{37})$		
26	$\max(S_r)$		
27	$\bar{x}_{harm}(S_r)$		

Tabelle D.2.: Features der adaptierten konzeptuellen Modelle. Für das auf Ebene 2 lokalisierte Modell L_2 sind statt den unmittelbar verwendeten Domänenmodellen der Ebene 1 diejenigen Features angegeben, die von diesen als Input-Features genutzt werden.

D.2. Prozedurales Wissen

Nach Abschluss der Adaption prozeduralen Wissens befinden sich acht maschinelle Modelle in der prozeduralen Wissensdomäne, die einen Zusammenhang zwischen der Datenbasis und der epithelialen Kapazität C^{epi} abbilden. Diese Modelle $C_1, C_2, C_3, C_4, C_5, C_6, C_7$ und C_8 sind Modelle der Ebene 1 und greifen somit unmittelbar auf Features der Datenbasis zurück. Modell der explorierten prozeduralen Domäne werden nicht verwendet.

Infolge von ΣZ -Dekonstruktionsprozessen wurden Modelle im Verlauf der Adaption temporal erweitert. Gemeinsam decken die adaptierten Modelle 168.257 Samples des Trainingsdatensatzes ($n=200.000$) und 63.246 Samples des Testdatensatzes ($n=75.000$) ab (Tab. D.3). Die einzelnen adaptierten Modelle decken zwischen 4,1 und 18,7 Prozent der Zeitspanne der gesamten Datenbasis ab (Abb. D.3). 17,2 Prozent der Zeitspanne der gesamten Datenbasis werden nicht abgedeckt.

Die adaptierten Modelle nutzen zwischen 16 und 41 Features (Tab. D.4). Zu unterscheiden ist dabei zwischen Features mit einem expliziten Bezug zu ein bzw. zwei Impedanzen und impliziten Features, die sich jeweils auf eine Gruppe von Impedanzen beziehen (vgl. Anhang C.1). Die expliziten Features ergeben sich unmittelbar aus den ineinander überführbaren Darstellungen in kartesischen und Polarkoordinaten, die Berechnungsgrundlage der impliziten Features kann Tab. C.1 entnommen werden. Die Zuordnung der Impedanz Z_i zur Frequenz f_i bzw. ω_i kann Tab. B.1 entnommen werden.

Modell	T	Σ	Z	Anzahl Features	Trainings-samples	Test-samples
C_1	[1473969552,1478259559]	{ANN, RF}	C^{epi}	22	19.757	7.298
C_2	[1480261571,1486132752]	{ANN, RF}	C^{epi}	16	29.700	11.119
C_3	[1486132783,1492081993]	{ANN, RF}	C^{epi}	17	29.700	11.260
C_4	[1478259751,1480261536]	{ANN, RF}	C^{epi}	18	9.900	3.613
C_5	[1498038829,1502260170]	{ANN, RF}	C^{epi}	18	19.800	7.616
C_6	[1502260177,1506646914]	{ANN, RF}	C^{epi}	22	19.800	7.316
C_7	[1506647917,1515811855]	{ANN, RF}	C^{epi}	17	29.700	11.207
C_8	[1515811915,1520460083]	{ANN, RF}	C^{epi}	41	9.900	3.817

Tabelle D.3.: Überblick über die Eigenschaften der adaptierten prozeduralen Modelle

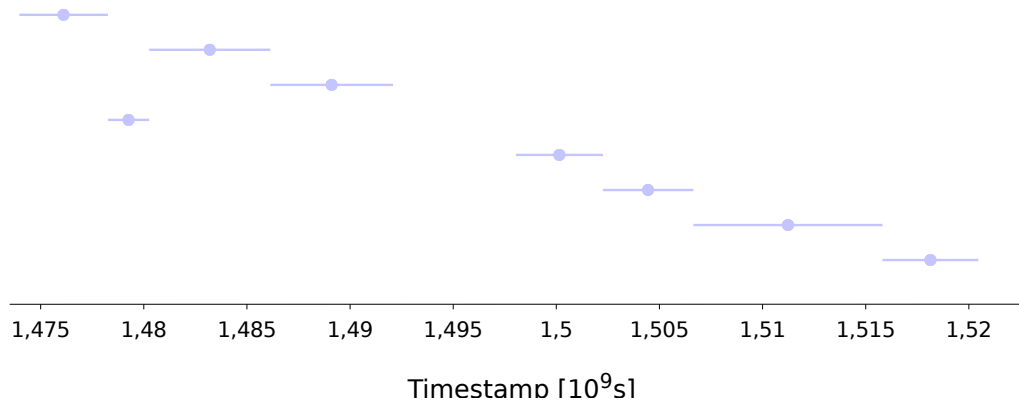
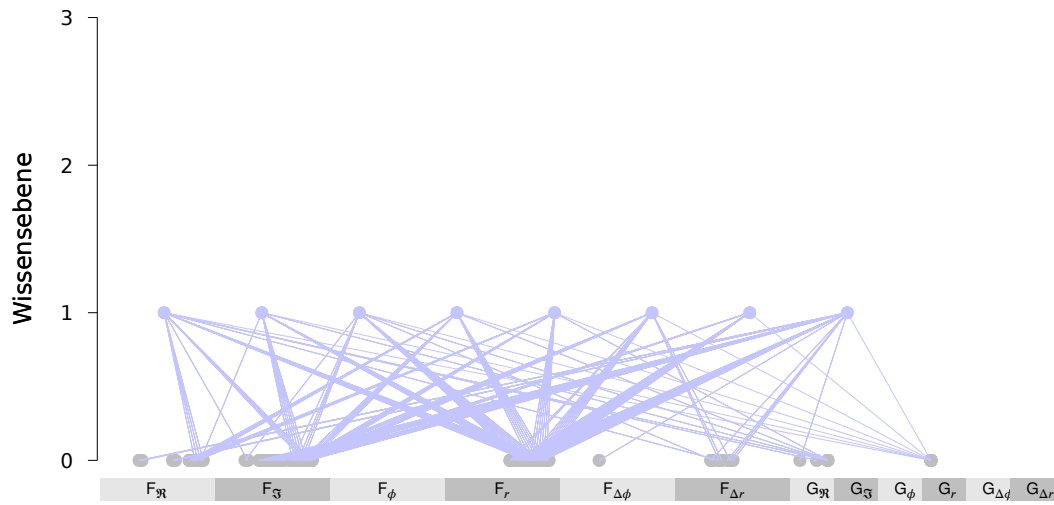
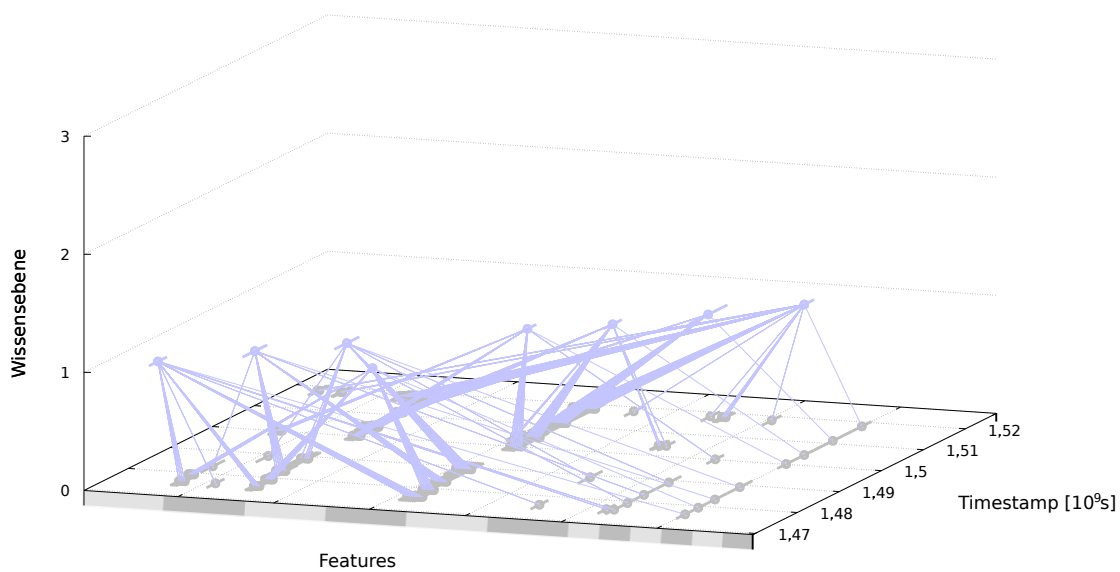


Abbildung D.3.: Temporale Gültigkeit des adaptierten prozeduralen Wissens. Dargestellt ist die zeitliche Ausdehnung der acht Modelle der Ebene 1. Die Modelle sind dabei entsprechend ihrer ID in 1er-Schritten auf einer einheitslosen y-Achse in absteigender Reihenfolge angeordnet.



a) Modellhierarchie



b) Dreidimensionale Aufsicht

Abbildung D.4.: Überblick über die adaptierte prozedurale Wissensdomäne. Modelle der gleichen Wissensebene sind jeweils im gleichen Farbton dargestellt. Die Feature-Sets der Trainingsdaten, die sich auf Ebene 0 befinden, sind entlang der x-Achse gekennzeichnet. a) Modellhierarchie. In dieser Frontal-Aufsicht werden Verknüpfungen zwischen den Domänenmodellen durch vertikale Geraden repräsentiert; die Zeitachse wird in dieser Darstellung vernachlässigt. b) Dreidimensionale Aufsicht. Zusätzlich zur Darstellung als Modellhierarchie ist die zeitliche Ausdehnung jedes Modells als Parallele zur z-Achse repräsentiert.

D. Charakterisierung des adaptierten Wissens

Feature	C_1	C_2	C_3	C_4	C_5	C_6	C_7	C_8
1	$\Re(Z_{35})$	$\Re(Z_{37})$	$\Im(Z_{11})$	$\Re(Z_{33})$	$\Re(Z_{33})$	$\Im(Z_{27})$	$\Im(Z_{29})$	$\Re(Z_{15})$
2	$\Re(Z_{36})$	$\Im(Z_{28})$	$\Im(Z_{28})$	$\Re(Z_{34})$	$\Re(Z_{34})$	$\Im(Z_{28})$	$\Im(Z_{30})$	$\Re(Z_{16})$
3	$\Re(Z_{37})$	$\Im(Z_{29})$	$\Im(Z_{30})$	$\Re(Z_{35})$	$\Re(Z_{35})$	$\Im(Z_{29})$	$\Im(Z_{31})$	$\Re(Z_{27})$
4	$\Re(Z_{38})$	$\Im(Z_{30})$	$\Im(Z_{31})$	$\Re(Z_{36})$	$\Re(Z_{36})$	$\Im(Z_{30})$	$\Im(Z_{32})$	$\Re(Z_{28})$
5	$\Im(Z_{12})$	$\Im(Z_{31})$	$\Im(Z_{32})$	$\Im(Z_{31})$	$\Im(Z_{31})$	$\Im(Z_{32})$	$\Im(Z_{33})$	$\Im(Z_{16})$
6	$\Im(Z_{32})$	$\Im(Z_{32})$	$r(Z_{25})$	$\Im(Z_{32})$	$\Im(Z_{32})$	$\Im(Z_{33})$	$r(Z_{24})$	$\Im(Z_{17})$
7	$\Im(Z_{33})$	$\Im(Z_{33})$	$r(Z_{26})$	$\Im(Z_{33})$	$\Im(Z_{33})$	$\Im(Z_{34})$	$r(Z_{25})$	$\Im(Z_{18})$
8	$\Im(Z_{34})$	$r(Z_{27})$	$r(Z_{27})$	$\Im(Z_{34})$	$\Im(Z_{34})$	$\Im(Z_{35})$	$r(Z_{26})$	$\Im(Z_{19})$
9	$\Im(Z_{35})$	$r(Z_{28})$	$r(Z_{28})$	$\Im(Z_{35})$	$\Im(Z_{35})$	$r(Z_{22})$	$r(Z_{27})$	$\Im(Z_{20})$
10	$r(Z_{28})$	$r(Z_{29})$	$r(Z_{29})$	$r(Z_{29})$	$r(Z_{29})$	$r(Z_{23})$	$r(Z_{28})$	$\Im(Z_{21})$
11	$r(Z_{29})$	$r(Z_{30})$	$r(Z_{30})$	$r(Z_{30})$	$r(Z_{30})$	$r(Z_{28})$	$r(Z_{29})$	$\Im(Z_{22})$
12	$r(Z_{30})$	$r(Z_{31})$	$r(Z_{31})$	$r(Z_{31})$	$r(Z_{31})$	$r(Z_{29})$	$r(Z_{30})$	$\Im(Z_{23})$
13	$r(Z_{31})$	$r(Z_{32})$	$r(Z_{32})$	$r(Z_{32})$	$r(Z_{32})$	$r(Z_{30})$	$r(Z_{31})$	$\Im(Z_{24})$
14	$r(Z_{32})$	$r(Z_{33})$	$r(Z_{33})$	$r(Z_{33})$	$r(Z_{33})$	$r(Z_{31})$	$r(Z_{32})$	$\Im(Z_{25})$
15	$r(Z_{33})$	$\bar{x}_{harm}(S_{\mathfrak{R}})$	$r(Z_{13}) - r(Z_{12})$	$r(Z_{34})$	$r(Z_{34})$	$r(Z_{32})$	$r(Z_{33})$	$\Im(Z_{28})$
16	$r(Z_{34})$	$max(S_r)$	$\bar{x}_{harm}(S_{\mathfrak{R}})$	$r(Z_{35})$	$r(Z_{35})$	$r(Z_{33})$	$r(Z_{34})$	$\Im(Z_{29})$
17	$r(Z_{35})$	$\bar{x}_{harm}(S_{\mathfrak{R}})$	$max(S_r)$	$\bar{x}_{harm}(S_{\mathfrak{R}})$	$\bar{x}_{harm}(S_{\mathfrak{R}})$	$r(Z_{34})$	$max(S_r)$	$\Im(Z_{30})$
18	$r(Z_{36})$			$max(S_r)$	$max(S_r)$	$r(Z_{35})$		$\Im(Z_{31})$
19	$r(Z_{15}) - r(Z_{14})$					$r(Z_{15}) - r(Z_{14})$		$\Im(Z_{32})$
20	$R_I P(S_{\mathfrak{R}}^*)$					$r(Z_{16}) - r(Z_{15})$		$\Im(Z_{33})$
21	$\bar{x}_{harm}(S_{\mathfrak{R}})$					$r(Z_{19}) - r(Z_{18})$		$\Im(Z_{34})$
22	$max(S_r)$					$max(S_r)$		$r(Z_{22})$
23								$r(Z_{23})$
24								$r(Z_{24})$
25								$r(Z_{25})$
26								$r(Z_{26})$
27								$r(Z_{27})$
28								$r(Z_{28})$
29								$r(Z_{29})$
30								$r(Z_{30})$
31								$r(Z_{31})$
32								$r(Z_{32})$
33								$r(Z_{33})$
34								$r(Z_{34})$
35								$\phi(Z_{13}) - \phi(Z_{12})$
36								$r(Z_{12}) - r(Z_{11})$
37								$r(Z_{18}) - r(Z_{17})$
38								$r(Z_{19}) - r(Z_{18})$
39								$r(Z_{20}) - r(Z_{19})$
40								$R(S_{\mathfrak{R}})$
41								$max(S_r)$

Tabelle D.4.: Features der adaptierten prozeduralen Modelle.

Literaturverzeichnis

- [1] E. Aarts et al. Simulated Annealing. In: E. K. Burke & G. Kendall (Hrsg.), *Search Methodologies. Introductory Tutorials in Optimization and Decision Support Techniques*, Kap. 7, S. 187–210. Springer, 2005.
- [2] Q. Al-Awqati. Terminal differentiation in epithelia: the role of integrins in hensin polymerization. *Annual Review of Physiology*, 73:401–412, 2011.
- [3] S. Alelyani et al. Feature selection for clustering: a review. In: C. C. Aggarwal & C. K. Reddy (Hrsg.), *Data Clustering: Algorithms and Applications*, S. 29–60. CRC Press, 2013.
- [4] D. G. Altman & P. Royston. What do we mean by validating a prognostic model? *Statistics in Medicine*, 19(4):453–473, 2000.
- [5] M. Amasheh et al. TNF α -induced and berberine-antagonized tight junction barrier impairment via tyrosine kinase, Akt and NF κ B signaling. *Journal of Cell Science*, 123(23):4145–4155, 2010.
- [6] S. Amasheh et al. Na⁺ absorption defends from paracellular back-leakage by claudin-8 upregulation. *Biochemical and Biophysical Research Communications*, 378(1):45–50, 2009.
- [7] G. M. Amdahl. Validity of the single processor approach to achieving large-scale computing capabilities. In: *AFIPS Spring Joint Computing Conference*, Vol. 30 of *AFIPS Conference Proceedings*, S. 483–485. ACM, 1967.
- [8] G. G. Anderson et al. Intracellular bacterial biofilm-like pods in urinary tract infections. *Science*, 301(5629):105–107, 2003.
- [9] L. W. Anderson & D. R. Krathwohl. *A taxonomy for learning, teaching and assessing: a revision of Bloom's taxonomy of educational objectives*. Longman, New York, 2001.
- [10] M. J. Apter. *Cybernetics and Development*. Pergamon Press, Oxford, 1966.
- [11] A. Arauzo-Azofra et al. Consistency measures for feature selection. *Journal of Intelligent Information Systems*, 30(3):273–292, 2007.
- [12] F. Baçao et al. Self-organizing maps as substitutes for k-means clustering. In: *Computational Science–ICCS 2005*, S. 476–483. Springer, 2005.
- [13] E. Bair & R. Tibshirani. Semi-supervised methods to predict patient survival from gene expression data. *PLoS Biol*, 2(4):e108, 2004.
- [14] G. Barker & N. L. Simmons. Identification of two strains of cultured canine renal epithelial cells (MDCK cells) which display entirely different physiological properties. *Quarterly Journal of Experimental Physiology*, 66(1):61–72, 1981.

Literaturverzeichnis

- [15] E. Bauer & R. Kohavi. An empirical comparison of voting classification algorithms: bagging, boosting, and variants. *Machine Learning*, 36(1):105–139, 1999.
- [16] D. C. Baumgart & S. R. Carding. Inflammatory bowel disease: cause and immunobiology. *The Lancet*, 369(9573):1627–1640, 2007.
- [17] M. F. Bear et al. (Hrsg.). *Neuroscience. Exploring the Brain*, Kap. 1 (Neuroscience: past, present, and future), S. 3–22. Lippincott Williams & Wilkins, Philadelphia, 3. Aufl., 2007.
- [18] D. A. Bell & H. Wang. A formalism for relevance and its application in feature subset selection. *Machine Learning*, 41(2):175–195, 2000.
- [19] Y. Bengio et al. Representation learning: a review and new perspectives. *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, 35(8):1798–1828, 2013.
- [20] E. M. Bennett et al. Communications through limited-response questioning. *Public Opinion Quarterly*, 18(3):303–308, 1954.
- [21] H. Berschneider. Development of normal cultured small intestinal epithelial cell lines which transport Na and Cl. *Gastroenterology*, 96(Suppl. Pt 2):A41, 1989.
- [22] C. A. Bertrand et al. System for dynamic measurements of membrane capacitance in intact epithelial monolayers. *Biophysical Journal*, 75(6):2743–2756, 1998.
- [23] V. Best et al. The role of high frequencies in speech localization. *Journal of the Acoustical Society of America*, 118(1):353–363, 2005.
- [24] W. Bialek et al. Reading a neural code. *Science*, 252(5014):1854–1857, 1991.
- [25] G. Biau. Analysis of a random forests model. *The Journal of Machine Learning Research*, 13(1):1063–1095, 2012.
- [26] C. M. Bishop. *Pattern Recognition and Machine Learning*, Kap. 1 (Introduction), S. 1–66. Springer, 2006.
- [27] B. S. Bloom. *Taxonomy of educational objectives. Vol. 1: cognitive domain*. McKay, New York, 1956.
- [28] A. P. Bradley. The use of the area under the ROC curve in the evaluation of machine learning algorithms. *Pattern Recognition*, 30(7):1145–1159, 1997.
- [29] L. Breiman. Bagging predictors. *Machine Learning*, 24(2):123–140, 1996.
- [30] L. Breiman. Random forests. *Machine learning*, 45(1):5–32, 2001.
- [31] R. Brette et al. Simulation of networks of spiking neurons: a review of tools and strategies. *Journal of Computational Neuroscience*, 23(3):349–398, 2007.
- [32] A. J. Brosnahan & D. R. Brown. Porcine IPEC-J2 intestinal epithelial cells in microbiological investigations. *Veterinary Microbiology*, 156(3):229–237, 2012.
- [33] P. Bühlmann. Bagging, subbagging and bragging for improving some prediction algorithms. In: M. G. Akritas & D. N. Politis (Hrsg.), *Recent Advances and Trends in Nonparametric Statistics*, S. 19–34. Elsevier, Amsterdam, 2003.
- [34] G. Carpenter et al. The ART of adaptive pattern recognition by a self-organizing neural network. *Computer*, 21(3):77–88, 1988.

- [35] R. G. Carraher & J. B. Thurston. *Optical Illusions and the Visual Arts*. Reinhold, New York, 1966.
- [36] M. A. Carreira-Perpinan & Z. Lu. Dimensionality reduction by unsupervised regression. *IEEE Conference on Computer Vision and Pattern Recognition*, 2008.
- [37] R. Caruana et al. Overfitting in neural nets: backpropagation, conjugate gradient, and early stopping. In: T. K. Leen et al. (Hrsg.), *Advances in Neural Information Processing Systems*, Vol. 13, S. 402–408. MIT Press, 2001.
- [38] A. V. Chadwick. Electrical conductivity measurements of ionic solids. *Philosophical Magazine A*, 64(5):983–998, 1991.
- [39] M. Chavent et al. ClustOfVar: an R package for the clustering of variables. *Journal of Statistical Software*, 50(13), 2012.
- [40] B. Cheng & D. Titterton. Neural networks: a review from a statistical perspective. *Statistical Science*, 9(1):2–30, 1994.
- [41] K. J. Cherkauer. Human expert-level performance on a scientific image analysis task by a system using combined artificial neural networks. In: *Proceedings of the Thirteenth National Conference on Artificial Intelligence*, S. 15–21. 1996.
- [42] C. Clausen et al. Impedance analysis of a tight epithelium using a distributed resistance model. *Biophysical Journal*, 26(2):291–317, 1979.
- [43] J. Cohen et al. A coefficient of agreement for nominal scales. *Educational and Psychological Measurement*, 20(1):37–46, 1960.
- [44] K. S. Cole. *Membranes, Ions, and Impulses*, Kap. Membrane capacity, S. 12. University of California Press, Berkley, 1972.
- [45] K. S. Cole & R. H. Cole. Dispersion and absorption in dielectrics. I. alternating current characteristics. *Journal of Chemical Physics*, 9:341–351, 1941.
- [46] C. Cortes & V. Vapnik. Support-vector networks. *Machine Learning*, 20(3):273–297, 1995.
- [47] T. Cover & P. Hart. Nearest neighbor pattern classification. *IEEE Transactions on Information Theory*, 13(1):21–27, 1967.
- [48] K. Craik. *The Nature of Explanation*. Cambridge University Press, 1943.
- [49] N. Cristianini & J. Shawe-Taylor. *An introduction to support vector machines and other kernel-based learning methods*, Kap. 6, S. 93–124. Cambridge University Press, 2000.
- [50] L. J. Cronbach. Coefficient alpha and the internal structure of tests. *Psychometrika*, 16(3):297–334, 1951.
- [51] Y. Dan. Spike timing-dependent plasticity: from synapse to perception. *Physiological Reviews*, 86(3):1033–1048, 2006.
- [52] R. D. De Veaux et al. A comparison of two nonparametric estimation schemes: MARS and neural networks. *Computers & Chemical Engineering*, 17(8):819–837, 1993.

- [53] J. M. Diamond. The mechanism of water transport by the gall-bladder. *Journal of Physiology*, 161:503–527, 1962.
- [54] J. M. Dietschy. Water and solute movement across the wall of the everted rabbit gall bladder. *Gastroenterology*, 47(4):395–408, 1964.
- [55] T. G. Dietterich. An experimental comparison of three methods for constructing ensembles of decision trees: bagging, boosting, and randomization. *Machine Learning*, 40(2):139–157, 2000.
- [56] A. C. Durham. A survey of readily available chelators for buffering calcium ion concentrations in physiological solutions. *Cell Calcium*, 4(1):33–46, 1983.
- [57] S. R. Eddy. What is a hidden Markov model? *Nature Biotechnology*, 22(10):1315–1316, 2004.
- [58] R. Eisenberg. Impedance measurement of the electrical structure of skeletal muscle. In: J. Field (Hrsg.), *Handbook of Physiology*, Kap. 11, S. 301–323. 1983.
- [59] K. El-Laithy. *Towards a Brain-inspired Information Processing System: Modelling and Analysis of Synaptic Dynamics*. Dissertation, Universität Leipzig, 2011.
- [60] L. E. Ericson & M. Nilsson. Structural and functional aspects of the thyroid follicular epithelium. *Toxicology Letters*, 64/65:365–373, 1992.
- [61] S. E. Fahlman & C. Lebiere. The cascade-correlation learning architecture. In: D. S. Touretzky (Hrsg.), *Advances in Neural Information Processing Systems*, Vol. 2, S. 524–532. Morgan Kaufmann, 1990.
- [62] R. Fergus et al. Semi-supervised learning in gigantic image collections. In: *Advances in Neural Information Processing Systems*, S. 522–530. 2009.
- [63] J. L. Fleiss. Measuring nominal scale agreement among many raters. *Psychological Bulletin*, 76(5):378–382, 1971.
- [64] M. J. Flores & J. A. Gámez. Breeding value classification in Manchego sheep: a Study of attribute selection and construction. In: R. Khosla et al. (Hrsg.), *Knowledge-Based Intelligent Information and Engineering Systems*, Vol. 3682 of *Lecture Notes in Computer Science*, S. 1338–1346. Springer, 2005.
- [65] M. J. Flynn. Some computer organizations and their effectiveness. *IEEE Transactions on Computers*, C-21(9):948–960, 1972.
- [66] R. Fox. Constructivism examined. *Oxford Review of Education*, 27(1):23–35, 2001.
- [67] M. J. Frank et al. By carrot or by stick: cognitive reinforcement learning in Parkinsonism. *Science*, 306(5703):1940–1943, 2004.
- [68] W. J. Frawley et al. Knowledge discovery in databases: an overview. *AI Magazine*, 13(3):57, 1992.
- [69] E. Fuchs. Scratching the surface of skin development. *Nature*, 445(7130):834–842, 2007.
- [70] S. D. Fuller & K. Simons. Transferrin receptor polarity and recycling accuracy in “tight” and “leaky” strains of Madin-Darby canine kidney cells. *The Journal of Cell Biology*, 103(5):1767–1779, 1986.

- [71] E. J. Furst. Bloom's taxonomy of educational objectives for the cognitive domain: philosophical and educational issues. *Review of Educational Research*, 51(4):441–453, 1981.
- [72] M. Furuse et al. Conversion of zonulae occludentes from tight to leaky strand type by introducing claudin-2 into Madin-Darby canine kidney I cells. *Journal of Cell Biology*, 153(2):263–272, 2001.
- [73] R. M. Gagné. Domains of learning. *Interchange*, 3(1):1–8, 1972.
- [74] M. Galar et al. A review on ensembles for the class imbalance problem: bagging-, boosting-, and hybrid-based approaches. *IEEE Transactions on Systems, Man, and Cybernetics, Part C*, 42(4):463–484, 2012.
- [75] C. R. Gaush et al. Characterization of an established line of canine kidney cells (MDCK). *Experimental Biology and Medicine*, 122(3):931–935, 1966.
- [76] S. Geisser. The predictive sample reuse method with applications. *Journal of the American Statistical Association*, 70(350):320–328, 1975.
- [77] A. H. Gitter et al. Impedance analysis for the determination of epithelial and subepithelial resistance in intestinal tissues. *Journal of Biochemical and Biophysical Methods*, 37(1-2):35–46, 1998.
- [78] A. H. Gitter et al. Trans/paracellular, surface/crypt, and epithelial/subepithelial resistances of mammalian colonic epithelia. *Pflügers Archiv*, 439(4):477–482, 2000.
- [79] A. B. Goldberg & X. Zhu. Seeing stars when there aren't many stars: graph-based semi-supervised learning for sentiment categorization. In: *Proceedings of the First Workshop on Graph Based Methods for Natural Language Processing*, S. 45–52. Association for Computational Linguistics, 2006.
- [80] D. A. Goodenough & J. P. Revel. A fine structural analysis of intercellular junctions in the mouse liver. *Journal of Cell Biology*, 45(2):272–290, 1970.
- [81] L. A. Goodman & W. H. Kruskal. Measures of association for cross classifications. *Journal of the American Statistical Association*, 49(268):732–764, 1954.
- [82] J. Griggs. One Minute with... Henry Markram. *New Scientist*, 217(2903):29, 2013.
- [83] K. R. Groschwitz & S. P. Hogan. Intestinal barrier function: molecular regulation and disease pathogenesis. *Journal of Allergy and Clinical Immunology*, 124(1):3–20, 2009.
- [84] D. Günzel et al. From TER to trans- and paracellular resistance: lessons from impedance spectroscopy. *Annals of the New York Academy of Sciences*, 1257:142–151, 2012.
- [85] I. Guyon & A. Elisseeff. An introduction to variable and feature selection. *Journal of Machine Learning Research*, 3:1157–1182, 2003.
- [86] I. Guyon et al. Gene selection for cancer classification using support vector machines. *Machine Learning*, 46(1/3):389–422, 2002.
- [87] M. A. Hall. *Correlation-based Feature Selection for Machine Learning*. Dissertation, University of Waikato, Hamilton, New Zealand, 1999.
- [88] A. G. Hamilton. *Numbers, Sets and Axioms: The Apparatus of Mathematics*. Cambridge University Press, 1982.

Literaturverzeichnis

- [89] J. S. Handler. Overview of epithelial polarity. *Annual Review of Physiology*, 51:729–740, 1989.
- [90] L. K. Hansen & P. Salamon. Neural network ensembles. *IEEE Transactions on Pattern Analysis & Machine Intelligence*, (10):993–1001, 1990.
- [91] J. R. Harkema et al. The nose revisited: a brief review of the comparative structure, function, and toxicologic pathology of the nasal epithelium. *Toxicologic Pathology*, 34(3):252–269, 2006.
- [92] J. Hartley. Programmed instruction 1954-1974 – a review. *Innovations in Education & Training International*, 11(6):278–291, 1974.
- [93] A. F. Hayes & K. Krippendorff. Answering the call for a standard reliability measure for coding data. *Communication Methods and Measures*, 1(1):77–89, 2007.
- [94] D. O. Hebb (Hrsg.). *The Organization of Behavior: A Neuropsychological Theory*. Wiley, New York, 1949.
- [95] B. R. Hergenhahn. *An Introduction to the History of Psychology*. Wadsworth, 4. Aufl., 2013.
- [96] B. R. Hergenhahn & M. H. Olson (Hrsg.). *An Introduction to Theories of Learning*, Kap. 3 (Early notions about learning), S. 51–52. Prentice-Hall, New Jersey, 4. Aufl., 1993.
- [97] H. Hertz. Die Prinzipien der Mechanik in neuem Zusammenhange dargestellt. In: H. Hertz (Hrsg.), *Gesammelte Werke*, Vol. 3. Barth, Leipzig, 1894.
- [98] G. E. Hinton & R. R. Salakhutdinov. Reducing the dimensionality of data with neural networks. *Science*, 313(5786):504–507, 2006.
- [99] Y. Hirose et al. A back-propagation algorithm which varies the number of hidden units. *Neural Networks*, 4(1):61–66, 1991.
- [100] D. Hofmann et al. Learning interpretable kernelized prototype-based models. *Neurocomputing*, 141:84–96, 2014.
- [101] O. R. Holsti. *Content Analysis for the Social Sciences and Humanities*. Addison-Wesley, Reading, MA, 1969.
- [102] R. W. Holz. Polyene antibiotics: nystatin, amphotericin B, and filipin. In: F. E. Hahn (Hrsg.), *Mechanism of Action of Antieukaryotic and Antiviral Compounds*, S. 313–340. Springer, 1979.
- [103] J. J. Hopfield. Neural networks and physical systems with emergent collective computational abilities. *Proceedings of the National Academy of Sciences*, 79(8):2554–2558, 1982.
- [104] K. Hornik et al. Multilayer feedforward networks are universal approximators. *Neural networks*, 2(5):359–366, 1989.
- [105] M. A. Hughes & D. E. Garrett. Intercoder reliability estimation approaches in marketing: a generalizability theory framework for quantitative data. *Journal of Marketing Research*, S. 185–195, 1990.
- [106] P. Hunter. Simulating the human brain. *EMBO Reports*, 16(6):685–688, 2015.

- [107] C. Igel & M. Hüsken. Empirical evaluation of the improved Rprop learning algorithms. *Neurocomputing*, 50:105–123, 2003.
- [108] C. M. V. Itallie & J. M. Anderson. Claudins and epithelial paracellular transport. *Annual Review of Physiology*, 68:403–429, 2006.
- [109] I. H. Iversen. Skinner’s early research: from reflexology to operant conditioning. *American Psychologist*, 47(11):1318–1328, 1992.
- [110] E. M. Izhikevich. Which model to use for cortical spiking neurons? *IEEE Transactions on Neural Networks*, 15(5):1063–1070, 2004.
- [111] R. A. Jacobs et al. Adaptive mixtures of local experts. *Neural Computation*, 3(1):79–87, 1991.
- [112] R. J. Jagacinski & R. A. Miller. Describing the human operator’s internal model of a dynamic system. *Human Factors: The Journal of the Human Factors and Ergonomics Society*, 20(4):425–433, 1978.
- [113] A. K. Jain. Data clustering: 50 years beyond k-means. *Pattern Recognition Letters*, 31(8):651–666, 2010.
- [114] T. J. Jentsch et al. Molecular structure and physiological function of chloride channels. *Physiological Reviews*, 82(2):503–568, 2002.
- [115] X. Jiang & A. H. K. S. Wah. Constructing and training feed-forward neural networks for pattern classification. *Pattern Recognition*, 36(4):853–867, 2003.
- [116] H. F. Kaiser. The application of electronic computers to factor analysis. *Educational and Psychological Measurement*, 20(1):141–151, 1960.
- [117] I. Kasa. A circle fitting procedure and its error analysis. *IEEE Transactions on Instrumentation and Measurement*, 1001(1):8–14, 1976.
- [118] T. Kato & R. L. Owen. Structure and function of intestinal mucosal epithelium. In: J. Mestecky et al. (Hrsg.), *Mucosal Immunology*, Kap. 8, S. 131–151. Academic Press, San Diego, 3. Aufl., 2005.
- [119] M. Kawato & K. Samejima. Efficient reinforcement learning: computational theories, neuroscience and robotics. *Current Opinion in Neurobiology*, 17(2):205–212, 2007.
- [120] B. S. Kim et al. A method for analyzing electrical impedance spectroscopy data from breast cancer patients. *Physiological Measurement*, 28(7):S237, 2007.
- [121] H.-C. Kim et al. Constructing support vector machine ensemble. *Pattern Recognition*, 36(12):2757–2767, 2003.
- [122] J. Kittler et al. On combining classifiers. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 20(3):226–239, 1998.
- [123] G. Klaus. *Wörterbuch der Kybernetik*. Dietz, Berlin, 1967.
- [124] J. Kober & J. Peters. Reinforcement learning in robotics: a survey. In: M. Wiering & M. van Otterlo (Hrsg.), *Reinforcement Learning*, Kap. 18, S. 579–610. Springer, 2012.

- [125] V. Koefoed-Johnsen & H. H. Ussing. The nature of the frog skin potential. *Acta Physiologica Scandinavica*, 42(3-4):298–308, 1958.
- [126] R. Kohavi & G. H. John. Wrappers for feature subset selection. *Artificial Intelligence*, 97(1-2):273–324, 1997.
- [127] R. Kohavi & F. Provost. Glossary of terms. *Machine Learning*, 30(2-3):271–274, 1998.
- [128] T. Kohonen. The self-organizing map. *Proceedings of the IEEE*, 78(9):1464–1480, 1990.
- [129] T. Kohonen. Physiological interpretation of the self-organizing map algorithm. *Neural Networks*, 6(7):895–905, 1993.
- [130] T. Kohonen. *Self-organizing maps*, Vol. 30 of *Springer Series in Information Sciences*. Springer, 3. Aufl., 2001.
- [131] R. H. Kolbe & M. S. Burnett. Content-analysis research: an examination of applications with directives for improving research reliability and objectivity. *Journal of Consumer Research*, 18:243–250, 1991.
- [132] J. Kornmeier & M. Bach. Ambiguous figures – what happens in the brain when perception changes but not the stimulus. *Frontiers in Human Neuroscience*, 6, 2012.
- [133] S. Kotsiantis & P. Pintelas. Recent advances in clustering: a brief survey. *WSEAS Transactions on Information Science and Applications*, 1(1):73–81, 2004.
- [134] O. Kramer. Dimensionality reduction by unsupervised k-nearest neighbor regression. *10th International Conference on Machine Learning and Applications and Workshops*, 2011.
- [135] D. R. Krathwohl. A revision of Bloom’s taxonomy: an overview. *Theory into Practice*, 41(4):212–218, 2002.
- [136] D. R. Krathwohl et al. *Taxonomy of Educational Objectives, Handbook II, Affective Domain*. David McKay Co, 1964.
- [137] K.-M. Kreusel et al. Cl- secretion in epithelial monolayers of mucus-forming human colon cells (HT-29/B6). *American Journal of Physiology-Cell Physiology*, 261(4):C574–C582, 1991.
- [138] K. Krippendorff. Bivariate agreement coefficients for reliability of data. In: Bortatta (Hrsg.), *Sociological Methodology*, S. 139–150. San Francisco, 1970.
- [139] K. Krippendorff. Reliability in content analysis: some common misconceptions and recommendations. *Human Communication Research*, 30(3):411–433, 2004.
- [140] J. Kriz et al. *Wissenschafts- und Erkenntnistheorie. Eine Einführung für Psychologen und Humanwissenschaftler*. VS Verlag für Sozialwissenschaften, 1987.
- [141] A. Krogh. What are artificial neural networks? *Nature Biotechnology*, 26(2):195–197, 2008.
- [142] A. Krogh et al. Neural network ensembles, cross validation, and active learning. *Advances in Neural Information Processing Systems*, 7:231–238, 1995.
- [143] S. M. Krug et al. Two-path impedance spectroscopy for measuring paracellular and transcellular epithelial resistance. *Biophysical Journal*, 97(8):2202–2211, 2009.

- [144] L. I. Kuncheva & C. J. Whitaker. Measures of diversity in classifier ensembles and their relationship with the ensemble accuracy. *Machine Learning*, 51(2):181–207, 2003.
- [145] S. W. Kwok & C. Carter. Multiple decision trees. In: *Proceedings of the Fourth Annual Conference on Uncertainty in Artificial Intelligence*, UAI '88, S. 327–338. North-Holland Publishing, Amsterdam, 1990.
- [146] U. G. Kyle et al. Bioelectrical impedance analysis – part I: review of principles and methods. *Clinical Nutrition*, 23(5):1226–1243, 2004.
- [147] U. G. Kyle et al. Bioelectrical impedance analysis – part II: utilization in clinical practice. *Clinical Nutrition*, 23(6):1430–1453, 2004.
- [148] P. Langley. Editorial: on machine learning. *Machine Learning*, 1(1):5–10, 1986.
- [149] Q. V. Le. Building high-level features using large scale unsupervised learning. In: *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, S. 8595–8598. IEEE, 2013.
- [150] C.-C. J. Lee et al. A kinetic model unifying presynaptic short-term facilitation and depression. *Journal of Computational Neuroscience*, 26(3):459–473, 2009.
- [151] J. Lee. Richard Wesley Hamming: 1915-1998. *IEEE Annals of the History of Computing*, 20(2):60–62, 1998.
- [152] J. Lee et al. Unsupervised dimensionality reduction: overview and recent advances. In: *The 2010 International Joint Conference on Neural Networks (IJCNN)*, S. 1–8. IEEE, 2010.
- [153] D. A. Leopold & N. K. Logothetis. Multistable phenomena: changing views in perception. *Trends in Cognitive Sciences*, 3(7):254–264, 1999.
- [154] R. E. Levien & M. Maron. A computer system for inference execution and data retrieval. *Communications of the ACM*, 10(11):715–721, 1967.
- [155] D. D. Lewis. Naive (Bayes) at forty: the independence assumption in information retrieval. In: C. Nédellec & C. Rouveirol (Hrsg.), *Machine Learning: ECML-98*, Vol. 1398 of *Lecture Notes in Computer Science*, S. 4–15. Springer, 1998.
- [156] F. Lewis et al. Reinforcement learning and feedback control: using natural decision methods to design optimal adaptive controllers. *IEEE Control Systems*, 32(6):76–105, 2012.
- [157] H. Li et al. Transepithelial electrical measurements with the Ussing chamber. *Journal of Cystic Fibrosis*, 3 Suppl 2:123–126, 2004.
- [158] T.-S. Lim et al. A comparison of prediction accuracy, complexity, and training time of thirty-three old and new classification algorithms. *Machine Learning*, 40(3):203–228, 2000.
- [159] W. P. Lincoln & J. Skrzypek. Synergy of clustering multiple back propagation networks. In: D. S. Touretzky (Hrsg.), *Advances in Neural Information Processing Systems*, Vol. 2, S. 650–657. Morgan Kaufmann, 1990.
- [160] M. Lombard et al. Content analysis in mass communication. *Human Communication Research*, 28(4):587–604, 2002.

Literaturverzeichnis

- [161] S.-Y. Lu & K. S. Fu. A sentence-to-sentence clustering procedure for pattern analysis. *IEEE Transactions on Systems, Man, and Cybernetics*, 8(5):381–389, 1978.
- [162] N. Luhmann. Sinn als Grundbegriff der Soziologie. In: J. Habermas & N. Luhmann (Hrsg.), *Theorie der Gesellschaft oder Sozialtechnologie*. Suhrkamp, Frankfurt a. M., 1971.
- [163] W. Maass. Networks of spiking neurons: the third generation of neural network models. *Neural Networks*, 10(9):1659–1671, 1997.
- [164] W. Maass & A. M. Zador. Dynamic stochastic synapses as computational units. *Neural Computation*, 11(4):903–917, 1999.
- [165] J. R. Macdonald. Impedance spectroscopy. *Annals of Biomedical Engineering*, 20(3):289–305, 1992.
- [166] H. Markram et al. Differential signaling via the same axon of neocortical pyramidal neurons. *Proceedings of the National Academy of Sciences*, 95(9):5323–5328, 1998.
- [167] W. S. McCulloch & W. Pitts. A logical calculus of the ideas immanent in nervous activity. *Bulletin of Mathematical Biology*, 5(4):115–113, 1943.
- [168] J. A. McRoberts et al. The Madin-Darby canine kidney (MDCK) cell line. In: G. Sato (Hrsg.), *Functionally Differentiated Cell Lines*, S. 117–139. Alan Liss, New York, 1981.
- [169] P. A. Merolla et al. A million spiking-neuron integrated circuit with a scalable communication network and interface. *Science*, 345(6197):668–673, 2014.
- [170] T. M. Mitchell. The discipline of machine learning. Report, Carnegie Mellon University, School of Computer Science, Machine Learning Department, 2006.
- [171] K. Mohraz & P. Protzel. FlexNet - a flexible neural network construction algorithm. In: M. Verleysen (Hrsg.), *Proceedings of the 4th European Symposium on Artificial Neural Networks (ESANN)*, S. 111–116. 1996.
- [172] P. R. Montague et al. A framework for mesencephalic dopamine systems based on predictive Hebbian learning. *The Journal of Neuroscience*, 16(5):1936–1947, 1996.
- [173] R. Müller. Zur Geschichte des Modelldenkens und des Modellbegriffs. In: H. Stachowiak (Hrsg.), *Modelle - Konstruktion der Wirklichkeit*, S. 17–86. Fink, München, 1983.
- [174] F. Murtagh. A survey of recent advances in hierarchical clustering algorithms. *The Computer Journal*, 26(4):354–359, 1983.
- [175] F. Murtagh & P. Contreras. Algorithms for hierarchical clustering: an overview. *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery*, 2(1):86–97, 2012.
- [176] Y. Nakazato et al. Characterization of subclones of Madin-Darby canine kidney renal epithelial cell line. *Biochimica et Biophysica Acta (BBA) - Molecular Cell Research*, 1014(1):57–65, 1989.
- [177] A. Navot et al. Nearest neighbor based feature selection for regression and its application to neural activity. In: *NIPS 2005*, Vol. 18 of *Advances in Neural Information Processing Systems*. 2005.

- [178] L. A. Necker. Observations on some remarkable optical phænomena seen in Switzerland; and on an optical phænomenon which occurs on viewing a figure of a crystal or geometrical solid. *Philosophical Magazine*, 1(5):329–337, 1832.
- [179] C. Nossol et al. Air–liquid interface cultures enhance the oxygen supply and trigger the structural and functional differentiation of intestinal porcine epithelial cells (IPEC). *Histochemistry and Cell Biology*, 136(1):103–115, 2011.
- [180] K. Oatley. Representations of the physical and social world. In: D. A. Oatley (Hrsg.), *Brain and Mind*, S. 32–58. Methuen, London, 1985.
- [181] S. Openshaw & C. Openshaw. *Artificial Intelligence in Geography*. John Wiley & Sons, Inc., 1997.
- [182] D. Opitz & R. Maclin. Popular ensemble methods: an empirical study. *Journal of Artificial Intelligence Research*, 11:169–198, 1999.
- [183] M. E. Orazem & B. Tribollet. Electrical circuits. In: *Electrochemical Impedance Spectroscopy*, Kap. 4, S. 61–72. Wiley, 2008.
- [184] M. E. Orazem & B. Tribollet. *Electrochemical Impedance Spectroscopy*. Wiley, 2008.
- [185] M. E. Orazem & B. Tribollet. Error structure of impedance measurements. In: *Electrochemical Impedance Spectroscopy*, Kap. 21, S. 407–425. Wiley, 2008.
- [186] M. E. Orazem & B. Tribollet. Methods for representing impedance. In: *Electrochemical Impedance Spectroscopy*, Kap. 16, S. 309–331. Wiley, 2008.
- [187] J. Park & I. W. Sandberg. Universal approximation using radial-basis-function networks. *Neural Computation*, 3(2):246–257, 1991.
- [188] I. P. Pavlov. *Die Arbeit der Verdauungsdrüsen*. Bergmann, Wiesbaden, 1898.
- [189] C. S. Peirce. *Collected Papers of Charles Sanders Peirce*, Vol. 5. Harvard University Press, Cambridge, MA, 1974.
- [190] C.-Y. J. Peng et al. An introduction to logistic regression analysis and reporting. *The Journal of Educational Research*, 96(1):3–14, 2002.
- [191] M. Pinto et al. Enterocytic differentiation of cultured human-colon cancer-cells by replacement of glucose by galactose in the medium. *Biology of the Cell*, 44(2):193–196, 1982.
- [192] K. R. Popper. *Logik der Forschung: zur Erkenntnistheorie der moderner Naturwissenschaft*. Springer, 1935.
- [193] K. R. Popper & H. Vetter. *Objektive Erkenntnis*. Hoffmann und Campe Hamburg, 1974.
- [194] R. Popping. On agreement indices for nominal data. In: W. E. Saris & I. N. Gallhofer (Hrsg.), *Sociometric Research: Data Collection and Scaling*, Vol. 1, S. 90–105. St. Martin's Press, New York, 1988.
- [195] W. J. Potter & D. Levine-Donnerstein. Rethinking validity and reliability in content analysis. *Journal of Applied Communication Research*, 27(3):258–284, 1999.

- [196] S. Prabhakaran & C. R. Sullivan. Impedance-analyzer measurements of high-frequency power passives: techniques for high power and low impedance. In: *Industry Applications Conference, 2002. Conference Record of the 37th IAS Annual Meeting*, Vol. 2, S. 1360–1367. IEEE, 2002.
- [197] L. Prechelt. Early stopping - but when? In: G. Orr & K.-R. Müller (Hrsg.), *Neural Networks: Tricks of the Trade*, Vol. 1524 of *Lecture Notes in Computer Science*, S. 53–67. Springer, 1998.
- [198] R. Pring. Bloom's taxonomy: a philosophical critique (2). *Cambridge Journal of Education*, 1(2):83–91, 1971.
- [199] V. Raghavan et al. A critical investigation of recall and precision as measures of retrieval system performance. *ACM Transactions on Information Systems*, 7(3):205–229, 1989.
- [200] I. D. Raistrick. Application of impedance spectroscopy to materials science. *Annual Review of Materials Science*, 16(1):343–370, 1986.
- [201] J. Rasmussen. On the structure of knowledge – a morphology of metal models in a man-machine system context. Report Riso-M-2192, Riso National Laboratory, Roskilde, Denmark, 1979.
- [202] A. D. Redish. Addiction as a computational process gone awry. *Science*, 306(5703):1944–1947, 2004.
- [203] K. Reich. Systemisch-konstruktivistische Didaktik. Eine allgemeine Zielbestimmung. *Die Schule neu erfinden*, S. 70–91, 1996.
- [204] K. Reich (Hrsg.). *Systemisch-konstruktivistische Pädagogik - Einführung in Grundlagen einer interaktionistisch-konstruktivistischen Pädagogik*. Luchterhand, Neuwied, 3. Aufl., 2000.
- [205] K. Reich. *Konstruktivistische Didaktik. Lehren und Lernen aus interaktionistischer Sicht*. Luchterhan, München, 2. Aufl., 2004.
- [206] J. B. Reitsma et al. Bivariate analysis of sensitivity and specificity produces informative summary measures in diagnostic reviews. *Journal of Clinical Epidemiology*, 58(10):982–990, 2005.
- [207] J. M. Rhoads et al. L-glutamine and L-asparagine stimulate $\text{Na}^+\text{-H}^+$ exchange in porcine jejunal enterocytes. *American Journal of Physiology-Gastrointestinal and Liver Physiology*, 266(5):G828–G838, 1994.
- [208] J. Richards & E. v. Glasersfeld. Die Kontrolle von Wahrnehmung und die Konstruktion von Realität. Erkenntnistheoretische Aspekte des Rückkopplungs-Kontroll-Systems. In: S. J. Schmidt (Hrsg.), *Der Diskurs des Radikalen Konstruktivismus*, S. 192–228. Frankfurt, 6. Aufl., 1994.
- [209] J. Richardson et al. Identification of two strains of MDCK cells which resemble separate nephron tubule segments. *Biochimica et Biophysica Acta (BBA) - General Subjects*, 673:26–36, 1981.
- [210] M. Riedmiller. Advanced supervised learning in multi-layer perceptrons – from backpropagation to adaptive learning algorithms. *Computer Standards & Interfaces*, 16(3):265–278, 1994.

- [211] M. Riedmiller & H. Braun. A direct adaptive method for faster backpropagation learning: the RPROP algorithm. *IEEE International Conference on Neural Networks*, S. 586–591, 1993.
- [212] S. Rimmon. *The Concept of Ambiguity. The Example of James Henry*. University of Chicago Press, 1977.
- [213] M. Ringnér. What is principal component analysis? *Nature Biotechnology*, 26(3):303–304, 2008.
- [214] Y. Rojanasakul et al. The transport barrier of epithelia: a comparative study on membrane permeability and charge selectivity in the rabbit. *Pharmaceutical Research*, 9(8):1029–1034, 1992.
- [215] L. Rokach. Ensemble-based classifiers. *Artificial Intelligence Review*, 33(1-2):1–39, 2009.
- [216] J. Rose. The early years: some comments on the origins and concepts of cybernetics. *Kybernetes*, 38(1/2):20–24, 2009.
- [217] F. Rosenblatt. The perceptron: a probabilistic model for information storage and organization in the brain. *Psychological Review*, 65(6):386–408, 1958.
- [218] R. Rosenthal et al. The effect of chitosan on transcellular and paracellular mechanisms in the intestinal epithelial barrier. *Biomaterials*, 33(9):2791–2800, 2012.
- [219] W. B. Rouse & N. M. Morris. On looking into the black box: prospects and limits in the search for mental models. *Psychological Bulletin*, 100(3):349–363, 1986.
- [220] E. J. Rubin. *Visuell wahrgenommene Figuren. Studien in psychologischer Analyse*. Dissertation, Kopenhagen/Berlin, 1921. Originaltitel: „Synsoplevede Figurer“.
- [221] D. W. Ruck et al. Feature selection using a multilayer perceptron. *Journal of Neural Network Computing*, 2(2):40–48, 1990.
- [222] R. Ruffo et al. Impedance analysis of silicon nanowire lithium ion battery anodes. *The Journal of Physical Chemistry C*, 113(26):11390–11398, 2009.
- [223] R. Ruiz et al. Incremental wrapper-based gene selection from microarray data for cancer classification. *Pattern Recognition*, 39(12):2383–2392, 2006.
- [224] G. Rusch. *Erkenntnis, Wissenschaft, Geschichte: von einem konstruktivistischen Standpunkt*. Suhrkamp, 1987.
- [225] D. Rustemeyer. Konstruktivismus in der Erziehungswissenschaft. In: *Stichwort: Zeitschrift für Erziehungswissenschaft*, Vol. 2, S. 125–144. Springer, 2013.
- [226] W. Samek et al. Explainable artificial intelligence: understanding, visualizing and interpreting deep learning models. *arXiv preprint*, (arXiv:1708.08296v1), 2017.
- [227] K. Sastry et al. *Search Methodologies. Introductory Tutorials in Optimization and Decision Support Techniques*, Kap. 4 (Genetic algorithms), S. 97–125. Springer, 2005.
- [228] P. Schierack et al. Characterization of a porcine intestinal epithelial cell line for in vitro studies of microbial pathogenesis in swine. *Histochemistry and Cell Biology*, 125(3):293–305, 2006.

- [229] T. Schmid et al. Using an artificial neural network to determine electrical properties of epithelia. *Artificial Neural Networks–ICANN 2010*, S. 211–216, 2010.
- [230] T. Schmid et al. Discerning apical and basolateral properties of HT-29/B6 and IPEC-J2 cell layers by impedance spectroscopy, mathematical modeling and machine learning. *PLOS ONE*, 8(7):e62913, 2013.
- [231] T. Schmid et al. Efficient prediction of x-axis intercepts of discrete impedance spectra. In: *Proceedings of the 21st European Symposium on Artificial Neural Networks, Computational Intelligence and Machine Learning (ESANN)*, S. 185–190. 2013.
- [232] T. Schmid et al. Automated quantification of the relation between resistor-capacitor subcircuits from an impedance spectrum. In: *Proceedings of the International Conference on Bio-Inspired Systems and Signal Processing*, S. 141–148. 2014.
- [233] T. Schmid et al. Automated quantification of the capacitance of epithelial cell layers from an impedance spectrum. In: A. Cheptsov & H. H. Ali (Hrsg.), *Proceedings of the Seventh International Conference on Bioinformatics, Biocomputational Systems and Biotechnologies*, S. 27–32. 2015.
- [234] H. Schröder. Ueber eine optische Inversion bei Betrachtung verkehrter, durch optische Vorrichtung entworfener, physischer Bilder. *Annalen der Physik*, 181(10):298–311, 1858.
- [235] J. D. Schulzke et al. Anti-diarrheal mechanism of the traditional remedy Uzara via reduction of active chloride secretion. *PloS one*, 6(3):e18107, 2011.
- [236] W. A. Scott. Reliability of content analysis: the case of nominal scale coding. *Public Opinion Quarterly*, 1955.
- [237] C. L. Sears & J. B. Kaper. Enteric bacterial toxins: mechanisms of action and linkage to intestinal secretion. *Microbiological Reviews*, 60(1):167, 1996.
- [238] M. Seeger. *Semi-supervised Learning*, Kap. 2 (A Taxonomy for semi-supervised learning methods), S. 17–32. MIT press, Cambridge, 2006.
- [239] L. B. Sheiner & S. L. Beal. Some suggestions for measuring predictive performance. *Journal of Pharmacokinetics and Biopharmaceutics*, 9(4):503–512, 1981.
- [240] D. Silver et al. Mastering the game of Go with deep neural networks and tree search. *Nature*, 529:484–489, 2016.
- [241] E. J. Simpson. The classification of educational objektives, psychomotor domain. In: C. J. Cotrell & E. F. Hauck (Hrsg.), *Educational Media in Vocational and Technical Education: A Report of a National Seminar*, S. 38–47. Ohio State University, 1967.
- [242] H. Sockett. Bloom’s taxonomy: a philosophical critique (I). *Cambridge Journal of Education*, 1(1):16–25, 1971.
- [243] H. Stachowiak. *Allgemeine Modelltheorie*. Springer, 1973.
- [244] J. Stephen. Pathogenesis of infectious diarrhea. *Canadian Journal of Gastroenterology*, 15(10):669–683, 2001.
- [245] C. Suhm. *Wissenschaftlicher Realismus. Eine Studie zur Realismus- Antirealismus-Debatte in der neueren Wissenschaftstheorie*. Ontos, Frankfurt am Main, 2005.

- [246] H. E. Tinsley & D. J. Weiss. Interrater reliability and agreement of subjective judgments. *Journal of Counseling Psychology*, 22(4):358–376, 1975.
- [247] E. C. Tolman. *Purposive behavior in animals and men*. The Century Company, New York, 1932.
- [248] E. C. Tolman. Cognitive maps in rats and men. *Psychological Review*, 55(4):189–208, 1948.
- [249] E. C. Tolman. There is more than one kind of learning. *Psychological Review*, 56(3):144–155, 1949.
- [250] H. Troeger et al. Effect of chronic *Giardia lamblia* infection on epithelial transport and barrier function in human duodenum. *Gut*, 56(3):328–335, 2007.
- [251] H. Troeger et al. Structural and functional changes of the duodenum in human norovirus infection. *Gut*, 58(8):1070–1077, 2009.
- [252] K. Tumer & J. Ghosh. Error correlation and error reduction in ensemble classifiers. *Connection Science*, 8(3-4):385–404, 1996.
- [253] A. M. Turing. Computing machinery and intelligence. *Mind*, 59(236):433–460, 1950.
- [254] J. R. Turner & J. L. Madara. Epithelia: biological principles of organization. In: T. Yamada (Hrsg.), *Textbook of Gastroenterology I*, S. 169–186. Wiley-Blackwell, 5. Aufl., 2009.
- [255] D. Umbach & K. N. Jones. A few methods for fitting circles to data. *IEEE Transactions on Instrumentation and Measurement*, 52(6):1881–1885, 2003.
- [256] L. J. van der Maaten et al. Dimensionality reduction: a comparative review. *Journal of Machine Learning Research*, 10(1-41):66–71, 2009.
- [257] W. Veldhuyzen & H. G. Stassen. The internal model concept: an application to modeling human control of large ships. *Human Factors: The Journal of the Human Factors and Ergonomics Society*, 19(4):367–380, 1977.
- [258] A. Vellido et al. Making machine learning models interpretable. In: *Proceedings of the European Symposium on Artificial Neural Networks, Computational Intelligence and Machine Learning*, S. 163–164. Bruges (Belgium), 2012.
- [259] E. Vigneau & E. M. Qannari. Clustering of variables around latent components. *Communications in Statistics - Simulation and Computation*, 32(4):1131–1150, 2003.
- [260] L. S. Vygotsky. *Problems of General Psychology*, Vol. 1 of *The Collected Works of L.S. Vygotsky*. Springer, 1987.
- [261] H. Wang & M. Song. Ckmeans.1d.dp: optimal k-means clustering in one dimension by dynamic programming. *The R Journal*, 3(2):29–33, 2011.
- [262] R. Wang. AdaBoost for feature Selection, classification and its relation with SVM: a review. *Physics Procedia*, 25:800–807, 2012.
- [263] T. Wardlaw et al. Diarrhoea: why children are still dying and what can be done. *The Lancet*, 375(9718):870–872, 2010.

Literaturverzeichnis

- [264] A. Whitney. A direct method of nonparametric measurement selection. *IEEE Transactions on Computers*, C-20(9):1100–1103, 1971.
- [265] C. D. Wickens. *Engineering Psychology and Human Performance*. Prentice Hall, New Jersey, 3. Aufl., 2000.
- [266] C. Yeaman et al. New perspectives on mechanisms involved in generating epithelial cell polarity. *Physiological Reviews*, 79(1):73–98, 1999.
- [267] J. You. Beyond the Turing test. *Science*, 347(6218):116–116, 2015.
- [268] X. Yuan et al. AC impedance technique in PEM fuel cell diagnosis – a review. *International Journal of Hydrogen Energy*, 32(17):4365–4380, 2007.
- [269] S. S. Zakrzewski et al. Improved cell line IPEC-J2, characterized as a model for porcine jejunal epithelium. *PloS one*, 8(11):e79643, 2013.
- [270] S. Zeissig et al. Changes in expression and distribution of claudin 2, 5 and 8 lead to discontinuous tight junctions and barrier dysfunction in active Crohn’s disease. *Gut*, 56(1):61–72, 2007.
- [271] T. Zhang & B. Yu. Boosting with early stopping: convergence and consistency. *The Annals of Statistics*, 33(4):1538–1579, 2005.
- [272] X. Zhu & A. B. Goldberg. Overview of Semi-Supervised Learning. In: R. J. Brachman & T. Dietterich (Hrsg.), *Introduction to Semi-supervised Learning*, Kap. 2, S. 9–20. Morgan & Claypool, 2009.
- [273] X. Zhu et al. Humans perform semi-supervised classification too. In: *Proceedings of the National Conference on Artificial Intelligence*, Vol. 22, S. 864–869. AAAI Press, 2007.
- [274] K. Zimmermann. Zur Morphologie der Epithelzellen der Säugetiere. *Archiv für Mikroskopische Anatomie*, 78:199–234, 1911.
- [275] W. Ziss & U. Hegel. Human colon cancer cells (HT-29) exhibit ion transport properties of intestinal crypt cells. *Pflügers Archiv European Journal of Physiology*, 411(S1):R81, 1988.
- [276] D. Zschocke. *Modellbildung in der Ökonomie: Modell-Information-Sprache*. Vahlen, München, 1995.

Abbildungsverzeichnis

1.1.	Gemeinfreie fotografische Reproduktion des Gemäldes “Blüte und Verwesung”	2
1.2.	Bekannte Kippbilder aus psychologischen Studien	4
1.3.	Grundbaustein des neurophysiologischen Paradigmas	5
2.1.	Schematische Darstellung der überarbeiteten Bloomschen Taxonomie	13
3.1.	Schematische Einordnung etablierter Lern- bzw. Cluster-Verfahren.	16
3.2.	Schematische Darstellung eines neuronalen Netzes	17
3.3.	Schematische Darstellung eines Ensembles aus drei neuronalen Netzen	19
4.1.	Darstellungsformen für Impedanzspektrn	25
4.2.	Schematische Darstellung eines einschichtigen Epithels	26
4.3.	Schematischer Aufbau einer Ussing-Kammer	28
5.1.	Grundstruktur konstruktivistischen maschinellen Lernens	35
5.2.	Ablauf des Lernprozesses	36
5.3.	Konstruktionsprozess für konzeptuelles Wissen	38
5.4.	Rekonstruktionsprozess	41
5.5.	Dekonstruktionsprozess	45
5.6.	Vollständige Dekonstruktion	48
6.1.	Elektrische Modelle für epitheliales Gewebe	50
6.2.	Frequenz- und widerstandsabhängiger Bias des Messaufbaus	55
7.1.	Verlauf der Exploration der konzeptuellen Wissensdomäne	61
7.2.	Überblick über die explorierte konzeptuelle Wissensdomäne	63
7.3.	Temporale Gültigkeit der Modelle der konzeptuellen Wissensdomäne	64
8.1.	Verlauf der Exploration der prozeduralen Wissensdomäne	71
8.2.	Wertebereiche der Modelle der explorierten prozeduralen Wissensdomäne	72
8.3.	Überblick über die explorierte prozedurale Wissensdomäne	73
8.4.	Temporale Gültigkeit der Modelle der prozeduralen Wissensdomäne	74
8.5.	Boxplot des relativen Fehlers der bestimmten epithelialen Kapazität [%].	77
9.1.	Charakteristische Parameter der Datenbasis	83
10.1.	Ergebnis einer Abstraktionskaskade	89
10.2.	Entwicklung des Speicherbedarfs der konzeptuellen Wissensdomäne im Verlauf der Exploration	91
10.3.	Rechenzeit für die Exploration der konzeptuellen Wissensdomäne	93
11.1.	Klassifikatoren-Konnektom der explorierten konzeptuellen Wissensdomäne	97
C.1.	Verteilung der epithelialen Kapazität über die Datenbasis	135
C.2.	Verteilung des zeitlichen Metadatumms über die Datenbasis	136

Abbildungsverzeichnis

D.1. Überblick über die adaptierte konzeptuelle Wissensdomäne	138
D.2. Temporale Gültigkeit des adaptierten konzeptuellen Wissens	139
D.3. Temporale Gültigkeit des adaptierten prozeduralen Wissens	140
D.4. Überblick über die adaptierte prozedurale Wissensdomäne	141

Tabellenverzeichnis

5.1.	Übersicht der Metadaten eines pragmatisch definierten maschinellen Modells.	34
5.2.	Übersicht der Blockverarbeitungsparameter im konstruktivistischen maschinellen Lernen.	37
6.1.	Publizierte und modellierte Parameter für die Zelllinie HT-29/B6 unter physiologischen Bedingungen	52
6.2.	Publizierte und modellierte Parameter für die Zelllinie IPEC-J2 unter physiologischen Bedingungen	52
6.3.	Modellierte Parameter für die Zelllinie MDCK I unter physiologischen Bedingungen	52
6.4.	Anzahl der modellierten Impedanzspektren pro Zellkulturlinie und funktionalem Zustand.	56
6.5.	κ_B -Werte pro Zellkulturlinie und funktionalem Zustand	57
7.1.	Übersicht der zur Exploration konzeptuellen Wissens verwendeten konstruktivistischen Lernparameter.	59
7.2.	Übersicht über die Modellzwecke der explorierten konzeptuellen Wissensdomäne.	62
7.3.	Übersicht der zur Adaption konzeptuellen Wissens verwendeten konstruktivistischen Lernparameter	65
7.4.	Verlauf der Adaption konzeptuellen Wissens	66
7.5.	Konfusionsmatrizen für die Klassifizierung nach zugehöriger Zelllinie	67
8.1.	Übersicht der zur Exploration prozeduralen Wissens verwendeten konstruktivistischen Lernparameter	69
8.2.	Übersicht der zur Adaption prozeduralen Wissens verwendeten konstruktivistischen Lernparameter.	75
8.3.	Verlauf der Adaption prozeduralen Wissens	76
8.4.	Relativer Fehler der bestimmten epithelialen Kapazität	77
A.1.	Übersicht der Trainingsparameter der unüberwachten Lernverfahren.	106
A.2.	Übersicht der Trainingsparameter der überwachten Lernverfahren	107
A.3.	Übersicht der konstruktivistischen Lernparameter	109
B.1.	Frequenzen der modellierten und gemessenen Impedanzspektren.	117
B.2.	Modellierte Parameter für die Zelllinie HT-29/B6 unter physiologischen Bedingungen.	121
B.3.	Modellierte Parameter für die Zelllinie HT-29/B6 unter Einfluss von Nystatin.	122
B.4.	Modellierte Parameter für die Zelllinie HT-29/B6 unter Einfluss von EGTA.	123
B.5.	Modellierte Parameter für die Zelllinie HT-29/B6 unter Einfluss von EGTA und Nystatin.	124
B.6.	Modellierte Parameter für die Zelllinie IPEC-J2 unter physiologischen Bedingungen.	125
B.7.	Modellierte Parameter für die Zelllinie IPEC-J2 unter Einfluss von Nystatin.	126
B.8.	Modellierte Parameter für die Zelllinie IPEC-J2 unter Einfluss von EGTA.	127
B.9.	Modellierte Parameter für die Zelllinie IPEC-J2 unter Einfluss von EGTA und Nystatin.	128

Tabellenverzeichnis

B.10. Modellierter Parameter für die Zelllinie MDCK I unter physiologischen Bedingungen.	129
B.11. Modellierter Parameter für die Zelllinie MDCK I unter Einfluss von Nystatin. . .	130
B.12. Modellierter Parameter für die Zelllinie MDCK I unter Einfluss von EGTA. . . .	131
C.1. Globale Features für das Feature-Set $S_{\mathfrak{R}}$	134
D.1. Überblick über die Eigenschaften der adaptierten konzeptuellen Modelle	137
D.2. Features der adaptierten konzeptuellen Modelle	139
D.3. Überblick über die Eigenschaften der adaptierten prozeduralen Modelle	140
D.4. Features der adaptierten prozeduralen Modelle.	142