

FH Aachen

**Fachbereich
Elektrotechnik und Informationstechnik**

Bachelorarbeit

**Prognose der Anwesenheit von Personen für die
Gebäudeautomatisierung mittels Umweltsensordaten**

**Alexander Loosen
Matr.-Nr.: 3167353**

Referent: Prof. Dr-Ing. Ingo Elsen

7. April 2022

Erklärung

Ich versichere hiermit, dass ich die vorliegende Arbeit selbstständig verfasst und keine anderen als die im Literaturverzeichnis angegebenen Quellen benutzt habe.

Stellen, die wörtlich oder sinngemäß aus veröffentlichten oder noch nicht veröffentlichten Quellen entnommen sind, sind als solche kenntlich gemacht.

Die Zeichnungen oder Abbildungen in dieser Arbeit sind von mir selbst erstellt worden oder mit einem entsprechenden Quellennachweis versehen.

Diese Arbeit ist in gleicher oder ähnlicher Form noch bei keiner anderen Prüfungsbehörde eingereicht worden.

Aachen, 7. April 2022

Geheimhaltung - Sperrvermerk

Die vorliegende Arbeit unterliegt bis [Datum] der Geheimhaltung. Sie darf vorher weder vollständig noch auszugsweise ohne schriftliche Zustimmung des Autors, des betreuenden Referenten bzw. der Firma [Firmenname und -sitz] vervielfältigt, veröffentlicht oder Dritten zugänglich gemacht werden.

Inhalt

1. Einleitung	5
1.1. Motivation und Aufgabenstellung	5
1.2. Vorgehensweise	6
2. Stand der Technik für Anwesenheitserkennung/-Prognose	7
3. Untersuchte Verfahren	8
3.1. Machine Learning	8
3.1.1. Random Forest Classifier	10
3.1.2. Gradient Boosting Classifier	10
3.1.3. Support Vector Classifier	11
3.1.4. Logistische Regression	12
3.1.5. K-Nearest-Neighbours	14
3.1.6. Neuronale Netzwerke	15
3.1.7. Long Short Term Memory	15
3.2. CO2 als Anwesenheitsindikator	17
3.3. Luftfeuchtigkeit und Temperatur	18
3.4. Sensordaten	18
4. Technische Umsetzung	19
4.1. Datenbeschaffung und -Vorbereitung	19
4.1.1. Gruppierung	19
4.1.2. Zyklische Codierung	19
4.1.3. Deltas und Shift-Werte	20
4.1.4. Outlier Detection	20
4.2. Validierung	21
4.2.1. Over- und Underfitting	23
4.2.2. Parameter Tuning	24
5. Anwendung und Ergebnisse	25
5.1. Feature Vektor	25
5.2. Ergebnisse	29
5.2.1. Decision Tree Modelle	29
5.2.2. Umgang mit fehlenden Präsenz-Labels	32
5.2.3. Neurale Netzwerke	33
5.2.4. Clustering Modelle	34
6. Zusammenfassung und Ausblick	36
6.1. CO2 als Anwesenheitsindikator	36
6.2. Mögliche Verbesserungen	36
Quellenverzeichnis	37
Abkürzungsverzeichnis	38

	Inhalt
Abbildungsverzeichnis	39
Tabellenverzeichnis	40
Anhang	40
A. Quellcode	41
B. Rohdatenvisualisierungen	42

1. Einleitung

Gebäudeautomatisierung bezeichnet die automatische Steuerung und Regelung von Gebäudetechnik wie Heizung, Lüftung oder Beleuchtung. Während sie bisher hauptsächlich für die Optimierung der Energieeffizienz von gewerblichen und öffentlichen Gebäuden genutzt wurde, welche in Zuge solcher Optimierungsschritte als „Smart Buildings“ bezeichnet werden, rückt sie in den letzten Jahren zunehmend unter dem Begriff „Smart Home“ auch in den privaten Bereich.

Die beiden Begriffe stehen in den letzten Jahren so im Vordergrund, weil eine Verbesserung der Energieeffizienz durch bauphysikalische Maßnahmen, wie verminderte Wärmeverluste durch bessere Isolation, an ihre Grenzen gestoßen sind.

In Deutschland haben private Haushalte einen Anteil am gesamten Energieverbrauch des Landes von etwa 29%¹. Davon nimmt die Erzeugung der Raumwärme einen Anteil von etwa 68%² ein. Vorallem bezogen auf den Energieverbrauch durch das Heizen stellt die Gebäudeautomatisierung also einen der größten Sektoren dar, in denen eine Optimierung der Prozesse entscheidend wäre. Beispielsweise können Wärmeverluste über die Nacht hinweg um 20% verringert werden, wenn Abends die Rolläden heruntergelassen werden.

Zur weiteren Steigerung der Energieeffizienz ist es also nötig, die Gebäudetechnik insofern automatisch anzusteuern, sodass sog. Performance-Gaps vermieden werden. Performance-Gaps stellen eine Diskrepanz im Energieverbrauch eines Gebäudes zwischen einem theoretischen Soll-Wert zu einem tatsächlichen Ist-Wert dar.

1.1. Motivation und Aufgabenstellung

Für nahezu alle Bereiche der Gebäudeautomatisierung stellt die Anwesenheit von Personen eine zentrale Variable dar. Es ist entscheidend bei einem Optimierungsprozess zu erkennen, zu welchen Zeitpunkten sich Menschen in einem Raum aufhalten. Beispielsweise werden in vielen gewerblich genutzten Gebäuden werden grundsätzlich alle Räume beheizt - unabhängig von tatsächlicher Anwesenheit von Personen. Eine Erkennung von menschlicher Anwesenheit ist bei der Optimierung also essenziell.

Da die direkte Messung von Anwesenheit über z.B. Infrarot- oder Kamerasensoren rechtlich problematisch ist, soll in dieser Arbeit untersucht werden, inwiefern Machine-Learning Algorithmen genutzt werden können, um eine genaue Erwartung über die Anwesenheit von Personen anhand von CO₂-Werten in der Raumluft zu treffen.

Die Motivation der Optimierung der Gebäudeautomatisierung existiert, da ein steigender CO₂-Gehalt der Raumluft nachweislich mit einer Abnahme der menschlich kognitiven Leistung einhergeht. Mehrere Studien konnten belegen, dass sowohl sprachliche als auch

¹[Umweltbundesamt, 2022]

²[Umweltbundesamt, 2020]

logisch- mathematische Fähigkeiten abnehmen, sobald der CO₂-Gehalt der Raumluft bestimmte Werte überschreitet.

Um eine angemessene Datengrundlage zu schaffen, wurden in diversen Büro-Räumen der FH Aachen Temperatur-, Luftfeuchtigkeits-, Infrarot- und CO₂-Sensoren angebracht, deren Messungen kontinuierlich auf einer Datenbank gespeichert wurden. Der Zeitraum der Messwerte begann Mitte 2021.

In allen Räumen sind täglich ein oder mehrere Personen im Rahmen eines ca. 8 stündigen Arbeitstages anwesend, weshalb die Temperatur-, Luftfeuchtigkeits- und CO₂-Messwerte über einen Tag hinweg Schwankungen aufweisen. Diese sollen von einem Machine Learning Modell mit der direkten Präsenzmessung des Infrarotsensors gegenübergestellt werden, sodass das fertig Trainierte Modell anhand der Messwerte von Temperatur, Luftfeuchtigkeit und CO₂ Aussagen über menschliche Präsenz treffen kann.

Die Frage, ob ein in einem Raum trainiertes Modell auch präzise Aussagen über die Anwesenheit in einem anderen Raum treffen kann, stellte eine besondere Bedeutung dar.

Es gab keine Einschränkungen hinsichtlich dessen, welche konkreten Machine-Learning Algorithmen benutzt werden sollen.

Als Programmiersprache für das Projekt wurde Python gewählt. Python ist wegen seiner umfassenden Machine-Learning Bibliotheken und einfachen Auswertungstechniken anhand von z.B. Graphen und Statistiken gut für diesen Anwendungsfall geeignet.

1.2. Vorgehensweise

Das Projekt beschäftigte sich im Schwerpunkt mit den folgenden Arbeitsschritten:

- Datenbeschaffung durch Datenbankzugriffe per SQL
- Analyse und Vorbereitung der Daten (Pre-Processing)
- Trainieren von Machine-Learning Modellen anhand der vorbereiteten Datensets
- Ergebnisauswertung durch Gegenüberstellung verschiedener Datensets und Modellen

Da es zwischen allen verfügbaren Datensets der einzelnen Räume und Machine-Learning-Modellen eine Vielzahl an Kombinationsmöglichkeiten gibt, war es ein Anspruch der Projektarbeit, ein möglichst übersichtliches, gut gekapseltes Python Programm zu erstellen, mit dem man einfach und schnell verschiedene Datensets verarbeiten und mit einer dem Forschungszweck angemessenen Anzahl von Machine-Learning Modellen auszuwerten.

Um einen Vergleich der Ergebnisse zu ermöglichen, sollen diese klar und verständlich dargestellt werden. Da nicht bei allen Algorithmen die gleichen Leistungsindikatoren genutzt werden, sollen hauptsächlich nur jene Indikatoren betrachtet werden, die bzgl. aller Algorithmen auch gleiche Bedeutung haben. Falls Modell- spezifische Leistungsindikatoren als besonders Erkenntnisreich erachtet werden, wird dies in dieser Arbeit angemerkt.

2. Stand der Technik für Anwesenheitserkennung/-Prognose

Man kann verschiedene Techniken zur Anwesenheitserkennung und -Prognose in Innenräumen in aktive und *passive* Verfahren unterteilen.

Aktive Verfahren messen die menschliche Anwesenheit etwa mit Kamera- oder Infrarotsensoren, welche den Menschen direkt im Raum erkennen können. Ein Kamerasensor nimmt dabei kontinuierlich Bilder des Raumes auf, während ein Algorithmus im Hintergrund versucht ein 3D-Modell eines menschlichen Skeletts auf sich bewegende Teile des Bildes zu legen. Bewegen sich über mehrere Bilder hinweg die Bereiche auf und unmittelbar neben dem 3D-Modell, registriert das System eine Anwesenheit.

Infrarotsensoren dagegen funktionieren eher wie ein klassischer Bewegungsmelder, welcher Unterschiede in der Infrarotstrahlung eines Raumes erkennt. Durch die generierte Körperwärme des menschlichen Körpers, sorgt ein Mensch, der sich durch den Sensor bewegt, für einen Ausschlag des Sensors.

Passive Verfahren orientieren sich an den physikalischen Eigenschaften des Raumes. Anstatt die Präsenz des Menschen direkt zu messen, wird versucht, diese von Veränderungen der Eigenschaften der Raumluft oder elektromagnetischer Strahlung des Raumes zu erkennen.

Verändert sich beispielsweise durch menschliche Anwesenheit die Luftfeuchtigkeit in einem Raum in sehr geringem Maß, wird die Laufzeit einer Ultraschallwelle die durch diesen Raum geschickt wird, leicht verringert, da die Schallgeschwindigkeit in dichteren Medien zunimmt.

Die Systeme unterscheiden sich neben der Art der Messung auch deutlich in der Komplexität der Implementierung. Während Ultraschallsensoren wegen der verhältnismäßig geringen Schallgeschwindigkeit keine besonders hohe Genauigkeit besitzen müssen, sind die an einen Mikrowellensensor gestellten Ansprüche wesentlich höher, da sich Mikrowellen mit $300 \cdot 10^6 \text{ m/s}$ gegenüber der Schallgeschwindigkeit mit lediglich etwa 330 m/s ausbreiten. Mit Mikrowellen können Bewegungen durch Ausnutzung des Doppler-Effektes einer elektromagnetischen Welle erkannt werden. Bewegen sich Objekte oder Personen auf den Sensor zu, wird das Echo der Welle gestaucht, wodurch die Bewegung erkannt wird.

Hier verdeutlicht sich nochmals die Motivation dieser Arbeit, da Sensoren für Temperatur, Luftfeuchtigkeit, oder CO₂ weder teuer sind, noch besondere Ansprüche an Positionierung im Raum oder Kalibrierung besitzen.

3. Untersuchte Verfahren

3.1. Machine Learning

Machine Learning ist ein Teilbereich der künstlichen Intelligenz und beschreibt das Entwickeln mathematischer Modelle zur statistischen Auswertung von Daten.

Dabei wird dem Modell anhand von Trainingsdaten beigebracht, Gesetzmäßigkeiten zu erkennen, um danach Erwartungen über bisher unbekannte Daten treffen zu können. Beispielsweise könnte ein solches Modell aus einem Datenset mit der aktuellen Jahreszeit, Uhrzeit und Position der Sonne am Himmel trainiert werden, sodass es auch schließlich in einem anderen Datenset aus Jahreszeit und Position der Sonne Rückschlüsse auf die Uhrzeit treffen kann.

Als Vorbild für diesen „Lernvorgang“ dient das menschliche Gehirn, welches ebenfalls versucht zwischen bestimmten Input-Parametern wie z.B. der Form und Farbe eines Gegenstandes eine Beziehung herzustellen, um das beobachtete Objekt in Zukunft schneller kategorisieren zu können.

Da eine Vielzahl von effektiven Machine Learning Algorithmen existiert, ist es essenziell, sich mit den Stärken und Schwächen einzelner Herangehensweisen zu befassen.

In dieser Arbeit sollen zwei Unterkategorien des Machine Learning genauer betrachtet werden:

- *Supervised Learning*
- *Unsupervised Learning*

Supervised Learning bedeutet zwischen bestimmten Feldern eines Datensets eine Beziehung zu einem sog. Label herzustellen, welches als eine Art Ergebnis aus den Eingabewerten gesehen werden kann. Ein so trainiertes Modell kann dann neue, ihm vorher unbekannte Datensets, mit einem Label versehen - etwa wie in dem o.g. Beispiel wo Jahreszeit und Sonnenposition die Eingabewerte und die Uhrzeit das Label darstellen. Der Begriff „*supervised*“ ergibt sich daraus, dass das Datenset, mit dem das Modell trainiert wird, diese Labels gegeben hat, sodass das Modell sich bei jedem Schritt des Lernvorgangs selbst korrigieren kann, falls eine Fehleinschätzung getroffen wurde. Bei einer sog. „*Klassifizierung*“ sind diese Labels fest vorgegeben, während sie in der „*Regression*“ kontinuierlicher Natur sind. Im Kontext dieser Arbeit wäre das Ergebnis einer Klassifizierung eine „1“ für Anwesenheit und eine „0“ für Abwesenheit, während das Ergebnis einer Regression eine Wahrscheinlichkeit auf Anwesenheit zwischen 0.0 und 1.0 darstellen würde.

Beim „*Unsupervised Learning*“ versucht das Modell ohne Referenz zu einem bestimmten Label, Zusammenhänge zwischen bestimmten Feldern des Datensets herzustellen. Solche Modelle arbeiten vorrangig mit „*Clustering*“ und „*Dimensionality Reduction*“.

„*Clustering*“-Algorithmen versuchen ein Datenset in kleinere Teilmengen einzuteilen und so aus den Feldern des Datensets bestimmte Abhängigkeiten abzuleiten.

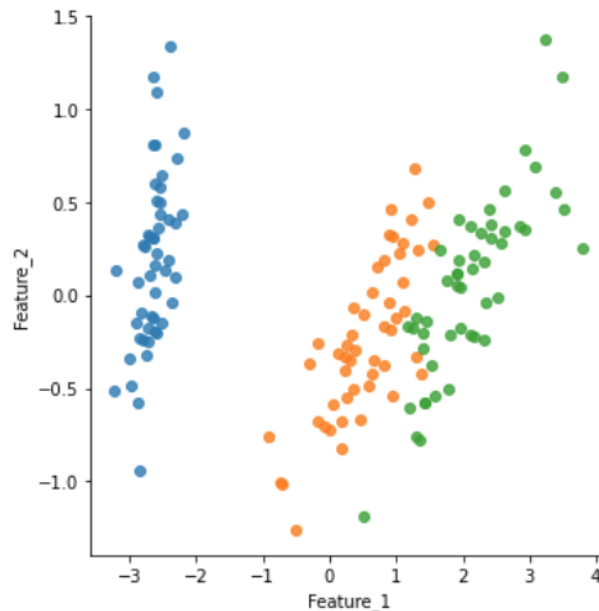


Abbildung 3.1.: Beispiel für Clustering

Bei der „*Dimensionality Reduction*“ versucht der Algorithmus das Datenset in seiner Dimensionalität, also seiner Anzahl an Feldern, zu reduzieren. Es wird also die Frage gestellt, ob sich in einem bestehenden Datenset auch mit weniger Feldern Abhängigkeiten feststellen lassen. Dieser Schritt wird vor allem für Modelle benutzt, die sensibel gegenüber hoher Dimensionalitäten sind, sodass das Datenset vor dem Training in seiner Dimensionalität heruntergebrochen werden kann.

Im Rahmen des Projektes wurden hauptsächlich Klassifizierungs-Algorithmen genutzt, da ein Großteil der Datensets Labels zur Überprüfung hatte.

Um einen Vergleich herzustellen werden später trotzdem noch einzelne Ergebnisse von Clustering und Dimensionality Reduction betrachtet. Im Folgenden sollen die genutzten Modelle erklärt werden.

3.1.1. Random Forest Classifier

Random Forest Classifier RFC stellen eine Unterkategorien der „*Decision Trees*“ dar. *Decision Trees* sind einfache Anordnungen von bestimmten Fragen, die über das Datenset gestellt werden, um eine Klassifikation zu erreichen.



Abbildung 3.2.: Beispiel eines Decision Trees

Erstellt man ein „*Ensemble*“ aus *Decision Trees* die Erwartungen über einen zufällig gewählten Teil des Datensets treffen können, entsteht ein *Random Forest*. Der *Random Forest Classifier* versucht, eine Menge einfacher Schätzfunktionen über einen komplexeren Sachverhalt „abstimmen“ zu lassen. Während sich in einem einzelnen Entscheidungsbaum Fehleinschätzungen entwickeln können, sinkt die Chance auf eine solche Fehleinschätzung, je mehr unabhängige Entscheidungsbäume man befragt.

3.1.2. Gradient Boosting Classifier

Der *Gradient Boosting Classifier* GBC ist eine Abwandlung eines RFC und versucht, seine Erwartungen aufgrund von Abweichungen eines Labels vom Durchschnitt dieses Labels zu treffen.

Erweitert man das Datenset im o.g. Beispiel um das Alter einer Maus, wird ein GBC als Ausgangswert den Durchschnitt aller Label-Werte, also dem Gewicht, berechnen. Danach werden die Abweichungen aller Label-Werte zu diesem Durchschnitt gebildet. Diese Abweichungen werden nun in Beziehung zu den anderen Spalten des Datensets gesetzt.

Beispielsweise könnte man so davon ausgehen, dass ausgewachsene Mäuse von einem bestimmten Alter über dem Durchschnittsgewicht liegen. Genauso liegen besonders junge Mäuse wahrscheinlich immer einen ähnlichen Wert unter dem Durchschnittsgewicht. So wurde zwischen dem Label *Gewicht* und der Spalte *Alter* eine Beziehung hergestellt. In weiteren Iterationen orientiert sich der GBC immer an der Abweichung zum Durchschnittswert des vorherigen Baumes. So werden die getroffenen Erwartungen über mehrere Iterationen immer präziser.

3.1.3. Support Vector Classifier

Der *Support Vector Classifier* SVC versucht in einem Datenset anhand von bestimmten Cut-Off-Values klare Grenzen zwischen Werten zu finden, sodass man alle Messwerte ober- und unterhalb der Grenze eindeutig Klassifizieren kann.

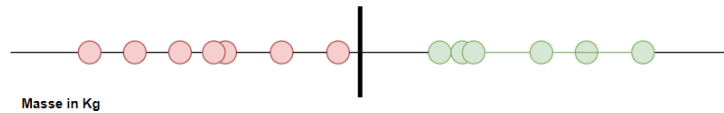


Abbildung 3.3.: Beispiel eines Support Vector Classifiers

In Abb. 3.3 ist der SVC ein Punkt auf einer eindimensionalen Linie, auf der das Gewicht in Kg von z.B. Mäusen in „Unter-“ und „Übergewichtig“ unterteilt wird. Dieser Punkt ist Ergebnis aller Verhältnisse der einzelnen Datenpunkte zueinander. Durch sog. „*Kernel Funktionen*“ versucht der Algorithmus nun Beziehungen in höheren Dimensionen zu finden, wie z.B. $Masse^2$, $Masse^3$ usw. . Der SVC stellt dann in diesen Dimensionen eine Linie in einem zwei-dimensionalen oder eine Ebene in einem drei-dimensionalen Koordinatensystem dar.

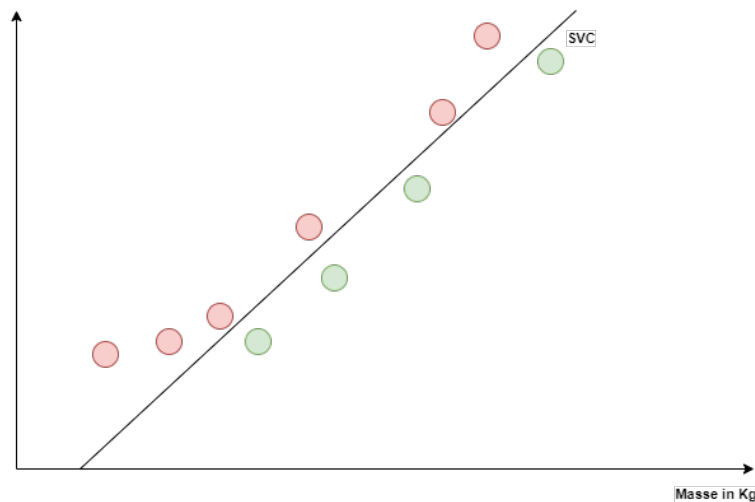


Abbildung 3.4.: Beispiel eines Support Vector Classifiers in der zweiten Dimension

Da der SVC die Verhältnisse aller Datenpunkte zueinander betrachtet, ist er sehr anfällig für Ausreißer in den Daten, was bei der Datenvorbereitung und der Auswertung beachtet werden muss.

3.1.4. Logistische Regression

Die logistische Regression ordnet Daten einem bestimmten Label anhand einer Wahrscheinlichkeit zu. Trotz dem Teilnamen „Regression“ handelt es sich um einen Klassifizierungsalgorithmus. Der Name kommt daher, dass die berechnete Wahrscheinlichkeit zu einem kontinuierlichen Label zwischen 0 und 1 führt, was durch die Form einer Sigmoid-Funktion ausgedrückt wird. Der Algorithmus eignet sich wegen dieser Eigenschaft besonders für das Klassifizierungsproblem der Projektarbeit, bei der ein Wahrheitswert, wie „Anwesenheit“ oder „Abwesenheit“ untersucht werden soll.

Nutzt man logistische Regression beispielsweise zur Abbildung einer Beziehung zwischen Gewicht von Mäusen in Gramm zu einer Wahrscheinlichkeit für Übergewicht, könnte das entstehende Modell wie folgt aussehen.

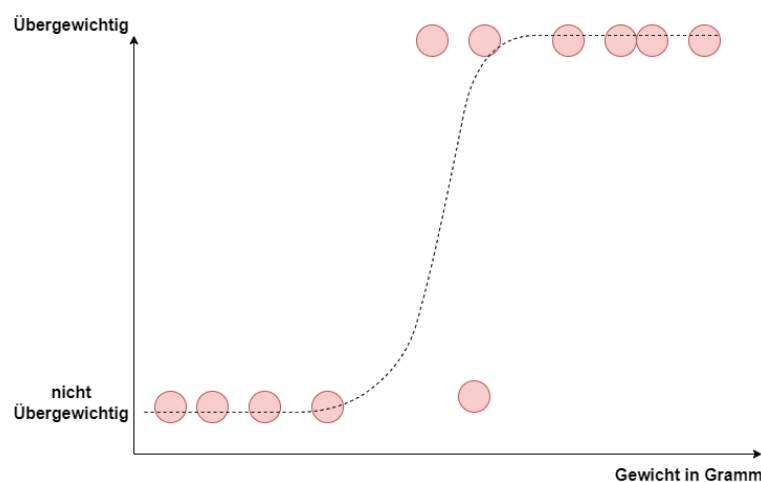


Abbildung 3.5.: Beispiel logistischer Regression

Die Sigmoid-Funktion wird grundsätzlich durch

$$p = \frac{1}{1 + e^{-(x-\mu)/s}} \quad (3.1)$$

beschrieben, wobei x der Eingabewert in Gramm und μ der Mittelpunkt der Kurve beschreibt, der sich aus dem Verhältnis von Eintrittswahrscheinlichkeit p_1 und Gegenwahrscheinlichkeit p_0 ergibt.

$$p(\mu) = 0.5 = \frac{p_1}{p_0} \quad (3.2)$$

Zusätzlich beschreibt s einen Skalierungsparameter, mit dem die Form der Kurve flacher oder steiler werden kann, was angibt, wie eindeutig die hereingegebenen Daten zugeordnet werden können. Wären im o.g. Beispiel also alle roten Punkte jeweils links und rechts von der Mitte einsortiert, wäre die Sigmoid-Funktion in der Mitte sehr steil, da alle Daten anhand des Mittelpunktes klar getrennt werden könnten.

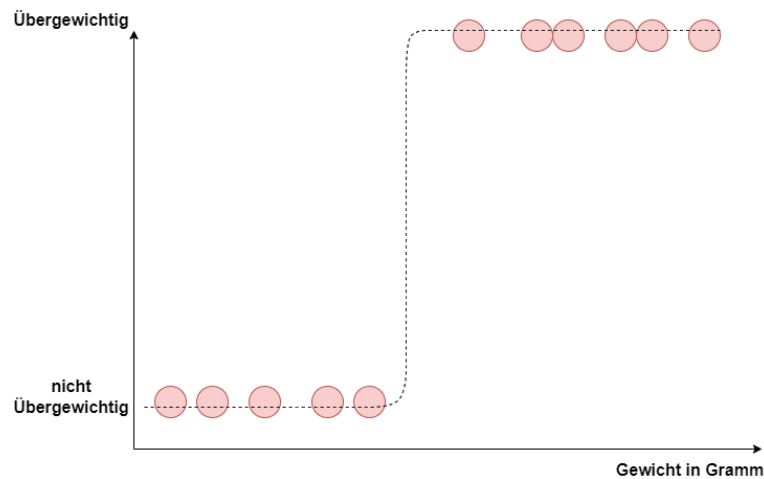


Abbildung 3.6.: Logistische Regression mit deutlicher Skalierung

Anhand der optimierten Sigmoid-Funktion können nun neue Daten klassifiziert werden, wobei die Wahrscheinlichkeit bei $p < 0.5$ auf 0 und $p \geq 0.5$ auf 1 gerundet wird.

3.1.5. K-Nearest-Neighbours

Der *K-Nearest-Neighbours*-Algorithmus KNN betrachtet den Abstand von einem gegebenen Punkt p zu einer Menge an Punkten in einem Datenset und ordnet ihn gemäß dieser Abstände einer bestimmten Kategorie zu. Die Zuordnung geschieht anhand der Kategorie der größten Menge an nächsten Nachbarn.

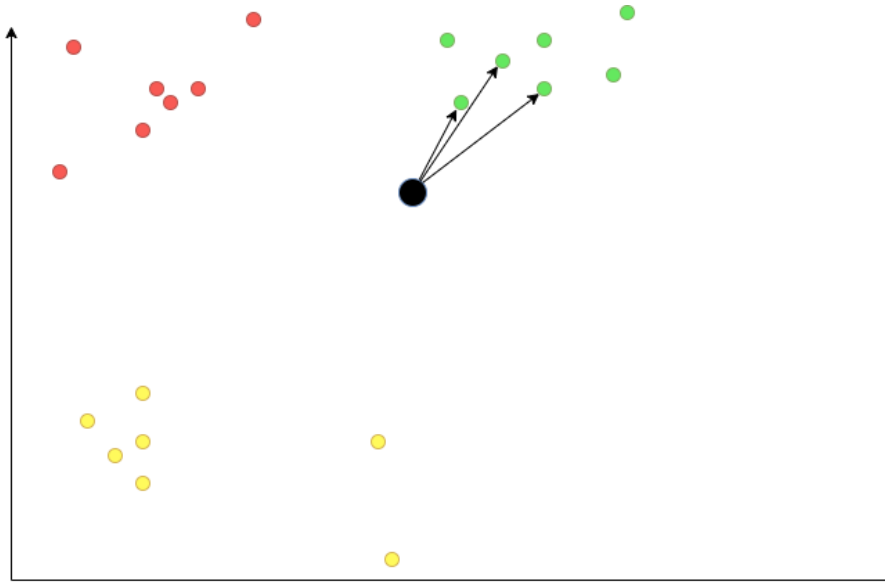


Abbildung 3.7.: Beispiel eines KNN

Im oben gezeigten Beispiel würde ein Punkt anhand seiner drei nächsten Nachbarn als "grün" kategorisiert. Die Anzahl der gewünschten Nachbarn die zu untersuchen sind, ist frei wählbar. Da die Abstandsberechnung über den *Euklidischen Abstand* berechnet wird, funktioniert dieser Algorithmus auch in höheren Dimensionen, da

$$d = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2 + (z_2 - z_1)^2 + \dots + (n_2 - n_1)^2} \quad (3.3)$$

sich in beliebig viele Dimensionen fortsetzen lässt.

3.1.6. Neuronale Netzwerke

Die Funktionsweise eines Neuronalen Netzwerks ist direkt angelehnt an die Funktionsweise des menschlichen Gehirns. Einzelne Knotenpunkten(Neuronen) werden mithilfe von Gewichteten Verbindungen verknüpft, sodass das Netzwerk versucht Eingabewerte bestimmten Ausgabewerten zuzuordnen. Diese Zuordnung der Ein- und Ausgabewerte im Input- und Output-Layer geschieht nicht direkt, sondern durch ein oder mehrere *Hidden Layer*, dessen Neuronenzahl üblicherweise über der Anzahl Neuronen im Input Layer liegt. Die Anzahl der Neuronen im Output-Layer entspricht der Anzahl an Ergebnissen, die sich aus dem Input ergeben können. Im Beispiel der Anwesenheitsanalyse entspräche das hier also zwei Neuronen für An- und Abwesenheit.

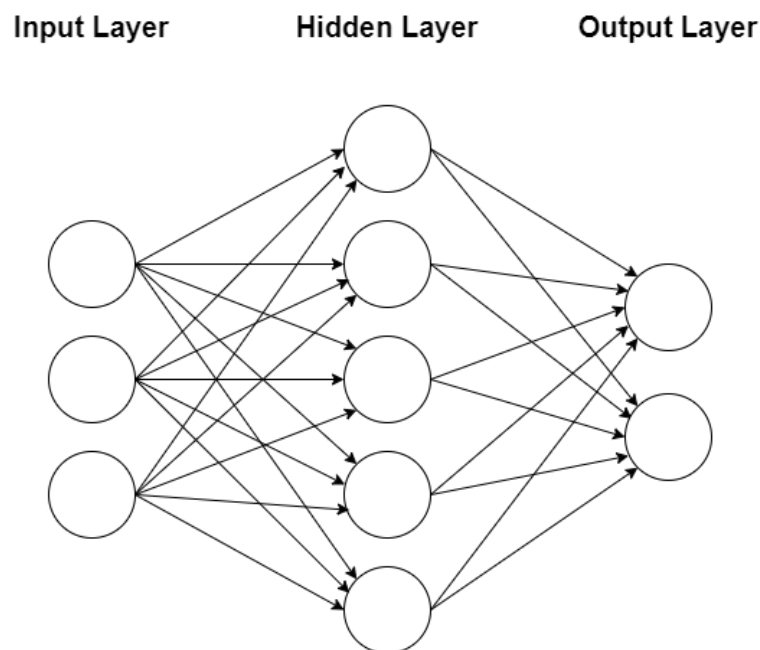


Abbildung 3.8.: Beispiel eines neuronalen Netzwerks

Liegt an einem Neuron eine Information an, wird diese als Eingabe einer Aktivierungsfunktion φ genutzt, die mithilfe eines bestimmten Schwellwertes bestimmt, ob dieses Neuron aufgrund der Eingabe aktiviert wird. Über eine bestimmte Anzahl von Iterationen werden die Gewichtungen zwischen den einzelnen Neuronen stärker oder schwächer.

3.1.7. Long Short Term Memory

Das *Long Short Term Memory* LSTM ist eine Abwandlung herkömmlicher Neuronalen Netzwerke. Es handelt sich um ein *rekurrentes* Neurales Netzwerk, was bedeutet, dass jedes Neuron seine Ausgabewerte auch wieder als Eingabewerte nutzt. Der Begriff *Memory* rührt daher, dass durch diese Rückkopplung eine Art Gedächtnis entsteht, durch die das Netzwerk bessere Rückschlüsse auf die Einordnung des aktuellen Input-Wertes ziehen kann.

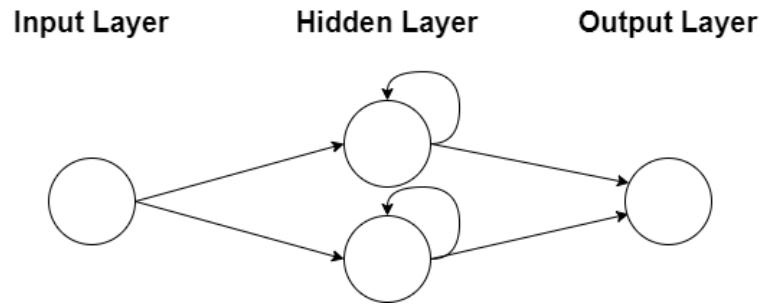


Abbildung 3.9.: Beispiel eines rekurrenten neuronalen Netzwerks

Bei einem LSTM verfügt jedes Neuron zusätzlich über einen Zell-Status C_t , welcher durch einen Input x_t verändert werden kann. Auf den Zell-Status nehmen *Input*-, *Forget*- und *Output*-Gates Einfluss. Diese drei Gates besitzen jeweils ein eigenes neuronales Netzwerk σ .

Das Input-Gate bestimmt, ob der Zell-Status anhand des anliegenden Inputs verändert werden darf. Das Forget-Gate gibt an, ob der Status der Zelle zurückgesetzt und somit das „Gedächtnis“ der Zelle gelöscht wird. Es stellt die zentrale Komponente des LSTMs dar. Das Output-Gate regelt, ähnlich wie Aktivierungsfunktion normaler NNs, ob der anliegende Input zu einem Output der Zelle führt.

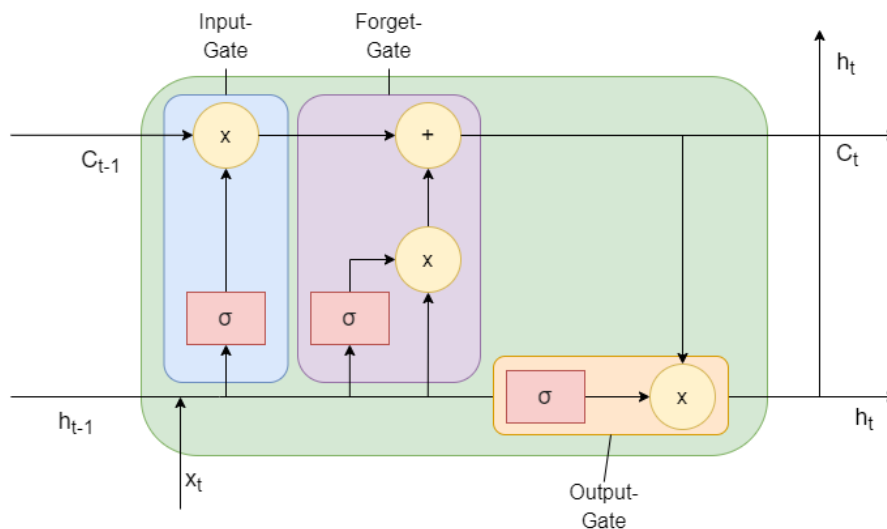


Abbildung 3.10.: Vereinfachter Aufbau einer LSTM-Zelle

Ein LSTM kann auf diese Weise eine Vielzahl von Zeitschritten zurückblicken und verlässt sich so nicht direkt auf einen gegebenen Input, sondern auf einen langen Verlauf von bereits verarbeiteten Input-Werten. LSTMs sind deshalb besonders interessant für Probleme bei denen Beziehungen zwischen kontinuierlichen Datenwerten gebildet werden müssen.

Geht man von einem LSTM aus, dass über eine Liste an Namen verfügt und anhand einem Datensatz gelernt hat, zwei Namen mit der Beziehung „sah“ zu z.B. „Jonas sah Jakob.“ zu verknüpfen.

Das LSTM hat nun intern einen Status, der angibt, dass auf einen Namen mit hoher Wahrscheinlichkeit das Wort „sah“ und mit niedriger Wahrscheinlichkeit ein zweiter Name folgt. Beginnt man nun einen neuen Satz mit „Jakob“, werden als mögliche nächste Worte „Jonas“, „Jakob“ und „sah“ ausgewählt. Bisher ähnelt dieser Verlauf einem normalen Neuronalen Netzwerks. Beim LSTM steht aber nun wegen dem vorherigen Durchlauf der Name „Jakob“ im Forget-Gate, sodass „Jakob“ aus der verfügbaren Wortauswahl für das nächste Wort gelöscht wird.

3.2. CO2 als Anwesenheitsindikator

Der Anteil von CO2 in frischer Atemluft beträgt zwischen 350 und 450 ppm. Es gibt in Deutschland und auch Europa keine grundsätzlich festgelegten Grenzwerte für akzeptable Raumluft, vielmehr raten Gesundheitsämter verschiedener Länder Grenzwerte zwischen 1200 und 1500 ppm einzuhalten. Bei der Obergrenze von 1500 ppm entstehen beim Menschen erste Müdigkeitserscheinungen, weshalb dieser Wert in der Literatur als maximaler Richtwert für Innenräume gilt.

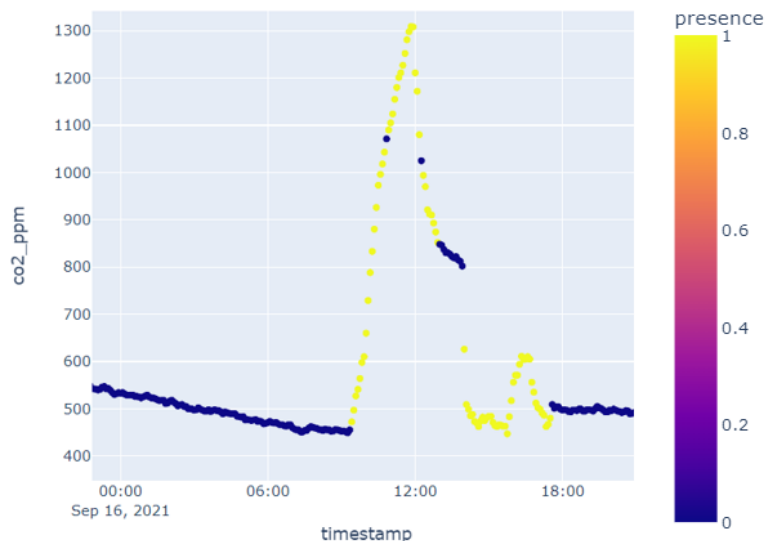


Abbildung 3.11.: CO2 Gehalt der Raumluft über einen Tag

Es ist zu erkennen, dass die CO2-Werte in einem normalen Büroraum innerhalb dieser empfohlenen Grenzen schwanken. Mit einigen kleinen Pausen sind fast durchgehend Personen anwesend, die täglich etwa um 12:00 den Raum lüften. Es ist auch klar zu erkennen, dass sich der CO2-Gehalt während Abwesenheit durch die passive Lüftung des Raumes (Tür-/Fensterspalten) langsam wieder gegen den Grundwert bewegt.

Tabelle 3.1.: Sensordaten

Sensordaten		
Name	Format	Beschreibung
timestamp	timestamp	Zeitpunkt der Messung
co2_ppm	integer	CO2-Wert
temperature_celsius	float	Temperatur in Grad Celsius
relative_humidity_percent	float	Luftfeuchtigkeit
presence	boolean	Aktivität Bewegungssensor

3.3. Luftfeuchtigkeit und Temperatur

Auch wenn Luftfeuchtigkeit und Temperatur Indikatoren für menschliche Präsenz in einem Raum sein können, unterliegen sie in ihrer Nützlichkeit dem CO₂-Gehalt der Raumluft in einem wichtigen Faktor. Sie sich beide stark beeinflusst, von äußeren Umweltfaktoren wie dem aktuellen Wetter oder der Jahreszeit. Beide Werte wurden zunächst in das Training aller Modelle miteinbezogen, es wird allerdings im Folgenden Kapitel gezeigt, dass diese beiden Werte zusammen mit dem CO₂-Werte nicht von besonderem Nutzen sind und deshalb später nicht weiter genutzt werden.

3.4. Sensordaten

Wie bereits beschrieben, wurden die Sensordaten in mehreren Räumen der FH Aachen kontinuierlich gesammelt. Durch das Filtern nach der Beziehung des Raumes ergab sich folgende Datenstruktur als Ausgangslage:

4. Technische Umsetzung

4.1. Datenbeschaffung und -Vorbereitung

Die Daten wurden lokal auf einem der FH-Server in Form eines *Hadoop Distributed File System* (HDFS) gespeichert. Mithilfe von Apache Drill konnten die Daten jederzeit mit einfachen SQL-Abfragen beschafft werden. Die Daten wurden so durch die gesamte Dauer des Projektes immer aktuell gehalten, damit alle Erkenntnisse immer auf der aktuellsten Datenlage basieren.

Die Datenvorbereitung oder Pre-Processing ist eine der wichtigsten Schritte bei der Anwendung von Machine Learning. Durch sie kann man beim Training des Modells durch Bearbeitung bestehender Spalten oder Hinzufügen von zusätzlichen Spalten im Datenset Schwerpunkte setzen, die es den Algorithmen beim Training zum einen erleichtern, ihre Erwartungen zu präzisieren, zum anderen aber auch die Leistung beim Verarbeiten bestimmter Spalten zu steigern.

4.1.1. Gruppierung

Die Sensordaten wurden alle sechs Sekunden erfasst. Da sich weder CO₂-Gehalt noch Feuchtigkeits- oder Temperaturwerte der Raumluft so schnell nicht verändert, wurden die Daten direkt beim Drill per SQL zu zwei-Minuten-Intervallen zusammengefasst. Dabei werden über alle Spalten hinweg Durchschnittswerte gebildet, die dann nachher zu einem Datensatz zusammengefasst werden. Dies steigert die Leistung aller Algorithmen erheblich, da sich die Zahl der Datensätze stark verringert. Da sich, wie oben erwähnt, der CO₂-Gehalt der Raumluft in einem Intervall von zwei Minuten kaum merklich verändert, verringert sich die Genauigkeit des gesamten Datensatz dadurch nicht maßgeblich.

4.1.2. Zyklische Codierung

Zyklische Codierung wird immer dort verwendet, wo Daten sich in wiederholenden Schemata bewegen. Diese Schemata, wie z.B. die Zahlenumbrüche bei einer Uhrzeit, sind für Algorithmen nicht direkt ersichtlich und sind zudem für Computer nicht leicht zu verarbeiten. Durch eine Encodierung in Sinus- und Cosinus-Werte können diese Zusammenhänge vereinfacht werden.

Hierzu wurde der Timestamp zuerst in Sekunden übersetzt, sodass sich ein bestimmter Zeitpunkt eines Tages immer zwischen 0 und 86400 Sekunden bewegt. Aus diesem Wert wurden dann zwei neue Datenspalten „*hour_sin*“ und „*hour_cos*“ in das Datenset eingefügt welche sich durch

$$hour_sin = \sin(2 * \pi * x / x_{max}) \quad (4.1)$$

$$hour_cos = \cos(2 * \pi * x / x_{max}) \quad (4.2)$$

ergeben. So kann jede Tageszeit einer eindeutigen Kombination aus Sinus- und Cosinus-Werten zwischen 0 und 1 zugeordnet werden.

4.1.3. Deltas und Shift-Werte

Desweiteren wurden von den Spalten „*co2_ppm*“, „*temperature_celsius*“ und „*relative_humidity_percent*“, die tatsächlich Rückschlüsse auf die Präsenz zulassen, zusätzliche Delta- und Shift-Spalten angelegt.

Ein Shift-Wert bedeutet lediglich, dass in einer Zeile x_n des Datensets zusätzlich, neben den aktuellen Werten, auch Werte von k Zeilen zuvor, also x_{n-k} stehen. So haben alle Algorithmen direkten Zugriff auf Vergangenheitswerte der ausgewählten Spalten.

Delta-Spalten stellen, dem Namen nach, Deltas zu vorherigen Werten dar:

$$\Delta x_k = x_n - x_{n-k} \quad (4.3)$$

Die Erwartung ist hier, dass die Änderung der CO₂-, Temperatur- und Luftfeuchtheitswerte ein wichtigerer Indikator sein könnte, als die tatsächlichen Werte. In einem schlecht klimatisierten Raum könnten Grundwerte von z.B. CO₂ höher sein, als in anderen Räumen. Durch die hinzugefügten Deltas werden diese Grundwerte ignoriert und Rückschlüsse auf die aktuelle Präsenz sind besser möglich.

Im Zuge der Projektarbeit wurden verschiedene Kombinationsmöglichkeiten von Delta- und Shiftwerten mit Zeitschritten zwischen zwei Minuten und einer Stunde mit Hinblick auf Verbesserungen der Modell-Genauigkeiten getestet.

4.1.4. Outlier Detection

Wie bereits erwähnt, spielen Datenausreißer für die Ergebnisse mancher Algorithmen eine große Rolle. Überall wo z.B. aus einer Reihe von Datenwerten Durchschnittswerte berechnet werden, würden Ausreißer in den Daten das Ergebnis verfälschen und die Leistung des Algorithmus deutlich senken. Um diese Ausreißer vor dem Training der Modelle zu beseitigen wurde das Verfahren des *Interquartilabstands* (IQR nach der englischen Bezeichnung *Interquartile Range*) gewählt.

Der IQR gibt die Intervallgröße an, die ein Wert vom Median einer Datenreihe abweichen darf. Bei einer der Größe nach sortierten Datenreihe $x = (x_0, x_1, \dots, x_n)$ bestimmt man die Mediane der unteren und oberen Hälfte des Datensets Q_1 und Q_2 . Der IQR ergibt sich nun aus

$$IQR = Q_2 - Q_1 \quad (4.4)$$

Mit diesem Wert kann man nun die erlaubten Ober- und Untergrenzen des Datensets mit

$$Limit_{upper} = Q_2 + 1.5 * IQR \quad (4.5)$$

$$Limit_{lower} = Q_1 - 1.5 * IQR \quad (4.6)$$

bestimmen. Alle Werte die außerhalb dieser Grenzen liegen, können als Ausreißer betrachtet werden.

Ausreißer zu entfernen, ist hier wichtig, da die Sensoren Messfehler erzeugen können, oder gelegentlich zur Demonstration von starken Veränderungen in der direkten Nähe des Sensors geatmet wurde und es so in vielen Datensets kurzfristige CO₂-Werte gibt, die natürlich in geschlossenen Räumen nicht vorkommen.

4.2. Validierung

Über das Projekt hinweg wurden verschiedene Validierungsmethoden verwendet. Grundsätzlich unterteilt man das Datenset in ein Trainings- und Testset mit einem Verhältnis von etwa 80-20. Das bedeutet, dass das Modell mit 80 Prozent der Daten trainiert wird, wobei die anderen 20 Prozent zurückgehalten werden, um daraufhin das fertig trainierte Modell daran zu testen. Aus dem Ergebnis dieses Tests ergeben sich verschiedene Werte, die die allgemeine Genauigkeit des Modells repräsentieren, d.h. wie verlässlich das Modell die An- und Abwesenheit des Testsets selbst berechnen kann, wenn man ihm den tatsächlichen Wert vorenthält.

Allgemeine Basis für die errechneten Genauigkeiten bildet die sog. *Confusion Matrix*.

Confusion Matrix		Predicted	
		Negative	Positive
Actual	Negative	True Negative	False Positive
	Positive	False Negative	True Positive

Abbildung 4.1.: Beispiel einer Confusion Matrix

Diese zeigt, inwiefern sich die Genauigkeit des Modells neben tatsächlich richtig erkannten Werten zusätzlich aus false negatives und false positives zusammensetzt.

Die *Accuracy Score* ist eine einfache Methode zur Bestimmung der Modellqualität und berechnet sich aus der Summe der richtigen Ergebnisse geteilt durch die Summe aller Ergebnisse.

$$AccuracyScore = \frac{TP + TN}{TP + TN + FP + FN}$$

Die *Recall Score* beschreibt die Genauigkeit bezogen auf alle erkannten positiven Ergebniswerte.

$$RecallScore = \frac{TP}{TP + FN}$$

Die *Precision Score* gibt die allgemeine Menge, an positiven Werten, die hätten erkannt werden müssen, an.

$$PrecisionScore = \frac{TP}{TP + FP}$$

Die *F1-Score* ist der Durchschnittswert aus Recall- und Precision-Score und gibt im Allgemeinen eine gute Auskunft über die Qualität eines Modells.

$$F1_Score = \frac{2}{\frac{1}{PrecisionScore} + \frac{1}{RecallScore}} = \frac{2 \cdot (PrecisionScore \cdot RecallScore)}{PrecisionScore + RecallScore}$$

In diesem Anwendungsfall ist es durchaus wichtig, Anwesenheitswerte genauer zu betrachten als Abwesenheitswerte, da ein Arbeitstag nur etwa ein Drittel eines tatsächlichen Tages darstellt. Das Datenset ist also mit einem Verhältnis von etwa 8/16 in Richtung der Abwesenheit unausgeglichen. Dies lässt die Erwartung zu, dass das trainierte Modell wesentlich besser darin sein wird, Abwesenheit zu erkennen, als Anwesenheit.

Hierzu wurde bei der Auswertung der Ergebnisse immer auch der *Classification Report* hinzugezogen, aus dem ersichtlich ist, wie genau das Modell beide möglichen Labelwerte berechnen konnte.

Desweiteren war es entscheidend zu erkennen, dass das Datenset über mehrere Monate hinweg nicht perfekt uniform ist, da es in bestimmten Monaten zu verhältnismäßig vielen Urlaubstagen (z.B. Weihnachten/Silvester) und damit einer Häufung an Abwesenheits-Werten kommt. Wenn die Daten nun im o.g. Verhältnis aufgeteilt werden und dann viele Daten des Testsets in solchen Monaten liegen, könnte es beim Ergebnis zu Verzerrungen kommen.

Hierfür wurde eine *Kreuzvalidierung* (engl. Cross Validation CV) implementiert. Bei einer CV wird das Datenset immernoch im gleichen Verhältnis aufgeteilt, allerdings mehrmals, sodass die Trainings- und Testsets jedes mal aus jeweils anderen Daten bestehen.



Abbildung 4.2.: Beispiel einer Kreuzvalidierung

Errechnet man nun die Durchschnittsgenauigkeit aller Iterationen der CV, ergibt sich eine Gesamtgenauigkeit, die das Modell besser repräsentiert, weil es an einer größeren Anzahl an verschiedenen Daten trainiert und getestet wurde.

Da es für manche Räume keine Datensätze mit Labelwerten gab, musste die Validierung dieser Datensätze augenscheinlich erfolgen. Hierfür wurde in ein Datenset eine Labelspalte eingefügt, die dann vom trainierten Modell selbst gefüllt werden sollte. Das Ergebnis musste dann in einem Graphen gezeichnet und per Hand validiert werden. Da sich das Ergebnis der Berechnung in diesem Anwendungsfall, wie oben gezeigt, sehr übersichtlich als Graph darstellen lässt, war diese Methode der Validierung zwar nicht perfekt, lieferte aber trotzdem einen ausreichenden Eindruck über die Qualität des trainierten Modells.

4.2.1. Over- und Underfitting

Over- und Underfitting sind Probleme, die bei allen Machine Learning Algorithmen auftreten können und bezeichnen im Allgemeinen einen falschen Umgang mit dem Datenset. *Overfitting* bedeutet, dass das Modell mit zu vielen, sich zu stark ähnelnden, Daten trainiert wurde, wodurch es nur in der Lage ist, neue Daten richtig zu kategorisieren, die dem Trainingsset besonders ähnlich sind. Overfitting kann während einer Kreuzvalidierung oder Validierung anhand eines anderen Datensets, als dem Trainingsset, erkannt werden.

Um Overfitting zu verhindern, kann das Datenset vereinfacht werden, um so die Beziehungen zwischen den für die Erwartungsberechnung relevanten Spalten zu stärken.

Underfitting bedeutet dagegen, dass es das Modell nicht geschafft hat, zwischen den gegebenen Daten Beziehungen herzustellen, die das Modell unbekannte Daten richtig kategorisieren lassen. Hier reicht es normalerweise aus, dem Modell mehr Daten zur Verfügung zu stellen und in der Datenvorbereitung Felder einzufügen, die dem Modell bestimmte Beziehungen zwischen Datenfeldern hervorheben sollen.

4.2.2. Parameter Tuning

Parameter Tuning beschreibt einen Schritt der Modell-Optimierung, der normalerweise stattfindet, nachdem ein funktionierendes Modell, das mit bereits verarbeiteten Daten gute Ergebnisse liefert, erstellt wurde. Jedes Modell besitzt Parameter, die bei der Erstellung festgesetzt werden. Diese Parameter haben großen Einfluss auf den Trainingsprozess, weshalb es sich anbietet, die Genauigkeiten mehrerer Modelle mit verschiedenen Parameter-Kombinationen zu testen.

Das im Scikit-Learn enthaltene *GridSearchCV*-Modul bietet die Möglichkeit einem Modell eine Vielzahl an verschiedenen Parameter-Optionen zu übergeben. Mit diesen werden automatisch Modelle trainiert und per Kreuzvalidierung verglichen. Am Ende liefert das Modul die Parameter-Kombination, mit der das beste Ergebnis erzielt wurde.

```
param_test = {
    'n_estimators':[25, 50, 100, 250],
    'max_depth': [2, 5, 10, 15],
    'min_samples_split':[25, 150, 500],
    'min_samples_leaf':[50, 75, 125],
    'max_features':[5, 6, 7, 8, 9],
    'subsample':[0.75, 0.9, 0.95]
}
```

Abbildung 4.3.: Beispiel einer Sammlung von Parameter-Optionen

Im obenstehenden Beispiel kann man erkennen, dass für bestimmte Parameter (rot) eine Sammlung an erlaubten Werten (grün) übergeben werden. Alle Kombinationen aus diesen Parameter-Optionen werden dann vom *GridSearchCV*-Modul ausgewertet.

Die Verbesserungen gegenüber Standard-Parametern bewegen sich normalerweise im niedrigen einstelligen Prozent-Bereich und kann auch zu einer leichten Verschlechterung führen, wenn die Standard-Parameter der einzelnen Modelle bereits gut gewählt waren.

5. Anwendung und Ergebnisse

5.1. Feature Vektor

Der Feature Vektor FV beschreibt das Ergebnis des Datensets nach der Vorverarbeitung. Durch die o.g. Schritte ergibt sich eine Vielzahl an möglichen zusätzlichen Features, mit denen man das Datenset versehen könnte. Es sollte zunächst ermittelt werden, welche diese Features wirklich aussagekräftig sind, um den FV nicht unnötig zu überladen. Einige Algorithmen geben nach dem Training Auskunft darüber, in welchem Maß ein Feature in die Berechnung einfließt. Zur Wahl eines FVs der nach dem Training mit einem Modell eine möglichst hohe Genauigkeit liefert, wurden verschiedene Kombinationen durch Verknüpfung von Temperatur, Luftfeuchtigkeit, CO2 mit jeweiligen Delta- und Shiftwerten vorgenommen. Ziel war es, mit einem möglichst kleinen Vektor eine möglichst hohe Genauigkeit zu erreichen, da die meisten Algorithmen eine direkte Abhängigkeit zwischen Dauer der Berechnung und Dimensionalität des FV haben.

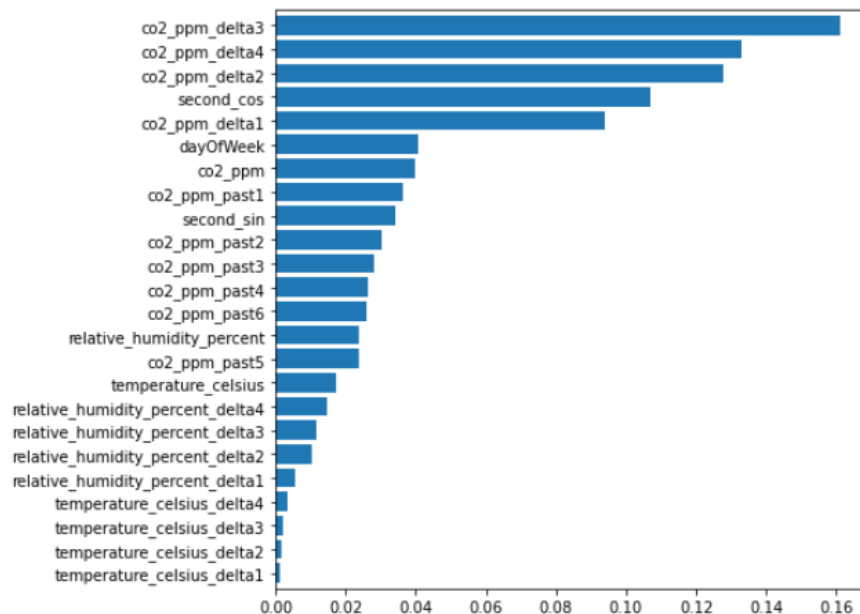


Abbildung 5.1.: Feature Importances eines Random Forest

In der o.g. Grafik sind neben den ursprünglichen Datenwerten nun auch Delta- und Shiftwerte zu sehen. Wie bereits im vorherigen Kapitel angedeutet, sieht man hier klar, dass die Temperatur- und Luftfeuchtigkeitswerte für die Errechnung der Labels kaum zu Rate gezogen werden. Weder die Grund- noch Deltawerte weisen eine hohe Wichtigkeit für das Modell auf.

Nach Exkludierung von Temperatur und Luftfeuchtigkeit wurden neue Deltas und Shifts für die CO₂-Werte eingefügt und die Feature Importances erneut betrachtet. Es ist zu erkennen, dass auch die CO₂-Werte von n Minuten zuvor ebenfalls nicht maßgeblich in die Berechnung einfließen.

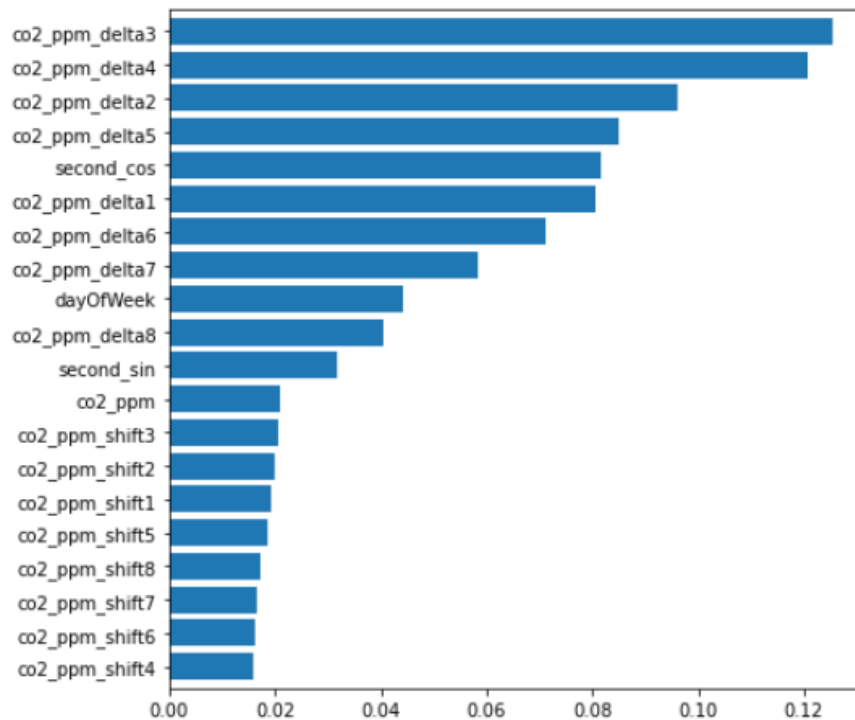


Abbildung 5.2.: Feature Importances eines Random Forest

Der Graph bestätigt die Annahme, dass beim CO₂ die Deltawerte deutlich ausschlaggebender sind, als die tatsächliche Messung zu einem bestimmten Moment.

Wenn diese Werte nur in solch geringem Maß zur Berechnung beitragen, muss ebenfalls die Schlussfolgerung daraus, dass die Genauigkeit auch ohne diese Werte innerhalb einer gewissen Toleranz konstant bleibt, überprüft werden.

Tabelle 5.1.: Vergleich Feature Vektoren

Feature Vektor	Beschreibung
FV1	CO2, Temperatur und Luftfeuchtigkeit mit Shift-Werten
FV2	CO2, Temperatur und Luftfeuchtigkeit ohne Shift-Werte
FV3	Nur CO2 mit Shift-Werten
FV4	Nur CO2 ohne Shift-Werte

Im Folgenden werden beispielhaft vier verschiedene Feature Vektoren beschrieben, aus denen der Vektor, der im weiteren Verlauf des Projektes benutzt werden sollte, abgeleitet wurde.

Mit diesen vier Vektoren wurden nun verschiedene Modelle trainiert und deren Genauigkeiten gegenübergestellt.

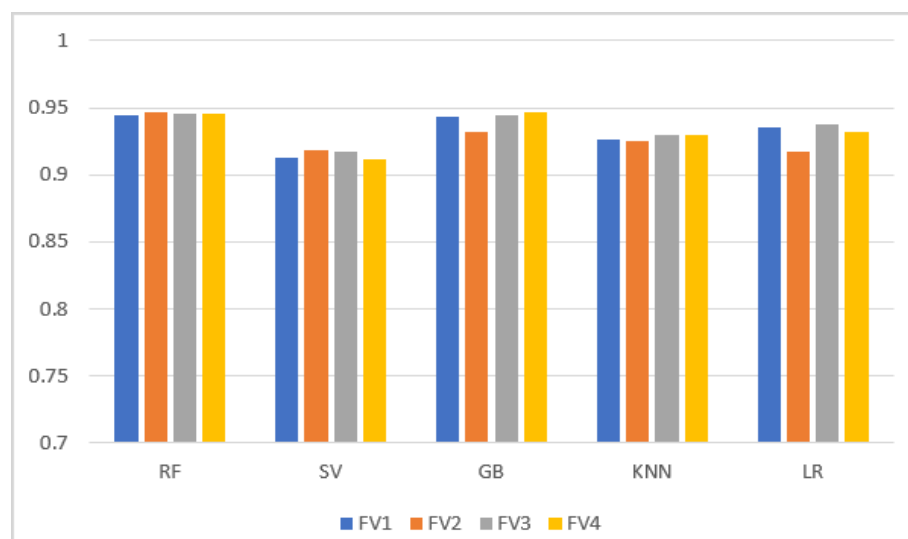


Abbildung 5.3.: Vergleich der Genauigkeiten

Wie deutlich sichtbar ist, sind die Genauigkeiten fast identisch. Diese Gegenüberstellung wurde über das Projekt hinweg mehrfach wiederholt und ergab immer ähnliche Unterschiede.

Dies war Anlass sowohl die Temperatur- und Luftfeuchtigkeitswerte, als auch die Shift-Werte der CO2-Messungen nicht weiter zu betrachten, um den FV für Modelle mit hoher Empfindlichkeit gegenüber Dimensionalität zu optimieren.

Tabelle 5.2.: Feature Vektor

Feld	Beschreibung
second_sin	timestamp-Sinusanteil
second_cos	timestamp-Cosinusanteil
day_of_week	Wochentag der Messung
co2_ppm	CO2-Wert
co2_ppm_delta1	Delta zum CO2-Wert vor 2 Minuten
co2_ppm_delta2	Delta zum CO2-Wert vor 4 Minuten
co2_ppm_delta3	Delta zum CO2-Wert vor 6 Minuten
co2_ppm_delta4	Delta zum CO2-Wert vor 8 Minuten
co2_ppm_delta5	Delta zum CO2-Wert vor 10 Minuten
co2_ppm_delta6	Delta zum CO2-Wert vor 12 Minuten
co2_ppm_delta7	Delta zum CO2-Wert vor 14 Minuten
co2_ppm_delta8	Delta zum CO2-Wert vor 16 Minuten

Aus diesen Betrachtungen ergab sich folgender FV der für den weiteren Verlauf des Projektes genutzt wurde:

5.2. Ergebnisse

Da nicht alle Algorithmen auf die gleiche Weise evaluiert und grafisch dargestellt werden können, sollen die Ergebnisse zu

- Clustering Modellen
- Decision Tree Modellen
- Neuralen Netzwerken

einzelnen betrachtet.

5.2.1. Decision Tree Modelle

Die ausgewählten Algorithmen konnten mit dem ausgewählten FV bei der Erwartungsberechnung über das gesamte Datenset eine Genauigkeit von etwa 94% erreichen. Eine *Confusion Matrix*, welche die errechneten Werte den tatsächlichen Labelwerten gegenüberstellt, zeigt, dass im hier benutzten Beispiel, ein Random Forest Classifier selbst bei der starken Unausgeglichenheit des Datensets ähnlich viele Fehler in sowohl An- als auch Abwesenheit macht.

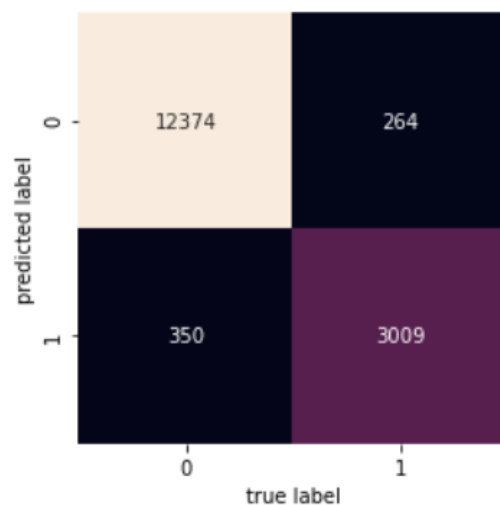


Abbildung 5.4.: Confusion Matrix eines RFC

In den Feldern oben links und unten rechts ist jeweils zu sehen, wann das Modell eine richtige Erwartung für jeweils Ab- und Anwesenheit errechnet hat, während die Felder oben rechts und unten links jeweils die falschen Erwartung für beide Werte zeigen.

Es ist ebenfalls aufschlussreich die erwarteten Labelwerte mit den tatsächlichen Labelwerten grafisch gegenüberzustellen, indem man das ursprüngliche Datenset mit timestamp und CO₂-Wert zeichnet. Dabei wird die gemessene Anwesenheit farblich als gelb für an- und blau für abwesend eingefärbt.

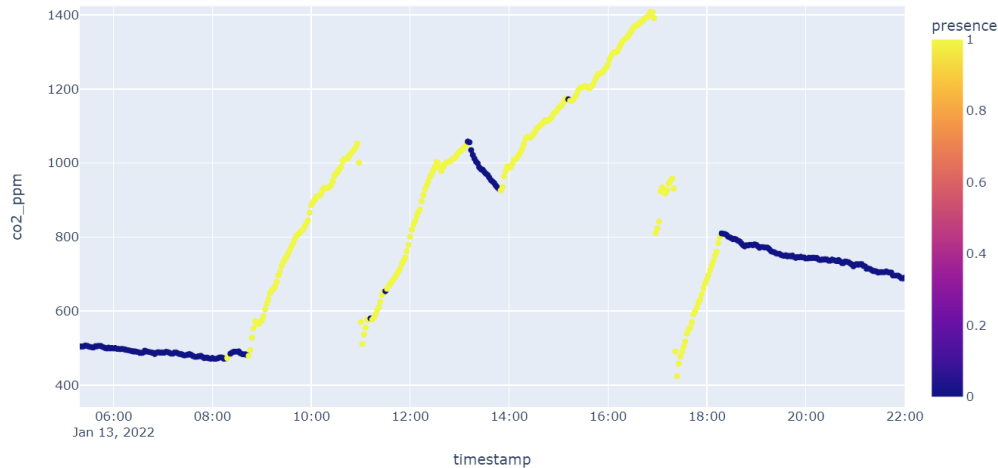


Abbildung 5.5.: Messwerte des 13. Januar

Löscht man aus dieser Datenreihe die tatsächlichen Labelwerte und fügt stattdessen die errechneten Anwesenheitswerte des Modells ein, zeigt sich die hohe Genauigkeit deutlich. Bis auf einige kleine Fehler trifft das Modell eine sehr präzise Aussage über die aktuelle Anwesenheit.

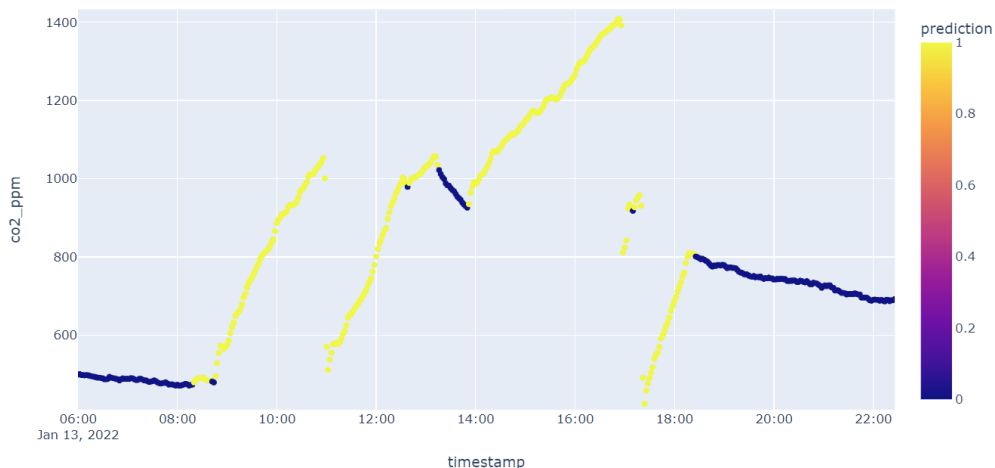


Abbildung 5.6.: Messwerte des 13. Januar

Dieses Ergebnis konnte über eine Vielzahl von Kombinationen verschiedener Räume und Modelle weiterhin bestätigt und repliziert werden. Die Tatsache, dass die Genauigkeiten aller Modelle so nah beieinander liegen, rührt daher, dass dieses Datenset ein typisches Klassifizierungsproblem darstellt, wodurch die sich wiederholenden Schemata in den CO₂-Werten von einer Vielzahl von Algorithmen schnell erkannt und verarbeitet werden können.

Beim Tuning der ausgewählten Modelle wurden als Parameter-Optionen ausschließlich Werte gewählt, die nicht den Standardwerten der Modelle entsprechen. Somit sollte das GridSearchCV-Modul ausschließlich in der Menge von Parametern nach Möglichkeiten suchen, die nicht den Standard-Einstellungen entsprechen.

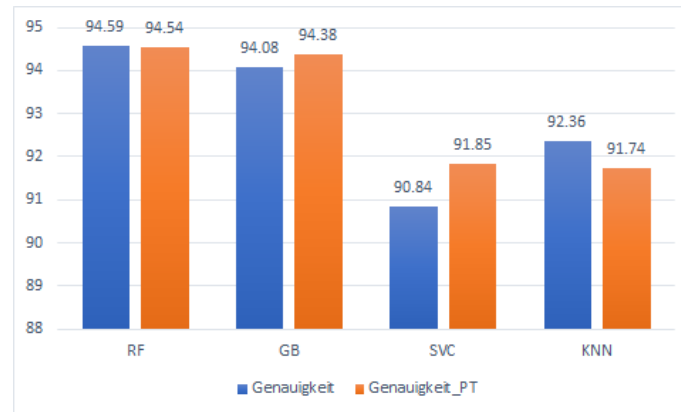


Abbildung 5.7.: !!PLACEHOLDER-GRAFIK!! Ergebnisse des Parameter Tuning

Es ist zu sehen, dass das Parameter Tuning die Genauigkeit der Modelle im erwarteten Rahmen erhöht oder verringert hat. Dies bestätigt die Annahme, dass dieser Anwendungsfall ein typisches Klassifizierungsproblem darstellt, weshalb die Modelle nicht mehr maßgeblich verbessert werden können. Der Schritt des Parameter Tunings sollte trotzdem grundsätzlich immer durchgeführt werden, um zu erkennen, ob man durch eine einfache Änderung der Parameter eine Verbesserung des Modells hervorbringen kann.

Zusätzlich wurden die *ROC-Curves* der True- und False-Positives gegeneinander gezeichnet um zu erkennen inwiefern die Ergebnisse zufällig oder errechnet sind.

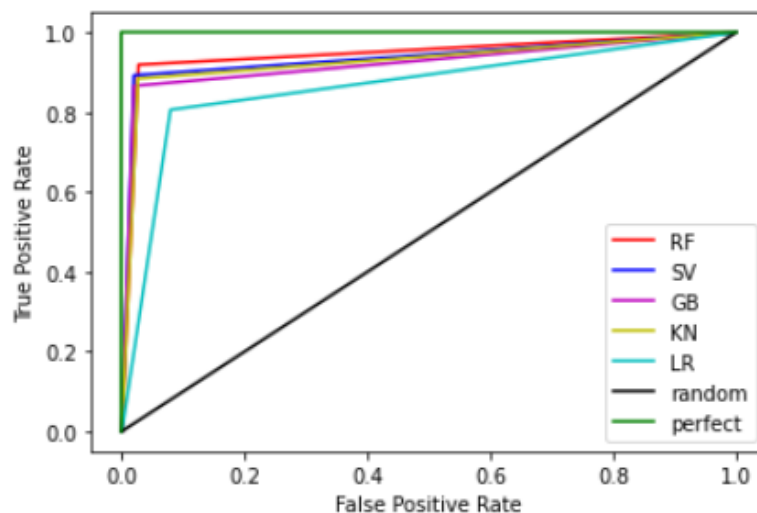


Abbildung 5.8.: ROC-Curves aller Modelle

Man sieht deutlich, wie alle Modelle sehr ähnliche Ergebnisse aufweisen. Keines der Modelle weist Anzeichen auf, dass Ergebnisse wegen Underfitting erraten werden.

5.2.2. Umgang mit fehlenden Präsenz-Labels

Da zur Auswertung auch eine große Menge an Datensätzen ohne Präsenz-Label zur Verfügung standen, wurden die trainierten Modelle zusätzlich auf diese Datensets angewandt und die Ergebnisse so gut wie möglich ausgewertet. Die Ergebnisse sind hier aufgrund des fehlenden Labels lediglich augenscheinlich zu bewerten, da eine mathematische Ermittlung der Genauigkeit nicht möglich ist.

Die Frage, die sich in diesem Kontext stellt, ist die, ob ein auf einen Büroraum trainiertes Modell auch richtige Erwartungen für einen Wohnraum treffen kann. Da das Modell unter anderen auf eine Erkennung der aktuellen Tageszeit trainiert wird, muss geprüft werden, inwiefern ein solches Modell auch Aussagen über CO₂-Werte treffen kann, dessen Veränderungen außerhalb der bisher trainierten Arbeitszeiten von etwa 7 Uhr morgens bis 17 Uhr Nachmittags befinden. Außerdem ist der Wohnraum deutlich größer, als ein Büro und hat in diesem Fall eine größere Anzahl anwesender Personen, wodurch sich die ermittelten Deltas der CO₂-Werte erheblich von den Büroräumen unterscheiden.

Eines der Datensets zeichnete kontinuierliche Werte aus einem Wohnzimmer auf, in dem sich über einen Tag hinweg maximal drei Personen befanden. Die Zeiten, zu denen in diesem Wohnzimmer anhand der CO₂-Werte Anwesenheit zu erkennen war, unterschied sich deutlich von den Datensets der Büroräume der FH-Aachen. Die starken Anstiege des CO₂-Wertes waren bei diesem Datenset viel mehr über den Tag verteilt, wodurch sich die Datenreihe für diesen Test besonders anbot.

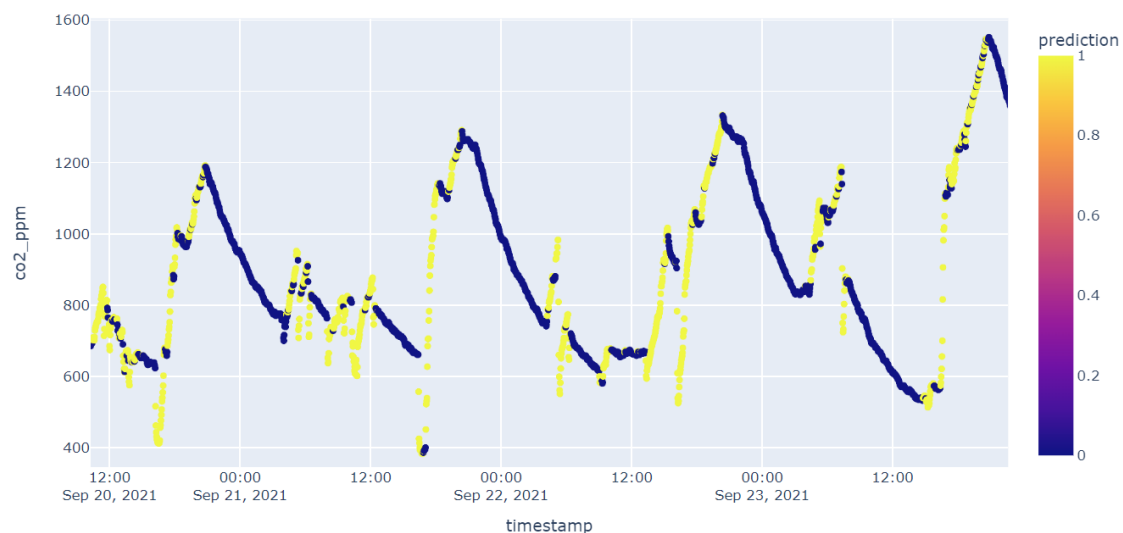


Abbildung 5.9.: Büroraum-Modell trifft Erwartungen für Wohnzimmer

Wie anhand der Einfärbung zu erkennen, ist, kann das Modell auch unabhängig von der Tageszeit und den bekannten Gegebenheiten eines normalen Arbeitstages Anwesenheit erkennen. Vorallem deutliche Anstiege und Abfälle werden deutlich jeweils Geld und Blau eingefärbt, während es auf einigen Strecken zwischen den Maxima und Minima nur gelegentlich zu Fehleinschätzungen kommt.

5.2.3. Neuronale Netzwerke

Die beiden implementierten Neuronale Netze konnten eine ähnliche Leistung, wie die bereits oben genannten Algorithmen erzielen. Beide Modelle wurden über so viele Iterationen (Epochen) trainiert, bis ersichtlich war, dass die Genauigkeit gegen einen Wert konvergiert. Zur Auswertung wurden hier die *Accuracy*- und *Loss*-Funktionen genutzt. Die Genauigkeit beschreibt, wie bei den anderen Modellen das Verhältnis zwischen richtigen und falschen Berechnungen.

Die Loss-Funktion hingegen gibt die Abweichung einer Schätzung zum tatsächlichen Wert an. Diese beiden Werte werden nach jeder Epoche ausgewertet, während zunächst der Erfolg während des Trainings als *train* angegeben wird und danach das Modell an einem zufälligen Teil des Validierungssets *val* getestet wird. Liegen diese beiden Kurven nah beieinander, ist dies ein gutes Anzeichen für ein funktionierendes Modell ohne Over- oder Underfitting, da sowohl beim Training als auch bei der Validierung die Genauigkeiten Ähnlichkeiten aufweisen, während zugleich die Loss-Funktion minimiert wird.

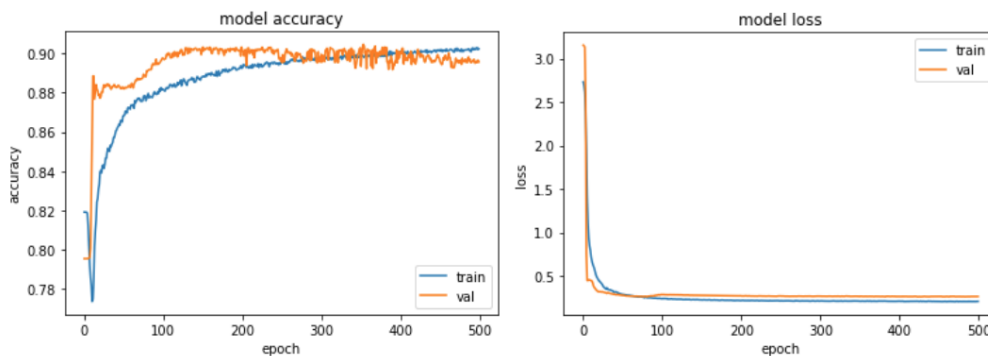


Abbildung 5.10.: Ergebnisse des NN

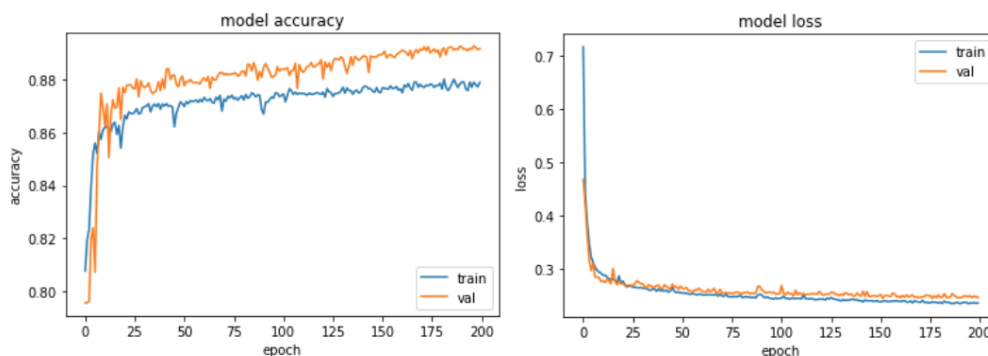


Abbildung 5.11.: Ergebnisse des LSTM

Wie zu sehen ist, weisen beide Arten neuronaler Netzwerke diese Eigenschaften auf. Die Genauigkeiten beider Modelle konvergierten auch nach verschiedenen Epochen-Werten etwa bei 89%.

Tabelle 5.3.: Genauigkeitsauswertung

Präsenz	Precision	Recall	F1
0	0.94	0.79	0.86
1	0.50	0.82	0.62

5.2.4. Clustering Modelle

Clustering Modelle hatten von allen benutzten Algorithmen die größten Schwierigkeiten bei der Klassifikation des Datensets. Während die Genauigkeit bei der Erwartungsbeurteilung von Abwesenheit bei 94% lag, konnte das K-Means-Modell Anwesenheiten nur zu 50% erkennen, was bedeutet, dass das Modell das Ergebnis errät, anstatt ihn zuverlässig zu berechnen.

Wegen der Unausgeglichenheit des Datensets sind ca. zwei Drittel des Datensets sehr einfach zu klassifizieren, da es im ganzen Datenset nie Anwesenheiten zwischen beispielsweise 20:00 Uhr und 06:00 Uhr gibt. Selbst wenn das Modell beim letzten Drittel, also den Anwesenheiten rät und somit nur zu 50% trifft, suggeriert eine daraus resultierende Genauigkeit von etwa 80% ein funktionierendes Modell, allerdings zeigt sich hier, warum auch die Gegenüberstellung von True und False Positives in der Modellauswertung von entscheidender Wichtigkeit sein kann.

Wie anhand der Precision-Score zu sehen ist, tritt hier genau der oben beschriebene Fall ein.

Zur Veranschaulichung der Ergebnisse hilft eine Dimensionalitätsreduktion, bei der man das Datenset auf zwei Dimensionen reduziert und die An- und Abwesenheiten in einem Graphen darstellt.

Die Menge der Anwesenheiten (Orange) ist nahezu identisch mit der Teilmenge der Abwesenheiten (Blau). Das heißt, dass alle Datenpunkte die in der linken Hälfte des blauen Datensets liegen, sehr einfach klassifiziert werden können.

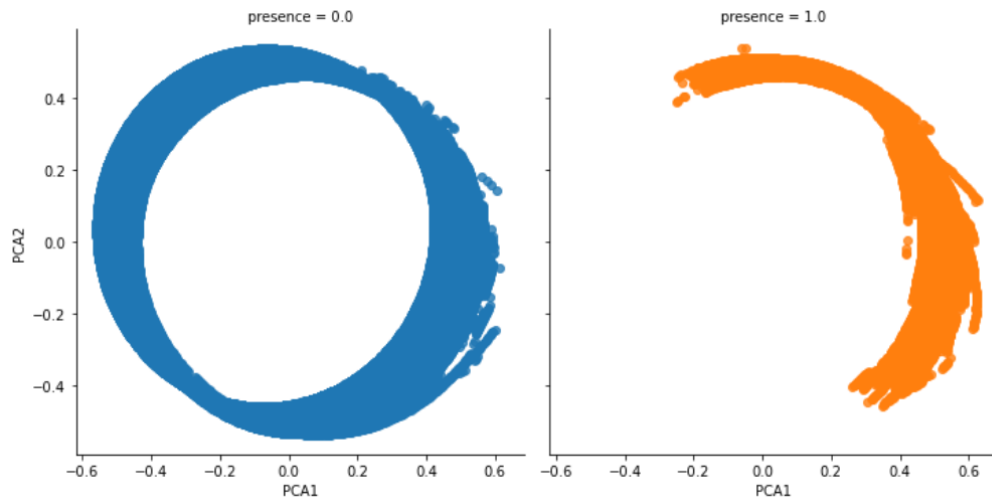


Abbildung 5.12.: Reduktion des Datensets auf zwei Dimensionen

Legt man beide Datensets übereinander ist es nun nahezu unmöglich zwischen An- und Abwesenheit zu unterscheiden.

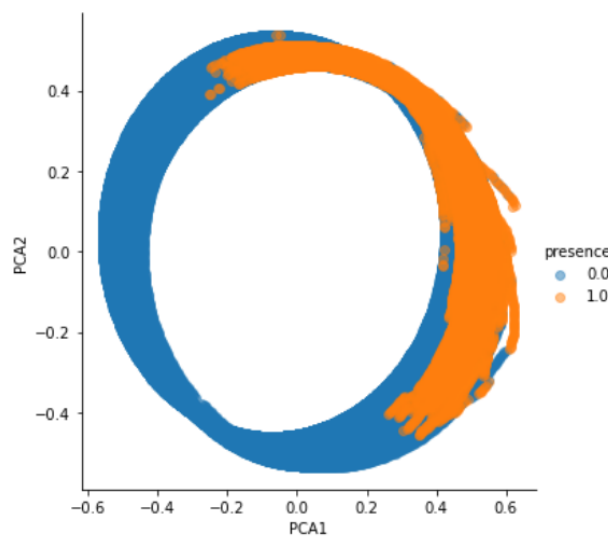


Abbildung 5.13.: Datenreduktion ohne Trennung

Das Erkennen von Abwesenheit im linken Teilbild ist nun erheblich einfacher, als die Trennung zwischen An- und Abwesenheit im rechten Teilbild.

6. Zusammenfassung und Ausblick

6.1. CO₂ als Anwesenheitsindikator

Die Genauigkeit aller Modelle bei der Präsenzerkennung während der gesamten Durchführungsphase des Projektes lag zwischen 88% bei neuronalen Netzwerken und 95% bei den Klassifizierungsalgorithmen. Die Beziehung zwischen der aktuellen Tageszeit und dem CO₂-Gehalt der Luft, konnte erfolgreich bestätigt werden. Dieses Ergebnis wird ebenfalls gestützt von der vorangegangenen Forschung¹. .

Weiterhin konnte gezeigt werden, dass alle Modelle zusätzlich in der Lage sind, menschliche Präsenz außerhalb der trainierten Uhrzeiten korrekt zu erkennen, wenn die Modelle mit Deltas zu Vergangenheitswerten des CO₂-Gehalts trainiert werden. Die Präsenzerkennung ist also vollständig unabhängig von den Umständen des Trainingssets und kann sowohl für andere Räume, als auch zu anderen Tageszeiten effektiv genutzt werden. CO₂ ist somit als geeigneter Indikator für menschliche Präsenz in Innenräumen anzusehen.

Die Zweckmäßigkeit von CO₂-Sensoren für Gebäudeautomatisierungssysteme konnte als kostengünstige und leicht zu implementierende Option gezeigt werden.

6.2. Mögliche Verbesserungen

Die Genauigkeit von Machine Learning Modellen kann allgemein durch die Hinzugabe von mehr Daten in ein Modell verbessert werden. Da der Datenbestand, auf die sich diese Arbeit stützt, über die Bearbeitungszeit des Projektes kontinuierlich vergrößert wurde, wuchs damit auch die Anzahl an Anhaltspunkten für Zusammenhänge zwischen den Datenfeldern stetig an. Besonders nützlich wäre das Sammeln von Messdaten über ein oder mehrere Jahre hinweg. Somit hätten die Modelle eine Chance Gewohnheiten in z.B. Urlaubstagen, Feiertage oder auch wiederkehrende Meetings korrekt einzuordnen, sodass Gebäudeautomatisierungssysteme ein Maximum an Komfort und Energieeffizienz herstellen können.

Da Infrarotsensoren inhärent nicht in der Lage sind kontinuierliche Präsenz zu erkennen, wenn sich die Personen im Raum nicht bewegen, kam es innerhalb des Datensets immer wieder zu kleinen Messfehlern, die durch die Gruppierung und Durchschnittsberechnung der Präsenzwerte behoben werden mussten. Es ist klar davon auszugehen, dass eine genauere Datenlage auch die Qualität der trainierten Modelle steigern würde.

Zusätzlich existierten im Datenset durch seltene Probleme mit der Hardware längere Datenreihen mit falschen Labelwerten, welche während der Bearbeitung ausgeschlossen werden mussten. Da sich dadurch die für das Training geeignete Datenmenge verringerte, ist auch hier davon auszugehen, dass das Aussortieren dieser Werte mit einer Verringerung der Modellqualität einherging.

¹[Hanfstaengl et al., 2019]

Quellenverzeichnis

- [Hanfstaengl et al., 2019] Hanfstaengl, L., Parzinger, M., Wirnsberger, M., Spindler, U., and Wellisch, U. (2019). Identifying the presence of people in a room based on machine learning techniques using data of room control systems. Conference.
- [Umweltbundesamt, 2020] Umweltbundesamt (2020). Endenergieverbrauch der privaten Haushalte. Online. <https://www.umweltbundesamt.de/daten/private-haushalte-konsum/wohnen/energieverbrauch-privater-haushalte/>.
- [Umweltbundesamt, 2022] Umweltbundesamt (2022). Energieverbrauch nach Energieträgern und Sektoren. Online. <https://www.umweltbundesamt.de/daten/energie/energieverbrauch-nach-energetraegern-sektoren>.

Abkürzungsverzeichnis

g	Gravitation in Nähe der Erdoberfläche
Nu	Nußelt-Zahl
ν_{Luft}	Kinematische Viskosität von Luft
Pr	Prandtl-Zahl
\dot{Q}	Wärmestrom
Ra	Rayleigh-Zahl
ρ_{Luft}	Dichte von Luft
T	Temperatur
T_{∞}	Umgebungstemperatur

Abbildungsverzeichnis

3.1. Beispiel für Clustering	9
3.2. Beispiel eines Decision Trees	10
3.3. Beispiel eines Support Vector Classifiers	11
3.4. Beispiel eines Support Vector Classifiers in der zweiten Dimension.....	11
3.5. Beispiel logistischer Regression	12
3.6. Logistische Regression mit deutlicher Skalierung	13
3.7. Beispiel eines KNN	14
3.8. Beispiel eines neuronalen Netzwerks	15
3.9. Beispiel eines rekurrenten neuronalen Netzwerks.....	16
3.10. Vereinfachter Aufbau einer LSTM-Zelle	16
3.11. CO2 Gehalt der Raumluft über einen Tag	17
4.1. Beispiel einer Confusion Matrix	21
4.2. Beispiel einer Kreuzvalidierung	23
4.3. Beispiel einer Sammlung von Parameter-Optionen	24
5.1. Feature Importances eines Random Forest	25
5.2. Feature Importances eines Random Forest	26
5.3. Vergleich der Genauigkeiten	27
5.4. Confusion Matrix eines RFC	29
5.5. Messwerte des 13. Januar.....	30
5.6. Messwerte des 13. Januar.....	30
5.7. !!PLACEHOLDER-GRAFIK!! Ergebnisse des Parameter Tuning.....	31
5.8. ROC-Curves aller Modelle	31
5.9. Büroraum-Modell trifft Erwartungen für Wohnzimmer.....	32
5.10. Ergebnisse des NN	33
5.11. Ergebnisse des LSTM	33
5.12. Reduktion des Datensets auf zwei Dimensionen	35
5.13. Datenreduktion ohne Trennung	35

Tabellenverzeichnis

3.1. Sensordaten	18
5.1. Vergleich Feature Vektoren.....	27
5.2. Feature Vektor.....	28
5.3. Genauigkeitsauswertung.....	34

A. Quellcode

Hier steht der Quellcode

B. Rohdatenvisualisierungen

1. Graustufen
2. Verteilungen