



Duale Hochschule Baden-Württemberg Mannheim

## **Projektreport**

# **Bundesliga Match Predictions**

## **Studiengang Wirtschaftsinformatik**

**Studienrichtung Data Science**

Verfasser:	Max Bernauer, Philipp Dingfelder, Julius Könning
Matrikelnummer:	5763624, 8687786, 7305370
Firma:	SAP SE, Schaeffler
Modul:	Data Exploration
Kurs:	WWI20DSB
Dozent:	Simon Poll
Bearbeitungszeitraum:	Sommersemester 2022

# Ehrenwörtliche Erklärung

Ich versichere hiermit, dass ich die vorliegende Arbeit mit dem Titel “*Bundesliga Match Predictions*” selbstständig verfasst und keine anderen als die angegebenen Quellen und Hilfsmittel benutzt habe. Ich versichere zudem, dass die eingereichte elektronische Fassung mit der gedruckten Fassung übereinstimmt.

Ort, Datum

Max Bernauer, Philipp Dingfelder, Julius Könning

# Inhaltsverzeichnis

<b>Abbildungsverzeichnis</b>	<b>iii</b>
<b>Tabellenverzeichnis</b>	<b>iv</b>
<b>Abkürzungsverzeichnis</b>	<b>v</b>
<b>1 Einleitung</b>	<b>1</b>
1.1 Hintergrund und Motivation . . . . .	1
1.2 Business Use Case . . . . .	1
<b>2 Theoretische Grundlagen</b>	<b>2</b>
2.1 Related Work . . . . .	2
2.2 Verwendete Technologien und Bibliotheken . . . . .	2
<b>3 Praktischer Teil</b>	<b>3</b>
3.1 Umsetzung . . . . .	3
3.2 Ergebnisse . . . . .	4
3.3 Use Case Validierung durch simulierte Wetten . . . . .	4
<b>4 Fazit</b>	<b>5</b>
<b>Anhang</b>	

# Abbildungsverzeichnis

# Tabellenverzeichnis

# Abkürzungsverzeichnis

xxx	Description
-----	-------------

# 1 Einleitung

## 1.1 Hintergrund und Motivation

Dieser Projektreport ist im Rahmen des Fachs Data Exploration entstanden. Das Ziel des Moduls ist die „Anwendung von Methoden und Verfahren des maschinellen Lernens auf eine vorgegebene Datenbasis unter Laborbedingungen“[Modulhandbuch]. Zusätzlich soll neben der informatischen Betrachtung auch der betriebswirtschaftliche Nutzen erörtert werden [vgl. Modulhandbuch].

Auf Basis dieser Vorgaben wurde das Thema des Projekts gesucht. Dabei ging es primär darum ein Themengebiet zu finden, welches sowohl breite Möglichkeiten für die informatische als auch die betriebswirtschaftliche Betrachtung bietet. Aufgrund der Interessen innerhalb der Gruppe wurde sich für das Thema **Bundesliga Match Predictions** entschieden. Wir wollten der Fragestellung auf den Grund gehen, ob es tatsächlich möglich ist, diese unzählig erscheinenden Faktoren des Fußballspiels durch Data Science-Prozesse für eine Vorhersage nutzen zu können.

## 1.2 Business Use Case

Wie bereits erwähnt spielt die wirtschaftliche Betrachtung dieses Projekts neben der informatischen Arbeit eine primäre Rolle. Das Ziel einer Bundesliga Match Prediction liegt hier auf der Hand. Ist es tatsächlich möglich Anbieter wie Tipico und Bwin durch ein mathematisches Modell zu schlagen?

Durch eine solche Vorhersage kann man potenziell starken Profit bei etwaigen Wettanbietern erzielen. Andersrum kann man natürlich auch diese Software an Wettanbieter verkaufen, damit diese ihre Quoten noch effizienter und genauer berechnen können.

Die Herangehensweise an dieses Projekt beginnt mit der richtigen Datenbasis. Durch diese kann man algorithmisch ein Modell erstellen, dass den genannten Business Use Case ermöglichen kann.

## 2 Theoretische Grundlagen

### 2.1 Related Work

- quellen/inspirationen"nennen und auf uns beziehen - wie gehen diese ihr problem an (quellen, algorithmen, ...)

### 2.2 Verwendete Technologien und Bibliotheken

Für die Datenvorbereitung werden klassische Python-Module wie Pandas, Sklearn und Matplotlib genutzt.

Bei der Modellimplementierung wurde vorerst ein Decision Tree Classifier genutzt. Dieser stammt aus dem Sklearn-Modul tree. Nach Ausprobieren dieses Classifiers wurde sich für das Ensemble-Learning entschieden. Dabei nutzt man mehrere Algorithmen um sie zeitgleich zu vergleichen und potenziell zu vereinen. Die Algorithmen wie beispielsweise Logistic Regression oder Gaussian Naive Bayes stammen auch von Sklearn. Final wurde ein Neural Network implementiert. Dazu wurde Keras genutzt. Aus Keras wurde außerdem eine Random Parameter Search zur weiteren Optimierung angewandt. Durch Matplotlib wurden die Ergebnisse visualisiert. Außerdem wurde eine Poisson-Verteilung genutzt um einen mathematischen Ansatz zu implementieren.



# 3 Praktischer Teil

## 3.1 Umsetzung

Eine der wohl wichtigsten Entscheidungen bei der Implementation eines Machine Learning Modells ist die Auswahl der Rohdaten. Nach intensiver Recherche wurde sich für einen Datensatz entschieden welcher Informationen über alle Bundesliga-Spiele seit der Saison 2005/2006 enthält. Dieser Rohdatensatz besitzt 65+ Features. Während des Preprocessings wurde die Anzahl der Features durch das Löschen unvollständiger oder für das Modell unnötiger Features, zum Beispiel die Odds verschiedener Wettanbieter, deutlich reduziert. Diese Features wurden dann harmonisiert und mit einer Korrelationsmatrix auf ihre Aussagekraft geprüft. Da diese nach der ersten Featureauswahl noch nicht zufriedenstellend war wurden weitere Features berechnet und in den Datensatz aufgenommen. So wurden beispielsweise die aktuelle Form der Mannschaften (Elo) sowie Angriffs- und Verteidigungswerte errechnet. Um die Aussagekraft weiter zu erhöhen wurden außerdem Features eingeführt, die einen direkteren Vergleich der Mannschaften ermöglichen. Diese Features sind zum Beispiel Differenzen der Elo und Angriffs- sowie Verteidigungswerte oder die Anzahl der Punkte und Tore, die eine Mannschaft während der letzten Spiele erzielt hat.

(Hier Grafik von Datensatz)

Dadurch, dass die Zielvariablen bereits im Datensatz vorliegen handelt es sich um ein Supervised Learning-Problem. Dieses wurde wie im Theorieteil erwähnt mit drei verschiedenen Methoden angegangen. Dem Ensemble Learning, einem Neural Network sowie einem mathematischen Ansatz.

Im Preprocessing wurden einige unvollständige Spieldaten gelöscht und viele unnötige Features - wie die Betting Odds der Anbieter. Nach Harmonisierung der Daten wurden die Features auf die Anzahl 21 reduziert. (Hier Grafik von Datensatz). Dadurch, dass die Zielvariablen im Datensatz bereits vorliegen handelt es sich um ein Supervised Learning-Problem. Dieses wurde wie bereits im Theorieteil erwähnt mit drei verschiedenen Methoden angegangen. Dem Ensemble Learning, einem Neural Network und einem mathematischen Ansatz.

## 3.2 Ergebnisse

Anfangen beim Ensemble Learning kann man verschiedene Accuracys je Algorithmus erkennen (Grafik zeigen). Hierbei sind die Logistic Regression und Linear Discriminant Analysis am Besten. Das Stacking Modell welches die anderen Modelle vereint schneidet auch gut ab. Alle drei liegen bei einer Accuracy von etwas unter 50 Prozent.

Bei Betrachtung des Neural Networks fällt auf, dass die Accuracy schon nach wenigen Epochen auf etwas unter 50 Prozent konvergiert (Grafik). Nach einer Random Parameter Search konnte diese auf 53 Prozent angehoben werden. Zuletzt haben wir mit der Poisson-Verteilung des mathematischen Ansatzes auch eine Accuracy von etwa 50 Prozent erreicht (Grafik).

## 3.3 Use Case Validierung durch simulierte Wetten

Um zu testen, wie gut das Modell unter Realbedingungen funktioniert wurden im Verlauf der Saison 2021/22 Wetten simuliert. Für die Simulation wurden die Predictions des Modells mit den Predictions von vier verschiedenen Wettanbietern verglichen. Um die Predictions der Wettanbieter zu erhalten wurde zuerst der Durchschnitt der Odds der Wettanbieter für jedes Ereignis (Sieg Heimteam, Unentschieden, Sieg Auswärtsteam) ermittelt. Dieser Durchschnitt wurde dann in eine Prozentzahl umgewandelt und nach dem Maximum gefiltert um eine Prediction für jeden Spieltag zu erhalten. Um auf die Predictions zu setzen mussten diese einen gewissen Schwellwert (XX% Sicherheit) überschreiten. Bei Überschreitung des Schwellwerts wurde je nach Höhe der % Zahl ein bestimmter Betrag gewettet. (Grafik wie viel ab welcher % Zahl gewettet wird)

Auf die gesamte Saison 2021/22 betrachtet ergibt sich folgende Grafik. (Grafik mit den Wetten) Es lässt sich erkennen, dass das Modell am Ende der Saison zwar mehr Gewinn gemacht hat, die Wettanbieter aber im Laufe der Saison besser performen. Auf die Bedeutung der Grafik in Bezug auf den Business Use Case wird im Fazit genauer eingegangen.

## 4 Fazit

Abschließend lässt sich feststellen, dass sich die Accuracy des Neural Networks im Verlauf des Projekts Schritt für Schritt durch das Optimieren der Datenbasis und des Modells selber erhöhen ließ. Die größte Erhöhung dabei kam durch das Fokussieren auf den direkten Unterschied der Werte der Mannschaften sowie die Random Parameter Search des Modells. Final konnte eine 53% Accuracy erreicht werden.

Wie in der Grafik in Abschnitt drei von Kapitel 3 gezeigt war es möglich mit Hilfe des Modells profitable zu wetten. Vergleicht man die Profit Kurve des Modells jedoch mit der der Wettanbieter fällt auf, dass diese im Verlauf der Saison Intervalle mit höheren Gewinnen aufweisen, ihre Predictions also besser funktionieren. Das lässt sich darauf zurück führen, dass das Modell mit maximal 65% Sicherheit das Ergebnis eines Spiels vorhersagt. Die Modelle der Wettanbieter weisen eine durchschnittlich höhere Sicherheit in Bezug auf die Ergebnisse der Spiele aus, was dazu führt, dass diese in der Simulation höhere Beträge einsetzen können und insgesamt auch mehr wetten eingehen als unser Modell. Um also langfristig die Wettanbieter gewinnbringend zu schlagen müsste das Modell weiterverbessert werden um eine höhere Accuracy zu erhalten. Mögliche Ansätze hierfür sind:

Zusammenfassend lassen sich die Ergebnisse des Projekts folgendermaßen zusammenfassen: Das Projekt war erfolgreich. Durch die 53 prozentige Accuracy sowie das halbwegs profitable Wetten konnte gezeigt werden, dass sich Fußballergebnisse zumindest in Ansätzen algorithmisch mit Hilfe von Modellen beschreiben lassen. Damit ist der Use Case zum Großteil erfüllt. Es ist jedoch kritisch anzumerken, dass auch sehr viele nicht genau spezifizierbare Zufallsvariablen einen Einfluss auf den Ausgang eines Spiels haben, sodass sich nie mit 100%er Wahrscheinlichkeit der genaue Ausgang eines Fußballspiels vorhersagen lässt.

Die Modell Accuracy von 53% ist für sich genommen relativ vielversprechend. Sie zeigt, dass sich Fußballergebnisse zumindest in Ansätzen algorithmisch bestimmen lassen. Hierbei ist jedoch kritisch anzumerken, dass auch sehr viele nicht genau spezifizierbare Zufallsvariablen einen Einfluss auf den Ausgang eines Spiels haben.

Abschließend lässt sich feststellen, dass man durch die Optimierung des Neural Networks die höchste Accuracy erreicht. Nach der Random Parameter Search konnten wir die Accuracy erhöhen. Wie bereits erwähnt lag diese final bei 53 Prozent.

Mithilfe dieses Modells wurde das Betting bei Wettanbietern simuliert. Leider kann man sehen, dass unser Modell schnell negative Summen erreicht und einiges an Geld verliert. Der Fußball ist wohl doch ein Gebiet mit zu vielen Zufallsvariablen, sodass man keine Genauigkeit erzielen kann, die einen Wettanbieter schlägt.

Die Genauigkeit von 53 Prozent ist jedoch für sich genommen recht vielversprechend. Es zeigt, dass sich der Fußball in Ansätzen algorithmisch bestimmen lässt. Für unseren Business Use Case ist dieses Ergebnis natürlich nicht hoch genug und somit muss dieser verworfen werden. Es hat sich gezeigt, dass man schnell Geld verliert bei der Nutzung unseres Modells.