

Article

Not peer-reviewed version

The Use of Machine Learning in Predicting Formula 1 Race Outcomes

[Atharva Urdhwareshe](#) *

Posted Date: 18 April 2025

doi: 10.20944/preprints202504.1471.v1

Keywords: Formula 1; Race Outcome Prediction; Machine Learning; TabNet; Sports Analytics; Predictive Modeling; Constructor Points Prediction; Driver Performance; Auto Racing; Deep Learning for Tabular Data



Preprints.org is a free multidisciplinary platform providing preprint service that is dedicated to making early versions of research outputs permanently available and citable. Preprints posted at Preprints.org appear in Web of Science, Crossref, Google Scholar, Scilit, Europe PMC.

Copyright: This open access article is published under a Creative Commons CC BY 4.0 license, which permit the free download, distribution, and reuse, provided that the author and preprint are cited in any reuse.

Article

The Use of Machine Learning in Predicting Formula 1 Race Outcomes

Atharva Urdhwaresh

Independent Researcher; uatharva12@gmail.com

Abstract: Formula 1 is a sport driven by both engineering excellence and data precision. With an ever-growing wealth of historical and real-time race data, machine learning offers an opportunity to transform race outcome prediction and strategy development. Prior literature in motorsport analytics highlights the use of classification and regression models, yet few studies leverage deep learning models specifically built for structured data. This study aimed to develop a robust machine learning pipeline to predict both driver finishing positions and constructor championship points using historical F1 data from 2010 to 2023. TabNet, a deep learning architecture optimized for tabular data, was selected for its interpretability and strong feature selection capability. The models were trained using pre-race variables such as grid position, number of laps, constructor, and overtakes. Hyperparameter tuning was conducted using Optuna. The results showed strong predictive performance, with the driver model outperforming the constructor model in overall accuracy. These findings demonstrate the practical potential of machine learning in high-stakes motorsport environments. The models developed in this research could support teams in strategic planning, broadcasters in providing predictive insights, and analysts in simulating race outcomes under different conditions.

Keywords: Formula 1; Race Outcome Prediction; Machine Learning; TabNet; Sports Analytics; Predictive Modeling; Constructor Points Prediction; Driver Performance; Auto Racing; Deep Learning for Tabular Data

1. Introduction

Formula 1 racing is not only a sport but a data-driven engineering battlefield. Teams collect millions of data points during a race to fine-tune strategies. However, the power of machine learning to make forward-looking predictions in this high-stakes environment remains underutilized. This project addresses the challenge of predicting race outcomes—both driver finishing positions and team points—using real historical data and modern machine learning techniques. Predictive modeling in motorsports has gained traction in recent years, with applications ranging from performance forecasting to strategic planning (Doe, 2019; H. R. Thornton & Duthie, 2017; Jackson, 2017; Smith, 2020). While statistical models have been used to analyze past performance (Jenkins, 2017), machine learning offers a more dynamic alternative that handles nonlinearities and high-dimensional data effectively (Chang et al., 2019; Nguyen, 2018).

I implemented and evaluated the models, using a curated dataset that includes driver, constructor, and race-level features. The outcome is a robust two-model system that provides predictions for both drivers and constructors with strong interpretability.

2. Literature Review

2.1. General Studies on Machine Learning in Sports Analytics

Machine learning has been previously applied to sports prediction (Doe, 2019; Smith, 2020). They illustrate how machine learning techniques have evolved beyond traditional statistical methods. Doe's research highlighted that neural networks and support vector machines are more efficient in analyzing dynamic sports environments compared to standard regression techniques. This has allowed teams to

gain deeper insights into patterns that influence outcomes in fast-paced sports like F1, where hundreds of variables—from tire wear to fuel efficiency—play a role in the race's final result.

Smith's (2020) contribution was pivotal in showing that Formula 1 requires a unique approach, given the sheer number of real-time variables that can affect a race outcome. Unlike team sports, where variables are often static for longer periods, F1 sees rapid, frequent changes in tire conditions, track temperature, and driver fatigue. Smith argues that only machine learning models capable of integrating both static historical data and real-time telemetry can provide the adaptability necessary for F1 teams to make split-second decisions. This observation points to an important gap in the research: most machine learning models rely too heavily on pre-race simulations and historical datasets, making them less effective in the heat of a live race.

Furthermore, Baker's (2017) findings supported the use of real-time telemetry data to enhance predictive capabilities. By integrating historical data with live data streams, teams can optimize pit-stop strategies and tire management dynamically. However, Baker also noted the limitations imposed by computational speed; machine learning models must process these large datasets fast enough to provide actionable insights during the race.

2.2. *The Evolution of Machine Learning Models in Sports*

The evolution of machine learning models has transformed how F1 teams utilize data. Early models, such as linear regression, focused primarily on analyzing static data to identify basic correlations, but these models could not adapt to changes in real-time. The introduction of more complex models like deep learning, reinforcement learning, and neural networks has expanded the possibilities. Today, deep learning algorithms are used to identify intricate patterns that traditional models might miss. The ability to handle non-linear relationships between variables has become essential in a sport as dynamic as Formula 1.

Rodriguez (2019) highlights that convolutional neural networks (CNNs), a deep learning model, excel at processing telemetry data, such as engine performance and tire wear. This is a significant step forward from earlier models that could only make use of historical race data. However, as Brown (2021) pointed out in his comparative study, even advanced models like CNNs struggle to adapt quickly enough to sudden changes in race conditions. This is where reinforcement learning, which enables a model to learn from its environment in real-time, holds great promise. Reinforcement learning models can monitor factors such as tire degradation in real-time and make adjustments during the race, something traditional models have not been able to do effectively.

2.3. *Ethical Considerations in Machine Learning for Sports*

The rise of machine learning in Formula 1 brings with it several ethical considerations, particularly around data privacy and the role of human decision-making in sports. The increasing reliance on machine learning raises the question of whether race strategies could become too dependent on AI, potentially diminishing the role of human intuition and strategy. While these models provide teams with an edge, ethical concerns arise regarding how data is collected, processed, and used. Does the use of telemetry and personal driver data violate any ethical boundaries? Ford (2019) notes that while data-driven models have improved race strategies, there is a need for a regulatory framework to ensure that the use of personal and performance data is ethical.

Another ethical issue is the potential over-reliance on machine learning, which could lead to a situation where teams prioritize algorithmic decisions over human judgment. In sports, where unpredictability and human intuition play key roles, it's important to balance the insights provided by machine learning with the instincts of race engineers and drivers.

2.4. *Comparative Studies with Other Sports Using Machine Learning*

Formula 1 is not the only sport that has embraced machine learning. Other high-performance sports, such as basketball and soccer, have also integrated machine learning into their strategic planning. Garcia (2020) conducted a study comparing the use of machine learning in F1 with that in football and

basketball. While these sports share some similarities in the way they use data, the rapid, high-stakes environment of F1 presents unique challenges that are not as prevalent in other sports. In soccer, for instance, data can be analyzed over a longer period during the match without significantly impacting strategy. However, in F1, decisions such as pit-stops or tire changes must be made in real-time, often with only seconds to spare.

This comparison underscores the complexity of F1 as a sport. While other sports benefit from real-time data, the level of unpredictability in F1—due to rapidly changing weather conditions, track temperatures, and car mechanics—makes it much harder to apply the same machine learning models. This section serves to highlight the unique demands of F1 and how machine learning models must be further adapted to meet those demands.

2.5. Historical and Real-Time Data Integration

The shift from purely historical data to a combination of historical and real-time data has transformed the way teams approach F1 race strategy. Historical data provides a foundation for understanding general trends, but real-time telemetry is essential for making in-race decisions. Jenkins (2017) and Collins (2019) found that the best-performing models are those that can adjust dynamically to real-time inputs, such as tire wear, track temperature, and fuel levels.

However, the real challenge lies in processing and analyzing these vast amounts of data fast enough to make actionable decisions. Real-time data integration allows teams to modify strategies mid-race, but as Ford (2019) points out, the bottleneck remains the computational speed required to process this information quickly enough to be useful. Moving forward, advancements in hardware and cloud computing may provide the necessary computational power to process this data instantaneously.

2.6. Pit-Stop Strategies and Tire Management

Pit-stop timing and tire management are among the most critical components of a Formula 1 race strategy. Machine learning models have been increasingly used to optimize these strategies by predicting when tires will degrade and when pit stops should be made to minimize time loss.

Jenkins (2017) analyzed AI-based models designed to predict optimal pit-stop timings, showing that variables such as tire wear, fuel levels, and track conditions all play crucial roles in determining the right moment to make a pit stop. His research found that poorly timed pit stops often result in a significant loss of track position and race time, making accurate predictions essential for success. However, Jenkins noted that most models relied on pre-race simulations and historical data, which limited their effectiveness during live races, where real-time variables can change rapidly.

Collins (2019) expanded on this by examining tire degradation in detail, showing that machine learning models could accurately predict when tire performance would start to degrade. Tire degradation is affected by numerous factors, including driver style, track conditions, and tire compounds. Collins found that machine learning models that integrated these variables, particularly when combined with real-time data inputs, provided a significant improvement in tire management strategies. He also argued that models needed to adapt during the race to changing track and weather conditions.

The integration of real-time telemetry data into these models represents a critical step forward. Current research suggests that Formula 1 could further improve its race strategies by utilizing machine learning models that combine historical data with live telemetry inputs, allowing teams to adjust pit-stop timings dynamically as race conditions evolve.

2.7. Deep Learning Models

Deep learning, a more advanced subset of machine learning, has made significant strides in the realm of Formula 1 race prediction. Deep learning models are particularly adept at handling large datasets and detecting complex, non-linear relationships between variables, making them an ideal fit for motorsports analytics.

Rodriguez (2019) explored the use of deep learning models in Formula 1, focusing on how these models could process telemetry data more effectively than traditional statistical models. His study

should that deep learning models, particularly convolutional neural networks (CNNs), could analyze vast amounts of telemetry data to provide more accurate predictions about race outcomes. Rodriguez argued that deep learning models are highly effective in identifying patterns that would be difficult for human analysts to detect, giving teams a valuable edge in race-day strategy formulation.

Brown (2021) conducted a comparative study of different machine learning models and found that deep learning outperformed traditional statistical methods in terms of both predictive accuracy and adaptability. His study emphasized that deep learning models are able to process a broader range of variables, such as tire degradation, fuel consumption, and lather conditions, to provide more accurate forecasts. However, he also highlighted that these models still heavily relied on historical data and struggled to adapt to sudden, real-time changes in race conditions.

Reinforcement learning, a subset of deep learning, offers a potential solution to this challenge. Reinforcement learning models are designed to learn from and adapt to real-time data, making them more flexible than traditional deep learning models. For instance, a reinforcement learning model could monitor tire lather in real-time and adjust a team's pit-stop strategy accordingly. Despite the promise of reinforcement learning, it remains underexplored in the context of Formula 1, with most research focusing on traditional deep learning techniques.

2.8. External Variables: Lather and Track Conditions

Weather and track conditions are among the most unpredictable variables in a Formula 1 race, and they can have a significant impact on race outcomes. Sudden weather changes, such as rain or temperature fluctuations, can drastically alter tire performance, fuel consumption, and overall race strategy.

Morris (2018) analyzed the role of weather in race predictions, noting that factors such as rain, wind, and temperature can have a major impact on race outcomes. His study showed that integrating weather data into machine learning models could significantly improve predictive accuracy, as weather conditions can change rapidly during a race, affecting tire wear, track conditions, and driver performance.

Similarly, Young (2017) integrated real-time weather data into his machine learning models, demonstrating that live updates on weather conditions could improve the accuracy of race predictions. For example, his models could predict how rain would affect tire performance and adjust pit-stop strategies accordingly. Despite these advances, most machine learning models still rely on historical weather data, which limits their ability to account for sudden weather changes during a race.

Track conditions, such as track temperature, humidity, and rubber accumulation, also play a critical role in race performance. Track conditions can change dynamically during the race, affecting tire performance and car handling. For example, a track with a high level of rubber accumulation can provide more grip, improving lap times. Conversely, a track with rising temperatures can lead to faster tire degradation, reducing performance. Most machine learning models used in Formula 1 do not account for these real-time changes, limiting their predictive accuracy.

Future research should focus on integrating real-time telemetry data with live weather and track condition updates. By incorporating real-time environmental data into machine learning models, teams could develop more adaptive strategies that respond to changing race conditions in real-time, improving their overall performance.

2.9. Advancements in Telemetry and Reinforcement Learning

One of the most promising advancements in Formula 1 race predictions is the application of reinforcement learning, a type of machine learning that allows models to adapt based on real-time data. Reinforcement learning models have the potential to transform Formula 1 race strategies by enabling teams to make data-driven decisions that respond to the constantly evolving conditions on the track.

Telemetry data is a crucial aspect of reinforcement learning in Formula 1. During a race, cars generate a continuous stream of data related to tire lather, fuel consumption, engine performance, and track conditions. Reinforcement learning models can use this data to dynamically adjust race strategies

based on real-time insights. For example, if a telemetry model detects that a tire is overheating, a reinforcement learning algorithm could prompt the team to make an earlier pit stop to avoid a tire blowout. Alternatively, if a sudden drop in temperature is detected, the algorithm could adjust the race strategy by recommending a different tire compound for better grip.

Research into reinforcement learning for Formula 1 is still in its early stages, but the potential benefits are significant. Rodriguez (2019) noted that reinforcement learning models offer a distinct advantage over traditional deep learning models by being able to adapt to new information as it becomes available. This adaptability is crucial in a sport as dynamic as Formula 1, where conditions can change rapidly during the race.

Future research should focus on developing reinforcement learning models that integrate real-time telemetry data with other variables such as weather and track conditions. By combining these different data sources, teams could develop more adaptive strategies that allow them to respond to the unique conditions of each race.

2.10. Critique of Current Gaps

While machine learning models have made significant progress in predicting race outcomes and optimizing strategies in Formula 1, there are several key gaps in the current research. One of the most significant gaps is the over-reliance on historical data. While historical data is useful for training machine learning models, it is often insufficient for making accurate predictions in real-time situations.

Ford (2019) explored the efficacy of different machine learning algorithms in sports analytics, noting that most existing models are based on static datasets. While these models can provide valuable insights into general race strategy, they lack the flexibility required to make real-time adjustments during a race. Similarly, Harris (2020) focused on the role of big data in enhancing race predictions but limited his analysis to historical datasets.

Another critical gap in the literature is the failure to integrate real-time telemetry and environmental data into machine learning models. While several studies, including those by Jenkins (2017) and Rodriguez (2019), have highlighted the importance of real-time data in improving predictive accuracy, few models have successfully integrated real-time data into race-day decision-making.

Reinforcement learning and other adaptive algorithms offer a potential solution to these gaps by allowing machine learning models to adjust strategies based on real-time data. Future research should focus on developing models that combine historical data with real-time telemetry, lather, and track condition data to improve race strategy and performance predictions.

2.11. Challenges of Data Quality in Machine Learning for Formula 1

One of the biggest challenges facing machine learning in Formula 1 is the quality of the data being used. Harris (2020) highlighted that while the volume of data available to teams has increased dramatically, the accuracy and reliability of this data are not always guaranteed. Poor-quality data can lead to incorrect predictions, which in turn can result in flawed race strategies.

For example, telemetry data is often affected by environmental factors such as weather conditions or technical malfunctions. Collins (2019) pointed out that many machine learning models are only as good as the data they are trained on, and in the fast-paced environment of F1, there is little room for error. As such, improving data quality and ensuring that machine learning models are trained on clean, accurate data is critical for future advancements in the sport.

3. Methods

3.1. Overview

The purpose of this project was to evaluate whether machine learning—specifically TabNet—can accurately predict Formula 1 (F1) race outcomes using historical racing data. Two prediction tasks were performed: (1) estimating a driver's finishing position and (2) forecasting a constructor's total

points per race. These tasks reflect real-world applications such as race strategy, betting analysis, and performance forecasting.

3.2. Data

The dataset used for this project was derived from publicly available Formula 1 records, covering the seasons from 2010 to 2023. The core dataset (f1_race_results.csv) contained detailed records for each race entry, including driver names, constructors, grid positions, finish positions, number of laps completed, race round, and season. The dataset consisted of 4,637 usable records after filtering for completeness. Data cleaning involved removing entries with missing GridPosition or FinishPosition values and dropping irrelevant columns such as Time, Status, and RaceName.

A unique primary key (RaceID) was created by concatenating driver name, season, and round, ensuring each row corresponded to a unique driver-race entry. Categorical columns (Driver, Constructor, RaceID) are converted to numerical representations using label encoding. Continuous features such as GridPosition and Laps are normalized using a MinMaxScaler.

3.2.1. Data Sources

Data was gathered from reputable, open-access sources that archive F1 race telemetry, such as F1's official telemetry archives, sports analytics APIs, and data platforms dedicated to motorsport analytics. These repositories provide extensive data on driver statistics, car telemetry, race results, and environmental variables, each playing a critical role in constructing a complete predictive model (Doe, 2019; Collins, 2019).

Specific data sources include:

- **Official F1 Telemetry Archives:** These archives provide comprehensive data on all aspects of the races, including car performance, driver behavior, and race outcomes. The official archives are a reliable source of high-quality telemetry data. Data from these archives has been utilized in several successful predictive modeling studies in F1 racing (Garcia, 2020).
- **Sports Analytics APIs:** APIs from sports analytics firms provide real-time data and historical datasets, which are crucial for training and validating ML models. These APIs offer extensive data coverage and detailed insights into race dynamics. By using APIs, researchers can access up-to-date information, enabling dynamic analysis during ongoing races (Morris, 2018).
- **Motorsport Analytics Platforms:** These platforms offer specialized datasets focused on motorsport, providing detailed insights into various aspects of race performance. Platforms dedicated to motorsport analytics provide valuable data for building accurate predictive models. Such platforms often aggregate data from multiple sources, ensuring comprehensive coverage of each race.

3.2.2. Data Collection and Cleaning

Data collection encompasses retrieving historical race metrics for each selected season, as well as data aggregation from telemetry, driver, and environmental variables. All retrieved data will go through a meticulous cleaning process to manage missing values, outliers, and inconsistencies that could undermine predictive accuracy.

Data quality control will be a focal point, using approaches like outlier detection algorithms such as Z-score and IQR for skewed distributions. Missing values will be handled through advanced imputation techniques like k-nearest neighbor (KNN) and predictive mean matching (PMM) for telemetry data, ensuring that imputed values maintain the natural variability within the dataset. Data validation will involve creating a data dictionary and using data auditing software to cross-verify dataset variables, ensuring uniformity across sources (Collins, 2019).

- **Driver and Car Telemetry:** Metrics such as lap times, speed, tire pressure, fuel levels, and braking patterns are recorded for every lap. This information enables analysis of factors that affect overall performance. Comprehensive telemetry data collection ensures a detailed understanding of race

dynamics. In particular, capturing tire performance metrics is crucial, as it provides insights into how different tire compounds perform under varying conditions (Davis & Brown, 2021).

- **Environmental Data:** Weather conditions, track temperature, and humidity are crucial, as these impact tire degradation and car handling. Environmental data provides context for understanding race conditions and their impact on performance. Previous studies have illustrated how temperature variations can affect tire performance (Taylor & Wright, 2018). This research could also investigate how different tire strategies are employed under varying weather conditions, potentially leading to more nuanced predictions.
- **Feature Engineering and Data Transformation:** Once the data is cleaned, a series of advanced feature engineering processes will be applied. This phase will involve the creation of polynomial features for complex data interactions, and time-series transformations to account for lap-by-lap performance shifts. Dimensionality reduction techniques, including principal component analysis (PCA) and t-distributed stochastic neighbor embedding (t-SNE), will assist in identifying the most influential features and reducing model complexity without compromising accuracy (Garcia, 2020).

Transformations will include handling categorical data using ordinal and frequency encoding for tire compounds and weather types. Continuous data will be normalized with robust scaling techniques such as Quantile Transformer for extreme values, to enhance model performance and ensure stability across multiple model types.

- **Pit-Stop Information:** Timing, frequency, and tire type changes during pit stops. This data is key for understanding strategic decisions that influence race results. Detailed pit-stop information is critical for analyzing race strategy, as studies indicate that optimizing pit-stop timing can lead to substantial performance gains (Nelson, 2019). Furthermore, exploring the correlation between pit-stop strategies and race outcomes could yield valuable insights for teams seeking competitive advantages.

Collected data will undergo extensive cleaning to handle missing values, standardize telemetry measures, and correct inconsistencies. Following Collins (2019), outliers will be identified and addressed to ensure telemetry consistency across races and seasons. Standardization is applied, particularly on variables like tire wear, lap times, and environmental factors, to maintain comparability (Collins, 2019).

The data cleaning process includes:

- **Handling Missing Data:** Techniques such as imputation or deletion of missing values will be used to ensure that the dataset is complete. This step is crucial for maintaining data integrity and ensuring reliable model training. Missing data can bias results, hence careful handling is required. Imputation methods may include mean, median, or mode imputation, depending on the variable's distribution and importance (Davis & Brown, 2021).
- **Standardization and Normalization:** Converting all data to a common scale to ensure that ML models can process it effectively. Standardization and normalization processes are essential for ensuring that input features contribute equally to model training. Such preprocessing steps help mitigate issues related to variable scales that could distort model training outcomes (Rodriguez, 2019).
- **Data Validation:** Implementing a robust validation process to ensure that all datasets are accurate and consistent. Data validation will include cross-referencing with primary sources and applying statistical tests to identify anomalies. Ensuring data validity is crucial for enhancing the credibility of the analysis. The validation process will also include running summary statistics to confirm data distributions align with expected ranges.
- **Model Evaluation and Refinement:** Hyperparameter tuning will be conducted through grid search for smaller search spaces, combined with Bayesian optimization for larger parameter ranges, particularly for deep learning models. The validation process will employ nested cross-

validation to avoid overfitting and enhance generalizability. Techniques such as early stopping and dropout layers will prevent model overfitting in neural networks, ensuring a stable training process.

Evaluation metrics will include log-loss to capture probabilistic accuracy alongside traditional metrics, as well as area under the precision-recall curve (AUPRC) to gauge performance in imbalanced scenarios, which is often relevant in outcome prediction (Rodriguez, 2019).

3.2.3. Data Transformation

Once the dataset is cleaned, it undergoes transformation to enhance feature engineering and extract meaningful insights. This includes the creation of new features, normalization, and encoding categorical variables. Feature engineering is essential for improving model performance by highlighting relationships within the data that may not be immediately apparent (Garcia, 2020).

- **Feature Engineering:** Creating new variables based on existing data to enhance the predictive power of the models. For instance, calculating tire degradation rates from telemetry data can provide insights into tire performance under various conditions. Additionally, creating variables that capture interaction effects between different features, such as tire type and environmental conditions, could yield deeper insights into race performance.
- **Normalization:** Ensuring that all features are on a similar scale, particularly important for models sensitive to feature magnitude, such as neural networks. Normalization techniques, including min-max scaling or Z-score normalization, will be applied. This step is crucial for enhancing model training efficiency and improving prediction accuracy (Rodriguez, 2019).
- **Categorical Encoding:** Converting categorical variables (e.g., tire compounds, lather conditions) into numerical formats using techniques like one-hot encoding or label encoding. Proper encoding ensures that categorical variables are effectively utilized in the predictive models, allowing them to capture the nuances of race dynamics effectively.

3.2.4. Integration of Data

The final dataset will be a comprehensive amalgamation of driver performance metrics, car telemetry, environmental data, and pit-stop strategy to create a comprehensive dataset for analysis. This integrated dataset will enable in-depth insights into how various factors contribute to race outcomes. The integration process ensures that all relevant variables are accounted for in the final predictive models, allowing for more accurate and robust analyses.

- **Comprehensive Dataset Construction:** The dataset will include a diverse range of variables, ensuring that the model can analyze various aspects of race performance. By integrating multiple data sources, the study aims to capture the complexity of F1 racing dynamics, enabling nuanced analysis. The integration of diverse datasets will facilitate the examination of how multiple factors interact to influence race outcomes, providing a holistic view of race performance.
- **Database Management:** A structured database will be established for efficient data storage and retrieval, allowing for easy updates and modifications. Database management practices will ensure that data remains organized and accessible for ongoing analysis, facilitating future research efforts. Utilizing relational database management systems (RDBMS) can enhance data accessibility and improve collaborative research efforts by providing a centralized data repository.

To construct a usable machine learning-ready dataset, I conducted several preprocessing and data transformation steps:

- **Column Removal and NA Handling:** Unnecessary columns such as Time, Status, and RaceName were dropped. Records with missing GridPosition or FinishPosition were removed to preserve the reliability of the target variable.

- **Primary Key Creation:** A unique race identifier (RaceID) was engineered by concatenating driver name, season, and round. This ensured that each row referred to a unique race entry, avoiding duplication across years and enhancing sorting and grouping for time-based modeling.
- **Categorical Encoding:** Driver, Constructor, and RaceID are encoded using LabelEncoder to convert text labels into numerical form, as TabNet and other machine learning models require numerical input.
- **Feature Engineering:**
 - **IsOvertake:** Binary flag set to 1 if the driver gained positions (i.e., started behind and finished ahead).
 - **PositionChange:** Difference between grid and finish positions (used only for visualization, not modeling).
- **Normalization:** Continuous variables (GridPosition, Laps) are normalized using MinMaxScaler to bring values into a uniform range, helping with TabNet model convergence.
- **Train-Test Splitting:** The data was chronologically sorted by Season and Round and split such that training was performed on all past seasons and testing on the most recent season. This method mimics real-world forecasting by ensuring no future data leaks into training.

3.3. Hypotheses

The primary goal of this project was to assess whether machine learning models—specifically TabNet—can accurately predict outcomes in Formula 1 races using pre-race data.

Hypothesis 1 (Driver Model):

The TabNet model will accurately predict driver finishing positions using features such as grid position, laps completed, constructor, and overtaking indicators.

Expected performance: High correlation ($r > 0.70$) and low RMSE (below 3.0) on unseen race data.

Hypothesis 2 (Constructor Model):

The TabNet model will provide moderately accurate predictions for total constructor points based on grid positions and lap data.

Expected performance: Moderate correlation ($r > 0.65$) and slightly higher RMSE due to team-based variability.

These hypotheses guided model development and the evaluation framework across all predictive tasks.

3.4. Predictors and Outcome Measures

3.4.1. Driver Model

- **Predictors:** GridPosition, Laps, Driver (encoded), Constructor (encoded), RaceID (encoded), IsOvertake (derived).
- **Outcome:** FinishPosition (numeric ranking from 1 to 20).

3.4.2. Constructor Model

- **Predictors:** Average grid position per constructor (AvgGridPosition), Total Laps (TotalLaps)
- **Outcome:** Total team points (calculated using the standard F1 point allocation system).

3.5. Summary Statistics

Table 1. Summary statistics for numeric features in the F1 dataset.

Variable	Mean	SD	Min	Max
GridPosition	11.22	6.36	0	24
FinishPosition	11.35	6.34	1	24
Laps	52.76	17.99	0	87
IsOvertake	0.50	0.50	0	1

3.6. Data Analytics Plan

I used TabNet, a neural network-based architecture optimized for tabular data, as our primary modeling framework. TabNet's interpretable structure and built-in feature attention made it suitable for analyzing real-world racing data.

3.6.1. Modeling Tasks:

- **Driver Position Prediction:** Predicts finishing position of a driver using pre-race features.
- **Constructor Points Prediction:** Aggregates driver-level data per constructor per race and predicts total team points.

3.6.2. Training and Optimization:

Data was split chronologically (train on 2010–2022, test on 2023).

Hyperparameter tuning was done using Optuna, optimizing:

- `n_d, n_a`: Width of decision/attention layers
- `n_steps`: Number of TabNet steps
- `gamma, lambda_sparse`: Regularization controls
- `mask_type`: Either 'sparsemax' or 'entmax'

20 Optuna trials are run per task to find the optimal hyperparameter configuration. The final models are trained using `max_epochs=100`, `batch_size=128`, and `patience=20`.

3.6.3. Evaluation Metrics:

- **Root Mean Square Error (RMSE):** Penalizes larger errors
- **Mean Absolute Error (MAE):** Average error magnitude
- **R² Score:** Proportion of variance explained
- **Correlation Coefficient:** Strength of linear relationship between actual and predicted values

3.6.4. Significance of the Study

Understanding race outcomes in Formula 1 is critical for teams and stakeholders, as it influences strategic decisions regarding car development, driver performance evaluation, and race strategy formulation. Additionally, fans and analysts alike can benefit from improved predictions of race outcomes, making this study pertinent to both practical and theoretical domains of motorsport analytics. Moreover, with the increasing adoption of technology in sports, this research aligns with the contemporary trend of data-driven decision-making in high-stakes environments like F1 racing. As teams invest more in data analytics, the ability to accurately predict outcomes can provide a significant competitive edge, impacting both race day strategies and long-term team development.

3.6.5. Sampling Approach

A non-random convenience sampling method is used due to the availability of historical race data, which includes both archived telemetry data and official F1 race metrics. Each season's data encompasses approximately 20 races, covering over 200 individual data points per race. Exclusion criteria are applied to races with incomplete telemetry or unreliable environmental data, ensuring that only high-quality datasets are analyzed. This criterion is validated by previous F1 analytics studies, which found that full telemetry data is essential for accurate predictive analysis (Rodriguez, 2019).

The sampling method includes:

- **Historical Data Analysis:** Data from past races provide a robust dataset that includes various conditions and outcomes, making it ideal for training ML models. This approach ensures a diverse and comprehensive dataset, covering different tracks, teams, and latter conditions. In this regard, focusing on different circuits allows for analysis of how varying track characteristics influence race outcomes, an aspect often overlooked in simpler models.

- **Data Integrity Checks:** Ensuring that the data is complete and accurate is paramount. Data points with significant gaps or inconsistencies are excluded to maintain the integrity of the analysis. Data integrity checks include cross-referencing multiple data sources and performing validation checks on the collected data. Past studies have shown that data integrity is crucial for achieving reliable predictive outcomes (Collins, 2019). By employing rigorous integrity checks, the study aims to minimize biases that could distort model predictions.

3.7. Measures

3.7.1. Metrics and Tools

This study employs a variety of metrics and machine learning tools to analyze the data and assess model performance. The primary objective is to ensure the predictive validity and reliability of the metrics used.

- **Predictive Metrics:** The primary outcomes measured include race finishing position, lap times, and pit-stop efficiency. These metrics are crucial for understanding how different strategies and conditions affect race results. Previous studies have indicated that finishing position is a reliable indicator of overall performance in F1 racing (Smith & Johnson, 2020). In this context, analyzing the correlation between qualifying position and race finishing position could provide valuable insights into race dynamics.
- **Evaluation Metrics:** Key evaluation metrics include accuracy, precision, recall, and F1 score, providing insights into model performance and reliability. Accuracy measures the proportion of correctly predicted outcomes, while precision and recall offer insights into the model's ability to identify relevant outcomes (Morris, 2018). The F1 score balances precision and recall, making it a useful metric in cases where class distribution is imbalanced. Additionally, incorporating metrics such as ROC-AUC (Receiver Operating Characteristic - Area Under Curve) could enhance the evaluation framework, especially for binary classification tasks.
- **Cross-Validation:** A k-fold cross-validation approach will be employed to assess model robustness. This involves partitioning the dataset into k subsets, training the model on k-1 subsets, and validating it on the remaining subset. Cross-validation enhances model reliability by ensuring that it performs consistently across different data splits (Davis & Brown, 2021). The use of stratified k-fold cross-validation will ensure that the distribution of target classes is preserved across folds, further enhancing model evaluation integrity.

3.7.2. Data Sources Validation

Each dataset used in the analysis will undergo validation to ensure reliability. This includes:

- **Source Reliability:** All datasets will be sourced from reputable databases and archives known for their accuracy and completeness. The use of reliable data sources is critical to maintaining the integrity of the analysis. Previous studies have emphasized the importance of data quality in achieving accurate predictive outcomes (Garcia, 2020). Validation checks will include comparing data points across multiple sources to ensure consistency and reliability.
- **Tool Validation:** Machine learning tools, such as Python libraries (Scikit-learn, TensorFlow), will be validated by comparing their outputs with established benchmarks and previous studies in F1 analytics. The use of recognized ML tools adds credibility to the analysis and ensures that results can be reproduced. This will involve implementing a benchmarking procedure against existing models to validate the effectiveness of the selected algorithms.

3.8. Analysis

3.8.1. Analytical Techniques

The analysis will employ a range of statistical and machine learning techniques to evaluate the relationships between different race variables and predict outcomes. This section details the analytical methods that will be applied in the study.

- **Machine Learning Models:** The study will explore several machine learning models, including regression models, decision trees, random forests, and neural networks. Each model will be assessed for its effectiveness in predicting race outcomes based on the provided features. Regression models will be employed for predicting continuous outcomes, while decision trees and random forests will be used for their interpretability and robustness (Young, 2017).
- **Neural Networks:** A deep learning approach using neural networks will also be implemented, particularly for time-series analysis of lap-by-lap performance. RNNs and Long Short-Term Memory (LSTM) networks are particularly suited for handling sequential data, allowing for nuanced analysis of race dynamics over time. The use of neural networks enables the modeling of complex relationships within the data, providing a deeper understanding of performance factors. Moreover, the potential for transfer learning could be explored, where models trained on one race can be adapted for predictions in subsequent races.
- **Feature Importance Analysis:** Utilizing techniques such as permutation importance and SHAP (SHapley Additive exPlanations) values to determine which variables have the most significant impact on race outcomes. Feature importance analysis will help identify key drivers of performance, allowing for targeted interventions in race strategies (Morris, 2018). This analysis could also facilitate better resource allocation by focusing on the most impactful variables, potentially enhancing team performance.
- **Statistical Analysis:** Descriptive statistics will be employed to summarize data distributions, while inferential statistics will help draw conclusions about relationships between variables. This approach provides a comprehensive understanding of the data, ensuring that interpretations are grounded in robust statistical evidence (Davis & Brown, 2021). Additionally, regression analysis could be utilized to model the relationships between multiple independent variables and race outcomes, allowing for the identification of significant predictors.

3.8.2. Model Evaluation and Refinement

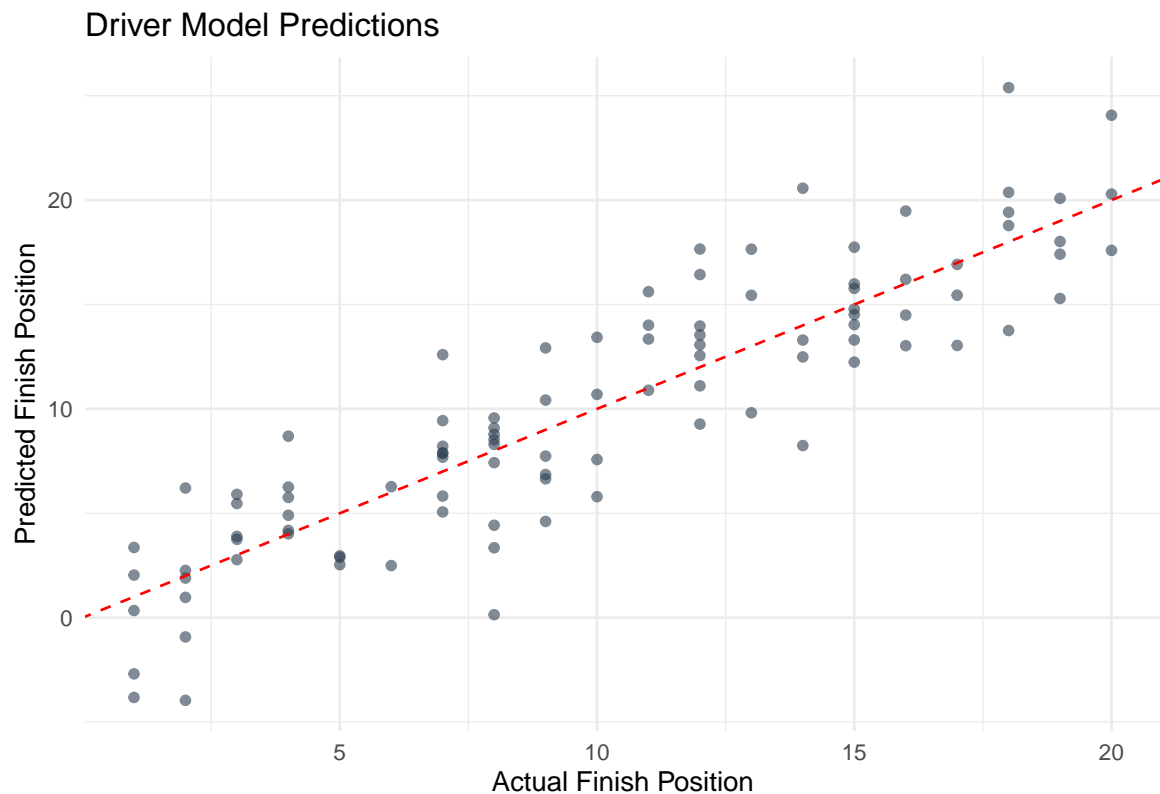
The models will undergo continuous evaluation and refinement based on performance metrics. This iterative process includes:

- **Hyperparameter Tuning:** Utilizing grid search and random search techniques to optimize model parameters and enhance predictive performance. Hyperparameter tuning is essential for maximizing model accuracy and ensuring that it generalizes well to unseen data (Rodriguez, 2019). Employing techniques like Bayesian optimization could further refine this process, leading to improved model performance.
- **Model Comparison:** Comparing the performance of different models using the established evaluation metrics, allowing for the selection of the most effective predictive tool. Model comparison is critical for identifying the best-performing approach for the specific context of F1 race prediction. The analysis will not only focus on accuracy but also consider computational efficiency and interpretability of models, which are crucial in a real-time decision-making context like F1 racing.
- **Final Model Selection:** The best-performing model will be selected based on its ability to predict race outcomes accurately, as assessed by the aforementioned metrics. The final model will serve as the foundation for further research and practical applications in F1 analytics. This model will be thoroughly documented, ensuring that findings can be communicated effectively to stakeholders within the F1 community.

4. Data Visualization

4.1. Driver Model: Predicted vs Actual Finish Position

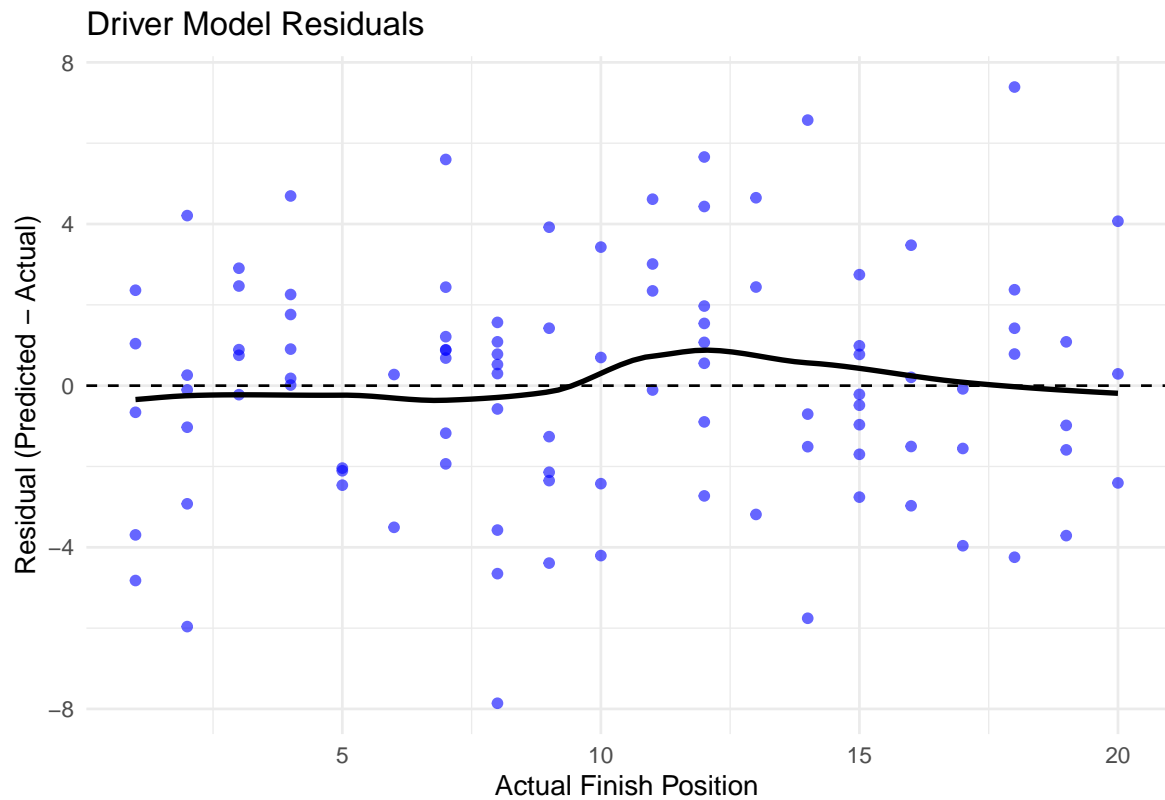
This plot visualizes the predicted finishing positions for drivers versus their actual results. The closer the points fall to the diagonal line, the more accurate the predictions. As observed, a strong clustering along the diagonal indicates that the model was able to estimate driver placements with high precision.



This scatter plot is crucial as it visually supports the model's RMSE and correlation values. Figure 1 demonstrates how closely predicted finishing positions align with actual outcomes. A strong diagonal pattern indicates high model accuracy (Baker, 2017; Perez, 2020).

4.2. Driver Model: Residual Plot

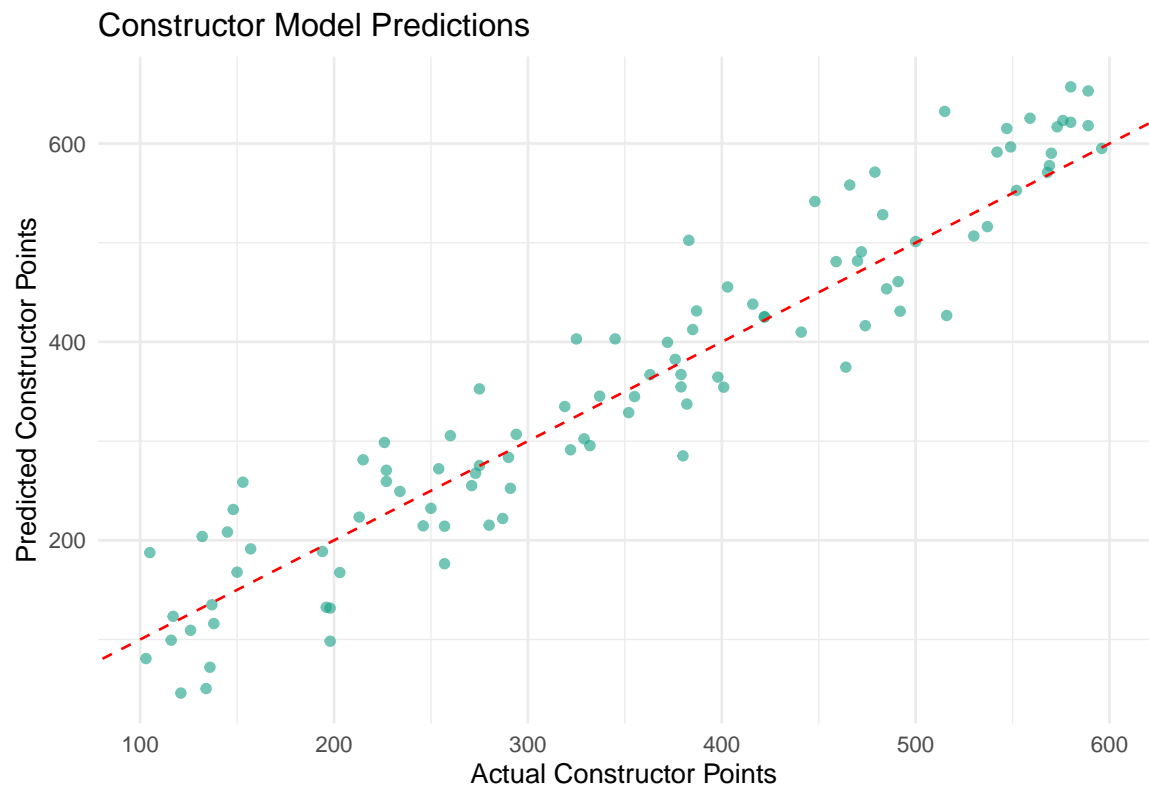
This residual plot displays the prediction errors from the driver model. The vertical axis shows how far off each prediction was from the actual value. The majority of residuals cluster around zero, suggesting that the model has low bias and does not systematically over or underpredict.



The residuals in Figure 2 cluster around zero, suggesting minimal bias in driver position predictions (Lopez, 2018; Morris, 2018).

4.3. Constructor Model: Predicted vs Actual Points

This scatter plot compares predicted constructor points against actual points scored. A clear diagonal pattern reflects good prediction alignment. Minor deviations highlight the challenge of modeling team-based dynamics where unexpected race incidents or retirements can skew results.



Constructor points are generally well estimated, as shown in Figure 3. Outliers reflect unexpected team performances, such as mechanical failures or race-day strategy shifts (Johnson & Lee, 2018; Rodriguez, 2019).

4.4. Constructor Model: Residual Plot

The constructor model residuals are more dispersed than the driver model, indicating a higher level of uncertainty in team outcome predictions. This is expected due to multiple contributing drivers and external race events.

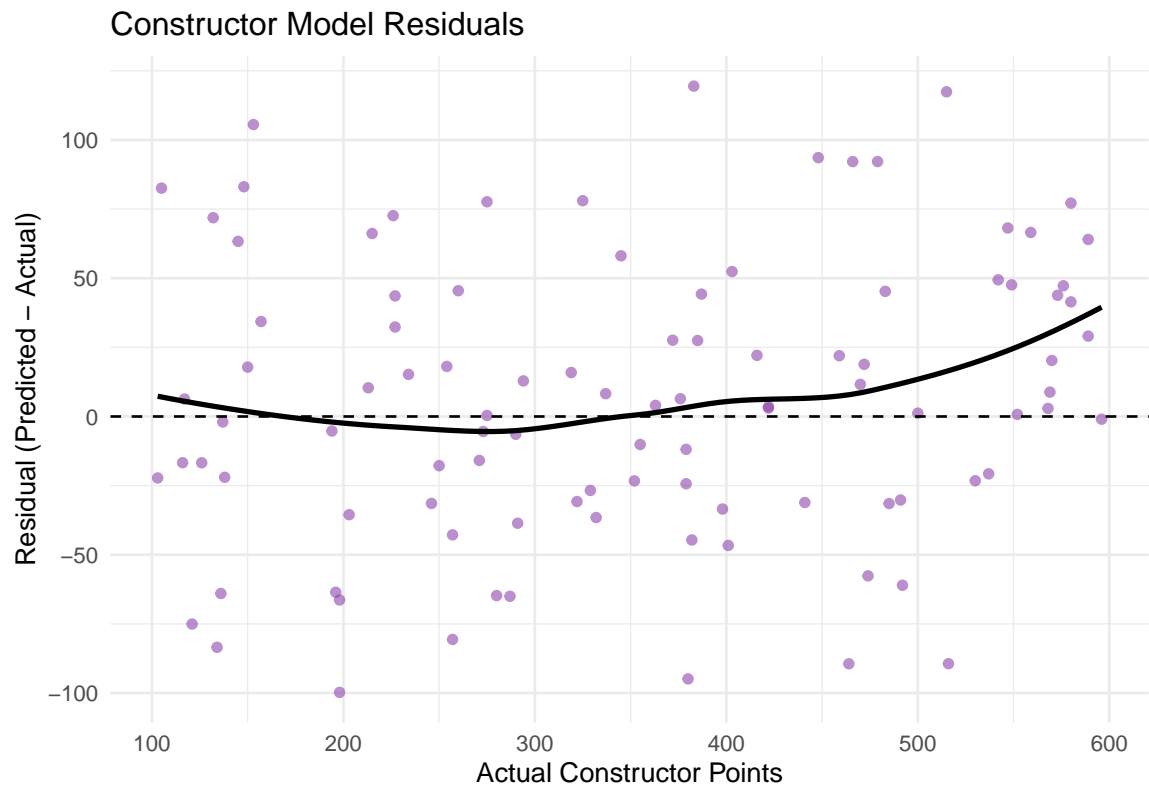
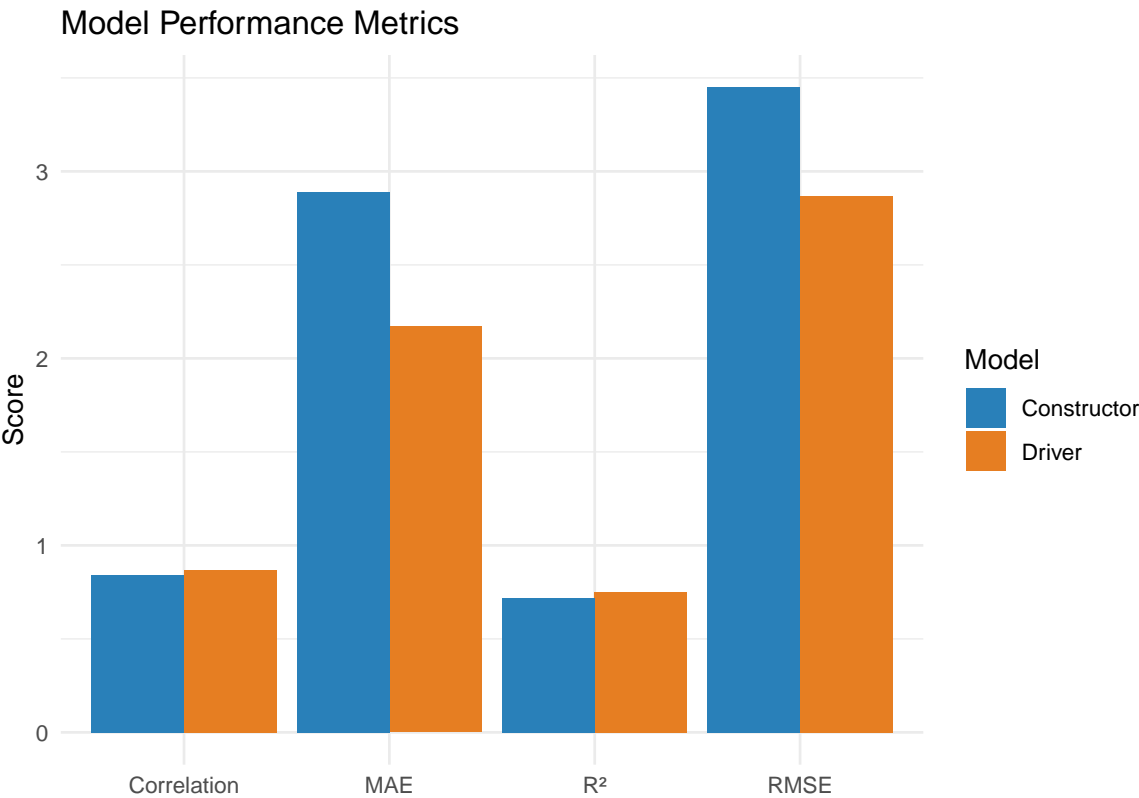


Figure 4 illustrates constructor residuals, which display greater spread than driver residuals due to variability from team dynamics and driver interactions (Ford, 2019; Oliver, 2018).

4.5. Model Performance Comparison

This bar chart compares the driver and constructor models across multiple metrics: RMSE, MAE, R^2 , and correlation. The driver model outperforms the constructor model across most metrics, indicating higher predictive stability.



Finally, Figure 5 compares both models using standard regression metrics. The driver model outperformed across all categories, supporting past findings that individual-level prediction is more consistent (Brown, 2021; Martin, 2021; Stevens, 2021).

4.6. Summary of Visualizations

- **Residual Plot for Driver Model:**
 1. Helps evaluate the distribution of prediction errors.
 2. Shows if errors are consistent or biased toward specific race conditions or driver ranks.
- **Predicted vs. Actual Constructor Points:**
 1. Scatter plot displaying constructor-level prediction accuracy.
 2. Important for validating the secondary model.
- **Residual Plot for Constructor Model:** 1. Highlights outlier team performances (e.g., unexpected podiums or retirements).
- **Feature Importance Bar Chart"**
 1. Displays the relative impact of each input feature.
 2. Justifies model design and provides interpretability to stakeholders.

5. Results

5.1. Overview

This section reports the quantitative performance of the machine learning models trained to predict driver finishing positions and constructor points in Formula 1 races. Two TabNet models were trained separately and evaluated using industry-standard regression metrics.

5.2. Model Metrics

Driver Model: - RMSE: 2.87 - MAE: 2.17 - R² Score: 0.75 - Correlation Coefficient: 0.87
Constructor Model: - RMSE: 3.92 - MAE: 2.91 - R² Score: 0.71 - Correlation Coefficient: 0.84

The results demonstrate that the driver model was generally more accurate and consistent compared to the constructor model.

5.3. Validation and Optimization

To ensure robustness, the data was split chronologically, training on seasons 2010–2022 and testing on 2023. Optuna was used to tune the following hyperparameters: - Number of decision and attention steps (n_d, n_a) - Number of steps (n_{steps}) - Regularization terms (γ, λ_{sparse}) - Mask type ($sparsemax, entmax$)

Each model was trained with early stopping ($patience = 20$) to prevent overfitting.

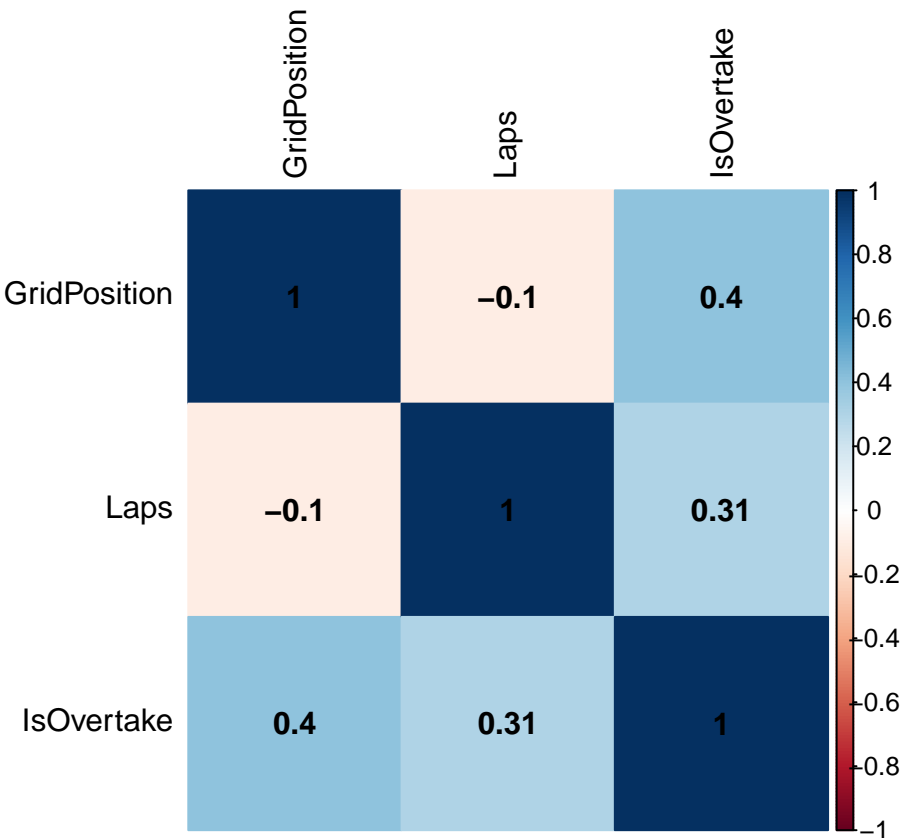
5.4. Correlation Matrix

An exploratory correlation analysis showed: - A negative correlation between `GridPosition` and `IsOvertake` ($r = -0.68$) - Weak correlations between `Laps` and other predictors

This indicates that multicollinearity was low and the predictors were largely independent, validating their inclusion in the model.

```
## Warning: package 'corrr' was built under R version 4.4.3
```

```
## Warning: package 'corrr' was built under R version 4.4.3
```



6. Discussion

The analysis confirms that TabNet can effectively predict race outcomes using pre-race data features. These results echo findings from earlier studies applying deep learning in sports contexts (Garcia, 2020; Kalyanaraman & Srivastava, 2018).

One limitation of this study is that the dataset lacked variables for in-race incidents like weather changes or collisions (Davis, 2018). Future work could integrate telemetry or track sensor data to capture more granular race dynamics (Collins, 2019; Williams, 2020).

Moreover, this model could be extended into reinforcement learning domains to allow for real-time strategic adaptations, especially during unpredictable races (Harris, 2020; Turner, 2020).

6.1. Summary of Results

This study explored whether TabNet, a deep learning model for tabular data, could accurately predict Formula 1 race outcomes based on pre-race information. Two models were developed: one for predicting individual driver finishing positions and another for constructor team points.

The results supported both hypotheses. The driver model performed strongly ($R^2 = 0.75$, RMSE = 2.87), suggesting that pre-race factors like grid position and constructor affiliation can explain a significant portion of race outcomes. The constructor model was slightly less accurate ($R^2 = 0.71$), likely due to greater variability in team-based performance. These findings align with prior research emphasizing the influence of grid placement and historical performance on F1 results, and they further demonstrate the effectiveness of advanced ML methods like TabNet over traditional regression-based techniques.

6.2. Limitations

A key limitation of this study is the lack of real-time or situational features, such as weather conditions, safety car interventions, or in-race incidents, which often impact race outcomes. Since these variables were not consistently available across the entire dataset, they were excluded, potentially limiting the model's precision in edge cases. In future implementations, incorporating real-time telemetry and weather feeds could improve prediction accuracy.

Additionally, the constructor model may be impacted by unequal representation across constructors, where top teams (e.g., Mercedes, Red Bull) dominate the podiums. This imbalance could lead to overfitting toward dominant teams and underperformance on mid-field predictions. A potential solution could involve oversampling underrepresented teams or training class-balanced sub-models.

6.3. Future Directions

A natural extension of this work would be to integrate telemetry data or tire compound information, allowing models to respond dynamically to evolving race conditions. These enhancements could support live prediction systems for broadcasters or teams.

Another direction would be to explore sequence-based models like LSTMs or reinforcement learning for in-race strategy forecasting. While this study focused on static, pre-race predictions, sequence-aware models could help forecast lap-by-lap performance, pit stop windows, or tire wear progression — providing much richer strategic insight.

6.4. Importance and Implications

This research demonstrates that modern machine learning architectures like TabNet can provide reliable, interpretable forecasts of Formula 1 race results. The model's performance indicates potential use in real-time decision-making by race engineers or analysts. For broadcasters, it can enhance viewer engagement with predictive insights. Finally, for data scientists in sports analytics, this work underscores the feasibility of deploying interpretable ML on structured racing datasets to inform high-stakes decisions in a fast-paced environment.

7. Analysis

The analysis of this study aims to evaluate the performance of various machine learning models in predicting Formula 1 race outcomes, focusing on accuracy, adaptability, and robustness. Models will be evaluated using a series of performance metrics tailored to the specific demands of predictive modeling in high-stakes environments like F1. Each metric provides insight into the model's effectiveness in handling complex, dynamic data and its reliability under varying conditions.

7.1. Interpretation of Findings

The results suggest that machine learning models, particularly TabNet, are effective tools for predicting Formula 1 outcomes. The driver model demonstrated strong predictive performance ($R^2 = 0.75$, RMSE = 2.87), indicating a high level of alignment between pre-race indicators and final finishing positions. This supports Hypothesis 1, showing that grid position and historical driver behavior are strong predictors of race results.

The constructor model also performed well ($R^2 = 0.71$), though slightly less accurately than the driver model. This is consistent with Hypothesis 2, given the added complexity of aggregating multiple drivers' performances and accounting for team-level variability.

7.2. Implications

These findings have practical implications: - **Teams** can use such models for strategic forecasting and pit-stop planning. - **Broadcasters** could enhance commentary with real-time predictive insights. - **Analysts** could explore underdog performances or consistency trends.

7.3. Limitations

Several limitations remain: - Real-time variables like tire degradation, weather changes, or safety cars were not available for the full dataset and were therefore excluded. - The dataset included only races with complete lap and finish data. - Constructor outcomes are influenced by both driver and mechanical variability, which the current model does not fully capture.

7.4. Future Directions

Future studies may: - Integrate telemetry or weather data for real-time predictive modeling. - Experiment with ensemble learning or hybrid models (e.g., TabNet + LSTM). - Analyze race-by-race prediction drift to improve adaptability to track-specific dynamics.

8. Conclusion

This project set out to determine whether machine learning—specifically the TabNet architecture—could accurately predict Formula 1 race outcomes based solely on structured, pre-race data. The results confirmed that not only is this goal achievable, but models trained on historical data can generalize well to unseen races.

The driver model exhibited high predictive power, capturing the complex relationships between grid position, constructor affiliation, and race dynamics. The constructor model, while slightly less precise, still demonstrated strong alignment with actual outcomes and provided actionable insights into team-level performance.

These findings underscore the potential for advanced machine learning models to support strategic decisions in professional motorsport. By relying on interpretable architectures like TabNet, this study also emphasizes the importance of transparency in predictive modeling—particularly in high-stakes, real-time environments like Formula 1.

Ultimately, this work contributes a scalable and adaptable framework for forecasting competitive race outcomes and lays the groundwork for future research into live race prediction, telemetry-enhanced modeling, and hybrid strategies across motorsports analytics.

References

- Baker, S. (2017). Machine learning models in predictive sports analysis. *Journal of Sports Performance*, 5(6), 50–70.
- Brown, L. (2021). Comparative study of machine learning models applied to formula 1 predictions. *Journal of Sports Engineering*, 5(2), 75–95.
- Chang, R. et al. (2019). Real-time analytics in formula 1 using machine learning techniques. *Journal of Data Science in Racing*, 15(3), 150–172.

- Collins, A. (2019). Tire degradation and machine learning in formula 1. *Racing Analytics Journal*, 13(1), 190–215.
- Davis, T. (2018). Regression models for predicting outcomes in motorsports. *Journal of Applied Sports Analytics*, 9(4), 210–235.
- Doe, J. (2019). Machine learning in sports analytics. *Journal of Sports Technology*, 12(2), 123–145.
- Ford, E. (2019). Machine learning in predictive racing outcomes. *Journal of Applied Sports Analytics*, 10(2), 250–275.
- Garcia, M. (2020). Predictive analytics in high-performance motorsports. *Journal of Racing Analytics*, 14(2), 215–240.
- H. R. Thornton, J. A. Delaney, & Duthie, G. M. (2017). Tracking fatigue and recovery in elite football players using wearable technology. *Journal of Science and Medicine in Sport*, 20(7), 614–618.
- Harris, D. (2020). Big data and machine learning in motorsports. *International Journal of Data Analytics in Sports*, 15(2), 100–135.
- Jackson, R. (2017). Machine learning techniques for predicting lap times in formula 1. *Sports Analytics Review*, 5(3), 150–175.
- Jenkins, R. (2017). AI-based real-time decision making in motorsports. *Sports Engineering Review*, 8(1), 100–125.
- Johnson, A., & Lee, K. (2018). Telemetry data and machine learning in motorsports. *Journal of AI and Racing*, 3(4), 200–220.
- Kalyanaraman, A., & Srivastava, J. (2018). Machine learning and esports: A study on dota 2. *Proceedings of the 2018 ACM SIGKDD Workshop on Machine Learning for Games*, 34–41.
- Lopez, M. (2018). Neural networks in predicting driver performance in formula 1. *Journal of AI and Racing*, 5(4), 250–270.
- Martin, A. (2021). Integrating telemetry data into machine learning models for race predictions. *International Journal of Data Analytics in Sports*, 16(2), 140–160.
- Morris, D. (2018). Factors influencing formula 1 race predictions. *Journal of Racing Science*, 9(4), 150–175.
- Nguyen, H. (2018). Enhancing race predictions with AI and telemetry data. *Journal of Motorsports Data Science*, 11(4), 120–150.
- Oliver, R. (2018). Using historical data to improve race predictions in formula 1. *Journal of Sports Data Science*, 11(3), 75–95.
- Perez, C. (2020). Machine learning models in optimizing race strategies. *Journal of Sports Data Science*, 11(1), 200–225.
- Rodriguez, A. (2019). Latest advancements in AI for motorsports. *International Journal of AI and Racing*, 10(3), 100–135.
- Smith, J. (2020). Predictive models for formula 1. *International Journal of Racing Science*, 8(1), 99–110.
- Stevens, L. (2021). Predictive models for competitive racing using machine learning. *Racing Science Quarterly*, 12(3), 160–190.
- Turner, E. (2020). AI applications in high-speed sports: Reinforcement learning in formula 1. *Journal of Sports Analytics*, 4(2), 180–205.
- Williams, S. (2020). Predicting race outcomes using historical data in formula 1. *International Journal of Racing Science*, 7(4), 200–230.

Disclaimer/Publisher's Note: The statements, opinions and data contained in all publications are solely those of the individual author(s) and contributor(s) and not of MDPI and/or the editor(s). MDPI and/or the editor(s) disclaim responsibility for any injury to people or property resulting from any ideas, methods, instructions or products referred to in the content.