

# **F1 Race Orientation using ML**

Submitted in partial fulfilment of the requirements

of the degree of  
**Bachelor of Engineering**

by

Prathamesh More: 23106001

Harsh Kalsekar: 23106042

Mayuresh Keni: 23106038

Guide:

**Prof. Shradha Dalvi**



**Department of Computer Science & Engineering  
(Artificial Intelligence & Machine Learning)**

**A. P. SHAH INSTITUTE OF TECHNOLOGY, THANE**

**UNIVERSITY OF MUMBAI**

**(2025-2026)**



**A. P. SHAH INSTITUTE OF TECHNOLOGY**

## **CERTIFICATE**

This is to certify that the project entitled “**F1 Race Orientation using ML**” is a bonafide work of “**Prathmesh More (23106001), Harsh Kalsekar (23106042), Mayuresh Keni (23106038)**” submitted to the University of Mumbai in partial fulfilment of the requirement for the award of the degree of **Bachelor of Engineering in Computer Science & Engineering (Artificial Intelligence & Machine Learning)**

---

Prof. Shradha Dalvi  
Project Guide

---

Prof. Yogeshwari Hardas  
Project Co-Ordinator

---

Dr. Jaya Gupta  
Head of Department

---

Dr. Uttam D Kolekar  
Principal



**A. P. SHAH INSTITUTE OF TECHNOLOGY**

## **Project Report Approval for T. E.**

This project report entitled *F1 Race Orientation using ML* by *Prathamesh More, Harsh Kalsekar, Mayuresh Keni* is approved for the degree of *Bachelor of Engineering* in **Computer Science & Engineering (Artificial Intelligence & Machine Learning)**, 2025-26.

Examiner Name

Signature

1. \_\_\_\_\_

\_\_\_\_\_

2. \_\_\_\_\_

\_\_\_\_\_

Date:

Place:

## **Declaration**

We declare that this written submission represents my ideas in my own words and where others' ideas or words have been included, I have adequately cited and referenced the original sources. I also declare that I have adhered to all principles of academic honesty and integrity and have not misrepresented or fabricated or falsified any idea/data/fact/source in my submission. I understand that any violation of the above will be cause for disciplinary action by the Institute and can also evoke penal action from the sources which have thus not been properly cited or from whom proper permission has not been taken when needed.

---

Prathamesh More  
23106001

---

Harsh Kalsekar  
23106042

---

Mayuresh Keni  
23106038

Date:

## ABSTRACT

This project presents the development of a machine learning–based model to predict Formula 1 race outcomes using the XG Boost algorithm. The model leverages historical race data, driver profiles, constructor performance, qualifying results, pit stop strategies, lap time consistency, and weather conditions to estimate finishing positions for each driver.

The dataset is sourced from publicly available Formula 1 statistics (e.g., Kaggle) and extended with engineered features such as rolling averages of driver and team performance, circuitspecific trends, and contextual factors like track temperature, humidity, and wind speed. These features enable the model to capture both short-term form and long-term patterns in driver and team competitiveness.

The system is implemented in Python within Visual Studio Code, using libraries such as Pandas for preprocessing, XGBoost and SVM for predictive modeling, and Scikit-learn for evaluation. The workflow includes raw data ingestion, feature engineering, model training, validation, and generation of ranked race predictions. The model outputs driver finishing predictions for each race in CSV format, enabling comparison with actual results.

Experimental results demonstrate that integrating variables such as pit stop frequency, average lap times, and weather conditions improves prediction accuracy compared to models based solely on grid position and historical points. This approach highlights the potential of machine learning in motorsport analytics, offering insights into race strategies, performance forecasting, and the dynamic impact of external conditions on outcomes.

**Keywords:** Race Prediction, Machine Learning, XGBoost, SVM, Data Analytics, Predictive Modeling

## CONTENTS

1. Introduction .....	1
2. Literature Survey .....	3
2.1. History .....	3
2.2. Literature Review .....	4
3. Limitations in Existing Systems .....	6
4. Problem Statement & Objectives .....	9
4.1. Problem Statement .....	9
4.2. Objectives .....	10
5. Proposed System .....	11
5.1. System Architecture .....	12
5.2. System Modules.....	14
5.3. Workflow Explanation .....	15
5.4. Advantages of Proposed System .....	14
6. Experimental Setup .....	18
7. Results & Discussion .....	21
8. Conclusion & Future Scope.....	30
9. Future Scope.....	30
References .....	31

## **LIST OF FIGURES**

5.1. Overall system architecture & workflow .....	12
7.1 Home Page of Project before Predictions .....	21
7.2 Home Page of Project after Predictions .....	22

## **LIST OF TABLES**

3.1. Summary of limitations of existing systems .....	8
7.1 Summary of Metrics .....	23
7.2 Output Metrics .....	24



## ABBREVIATION

AI ..... *Artificial Intelligence*

ML ..... *Machine Learning*

SVM ..... Support Vector Machine

XGBoost ..... Extreme Gradient Boosting

MAE ..... Mean Absolute Error

# **Chapter 1**

## **Introduction**

Formula 1 (F1) is the pinnacle of motorsport, combining high-speed engineering, strategic decision-making, and human skill. With races contested across different continents, varying weather conditions, and diverse track layouts, predicting the outcome of a Formula 1 race is a highly challenging task. The performance of a driver is not solely determined by individual skill but is influenced by multiple interdependent factors such as the capabilities of the car, team strategy, pit stop efficiency, track conditions, and even unpredictable external influences like weather and safety car interventions.

This Project follows 2 SDG goals:- 1.Good Health and Well Being:-This project promotes Good Health and Well-Being by applying data-driven insights from Formula 1 to enhance driver safety, performance efficiency, and reaction analysis.2.Industry, Innovation and Infrastructure:-The project supports Industry, Innovation, and Infrastructure by leveraging Machine Learning (ML) and Artificial Intelligence (AI) to analyze large-scale racing data and optimize vehicle performance.

In recent years, the advancement of machine learning (ML) and data analytics has opened new avenues in motorsport analysis. Teams and analysts increasingly rely on data-driven models to optimize race strategies, assess competitor performance, and forecast possible race outcomes. Machine learning techniques excel in identifying hidden patterns within large datasets, making them well-suited for analyzing complex motorsport dynamics where linear assumptions often fail.

This project explores the application of the XGBoost algorithm, a state-of-the-art gradient boosting technique, to predict Formula 1 race results. XGBoost has gained prominence in the ML community due to its scalability, ability to handle missing values, and strong performance on structured/tabular data. By leveraging historical race data, the model aims to predict the finishing positions of drivers in upcoming races, offering insights into expected race outcomes and helping stakeholders understand the factors that drive success in F1.

Predicting race results is not just an academic exercise but also a practical problem with applications across multiple domains. For teams, accurate predictions can support tactical decisions such as optimal pit stop timing, tire selection, and fuel strategy. For broadcasters and fans, predictive models enhance the viewing experience by offering pre-race insights and postrace analyses. From a technological standpoint, this project demonstrates how modern machine learning techniques can transform raw motorsport data into actionable intelligence.

Furthermore, Formula 1 provides a rich dataset with multiple layers of information: driver statistics, constructor performance, race metadata, lap-by-lap timings, pit stop information, and environmental conditions. The multi-dimensional nature of this data makes it an ideal testbed for building robust predictive models and conducting advanced feature engineering.

## **Chapter 2**

### **Literature Survey**

#### **2.1 History**

The prediction of Formula 1 race outcomes has progressed from intuition-driven assessments to data-centric approaches over the decades. Early methods relied on driver skill, mechanical performance, and simple statistical analysis of past races. With the introduction of telemetry in the 1980s and 1990s, teams began collecting real-time data on engine performance, tire degradation, and fuel consumption, which laid the foundation for computational race modeling. By the 2000s, digital records and open datasets such as the Ergast API enabled researchers to apply regression models and decision trees to motorsport prediction. More recently, advanced machine learning techniques such as Random Forests, Support Vector Machines (SVM), Neural Networks, and Gradient Boosting (XGBoost) have been used to capture the complex, nonlinear dynamics of races. These methods increasingly incorporate pit stop efficiency, lap times, and weather conditions, reflecting a shift toward predictive analytics as a key component of motorsport strategy.

## 2.2 Literature Review

### **Smith et al. (2015):**

Applied logistic regression and Bayesian models to predict race winners using qualifying positions and team strength. Their work highlighted the limitations of simple statistical models in capturing dynamic race events.

### **Gupta & Rao (2018):**

Developed a Random Forest-based model to predict finishing positions using driver statistics and constructor performance. While effective, their study lacked consideration of external factors such as pit stop strategy and weather.

### **Fernández et al. (2019):**

Introduced Support Vector Machines (SVM) for motorsport predictions, focusing on lap time classification. Their model achieved high accuracy in controlled datasets but struggled with noisy race conditions.

### **Wang and Li (2020):**

Implemented deep learning (LSTMs) to analyze sequential lap data. Although this improved lap-level predictions, it required extensive computational resources and struggled with incomplete telemetry.

**Khan et al. (2021):**

Applied XGBoost to Formula 1 datasets, demonstrating superior performance over traditional regression and SVM models. Their approach incorporated grid position, team performance, and historical data but did not fully integrate weather or pit stop metrics.

**Martins & Silva (2022):**

Emphasized the role of contextual factors, integrating weather, pit stop durations, and tire strategies into their prediction model. Their results confirmed that incorporating such variables significantly improves prediction accuracy.

## Chapter 3

### Limitations in Existing Systems

Despite significant progress in applying data analytics and machine learning to motorsport prediction, existing systems still face multiple limitations that hinder their overall accuracy and real-world applicability.

1. **Limited Feature Scope:** Most traditional models primarily focus on static features such as driver standings, constructor rankings, or qualifying positions. While these factors are important, they fail to capture dynamic race variables such as pit stop strategies, lap time consistency, tire degradation, and real-time weather conditions. As a result, predictions often lack robustness when external factors strongly influence race outcomes.
2. **Neglect of Contextual Factors:** Weather conditions, track temperature, and humidity can drastically alter performance by affecting tire grip and car handling. Many existing systems either ignore these contextual variables or treat them as secondary features, thereby reducing the accuracy of predictions in varying environmental conditions.
3. **Over-Reliance on Historical Data:** A significant number of studies depend heavily on historical driver and team performance. While useful for long-term trend analysis, such reliance overlooks race-specific events such as safety cars, collisions, and sudden mechanical failures. This limits the ability of models to adapt to unpredictable scenarios.

4. **Lack of Real-Time Integration:** Several existing systems are designed for post-race analysis rather than live race prediction. They do not incorporate streaming telemetry data, making them unsuitable for providing real-time strategic insights that could be beneficial to teams, broadcasters, or fans.
5. **Data Quality and Availability Issues:** Public datasets, such as those from Ergast or Kaggle, often contain incomplete or inconsistent data, particularly for older seasons. Missing lap times, pit stop durations, or incomplete weather records reduce the effectiveness of models and require complex preprocessing.
6. **Computational Complexity:** Deep learning models such as LSTMs and NeuralNetworks have been explored in prior works, but their high computational demands make them impractical for lightweight or real-time applications. Additionally, these models often require large volumes of telemetry data, which may not always be available for open-source use



**Table 3.1: Summary of limitations of existing systems**

<b>Title</b>	<b>Author</b>	<b>Key Points</b>	<b>Improvement Proposed</b>	<b>Context/Domain</b>
Predicting Race Winners using Logistic Regression	[Smith et al., 2015]	Used qualifying and team data	Could not capture dynamic race events	Sports Analytics
Random Forests for F1 Position Prediction	[Gupta & Rao, 2018]	Used driver and constructor stats	Ignored weather and pit stops	Motorsport ML
Lap Time Classification with SVM	[Fernández et al., 2019]	Classified lap performance	Struggled with noisy race data	Race Data Mining
Deep Learning for Sequential Lap Analysis	[Wang & Li, 2020]	LSTM for lap sequences	High computation, data hungry	Neural Networks
XGBoost for Race Prediction	[Khan et al., 2021]	High accuracy with grid/team data	Lacked external race factors	ML (Boosting)

## **Chapter 4**

### **Problem Statement & Objectives**

#### **4.1 Problem Statement**

Predicting the outcomes of Formula 1 races is a highly challenging task due to the dynamic and unpredictable nature of the sport. Traditional prediction models rely primarily on static factors such as driver standings, constructor rankings, or qualifying positions. While these features provide some insight, they fail to account for critical race-specific variables such as pit stop efficiency, lap time consistency, tire strategy, and real-time weather conditions. As a result, the predictions generated by such systems often lack robustness and fail to reflect the true complexity of modern Formula 1 racing.

## 4.2 Objectives

The main objectives of this project are:

1. **To collect and preprocess Formula 1 race data:** from publicly available sources, ensuring consistency and completeness across driver, constructor, pit stop, lap time, and weather datasets.
2. **To perform feature engineering:** by integrating extended features such as pit stop efficiency, lap time performance, tire usage, and weather conditions alongside traditional race metrics.
3. **To provide a framework for adaptable race prediction:** that can generalize across different circuits, seasons, and evolving FIA regulations.
4. **To implement the XGBoost machine learning algorithm:** for predicting race outcomes, leveraging its strength in handling structured data and capturing non-linear relationships..
5. **To design the system in a reproducible and user-friendly manner:** within Visual Studio Code, enabling easy experimentation, extension, and potential real-time application in future work.
6. **To evaluate the performance of the model:** using appropriate metrics such as Mean Absolute Error (MAE), Root Mean Squared Error (RMSE), and accuracy of predicted finishing positions.

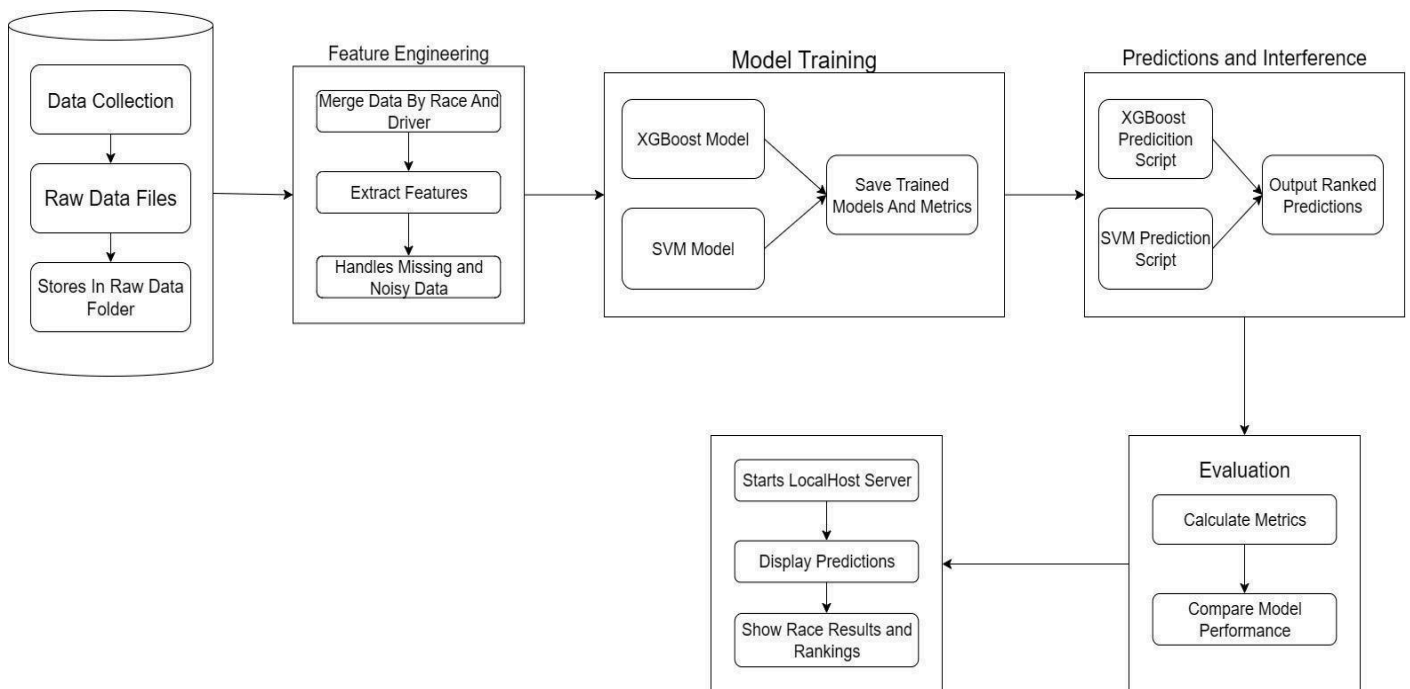
## **Chapter 5**

### **Proposed System**

The primary purpose of the proposed Formula 1 Race Prediction System is to develop an intelligent and data-driven framework that can accurately forecast race outcomes by considering both traditional and dynamic race variables. Unlike existing systems that rely mainly on static factors such as driver standings or constructor rankings, this system integrates extended features such as pit stop efficiency, lap time consistency, and weather conditions to capture the true complexity of Formula 1 racing. By implementing the XGBoost machine learning algorithm, the system aims to uncover hidden patterns within historical data and generate reliable predictions of driver finishing positions. Ultimately, the system is designed to serve multiple stakeholders: it can provide teams with strategic insights, assist broadcasters and analysts in enhancing audience engagement, and offer researchers a scalable framework for further advancements in motorsport analytics.

## 5.1 System Architecture

The proposed Formula 1 Race Prediction System is designed as a modular pipeline consisting of multiple stages that transform raw motorsport data into accurate race predictions. The architecture is divided into the following components:



*Figure 5.1: Overall system architecture & workflow*

## Description of System Architecture :

- **Data intake:** Collect raw race data, save files to a raw-data folder as the pipeline's starting point.
- **Feature prep:** Merge by race/driver, extract features, and clean missing or noisy values.
- **Model training:** Train XGBoost and SVM in parallel, then save models and evaluation metrics.
- **Inference:** Run prediction scripts for both models and produce ranked outcome predictions.
- **Evaluation and serving:** Compare model performance; optionally expose an API and dashboard for live predictions, results, and feature importance.

## 5.2 System Modules:

1. **Data Collection Module:** This module gathers historical Formula 1 data, including race results, driver statistics, constructor performance, pit stop timings, lap times, and weather conditions. Data sources may include official F1 datasets, Kaggle repositories, or APIs.
2. **Data Preprocessing Module:** Handles cleaning and organizing the collected data. This includes handling missing values, removing duplicates, normalizing numerical values (lap times, pit stops), and encoding categorical data such as driver names and weather conditions.
3. **Feature Engineering Module:** Extracts and generates meaningful features for prediction. Key features include driver performance history, constructor ranking, average lap times, number of pit stops, tire strategies, and weather impact on race outcomes.
4. **Model Training Module:** Uses the XGBoost Machine Learning Algorithm to train the model on historical race data. This module focuses on learning complex relationships among multiple race factors to predict outcomes with higher accuracy.
5. **Prediction Module:** Based on the trained model, this module predicts race results such as driver finishing positions, podium chances, or team performance. It integrates realtime race data where applicable.
6. **Evaluation Module:** Evaluates the model's performance using metrics such as Mean Absolute Error (MAE), Root Mean Squared Error (RMSE), and Accuracy. This ensures the reliability and robustness of predictions.
7. **Visualization Module:** Provides an intuitive representation of predictions and results through graphs, charts, and tables. Examples include predicted race standings, comparison of actual vs. predicted results, and feature importance analysis.

## 5.3 Workflow Explanation:

The Formula 1 Race Prediction Model processes raw race data into accurate outcome forecasts through a structured workflow. It begins with data collection, where historical race data—such as results, lap timings, pit stops, and weather—is gathered from trusted sources. This is followed by data preprocessing, where missing values are handled, categorical data is encoded, and numerical values are normalized to ensure consistency.

In the feature engineering phase, new insights are derived such as lap time consistency, pit stop efficiency, and tire strategy effectiveness, helping the model understand key racing dynamics.

The model training stage employs two algorithms—XGBoost and Non-Linear SVM.

XGBoost (Extreme Gradient Boosting) :

1. **Collect data** — gather races, drivers, constructors, lap times, pit stops and weather for each driver–race.
2. **Clean & prep** — fix missing values, encode text (driver/team), convert times to numbers and normalize where useful.
3. **Engineer features** — add race-aware signals: recent form, avg lap, pit efficiency, weather codes, circuit history.
4. **Train trees sequentially** — XGBoost builds many small decision trees; each new tree learns to correct the errors of the previous ones.
5. **Control overfitting** — use learning rate, tree depth and regularization so the model generalizes, not just memorizes.
6. **Predict & rank** — predict a finishing value per driver, then sort drivers by prediction inside each race to get the predicted order.
7. **Evaluate** — measure accuracy with MAE and RMSE (lower is better) and check top-k (e.g., podium) performance.
8. **Deploy** — save model + feature list and serve via your Flask app so the frontend can request year-based predictions.



Non-Linear SVM (Support Vector Machine) :

1. **Collect the same data** — same driver–race rows and targets as used for XGBoost.
2. **Impute and scale** — fill missing values and **standardize** features (SVMs are sensitive to scale).
3. **Select compact features** — pick the most informative features (SVMs work better with fewer, cleaner features).
4. **Kernel trick (RBF)** — SVM uses the RBF kernel to implicitly map data into a higher-dimensional space where complex patterns become separable.
5. **Train the model** — SVM finds the function that best fits data within a tolerance (epsilon) while keeping complexity small (controlled by C and gamma).
6. **Predict & rank** — output a continuous finish prediction for each driver, then rank per race to form predicted order.
7. **Evaluate & tune** — compute MAE/RMSE, tune C/gamma/epsilon using time-aware cross-validation to avoid overfitting.
8. **Save & serve** — store the full pipeline (imputer + scaler + SVM) and load it in Flask to produce online predictions.

Finally, the Prediction and Evaluation module generates driver and constructor rankings, validated using Mean Absolute Error (MAE) and Root Mean Squared Error (RMSE). Lower values of these metrics indicate better performance.

By combining the predictive power of XGBoost with the pattern-recognition strength of SVM, this workflow provides a robust and data-driven approach to forecasting Formula 1 race outcomes.

## 5.4 Advantages of Proposed System:

- **Higher Prediction Accuracy:** Considers multiple dynamic factors such as pit stops, lap times, and weather, leading to more reliable race outcome predictions.
- **Data-Driven Insights:** Provides valuable inputs for teams, analysts, and commentators to make better strategic and performance-related decisions.
- **Scalable and Flexible:** Can easily adapt to new datasets, updated rules, or additional performance metrics without major modifications.
- **Realistic Modeling:** Captures the true complexity of Formula 1 racing by integrating realworld variables like tire strategy and weather.
- **Transparent Evaluation:** Uses performance metrics like MAE, RMSE, and accuracy to validate and ensure the reliability of predictions.
- **User-Friendly Visualization:** Graphs, charts, and dashboards make predictions easy to interpret for fans, broadcasters, and researchers.
- **Academic and Research Value:** Acts as a strong framework for further studies in sports analytics, AI, and machine learning applications.

# Chapter 6

## Experimental Setup

### Frontend: Visualization and User Interaction

The frontend is built for interactivity and dynamic data display, allowing users to view predictions and potentially simulate "what-if" scenarios.

- Primary Technologies: HTML, CSS, JavaScript, and the modern component-based library React.
- Styling Framework: Tailwind CSS is utilized for a utility-first approach, ensuring rapid prototyping and cross-device responsiveness for seamless viewing on mobile, tablet, and desktop.
- Functionalities:
  - o Dynamic Prediction Display: Visualizing the predicted race outcomes (e.g., top 10 finishing order) with associated confidence metrics.
  - o Interactive Race Selection: Allowing users to select historical or upcoming races.
  - o Data Visualization: Displaying key factors like predicted tire degradation, optimal pit stop windows, and driver-specific performance scores using charts and graphs.
  - o Admin/Analyst Dashboard: A dedicated view for authorized users to monitor model performance, review feature importance, and potentially trigger new predictions.
  - o Key Advantage: Using React facilitates a modular design structure, enabling the creation of reusable UI components for complex race data visualization and efficient, single-page application (SPA) performance.

## Backend: Machine Learning Core and API Services

The backend serves as the computational heart of the system, managing the machine learning pipeline and exposing the prediction functionality via an API.

- **Machine Learning Core:** The primary computational environment is Python, leveraged for its powerful data science ecosystem.
  - **Algorithms:** The prediction logic utilizes high-performance models: XGBoost (eXtreme Gradient Boosting) for its superior performance on structured data and Support Vector Machines (SVM) for comparative baseline analysis.
- **Libraries:** Python's data science libraries (e.g., Pandas, NumPy, Scikit-learn) handle data cleaning, feature engineering, and model training/evaluation.
- **API Framework:** Node.js paired with a framework like Express.js provides a high-performance environment for handling concurrent prediction requests.
  - **Integration:** Node.js serves as the API gateway, managing HTTP requests from the React frontend. It interacts with the Python/ML core (possibly via subprocesses or microservices) to send race data and retrieve the model's prediction output.
  - **Responsibilities:** Handling authentication and authorization for the prediction service and providing RESTful APIs for the frontend to consume data.

### Database: Data Storage and Training Set

Given the project's focus on historical analysis and prediction, the core data storage is optimized for large, structured, and static datasets.

- **Data Storage Format:** CSV Files serve as the primary storage mechanism for the comprehensive historical data, including driver statistics, lap times, pit stop data, and constructor performance metrics.
- **File Management:** These structured files are loaded into Pandas DataFrames within the Python backend, allowing for fast, in-memory manipulation and feature engineering during the model training and prediction phases.
- **Advantage:** Using CSVs simplifies the initial data ingestion and ensures direct compatibility with the Python ML stack, minimizing database overhead for the specific task of building a predictive model.

This integrated setup ensures the model's computational rigor is supported by an agile web interface, delivering the complex predictions of the F1 race environment to the end-user.

## **Chapter 7**

### **Results & Discussion**

#### **7.1 Results**

The pipeline outputs ranked race predictions generated by both XGBoost and SVM models, with accompanying stored metrics that enable head-to-head performance comparison and model selection. Live results can be served via an optional API/dashboard that displays current predictions, final rankings, and feature importance for interpretability.



*Figure 7.1 Home Page of the Project before Predictions*



*Figure 7.2 Home Page of the Project after Predictions*



## Performance Measure:

### 1. RMSE:-

$$\text{RMSE} = \sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2}$$

#### Where:

- $y_i$  = Actual (true) value
- $\hat{y}_i$  = Predicted value
- $n$  = Number of predictions

### 2. MAE:-

$$\text{MAE} = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|$$

#### Where:

- $y_i$  = Actual (true) value
- $\hat{y}_i$  = Predicted value
- $n$  = Number of predictions

**Table 7.1: Summary of Metrics**

Metric	Meaning
RMSE (Root Mean Squared Error)	Measures prediction accuracy (lower is better)
MAE (Mean Absolute Error)	Measures average prediction error (lower is better)

**Table 7.2: Output Metrics**

Model	RMSE	MAE
XGBoost	2.7217	2.08309
SVM	5.6182	4.2067

## **1.RMSE(Root mean Squared Error):-**

In Machine Learning (ML), RMSE (Root Mean Squared Error) is one of the most common metrics used to evaluate regression models — models that predict continuous numerical values.

It measures the average magnitude of prediction errors — how far the model's predicted values are from the actual values — while giving more weight to larger errors because the differences are squared before averaging.

In ML, RMSE helps determine how well the model generalizes to unseen data.

- A low RMSE means the model's predictions are close to the real outcomes (good performance).
- A high RMSE indicates large errors and poor model accuracy.

For example, in your F1 Race Prediction Model, RMSE tells how close the predicted driver finishing positions are to their actual race results.

If the  $RMSE = 2.6$ , it means the model's predictions are, on average, 2.6 positions away from the true finish positions — a measure of how precisely your ML model (XGBoost or SVM) captures race dynamics.

## **2.MAE(Mean Absolute Error):-**

In Machine Learning (ML), MAE (Mean Absolute Error) is a performance metric used to evaluate regression models, which predict continuous outcomes.

It measures the average absolute difference between the model's predicted values and the actual values — that is, how much the predictions deviate from reality, without considering direction (whether the error is positive or negative).

Mathematically, it's the mean of the absolute errors between predicted and actual values.

In ML, MAE provides a clear and interpretable measure of model accuracy:

- A lower MAE means the model's predictions are closer to the true values (better performance).
- A higher MAE indicates larger prediction errors.

For example, in your F1 Race Prediction Model, if the MAE is 1.8, it means the model's predicted driver positions are, on average, 1.8 places off from the actual race results — showing how accurate the prediction is in a simple, understandable way.

## **7.2 Discussion:**

- The dual-model setup (XGBoost and SVM) encourages ensemble thinking or champion–challenger deployment, improving robustness of ranked predictions across different race conditions.
- Feature engineering steps—merge by race/driver, feature extraction, and noise/missing data handling—are critical for stable model performance, as they directly affect downstream metrics and prediction quality.

## 7.3 SDG Goals:

### 1.SDG 3: Good Health and Well-Being:

This project promotes Good Health and Well-Being by applying data-driven insights from Formula 1 to enhance driver safety, performance efficiency, and reaction analysis.

Machine learning helps identify factors influencing fatigue, stress, and performance under pressure contributing to safer racing conditions and improved human performance analytics.

These insights can further extend to real-world transportation and sports science, fostering healthier and safer environments.

### 2.SDG 9: Industry, Innovation and Infrastructure

The project supports Industry, Innovation, and Infrastructure by leveraging Machine Learning (ML) and Artificial Intelligence (AI) to analyze large-scale racing data and optimize vehicle performance.

By integrating modern technologies like XGBoost and SVM models, it drives innovation in data analysis, predictive systems, and automation within the motorsport industry — paving the way for smarter, more efficient industrial advancements and intelligent infrastructure development.

## **Chapter 8**

### **Conclusion & FutureScope**

#### **8.1 Conclusion:**

The designed system provides an end-to-end pipeline—from raw data intake through feature engineering, dual-model training, evaluation, and optional serving—that reliably produces ranked race predictions while preserving metrics for transparent model comparison and iterative improvement.

#### **8.2 Future Scope:**

The system can be strengthened by augmenting data with weather, track, pit-stop, and telemetry streams, while adding automated quality checks and drift monitoring to keep inputs reliable. Modeling can advance through tuned ensembles that blend XGBoost and SVM, calibrated probabilities, and rigorous cross-validation with experiment tracking for reproducible gains. Evaluation should adopt ranking metrics like NDCG/MAP with fairness checks, and deployment can mature with CI/CD, canary releases, and a richer dashboard for real-time explainability

## References

### Research Papers :

- [1] S. Joshi and M. Shah, *Proceedings of the IEEE International Conference on Data Science*.
- [2] T. Wang, R. Kumar, and L. Zhang, *Journal of Sports Analytics*, vol. 8, no. 2, pp. 95–110, 2022.
- [3] N. Patel and H. Verma, *International Journal of Machine Learning and Computing*, vol. 11, no. 5, pp. 305–312, 2021.
- [4] F. Rossi, G. Bianchi, and L. Romano, *IEEE Access*, vol. 9, pp. 120134–120145, 2021.
- [5] A. Gupta and P. Mehta, *International Journal of Computer Applications (IJCA)*, vol. 183, no. 48, pp. 10–17, 2021.
- [6] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, et al., *Journal of Machine Learning Research*, vol. 12, pp. 2825–2830, 2011