# Formula 1 Machine Learning Project

Stephan Karas



## TABLE OF CONTENTS





**Model Choice** 



Data Overview









# Problem Statement





The primary objective is to predict the finishing position of a Formula One driver in a race. We propose developing a predictive model incorporating time-sensitive qualifying data and weather data for short-term forecasts





# 02

# **Data Overview**



## Data Sources



developed by Chris Newell. It is a free service that provides a historical record of Formula One (F1) racing results and statistics.



#### **Visual Crossing**

The Visual Crossing Weather API provides historical, current, and forecast weather data for any location worldwide, enabling detailed weather insights for various applications.

## Dataset

To achieve this objective, we will use a comprehensive dataset that includes:

- **Driver Attributes:** Information such as age, nationality, career length, and average points per race.

- **Constructor Attributes:** Team references, nationality, and historical performance.

- **Race and Circuit Information:** Details like race year, round, and circuit specifics.

- Weather Conditions: Expected weather conditions on race day, sourced from the Visual Crossing API.

- **Circuit-Specific Historical Performance:** Past performance data of drivers on specific circuits.



## Dataset

Each row represents a driver's participation in a specific race.

	raceld	race_year	race_round	race_circuitId	circuitld	circuit_alt	tempmax	tempmin	humidity	precip	 circuit_country_Turkey	circuit_count
0	1	2009	1	1	1	10	76.7	47.1	50.7	0.000	0	
1	2	2009	2	2	2	18	89.7	75.5	84.8	1.075	0	
2	3	2009	3	17	17	5	71.4	60.9	87.3	0.370	0	
3	5	2009	5	4	4	109	68.2	59.4	78.7	0.009	0	
4	6	2009	6	6	6	7	75.8	62.3	77.6	0.000	0	
9043	939	2015	13	15	15	18	93.5	57.5	40.9	0.000	0	
9044	940	2015	14	22	22	45	86.7	66.0	73.7	0.315	0	
9045	944	2015	18	18	18	785	83.1	65.3	75.1	0.197	0	
9046	943	2015	17	32	32	2227	75.0	48.3	73.3	0.008	0	
9047	942	2015	16	69	69	161	63.2	57.8	89.1	0.688	0	

9048 rows × 90 columns



## Target Variable

The target variable for our model is the driver's finishing position, categorized as 'driver\_race\_result\_category', with the following classifications:

- 0: Finished outside the top 10.
- 1: Finished in the top 10 but outside the top 5.
- 2: Finished in the top 5 but did not win.
- 3: Won the race.



# Challenges

- Integrating multiple data sources, including historical race data from the Ergast API and weather data from Visual Crossing Weather API.
- Ensuring data consistency and synchronization between different sources.
  - Handling potential class imbalance in race outcome categories.
- Effectively incorporating real-time data (qualifying results) into the model.





# 03

# Data Preprocessing



#### DATA PRELOADING

import pandas as pd

drivers = pd.read\_csv('f1db\_csv\drivers.csv')
constructors = pd.read\_csv('f1db\_csv\constructors.csv')
races = pd.read\_csv(r'f1db\_csv\races.csv')
results = pd.read\_csv(r'f1db\_csv\results.csv')
circuits = pd.read\_csv('f1db\_csv\circuits.csv')
qualifying = pd.read\_csv('f1db\_csv\qualifying.csv')
sprint = pd.read\_csv('f1db\_csv\sprint\_results.csv')





### **Feature Engineering**

# Calculate historical points per race
results['points'] = results['points'].astype(float) # convert points to floaat

# Calculate the average points each driver has earned per race, grouped by their unique driver ID
avg\_points = results.groupby('driverId')['points'].mean().reset\_index() # Group by driver and calculate mean points

# Merge the average points data back into the drivers DataFrame # This adds the avg\_points column to the drivers DataFrame where the driverId matches drivers = drivers.merge(avg\_points, on='driverId', how='left')



# Merging DB Tables

Merging the relevant database tables, bringing together our dataset

- 1. Renaming the columns of each dataset to include its name
- 2. Merging Races & Circuits table (to initially incorporate date and location data)

races\_circuits = pd.merge(races, circuits, left\_on='race\_circuitId', right\_on='circuitId')



#### **Incorporating Weather Data**

- We used the Visual Crossing Weather API to include weather forecasts in our predictions.
- Relevant weather features included: tempmax, tempmin, humidity, precip, windspeed, and conditions
  - Example code to fetch and integrate weather data:

```
import requests
weather features = ['tempmax', 'tempmin', 'humidity', 'precip',
                                                                                       def get relevant weather features(location, date, api key):
'windspeed', 'conditions']
                                                                                         api url = f"https://weather.visualcrossing.com/VisualCrossingWebServices/rest/services/timeline/{location}/{date}"
for feature in weather features:
                                                                                         params = {
                                                                                           'unitGroup': 'us'.
  races circuits[feature] = None
                                                                                           'key': api key,
                                                                                           'include': 'obs'
races circuits['race date'] =
                                                                                         response = requests.get(api_url, params=params)
pd.to datetime(races circuits['race date']).dt.strftime('%Y-%m-%d'
                                                                                         if response.status code == 200:
                                                                                           data = response.json()
                                                                                           if 'days' in data and len(data['days']) > 0:
                                                                                             weather day = data['days'][0]
for index, row in races circuits.iterrows():
                                                                                             relevant weather features = {
                                                                                                'tempmax': weather day.get('tempmax'),
  weather data =
                                                                                                'tempmin': weather day.get('tempmin'),
get relevant weather features(str(row['circuit location']),
                                                                                               'humidity': weather day.get('humidity'),
                                                                                                'precip': weather day.get('precip').
str(row['race date']), "YOUR API KEY")
                                                                                               'windspeed': weather_day.get('windspeed'),
  if weather data:
                                                                                                'conditions': weather day.get('conditions')
      for feature in weather features:
                                                                                             return relevant weather features
         races circuits.at[index, feature] = weather data[feature]
                                                                                         else:
         print("Collected data for", str(row['circuit location']), ",",
                                                                                           print(f"Failed to retrieve data for {location} on {date}: {response.status code}")
                                                                                           return None
str(row['race date']))
```



# Merging DB Tables

Merging the relevant database tables to create a comprehensive dataset

1. Merging Results with Circuits and Races.

races\_circuits\_results = pd.merge(races\_circuits, results, on='raceId')

2. Creating the target variable, categorizing the results value





# Merging DB Tables

Merging the relevant database tables to create a comprehensive dataset

3. Merging the Drivers and Constructors, and Qualification tables

df quali = pd.merge(df cleaned, qualifying, on=['raceId', 'driverId'])



# **Cleaning weather data**

- Removed rows with NA weather data (older races) to ensure data quality

- Combined driver's forename and surname for a unique identifier.

Example code:

df\_filtered\_weather = races\_circuits.dropna(subset=['tempmax', 'tempmin',
'humidity', 'precip', 'windspeed', 'conditions'])
df\_filtered\_weather['fullname'] = df\_filtered\_weather['driver\_forename'] + '\_' +
df\_filtered\_weather['driver\_surname']
df\_filtered\_weather.to\_csv('data\_preview/filtered\_weather\_data.csv')

## **Removing Irrelevant Features**

Dropped columns that are not relevant to the prediction task to reduce dimensionality and improve model performance.

- Example code:

df\_cleaned = pd.read\_csv('data\_preview/filtered\_weather\_data.csv')
df\_quali = pd.merge(df\_cleaned, qualifying, on=['raceId', 'driverId'])

df\_cleaned = df\_quali.drop(columns=[
 'race\_url', 'circuit\_circuitRef', 'circuit\_url', 'result\_positionText',
 'driver\_forename', 'driver\_surname', 'driver\_number', 'driver\_driverRef',
 'driver\_code', 'driver\_age', 'driver\_url', 'constructor\_constructorRef',
 'constructor\_url', 'race\_fp1\_date', 'race\_fp1\_time', 'race\_fp2\_date',
 'race\_fp2\_time', 'race\_fp3\_date', 'race\_fp3\_time', 'race\_quali\_date',
 'race\_quali\_time', 'race\_sprint\_date', 'race\_sprint\_time', 'Unnamed: 0', 'quali\_qualifyId',
 'constructorId\_y', 'quali\_number', 'quali\_position', 'constructorId\_x'
])
df cleaned.to csv('data quali preview/cleaned data quali.csv')

# **Converting Qualifying Timings**

- Converted qualifying timings from string format to numerical seconds to facilitate model training
- Filled NaN values in qualifying columns with a large time value to handle missing data.

def time\_to\_seconds(time\_str):
 """Convert time string in format 'm:ss.mmm' to seconds."""
 if pd.isnull(time\_str):
 return np.nan
 mins, secs = time\_str.split(':')
 return int(mins) \* 60 + float(secs)

# Convert time strings to numerical format all\_quali\_columns = ['quali\_q1', 'quali\_q2', 'quali\_q3'] # Example column names for col in all\_quali\_columns: df\_na\_cleaned[col] = df\_na\_clean



- majority of drivers finished outside the top 10 (category 0)
- a smaller number finishing in the top 10 but not the top 5 (category 1)
- fewer still in the top 5 but not winning (category 2)
- the least number of drivers won the race (category 3).



# 04

# **Model Choice**



## RandomForest + Parameter Grid Tuning

- We chose the Random Forest algorithm due to its ability to handle large datasets and capture complex interactions between features.
- Initial parameter selection based on common starting values and domain knowledge.
  - Created a parameter grid for hyperparameter tuning.

param\_grid = {
 'n\_estimators': [100, 200, 300],
 'max\_depth': [10, 20, 30],
 'min\_samples\_split': [2, 5, 10],
 'min\_samples\_leaf': [1, 2, 4],
 'max\_features': ['sqrt', 'log2', None],
 'bootstrap': [True, False]

## How RandomForest Works for us

**Data**: We have a comprehensive dataset that includes our relevant features.

-

**Training**: Creates multiple decision trees, each trained on a different random subset of our dataset.

- allows each tree to learn from different parts of the data, such as specific race days, varying weather conditions, or different driver performance statistics.

**Prediction**: Each decision tree independently analyzes its subset of data and makes a prediction about the driver's finishing position.

I.e. one tree might consider how the driver performed in similar weather conditions, while another looks at qualifying results.

**Final Decision**: The Random Forest aggregates the predictions from all the trees. By using a majority vote system, it determines the most likely finishing position for the driver.



### Parameters choice

- Choosing the right parameters is crucial for the performance of our Random Forest model.
- n\_estimators: specifies the number of trees in the forest.

max\_depth: limits the number of nodes in the tree.

min\_samples\_split: sets minimum number of samples required in a node to split into two new nodes.

min\_samples\_leaf: min no. of samples required to be at a leaf node.

max\_features: no. of features to consider when looking for the best split.

bootstrap: boolean to use bootstrap samples when building trees. (Drawing random samples with replacement from the original dataset.)

## Ensuring Fair Train-Test Data Split

- We used the StratifiedShuffleSplit technique to split our dataset. This method ensures that both training and test sets maintain the same class distribution as the original dataset.
  - stratifying the data on driver\_race\_result\_category, we preserve the proportion of each race result category in both sets
  - provides a balanced and representative sample for training and evaluation.

- Applied SMOTE (Synthetic Minority Oversampling Technique) to the training set.
  - generating synthetic samples for the minority classes, enhancing their representation
  - improving the model's ability to learn from the underrepresented categories.
- set class\_weight='balanced' to ensure the model pays more attention to minority classes.

### RandomizedSearchCV Hyperparameter tuning

RandomizedSearchCV works by randomly sampling a fixed number of hyperparameter combinations from a specified range and evaluating each combination using cross-validation, then selecting the best combination based on performance.

- It uses K-fold cross-validation to evaluate the performance of each parameter combination.



## **Baseline Classifier: KNN**

- KNN classifies a data point by looking at the k closest data points (neighbors) in the dataset and assigning it to the most common class among those neighbors.
- Same Data preparation as RandomForest:
  - Grouped Train-Test Split (20% test size)
  - Standardized features with StandardScaler
- KNN (K-Nearest Neightbors):
  - Optimized k using RandomizedSearchCV
  - Tested odd values of k (1-29) to avoid ties (even k can result in equal votes from neighbors)
  - In KNN, a tie occurs when an equal number of neighbors vote for different classes.
- Results:
  - Best k: 27



# 05

# Model Evaluation





#### **Classification Report**

Best parameters: {'n\_estimators': 400, 'min\_samples\_split': 2, 'min\_samples\_leaf': 1, 'max\_features': 'sqrt', 'max\_depth': 100, 'bootstrap': False}

Accuracy	on te	est set: 0.59	944751381	21547	
Classific	catio	n Report on t	est set:		
		precision	recall	f1-score	support
	0	0.72	0.76	0.74	963
	1	0.40	0.30	0.35	423
	2	0.48	0.55	0.52	339
	3	0.44	0.44	0.44	85
accur	racy			0.60	1810
macro	avg	0.51	0.51	0.51	1810
weighted	avg	0.59	0.60	0.59	1810



#### **Confusion Matrix**



#### **Baseline Classifier: KNN**

Best paramete	ers: {'n_neig	hbors': 2	7}	
	precision	recall	f1-score	support
0	0.51	0.93	0.66	892
1	0.25	0.10	0.14	473
2	0.46	0.13	0.20	406
3	0.00	0.00	0.00	132
accuracy			0.49	1903
macro avg	0.31	0.29	0.25	1903
weighted avg	0.40	0.49	0.39	1903

# **Thank You!**

