



POLITECNICO
MILANO 1863

SCUOLA DI INGEGNERIA INDUSTRIALE
E DELL'INFORMAZIONE

Formula One Results Analysis with Documental Database approaches - SMBUD Project

Author(s): **Paolo Riva - 10938975**

Group Number: **15**

Academic Year: 2024-2025

Contents

Contents	i
-----------------	----------

1 Introduction	1
1.1 State of the Art	1
1.2 Project Specification	2
2 Dataset Structure	3
2.1 Results Database (results)	3
2.2 Drivers Database (drivers)	4
2.3 Circuits Database (circuits)	4
2.4 Races Database (races)	4
3 Commands and Queries	7
3.1 Dataset Format	7
3.2 Preprocessing	7
3.3 Interface and Import	9
3.4 Setup and Operators	9
3.4.1 Query syntax	9
3.4.2 Import Data	9
3.4.3 Specific operators	10
3.5 Queries	11
3.5.1 Query 1: Average Driver position	12
3.5.2 Query 2: Most Races in F1	14
3.5.3 Query 3: Most F1 Wins	15
3.5.4 Query 4: Most laps ever completed by a driver	17
3.5.5 Query 5: Most Points in F1	18
3.5.6 Query 6: Average Time of completion for each race	21
3.5.7 Query 7: Most Wins per Circuit	22

3.5.8	Query 8: Highest average speed in a race	25
3.5.9	Query 9: Most Podiums in F1	27
3.5.10	Query 10: Races finished disastrously	29
4	Conclusion	31
A	External Links	33

1 | Introduction

Formula One is widely regarded as the pinnacle of motorsport, as the major open-wheel competition available to compete in for all top car manufacturer, motorsport entities and drivers. Formula One nowadays involves more than one-thousand people per team, focused on building and updating the fastest track vehicles in the world in order to race at the highest of levels among the competition. Everything, from car design and race strategy, to brand reputation and economic return, comes down to the two single race results each driver of a considered team score in each race of the current season, so analyzing on-track performance and overall results allows teams to better identify and choose their drivers as the elite ambassadors of the team's heritage in the Formula One environment. With a rich history spanning decades, F1 has accumulated vast datasets documenting race results, driver statistics, and track performances that provide a foundation for analyzing trends, uncovering insights, and developing a deeper understanding of the sport.

In this chapter, the state-of-the-art and the hypotheses are presented, in order to provide an overview of the underlying context the project addresses with respect to Formula One history data, considering its recording and managing, and the technologies best suiting the task of data gathering and analysis among

1.1. State of the Art

The application of data analytics in sports has evolved rapidly in the past decade, and Formula One is no exception. Public databases of race results have opened the doors to explore patterns, assess driver performances, and evaluate track-specific trends.

This project wants to analyze results and available data in order to provide an overview of F1 driver performance in history, i.e. answering questions like:

- What drivers can be considered the greatest to ever take part in F1?
- Which races were the most challenging among history?
- What drivers have shown the highest consistency in their career?

Much of these questions' answers depend on the availability of clean, structured datasets, requiring assumptions about how results are stored, managed, and processed. Therefore, choosing a Documental Database technology to perform queries on data allows to gather several analytic perspectives by performing tailored queries on available driver results, due to flexibility in handling hierarchical and semi-structured data, as well as powerful aggregation operations.

1.2. Project Specification

This project aims at building and querying on a documental database handling sportive results. Having chosen Formula One as the sport of interest, the goal is to design and implement queries capable of analyzing available data on races, drivers, and track performances. MongoDB has been chosen as the database management system, allowing the import of CSV-formatted data as Documental data and allowing various aggregation and filtering operations to be performed in order to gather and analyze the data of interest.

The main objectives of the project include:

- Importing a properly structured database containing data about races, drivers, circuits, and results.
- Designing efficient queries to extract meaningful insights, such as driver performance metrics, track-specific trends, and comparative analysis.

2 | Dataset Structure

To analyze Formula One data effectively, we rely on a main database "results" and 3 secondary databases, all available publicly on Kaggle. While the first one contains all results in races achieved by any driver ever, the others come as a support to associate IDs of tracks and drivers with their details.

2.1. Results Database (results)

The **results** database is the backbone of our analysis, recording the detailed outcomes of every race. Each row represents a single driver's performance in a specific race. The key fields include:

- **resultId**: A unique identifier for each result entry.
- **raceId**, **driverId**, **constructorId**: Link to their respective races, driver corresponding to the result, and constructor corresponding to the result.
- **grid** and **position**: Driver's starting grid position and final finishing place.
- **points**: Points awarded to the driver based on their finishing position.
- **fastestLap**, **fastestLapTime**, **fastestLapSpeed**: Metrics for the driver's fastest lap during the race.
- **statusId**: Indicates whether the driver finished the race or faced an issue (e.g., mechanical failure or crash).
- **time**: Total time to complete the race (if not lapped).

The field has been **preprocessed** to allow a more consistent comparison with respect to MongoDB types.

2.2. Drivers Database (`drivers`)

The `drivers` database is about current and past drivers to ever compete in F1, with fields:

- **driverId**: A unique identifier for each driver.
- **forename** and **surname**: Real name of the driver.
- **dob** and **nationality**: Date of birth and National Flag.
- **code**: 3-letter broadcast identifier (e.g. "HAM")
- **number**: Driver's race number.
- **url**: Links to detailed Wiki profiles about each driver.

2.3. Circuits Database (`circuits`)

The `circuits` database is about current and past tracks to ever host a race in F1, with fields:

- **circuitId**: A unique identifier for each circuit.
- **name**: The official name of the track.
- **location** and **country**: Where in the world you'll find the circuit.
- **lat**, **lng**, and **alt**: Geographic coordinates and altitude.
- **url**: Links to more information about each circuit.

2.4. Races Database (`races`)

The `races` database about the details of a single *race* event, specifying each season's specific details about the event hosted at a specific circuit, with fields:

- **raceId**: The unique identifier for each race.
- **year** and **round**: Which season and which round of the season the race occurred in.
- **circuitId**: A link to the circuit where the race was held.
- **name**, **date**, and **time**: The race's official name and schedule.

- **fp1_date**, **fp2_time**, etc.: Dates and times for practice sessions, qualifying, and sprint races (if applicable).

3 | Commands and Queries

This section presents how the MongoDB environment was set up, the preprocessing on existing data and following MongoDB queries on the imported databases.

3.1. Dataset Format

The "Formula 1 World Championship (1950 - 2024)" dataset was retrieved on Kaggle as a publicly available dataset, including results updated up to the last 6 months of races held. Four CSV files were obtained: *results.csv*, *drivers.csv*, *races.csv*, *circuits.csv*, matching the importing description of *Section 2*.

3.2. Preprocessing

Preprocessing was required on the CSV *results.csv* in order to perform some of the queries more efficiently. When a Formula 1 race ends, the timing data usually shows the winner's total race duration (like "1 hour, 30 minutes, and 45.123 seconds"), while all other drivers' times are shown as gaps to the winner (like "+45.123 seconds" or "+1 minute, 30.456 seconds" behind). This preprocessing Python script converts the winner's time into plain seconds, then for each other driver, it takes their gap time, converts it to seconds too, and adds it to the winner's time. This gives us every driver's actual finish time in seconds, making it much easier to compare and analyze race performances between drivers in the same race.

```

1 def transform_time(input_csv, output_csv):
2     def parse_time_to_seconds(time_str):
3         if "+" in time_str: # Gap time, like "+minutes:seconds.
4             # milliseconds" or "+seconds.milliseconds"
5             time_str = time_str[1:] # Remove the '+' sign
6             parts = time_str.split(":")
7             if len(parts) == 2: # Format: "minutes:seconds.milliseconds"
2             minutes, seconds = map(float, parts)

```

```

8         return round(minutes * 60 + seconds, 3)
9     elif len(parts) == 1: # Format: "seconds.milliseconds"
10         return round(float(parts[0]), 3)
11     else:
12         raise ValueError(f"Unexpected time format: {time_str}")
13
14 def winner_time_to_seconds(winner_time):
15     parts = winner_time.split(":")
16     if len(parts) == 3: # Format: "hours:minutes:seconds.
17         milliseconds"
18         hours, minutes, seconds = map(float, parts)
19         return round(hours * 3600 + minutes * 60 + seconds, 3)
20     elif len(parts) == 2: # Format: "minutes:seconds.milliseconds"
21         minutes, seconds = map(float, parts)
22         return round(minutes * 60 + seconds, 3)
23     else:
24         raise ValueError(f"Unexpected winner time format: {
25             winner_time}")
26
27 with open(input_csv, mode="r", newline="", encoding="utf-8") as
28     infile, \
29     open(output_csv, mode="w", newline="", encoding="utf-8") as
30         outfile:
31     reader = csv.DictReader(infile)
32     fieldnames = reader.fieldnames
33     writer = csv.DictWriter(outfile, fieldnames=fieldnames)
34     writer.writeheader()
35
36     for row in reader:
37         if row["time"] == "\\N": # Skip rows where the driver has
38             retired or been lapped
39             writer.writerow(row)
40             continue
41
42         if row["positionOrder"] == "1": # For the winner
43             winner_time = row["time"]
44             winner_seconds = winner_time_to_seconds(winner_time)
45             row["time"] = winner_seconds
46         else: # For other drivers
47             gap_time = row["time"]
48             try:
49                 gap_seconds = parse_time_to_seconds(gap_time)
50                 row["time"] = round(winner_seconds + gap_seconds,
51                     3) # Add the gap to the winner's total time

```

```
46         except ValueError as e:
47             print(f"Skipping row with invalid time format: {row}")
48             continue
49
50     writer.writerow(row)
```

An outputted *f1_results_transformed.csv* file is used as preprocessed replacement of the original *results.csv* file, allowing more efficient queries overall on the **time** field.

3.3. Interface and Import

MongoDB Compass was chosen as the official GUI tool for MongoDB, from the official MongoDB website. After installation on Windows 11, the local server was configured to run on the default port 27017. The data import process used MongoDB Compass's import functionality, which allows direct CSV-to-document conversion while maintaining data (possibly mixed) types integrity. Each CSV file was mapped to a corresponding collection within the database, with appropriate data type conversions and indexing strategies applied to optimize query performance.

3.4. Setup and Operators

3.4.1. Query syntax

MongoDB Compass allows to perform aggregation queries by simply specifying aggregation operations inside branches ("[]"). The following queries are, for this reason, in such form. To perform such queries in a regular MongoDB Shell, it's enough to enclosure the following queries as:

```
1 db.results.aggregate([QUERY]);
```

3.4.2. Import Data

To import data and create Documents, it's enough to create a Local MongoDB connection, and import the 4 specified databases using the MongoDB Compass import function from JSON/CSV as:

- results.csv as *results*
- drivers.csv as *drivers*

- `circuits.csv` as *circuits*
- `racers.csv` as *racers*

3.4.3. Specific operators

\$lookup Operator

The `$lookup` operator performs a left outer join operation between collections in MongoDB, enabling the combination of documents from multiple collections within the same database. This operator requires four essential parameters:

- `from`: Specifies the target collection to join with
- `localField`: Defines the field from the input documents to match
- `foreignField`: Indicates the field to match against in the target collection
- `as`: Designates the name for the new array field that will contain the matching documents

When executed, `$lookup` adds a new array field to each input document. This array contains the matching documents from the foreign collection. If no matches are found, the array will be empty but still present in the result document. In our case, this allows to associate each result with the specific driver involved in it just from the specified ID, along with race and track details.

\$concat Operator

The `$concat` operator is a string manipulation tool that combines multiple strings into a single string value. It accepts an array of string expressions as input and concatenates them in the specified order. This operator proves particularly useful when:

- Combining multiple fields into a single string field
- Adding constant strings between field values
- Formatting output data for presentation or export

In our case, this allows to join both name and surname of the drivers.

\$switch Operator

The `$switch` operator provides conditional branching within MongoDB's aggregation pipeline, similar to switch-case statements in traditional programming languages. This

operator evaluates a series of cases and returns a value based on the first matching condition.

The operator requires two main components:

- **branches**: An array of possible cases, where each case contains:
 - **case**: The condition to evaluate
 - **then**: The value to return if the condition is true
- **default**: The value to return if no cases match (required to handle all possible scenarios)

The `$switch` operator evaluates each case in order until it finds a true condition, then returns the corresponding **then** value. In our case, it is used to evaluate the points achieved by the driver based on its finishing position.

`$first` Operator

The `$first` operator is an accumulator that returns the first document from a grouped result set in MongoDB's aggregation pipeline. It operates in two distinct contexts: group stages and array manipulation.

Within a `$group` stage, `$first` returns the first value it encounters in each group, making it particularly useful to retrieve initial documents in a list.

`$cond` Operator

The `$cond` operator provides if-then-else logical operations within MongoDB's aggregation pipeline. This operator can be expressed in two syntactical forms: a ternary expression or an object expression.

The operator requires three components:

- **if**: The boolean expression to evaluate
- **then**: The value to return when the condition is true
- **else**: The value to return when the condition is false

3.5. Queries

3.5.1. Query 1: Average Driver position

This query evaluates the average grid positions of all drivers with at least one result, returning the drivers in descending order.

```
1 [
2   {
3     // Grouping by driverId to calculate total wins and average grid
      position
4     $group: {
5       _id: "$driverId",
6       avgGridPosition: { $avg: "$grid" }
7     }
8   },
9   {
10    $match: {
11      avgGridPosition: { $ne: 0 }
12    }
13  },
14  {
15    // Join with the drivers_id collection to get driver details
16    $lookup: {
17      from: "drivers",
18      localField: "_id",
19      foreignField: "driverId",
20      as: "driverDetails"
21    }
22  },
23  {
24    // Unwind the driverDetails array to flatten it
25    $unwind: "$driverDetails"
26  },
27  {
28    //Project final fields: driver name, total wins, and avg grid
      position
29    $project: {
30      _id: 0,
31      driverName: {
32        $concat: [
33          "$driverDetails.forename",
34          " ",
35          "$driverDetails.surname"
36        ]
37      },
38      totalWins: 1,
```



```
39     avgGridPosition: {
40       $round: ["$avgGridPosition", 2]
41     } // Round to 2 decimals
42   }
43 },
44
45 {
46   // Sort by total wins in descending order
47   $sort: { avgGridPosition: 1 }
48 }
49 ]
```

Returned:

```
1  [{
2    "driverName": "Pierre-Henri Raphanel",
3    "avgGridPosition": 1.06
4  },
5  {
6    "driverName": "Juan Fangio",
7    "avgGridPosition": 2.48
8  },
9  {
10   "driverName": "Ayrton Senna",
11   "avgGridPosition": 3.13
12 },
13 {
14   "driverName": "Jack McGrath",
15   "avgGridPosition": 3.17
16 },
17 {
18   "driverName": "Nino Farina",
19   "avgGridPosition": 3.19
20 },
21 {
22   "driverName": "Neville Lederle",
23   "avgGridPosition": 3.33
24 },
25 {
26   "driverName": "Emilio de Villota",
27   "avgGridPosition": 3.5
28 },
29 {
30   "driverName": "Jim Clark",
31   "avgGridPosition": 3.6
```

```
32 }  
33 ...
```

3.5.2. Query 2: Most Races in F1

This query computes the total races any driver has ever took part in, sorting the results in descending order.

```
1 [  
2   {  
3     // Group by driverId and count the total races  
4     $group: {  
5       _id: "$driverId", // Group by driverId  
6       totalRaces: { $sum: 1 } // Count the total number of races  
7     }  
8   },  
9   {  
10    // Lookup to join with the drivers_id collection  
11    $lookup: {  
12      from: "drivers", // The collection with driver details  
13      localField: "_id", // driverId from the grouped data  
14      foreignField: "driverId", // driverId field in drivers_id  
15      as: "driverDetails" // Alias for the joined data  
16    }  
17  },  
18  {  
19    // Unwind the driverDetails array to normalize  
20    $unwind: "$driverDetails"  
21  },  
22  {  
23    // Project the desired fields: name and total races  
24    $project: {  
25      _id: 0,  
26      driverName: {  
27        $concat: [  
28          "$driverDetails.forename",  
29          " ",  
30          "$driverDetails.surname"  
31        ]  
32      }, // Combine forename and surname  
33      totalRaces: 1 // Include the total races  
34    }  
35  },  
36  {
```

```
37 // Sort the results by totalRaces in descending order
38 $sort: { totalRaces: -1 }
39 }
40 ]
```

Returned:

```
1 [{
2   "totalRaces": 392,
3   "driverName": "Fernando Alonso"
4 },
5 {
6   "totalRaces": 352,
7   "driverName": "Kimi Raikkonen"
8 },
9 {
10  "totalRaces": 344,
11  "driverName": "Lewis Hamilton"
12 },
13 {
14  "totalRaces": 326,
15  "driverName": "Rubens Barrichello"
16 },
17 {
18  "totalRaces": 309,
19  "driverName": "Jenson Button"
20 },
21 {
22  "totalRaces": 308,
23  "driverName": "Michael Schumacher"
24 },
25 {
26  "totalRaces": 300,
27  "driverName": "Sebastian Vettel"
28 },
29 ...
```

3.5.3. Query 3: Most F1 Wins

This query computes the total number of wins for each driver, and returns the most winning drivers in history in descending order.

```
1 [
2   {
```

```

3      // Filter race results to include only wins (position: 1)
4      $match: { position: 1 }
5  },
6  {
7      // Group by driverId and count the number of wins
8      $group: {
9          _id: "$driverId", // Group by driverId
10         totalWins: { $sum: 1 } // Count total wins
11     }
12 },
13 {
14     // Lookup to join with the drivers_id collection
15     $lookup: {
16         from: "drivers", // The collection with driver details
17         localField: "_id", // driverId from the grouped data
18         foreignField: "driverId", // driverId field in drivers_id
19         as: "driverDetails" // Alias for the joined data
20     }
21 },
22 {
23     // Unwind the driverDetails array to normalize
24     $unwind: "$driverDetails"
25 },
26 {
27     // Step 5: Project the desired fields: name and total wins
28     $project: {
29         _id: 0, // Hide the default _id
30         driverName: {
31             $concat: [
32                 "$driverDetails.forename",
33                 " ",
34                 "$driverDetails.surname"
35             ]
36         }, // Combine forename and surname
37         totalWins: 1 // Include the total wins
38     }
39 },
40 {
41     // Step 6: Sort the results by totalWins in descending order
42     $sort: { totalWins: -1 }
43 }
44 ]

```

Returning:

```
1 [{
2   "totalWins": 104,
3   "driverName": "Lewis Hamilton"
4 },
5 {
6   "totalWins": 91,
7   "driverName": "Michael Schumacher"
8 },
9 {
10  "totalWins": 61,
11  "driverName": "Max Verstappen"
12 },
13 {
14  "totalWins": 53,
15  "driverName": "Sebastian Vettel"
16 },
17 {
18  "totalWins": 51,
19  "driverName": "Alain Prost"
20 },
21 {
22  "totalWins": 41,
23  "driverName": "Ayrton Senna"
24 },
25 {
26  "totalWins": 32,
27  "driverName": "Fernando Alonso"
28 },
29 ...
```

3.5.4. Query 4: Most laps ever completed by a driver

This query computes the maximum laps ever completed by a driver in F1 during his career.

```
1 [
2   {
3     $group: {
4       _id: "$driverId", // Group by driverId
5       totalLaps: { $sum: "$laps" }
6     }
7   },
8   {
```

```

9      $sort: { totalLaps: -1 } // Sort descending
10    },
11    {
12      $limit: 1 // Show the highest result
13    }
14  ]

```

Returning:

```

1  [{
2    "_id": 4,
3    "totalLaps": 21173
4  }]

```

3.5.5. Query 5: Most Points in F1

This query shows the drivers who achieved the most points in F1. Although the *points* field exists, we want to compare all drivers in history equally, since the point system in F1 actually changed between 2009 and 2011 assigning much more points to the top 10 grid.

```

1  [
2    {
3      // Add points based on the modern scoring system
4      $project: {
5        driverId: 1,
6        points: {
7          $switch: {
8            branches: [
9              {
10                 case: { $eq: ["$position", 1] },
11                 then: 25
12               },
13               {
14                 case: { $eq: ["$position", 2] },
15                 then: 18
16               },
17               {
18                 case: { $eq: ["$position", 3] },
19                 then: 15
20               },
21               {
22                 case: { $eq: ["$position", 4] },
23                 then: 12

```

```

24         },
25         {
26             case: { $eq: ["$position", 5] },
27             then: 10
28         },
29         {
30             case: { $eq: ["$position", 6] },
31             then: 8
32         },
33         {
34             case: { $eq: ["$position", 7] },
35             then: 6
36         },
37         {
38             case: { $eq: ["$position", 8] },
39             then: 4
40         },
41         {
42             case: { $eq: ["$position", 9] },
43             then: 2
44         },
45         {
46             case: { $eq: ["$position", 10] },
47             then: 1
48         }
49     ],
50     default: 0 // Positions outside top 10 get 0 points
51 }
52 }
53 }
54 },
55 {
56     // Group by driverId to calculate total points
57     $group: {
58         _id: "$driverId",
59         totalPoints: { $sum: "$points" }
60     }
61 },
62 {
63     // Join with the drivers_id collection to get driver names
64     $lookup: {
65         from: "drivers",
66         localField: "_id",
67         foreignField: "driverId",

```

```
68     as: "driverDetails"
69   }
70 },
71 {
72   // Unwind the driverDetails array
73   $unwind: "$driverDetails"
74 },
75 {
76   // Project the driver name and total points
77   $project: {
78     _id: 0,
79     driverName: {
80       $concat: [
81         "$driverDetails.forename",
82         " ",
83         "$driverDetails.surname"
84       ]
85     },
86     totalPoints: 1
87   }
88 },
89 {
90   // Sort by total points in descending order
91   $sort: { totalPoints: -1 }
92 }
93 ]
```

Returning:

```
1 [{
2   "totalPoints": 5048,
3   "driverName": "Lewis Hamilton"
4 },
5 {
6   "totalPoints": 3890,
7   "driverName": "Michael Schumacher"
8 },
9 {
10  "totalPoints": 3287,
11  "driverName": "Sebastian Vettel"
12 },
13 {
14  "totalPoints": 3143,
15  "driverName": "Fernando Alonso"
16 },
```



```
17 {
18   "totalPoints": 2791,
19   "driverName": "Kimi Raikkonen"
20 },
21 {
22   "totalPoints": 2730,
23   "driverName": "Max Verstappen"
24 },
25 {
26   "totalPoints": 2483,
27   "driverName": "Alain Prost"
28 },
29 ...
```

3.5.6. Query 6: Average Time of completion for each race

This query evaluates the average time of completion for each race ever held in history, to highlight the longest races ever completed by drivers in terms of total time.

```
1 [
2   {
3     $match: {
4       time: { //Match only drivers with full laps, not DNF, not Lapped
5         $exists: true,
6         $ne: null
7       }
8     }
9   },
10  {
11    $group: { // Evaluate Average time
12      _id: "$raceId",
13      averageTime: {
14        $avg: "$time"
15      }
16    }
17  },
18  {
19    $sort: {
20      averageTime: 1
21    }
22  },
23  {
24    $project: {
25      _id: 0,
```

```
26     raceId: "$_id",
27     averageTime: 1
28   }
29 }
30 ]
```

Returning:

```
1 [{
2   "averageTime": 226.323950000000002,
3   "raceId": 1063
4 },
5 {
6   "averageTime": 1527.72299999999997,
7   "raceId": 320
8 },
9 {
10    "averageTime": 2574.25,
11    "raceId": 579
12 },
13 {
14    "averageTime": 3534.8642857142854,
15    "raceId": 587
16 },
17 {
18    "averageTime": 3690.9816,
19    "raceId": 441
20 },
21 {
22    "averageTime": 4066.579090909091,
23    "raceId": 540
24 },
25 {
26    "averageTime": 4124.138,
27    "raceId": 154
28 },
29 ...
```

3.5.7. Query 7: Most Wins per Circuit

This query joins both circuit and race details to the result documents in order to compute the drivers with the highest number of wins on each track, showing for each circuit the driver that achieved the highest number of wins in such track.

```
1  [
2    {
3      // Filter for race winners
4      $match: { position: 1 }
5    },
6    {
7      // Lookup the race details to get the circuitId
8      $lookup: {
9        from: "races",
10       localField: "raceId",
11       foreignField: "raceId",
12       as: "raceDetails"
13     }
14   },
15   {
16     //Unwind the raceDetails array to access circuitId
17     $unwind: "$raceDetails"
18   },
19   {
20     // Lookup the circuit details to get circuit data
21     $lookup: {
22       from: "circuits",
23       localField: "raceDetails.circuitId",
24       foreignField: "circuitId",
25       as: "circuitDetails"
26     }
27   },
28   {
29     // Unwind the circuitDetails array
30     $unwind: "$circuitDetails"
31   },
32   {
33     // Group by circuit and driver, and count wins
34     $group: {
35       _id: {
36         circuitId: "$circuitDetails.circuitId",
37         driverId: "$driverId"
38       },
39       circuitName: {
40         $first: "$circuitDetails.name"
41       }, // Get circuit name
42       driverId: { $first: "$driverId" }, // Get driverId
43       wins: { $sum: 1 } // Count wins
44     }
45   }
46 ]
```

```
45 },
46 {
47     // Group by circuit to find the driver with the most wins
48     $group: {
49         _id: "$_id.circuitId", // Group by circuitId
50         circuitName: { $first: "$circuitName" }, // Include circuit name
51         driverId: { $first: "$driverId" }, // Driver with the most wins
52         wins: { $max: "$wins" } // Maximum wins
53     }
54 },
55 {
56     // Lookup driver details to get the driver's name
57     $lookup: {
58         from: "drivers",
59         localField: "driverId",
60         foreignField: "driverId",
61         as: "driverDetails"
62     }
63 },
64 {
65     // Unwind the driverDetails array
66     $unwind: "$driverDetails"
67 },
68 {
69     // Project the final output
70     $project: {
71         _id: 0,
72         circuitName: 1,
73         driverName: {
74             $concat: [
75                 "$driverDetails.forename",
76                 " ",
77                 "$driverDetails.surname"
78             ]
79         },
80         wins: 1
81     }
82 },
83 {
84     // Sort by circuit name (alphabetic order)
85     $sort: { circuitName: 1 }
86 }
87 ]
```

Returning:

```
1 [{
2   "circuitName": "Aintree",
3   "wins": 2,
4   "driverName": "Stirling Moss"
5 },
6 {
7   "circuitName": "Albert Park Grand Prix Circuit",
8   "wins": 4,
9   "driverName": "David Coulthard"
10 },
11 {
12   "circuitName": "Autodromo Enzo e Dino Ferrari",
13   "wins": 7,
14   "driverName": "Damon Hill"
15 },
16 {
17   "circuitName": "Autodromo Internazionale del Mugello",
18   "wins": 1,
19   "driverName": "Lewis Hamilton"
20 },
21 {
22   "circuitName": "Autodromo Nazionale di Monza",
23   "wins": 5,
24   "driverName": "Graham Hill"
25 },
26 }
```

3.5.8. Query 8: Highest average speed in a race

This query averages the drivers' fastest lap on each race event to average the highest average speed ever observed in a race. This query aims at showing the event where there raced the fastest generation of cars on a single track in history.

```
1 [
2   {
3     // Filter records where fastestLapSpeed exists
4     $match: { fastestLapSpeed: { $ne: null } }
5   },
6   {
7     // Group by raceId and calculate the average fastestLapSpeed
8     $group: {
9       _id: "$raceId",
10      avgFastestLapSpeed: {
```

```

11     $avg: "$fastestLapSpeed"
12   }
13 }
14 },
15 {
16   // Sort by the average fastestLapSpeed in descending order
17   $sort: { avgFastestLapSpeed: -1 }
18 },
19 {
20   // Limit to the race with the highest average fastestLapSpeed
21   $limit: 1
22 },
23 {
24   // Lookup race details for the corresponding raceId
25   $lookup: {
26     from: "races",
27     localField: "_id",
28     foreignField: "raceId",
29     as: "raceDetails"
30   }
31 },
32 {
33   // Unwind the race details
34   $unwind: "$raceDetails"
35 },
36 {
37   // Project the final result with race details and average speed
38   $project: {
39     _id: 0,
40     raceName: "$raceDetails.name",
41     raceDate: "$raceDetails.date",
42     circuitId: "$raceDetails.circuitId",
43     avgFastestLapSpeed: 1
44   }
45 }
46 ]

```

Returning:

```

1 [{
2   "avgFastestLapSpeed": 250.63199999999998,
3   "raceName": "Italian Grand Prix",
4   "raceDate": {
5     "$date": "2004-09-12T00:00:00.000Z"
6   },

```

```
7   "circuitId": 14
8 }}
```

3.5.9. Query 9: Most Podiums in F1

This query aims at showing the 10 drivers with the most number of podiums in F1.

```
1  [
2    {
3      $match: { position: { $lte: 3 } } // Positions 1, 2, 3
4    },
5    {
6      $group: {
7        _id: "$driverId", // Group by driverId
8        podiums: { $sum: 1 }
9      }
10   },
11
12   {
13     // Join with the drivers_id collection to get driver details
14     $lookup: {
15       from: "drivers",
16       localField: "_id",
17       foreignField: "driverId",
18       as: "driverDetails"
19     }
20   },
21   {
22     // Unwind the driverDetails array to flatten it
23     $unwind: "$driverDetails"
24   },
25   {
26     // Project final fields: driver name, total wins, and avg grid
27     // position
28     $project: {
29       _id: 0,
30       driverName: {
31         $concat: [
32           "$driverDetails.forename",
33           " ",
34           "$driverDetails.surname"
35         ]
36       },
37       podiums: 1
38     }
39   }
40 ]
```

```
37     }
38   },
39
40   {
41     $sort: { podiums: -1 }
42   },
43   { $limit: 10 }
44 ]
```

Returning:

```
1 [{
2   "podiums": 199,
3   "driverName": "Lewis Hamilton"
4 },
5 {
6   "podiums": 155,
7   "driverName": "Michael Schumacher"
8 },
9 {
10    "podiums": 122,
11    "driverName": "Sebastian Vettel"
12 },
13 {
14    "podiums": 107,
15    "driverName": "Max Verstappen"
16 },
17 {
18    "podiums": 106,
19    "driverName": "Alain Prost"
20 },
21 {
22    "podiums": 106,
23    "driverName": "Fernando Alonso"
24 },
25 {
26    "podiums": 103,
27    "driverName": "Kimi Raikkonen"
28 },
29 {
30    "podiums": 80,
31    "driverName": "Ayrton Senna"
32 },
33 {
34    "podiums": 68,
```



```
35     "driverName": "Rubens Barrichello"
36 },
37 {
38     "podiums": 67,
39     "driverName": "Valtteri Bottas"
40 }]
```

3.5.10. Query 10: Races finished disastrously

This query returns the number of races with the highest number of crashes, technical DNFs and mistakes in F1 history, specifically counting all the races where the number of drivers that finished the race (`statusId == 1`) was lower than the ones that didn't.

```
1  [
2    {
3      // Group results by raceId and calculate counts for statusId == 1
4      // and statusId > 1
5      $group: {
6        _id: "$raceId",
7        totalFinishedRegularly: {
8          $sum: {
9            $cond: [{ $eq: ["$statusId", 1] }, 1, 0]
10         }
11      },
12      totalNotFinishedRegularly: {
13        $sum: {
14          $cond: [{ $gt: ["$statusId", 1] }, 1, 0]
15        }
16      }
17    },
18    {
19      // Filter only races where not finished > finished regularly
20      $match: {
21        $expr: {
22          $gt: [
23            "$totalNotFinishedRegularly",
24            "$totalFinishedRegularly"
25          ]
26        }
27      }
28    },
29    {
30      //Count the number of races
```

```
31     $count: "racesWithMoreNonFinishers"  
32   }  
33 ]
```

Returning:

```
1 [{  
2   "racesWithMoreNonFinishers": 893  
3 }]
```

4 | Conclusion

This project aimed at showing different analytics regarding Formula One over the history of the sport demonstrating the use of MongoDB in analyzing several data entries. By leveraging advanced aggregation techniques, such as grouping, filtering, and lookups, we were able to meaningful insights into drivers' performances, race trends, and circuit statistics. The flexibility of MongoDB allowed for efficient handling of complex relationships between datasets, such as driver details, race results, and circuit information. Although an agglomerated grouping of the four databases involved could have also be considered by utilizing subdocument structures inside Results document, which was not chosen to avoid excessive redundancy, this analysis highlighted the potential of NoSQL databases for exploring and uncovering patterns in large, interconnected datasets, providing valuable perspectives on the evolution of Formula One and its entities.

A | External Links

- Formula 1 World Championship (1950 - 2024) Dataset:
<https://www.kaggle.com/datasets/rohanrao/formula-1-world-championship-1950-2020>
- MongoDB Compass: <https://www.mongodb.com/products/tools/compass>

