# DATA ANALYSIS AND VISUALIZATION IN AWS

Learn how to perform serverless analysis in AWS using visual data analysis tools- AWS Glue DataBrew and Amazon QuickSight. A data analyst will be able to build an end-to-end data analysis tool following this workshop.

*Saunak Chandra & Arturo Duarte, AWS PSA Analytics*
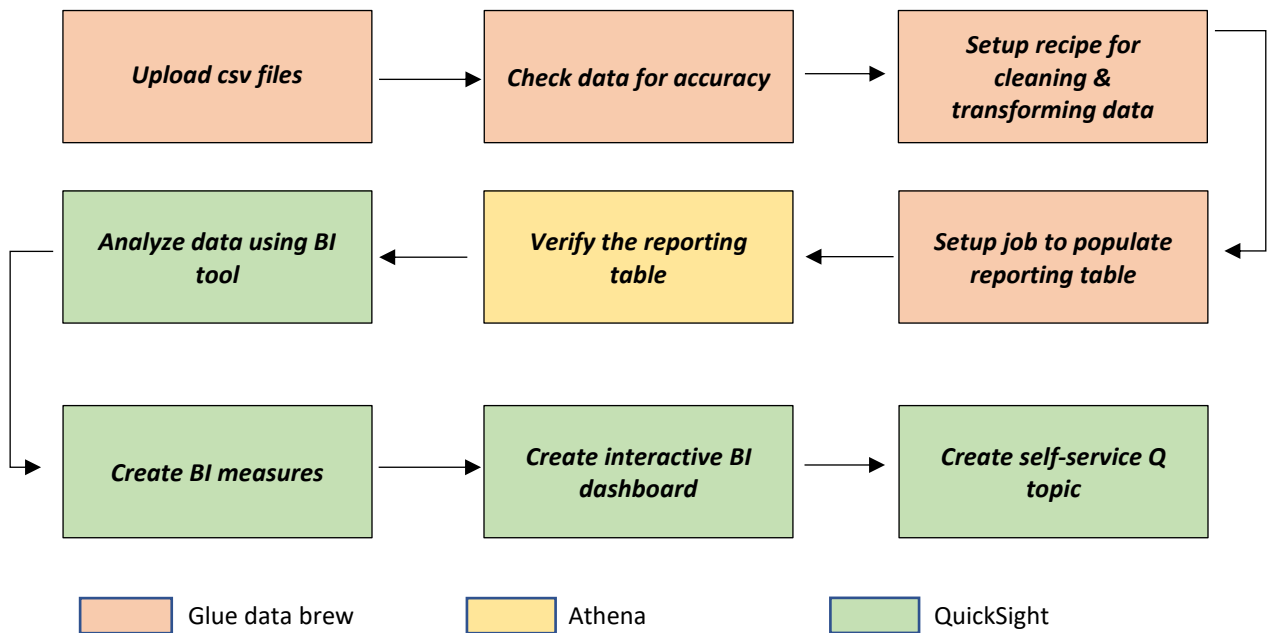
# Problem statement:

You are a data analyst for Formula 1 business group. You are given several CSV files containing historical race results. Your job is to perform a data analysis- look for any inconsistencies that may exist in the data and provide a self-service BI tool- which a formula 1 enthusiast can ask various questions to get answers from.

The challenge is this data is not ready for data analysis such as 1/ this data is broken into several csv files containing **race, results, constructors, circuits** etc., 2/ there are data quality issues such as unknown values, 3/ some attributes in the data are too granular to perform meaningful analysis; an example is age of the driver.

You are using AWS services to perform these jobs. The objective here is to build a data analysis and visualization tool-

1.  That can perform initial analysis of the data to check for consistency and accuracy. In case of missing or invalid data for certain fields you should change those values using mathematical functions such as mean value, most frequent values.
2.  Build BI dashboard to interactively explore formula 1 facts, trends and KPIs. The end users also can to ask ad-hoc questions in natural language to get answers from this tool.
3.  Setup an end-to-end process to automate data ingestion, data cleaning, data transformation; then populate your final tables with this transformed data and finally refresh the reports.

In AWS you can perform the above steps using AWS Glue Databrew, Amazon Athena and Amazon QuickSight. The logical steps to automate the data analysis process is captured in the below diagram-

| Glue data brew | | Athena | | QuickSight |

## Pre-requisite:

Create an S3 bucket

1. Create an s3 bucket named **f1-data-<AWS-account_number>** to store the CSV files for this workshop.

Amazon S3 > Create bucket

# Create bucket Info

Buckets are containers for data stored in S3. Learn more 🗗

## General configuration

**Bucket name**

f1-data-72⋯

Bucket name must be unique and must not contain spaces or uppercase letters. **See rules for bucket naming** 🗗

**AWS Region**

US East (N. Virginia) us-east-1 ▼

Copy settings from existing bucket - *optional*

Only the bucket settings in the following configuration are copied.

Choose bucket

## Object Ownership Info

Control ownership of objects written to this bucket from other AWS accounts and the use of access control lists (ACLs). Object ownership determines who can specify access to objects.

○ **ACLs disabled (recommended)**
All objects in this bucket are owned by this account. Access to this bucket and its objects is specified using only policies.

○ **ACLs enabled**
Objects in this bucket can be owned by other AWS accounts. Access to this bucket and its objects can be specified using ACLs.

Object Ownership

Bucket owner enforced

2. Create an S3 folder *input-files*.

Amazon S3 > f1-data > Create folder

## Create folder Info

Use folders to group objects in buckets. When you create a folder, S3 creates an object using the name that you specify followed by a slash (/). This object then appears as folder on the console. Learn more [Z]

(i) **Your bucket policy might block folder creation**
If your bucket policy prevents uploading objects without specific tags, metadata, or access control list (ACL) grantees, you will not be able to create a folder using this configuration. Instead, you can use the upload configuration to upload an empty folder and specify the appropriate settings.

### Folder

Folder name

input-files

Folder names can't contain "/". See rules for naming [Z]

### Server-side encryption

(i) The following settings apply only to the new folder object and not to the objects contained within it.

Server-side encryption
O  Disable
O  Enable

Cancel    **Create folder**

3. Upload the below CSV files into the above S3 folder.

https://awspsa-sampledata.s3.amazonaws.com/f1data/results/results.csv
https://awspsa-sampledata.s3.amazonaws.com/f1data/races/races.csv
https://awspsa-sampledata.s3.amazonaws.com/f1data/constructors/constructors.csv
https://awspsa-sampledata.s3.amazonaws.com/f1data/circuits/circuits.csv
https://awspsa-sampledata.s3.amazonaws.com/f1data/driver/drivers.csv

# 1. Data Preparation

## 1.1 Create Glue Databrew Datasets

1. Navigate to **AWS Glue DataBrew** service by logging into your AWS account. Select the AWS region of your preference; if none mentioned use US East (N. Virginia).
2. Click **DATASETS** from the left pane. Click on **Connect new dataset**.
3. Enter *f1-race-results* as **Dataset name**.
4. Select **Amazon S3** under **Connect to new dataset.**
5. Enter the below S3 path for **Enter your source from S3**.

```
s3://f1-data-<AWS-Account_number>/input-files/results.csv
```

6. Scroll down the page until the **Preview data** is visible. Note: you may need to expand the **Additional configurations** if Preview data is not visible.
7. Click **Create dataset** in the New dataset details page.



8. (Optional)

Click on Preview data and examine the various fields in the f1-race-results dataset in the Grid format. Make yourself aware of the fields available. Note the below-

    I.    There are various id fields- *resultId, raceId, driverId, constructorId*
    II.    Numerical field names are prefixed with "#" and text field names with "ABC".
    III.    Click Done to close the preview screen.

## Dataset preview                                                      ✕

| Dataset name | Data source | S3 location | Created by |
|---|---|---|---|
| - | - | s3://▪▪▪▪▪/f1d ata/results/results.csv | - |

| Size | Rows | Columns | Last modified by |
|---|---|---|---|
| 1.6 MB | 284 | 18 | - |

**Selected file type**
Format of the selected file

[ CSV ▼ ]

**CSV delimiter**

[ Comm... ▼ ]

**Column header values**

● **Treat first row as header**
The first row in your dataset will be treated as column header values

○ **Add default header**
Default headers will be added with values Column_1, Column_2 ...

[ ▦ Grid ]  [ ▥ Schema ]  [ T Text ]  [ ⊨ Tree ]

| # resultId | # raceId | # driverId | # |
|---|---|---|---|
| 1 | 18 | 1 | 1 |
| 2 | 18 | 2 | 2 |
| 3 | 18 | 3 | 3 |
| 4 | 18 | 4 | 4 |
| 5 | 18 | 5 | 1 |
| 6 | 18 | 6 | 3 |
| 7 | 18 | 7 | 5 |
| 8 | 18 | 8 | 6 |
| 9 | 18 | 9 | 2 |
| 10 | 18 | 10 | 7 |
| 11 | 18 | 11 | 8 |
| 12 | 18 | 12 | 4 |
| 13 | 18 | 13 | 6 |
| 14 | 18 | 14 | 9 |
| 15 | 18 | 15 | 7 |
| 16 | 18 | 16 | 10 |
| 17 | 18 | 17 | 9 |
| 18 | 18 | 18 | 11 |

[ Done ]

Similarly create datasets- ***f1-races, f1-circuits, f1-drivers and f1-constructors.*** The S3 locations for these datasets are
- f1-races
- f1-circuits
- f1-drivers
- f1-constructors

At the end of these steps you have 5 datasets created as shown below

## 1.2   Create Glue Databrew Project

Now the datasets are uploaded you are ready to check the data for accuracies and then perform data cleaning and transformation. This is essential step before your proceed further for BI analysis as any inaccuracies that may exist in the data may provide inconsistent view of the reports.

Data cleaning and transformation is done in Glue DataBrew project. A Glue DataBrew project is the workspace for data analysis and transformation. When you create a project in Glue DataBrew, it downloads a sample of your dataset(s) and provide you a ways to define steps for cleaning and transforming data, preview transformed data and finally create job that applies the recipe steps on the entire dataset.

In this exercise we will create *f1-races-analysis* project by connecting to the datasets we created earlier.

1. Click **PROJECTS** from the left pane. Click on **Connect project**.
2. Enter *f1-races-analysis* as **Project name**.
3. Check a new recipe is getting created as *f1-races-analysis-recipe.*
4. Under select a dataset click **My datasets** > **f1-race-results.**
5. Scroll down until Permissions. Under Role name select option **Create new IAM role.**
6. Enter *data-analyst* as **New IAM role suffix.**
7. Hit Create project.

## 1.3  Working in Glue Databrew Workspace

Now a sample data is loaded into you project workspace you are ready to perform data analysis and transformation. First note that the *f1-race-results* dataset alone have missing details on race, drivers, constructors and circuits. We need to join this with the rest of the datasets.

### 1.3.1 Join datasets

> I.  In the project workspace identify **JOIN** icon which is located between GROUP and UNION. It will look as two vertical bars and an arrow pointing to the left bar from the right bar.

II. Select *f1-races* from the **Select dataset** dropdown. Hit **Next**.

III. This will open up the join defining screen. Click **Left join** as the join type. For the Join keys select the below keys-

Table A f1-race-results: raceId = Table B f1-races: raceId



IV. Click **Joined table preview**. Note that the fields from *f1-races* dataset are included in the joined dataset. Hit **Finish**.

V.   Apply the similar joins for *f1-circuits*, *f1-drivers*, *f1-constructors*. When you define the join make sure the Table A always shows *f1-race-results*.



VI.   After the 3 joins defined you can verify the recipe steps have appeared.



At the end of these steps you have created a joined table in the project workspace combining data from 5 datasets. This table is logical since you have not selected a location where this table will be created. We will do that at later stage.

## 1.3.2 Clean data

Now let's inspect for any data issues. As we scroll towards right we find the **position** field has values marked as "\N". These are drivers who could not finish the race. Let's replace them as more subtle "DNF" (Did Not Finish) value. This is also evident from the **positionText** field which is marked as "R as retired.

### 1.3.3 Change Data type

Further scrolling down will reveal that the ***fastestLapSpeed*** field has been identified as String data type by Glue DataBrew. We will change that to Decimal field.
1. Click on "ABC" data type left of the field name ***fastestLapSpeed.***
2. Select **decimal**.
3. Hit **Apply** in the **Change data type** prompt.





4. After changing the decimal data type the fields is overflowed with spurious 0's upto 18 decimal places. This is not helpful to view results and storing the field values. We will change the decimal precision to 3.

## Format column

Format column to

Decimal precision ▼

Decimal precision symbol

. (Dot) ▼

Decimal places

3

Fill in a value from 0 to 18. If blank, decimal places stay unchanged.

☐ Thousands separator
Specify symbol to separator thousands

Format preview
30000000.000

Apply transform to

◉ **All rows** (500 rows)
Transformation will be applied to all rows in the dataset

○ **Filtered rows - 0 filters applied** (500/500 rows)
Transformation will be applied to filtered rows in the grid

⚠ Column type will be changed to string after you apply this transformation

👁 Preview changes

Cancel     **Apply**

### 1.3.4 Split Column

### 1.3.5 Create Column

### 1.3.6 Merge Column

We notice there are 2 fields for driver first name and last name. We will merge these 2 fields into a single *DriverFullname* field.

## 1.3.7 Date difference Column

An interesting field is driver dob (date of birth), which we can use to extract driver age on the race date.

1. Click on **dob** field. This will highlight the field.
2. Select **FUNCTIONS** from the top menu options. Under FUNCTIONS we have many options of which we will select **Date functions** > **DATEDIFF**.

**3.** For Value 1 choose **Select source column** as *dob*. For Value 2 choose the field *date*. Select "Years" as **Unit** and *DriverAge* as **Destination column**.



4. Hit **Apply**.

## 1.3.8 Binning

Now we have the driver age which is numerical values between 20 and 37. However in many analyses we would like to bucket the ages for the purpose of histogram. This is done by binning function.

### 1.3.9 Rename columns

We will rename the below columns
name → grandPrixName
raceStartTime_1 → racteStartHourOfDay
TABLE_B_number → driverNumber
TABLE_B_name → constructorName
TABLE_B_nationality → constructorNationality
TABLE_B_name1 → circuitName

### 1.3.10 Publish Recipe

Now the transformation steps are defined we are ready to publish this recipe for future use.

## 1.4 Create Job

The above transformation steps are used to create a transform job that can be run as and when any new race data comes. The job will output the transformed data into another s3 folder- this folder will be registered as a Data Catalog S3 table by Glue DataBrew.

First you need to create a database in Glue if you are using Glue for the first time. If you have Glue database "default" already exists you can skip the below step.

Navigate to service **AWS Glue > Glue console > Glue catalog > Databases > Add database**



Note: You may also need to attach the "S3FullAccess" policy to the "AWSGlueDataBrewServiceRole-data-analyst" role that was created during Project creation.

# Create job Info

## Job details

Job name
Identifier for the jobs

f1-data-analysis-job

The job name must contain 1-240 characters. Valid characters are alphanumeric (A-Z, a-z, 0-9), hyphen (-), period (.), and space.

## Job type

Recipe job
A recipe job runs the transformation from the associated recipe on the population of the associated dataset.

Associated dataset
f1-race-results
S3 | s3://awspsa-sampledata/f1data/results/results.csv

Associated recipe
f1-races-analysis-recipe
Working version

## Job output settings Info
Running a job generates output files at specified file destinations.

Output 1                                                    Settings

---

## Job output settings Info
Running a job generates output files at specified file destinations.

Output 1                                              7   Settings

| Output to | 2 | File type | 1 | Delimiter | Compression | Setting summary Info |
| Output location | | Output format | | CSV separator | Available types | |
| Data Catalog S3 tables ▼ | | PARQUET ▼ | | Comma (,) ▼ | None ▼ | File output storage |
| | | | | | | Create a new folder for each job run |

Database name
Name of database                3

Q  default                              ✕   Browse

Custom partition by column values
None

Table name
Name of table                    4

Create new table                         ▼

Catalog table name              5
Catalog table within the database instance. The catalog table name will be prefixed with "awsgluedatabrew_"

f1_race_results

S3 location
Format is: s3://bucket/folder/    6

s3://a ━━━━━━/f1data/results/           Browse

## Settings ✕

### File output storage

○ **Create a new folder for each job run**
Under specified S3 path, a new folder will be created for each job run and for each output file type. The output folder and file name contains job name and job run time. Example:
s3://bucket/myfolder/jobname_10may2020_timestamp/filetype_compression/jobname_10may2020_timestamp_part1.csv

● **Replace output files for each job run**
Flat output files will be created under the specified S3 path. For every job run, the previous output files will be replaced with files from the latest job run. You can enable bucket versioning to be able to restore previous file versions. Example:
s3://bucket/myfolder/jobname_part1.csv

### Custom partition by column values

Partition by unique values of columns. The file is partitioned and stored in a folder path based on the order of columns provided. Example: A file partitioned by Column A and Column B is stored at this S3 path: s3://output file path.../Column A/Column B/

| Enter a column name | | Add |

Cancel    Save

## 1.5 View output table

Once the Job succeeded you are ready to view the output table and run some query on.

# 2. Dashboarding in QuickSight

Create dataset

**Datasets > New dataset > Athena**

## Exclude fields

Organize Fields in Folders

Focus

All fields
driver

Select **All** | **None**

⬜ constructor ⋮

📁 **driver** ⋮

⬜ code ⋮
⬜ dob ⋮
● driverage ⋮
⬜ driveragebins ⋮
⬜ driverfullname ⋮
# driverid ⋮
⬜ drivernumber ⋮
⬜ driverref ⋮
⬜ nationality ⋮

⬜ race ⋮

**Excluded fields** 7 fields excluded ⌃

# table_b_raceid ⋮
# table_b_driverid ⋮

**Fields** 45 fields included ⌃

⬜ driver ⋮

📁 **race** ⋮

# alt ⋮
# circuitid ⋮
⬜ circuitname ⋮
⬜ circuitref ⋮
○ country
Country ⋮
📅 date ⋮
⬜ fastestlap ⋮
# fastestlapspeed ⋮
⬜ fastestlaptime ⋮
⬜ grandprixname ⋮
# grid ⋮
# laps ⋮
○ lat
Latitude ⋮
○ lng
Longitude ⋮
⬜ location ⋮
⬜ milliseconds ⋮
⬜ number ⋮
# points ⋮
⬜ position ⋮
# positionorder ⋮
⬜ positiontext ⋮

**Excluded fields** 7 fields excluded ⌄

## Create calculated fields

1. _race_win_

ifelse(positionorder = 1, 1, 0)



2. _driver_championship_points_

sumOver(points, [driverid, year], PRE_FILTER)

3. _driver_championship_rank_

rank([sum(points) DESC], [year])

4. _driver_champion_name_

firstValue(driverfullname, [sum(points) DESC], [year])

5. _constructor_championship_points_

sumOver(points, [constructorid, year], PRE_FILTER)

6. _constructor_championship_rank_

denseRank([{constructor_championship_points} DESC], [year], PRE_FILTER)

7. _constructor_champion_name_

firstValue(constructorname, [sum(points) DESC], [year])

8. _constructor_last_season_

maxOver(year, [constructorid], PRE_FILTER)

**Fields** 53 fields included

Add calculated field

Augment with SageMaker

Search fields    🔍

Focus

Calculated fields    ∨

Select  All | None

\#   race_win      ⋮

\#   constructor_championship_rank      ⋮

\#   driver_championship_points      ⋮

▭   constructor_champion_name      ⋮

▭   driver_champion_name      ⋮

\#   driver_championship_rank      ⋮

\#   constructor_championship_points      ⋮

\#   constructor_last_season      ⋮

Save & publish

Start analysis

*Sheet - Overview*

**Races over the years**

Fields to select: *date (YEAR), raceId(Count distinct)*

**Geo locations where races been played**
**Points on Map**
Fields to select:
lat, lng → Geospatial
*racedId(Count distinct)* →Size
*country* → Color

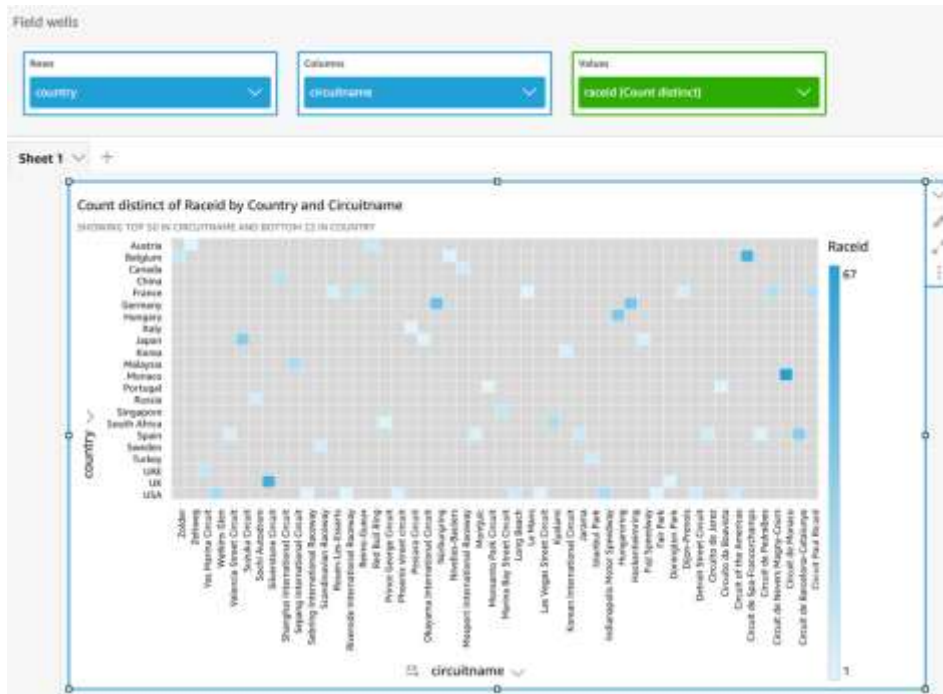**Races in Countries & Circuits**
**Heat map**
Fields to select:
*country* → Rows
*circuitname* →Columns
*racedId(Count distinct)* → Values

**Fastest LapSpeed**
**Clustered bar combo chart**
Fields to select:
*Country, grandprixname drill down* → X axis
*fastestlatspeed (P90)* →Bar
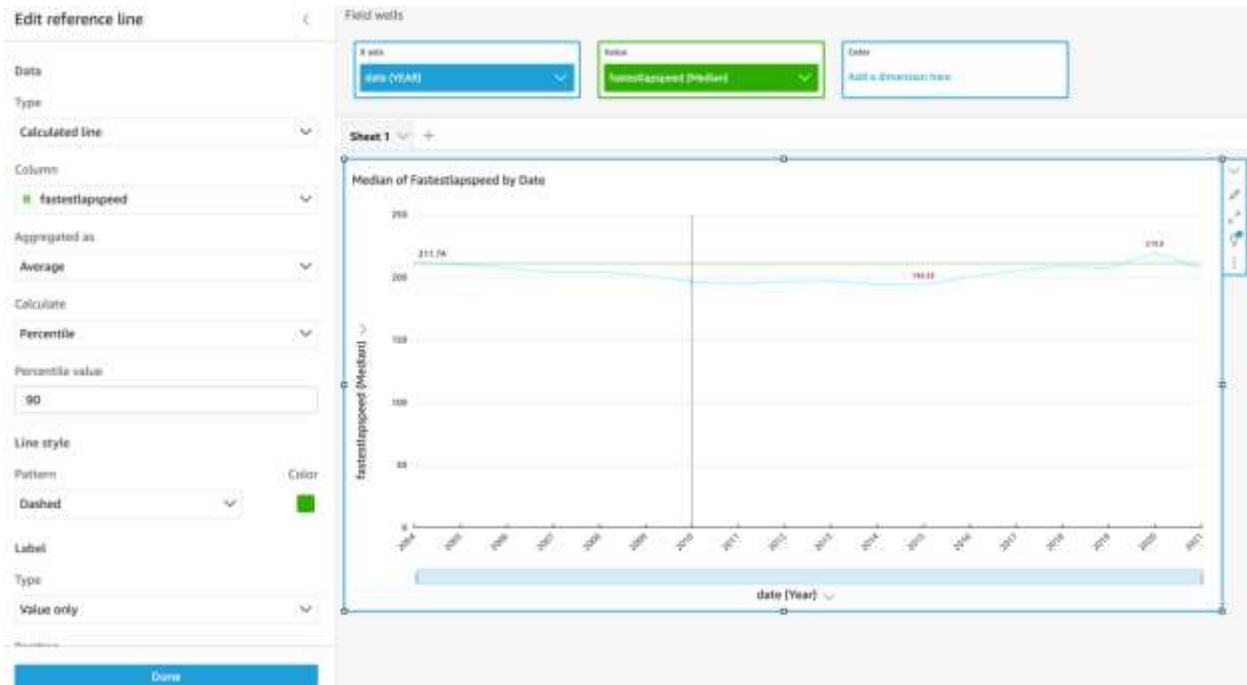*racedId(Count distinct)* → Lines
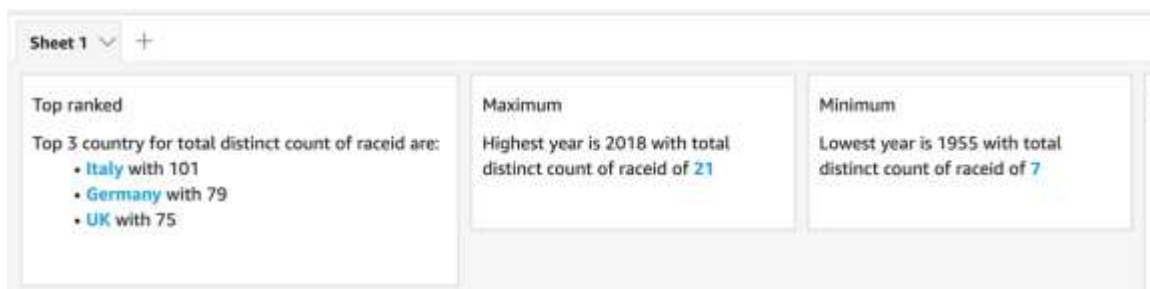


LapSpeed Overtime
**Line chart**
Fields to select:

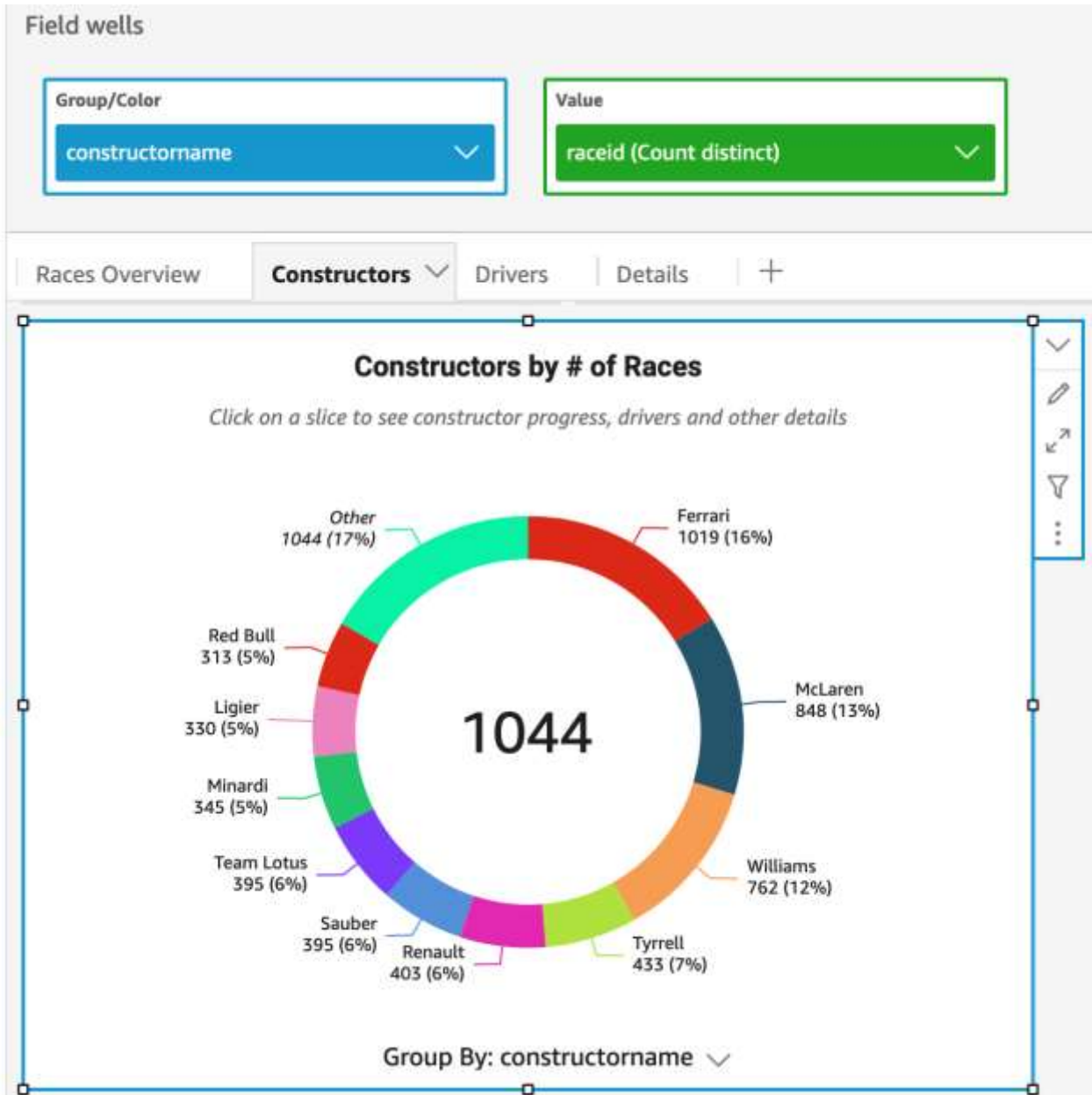*Date(Year)*→ X axis
*fastestlapspeed (Median)* →Value



**Insights**



**Extras:**
KPIs
- # of Races
- # of Circuits
- # of Drivers
- # of Championships
- # of Countries

Formatting of Insights, KPIs and Analysis

*Sheet – Constructors*
**Constructors by # of Races**

## Field wells

**Group/Color**
constructorname ∨

**Value**
raceid (Count distinct) ∨

| Races Overview | **Constructors** ∨ | Drivers | Details | + |

### Constructors by # of Races

*Click on a slice to see constructor progress, drivers and other details*

Other
1044 (17%)

Ferrari
1019 (16%)

Red Bull
313 (5%)

McLaren
848 (13%)

Ligier
330 (5%)

**1044**

Minardi
345 (5%)

Williams
762 (12%)

Team Lotus
395 (6%)

Sauber
395 (6%)

Renault
403 (6%)

Tyrrell
433 (7%)

Group By: constructorname ∨

**Constructor progress over the years by race wins**

**Constructors by # of race wins**

## Field wells

Group by
constructorname ∨

Value
race_win (Sum) ∨

Races Overview | **Constructors** ∨ | Drivers | Details | +

### Constructors by # of race wins



Ferrari 239
McLaren 178
Mercedes 118
Williams 114
Red Bull 70
Team Lotus 45
Renault 35
Other 33
Benetton 27
Brabham 23
Tyrrell 23
Lotus-Climax 22
BRM 17
Cooper-Climax 12
Alfa Romeo 11
Lotus-Ford 11
Vanwall 10
Ligier 9
Maserati 9
Matra-Ford 9
Brabham-Repco 8
Brawn 8
Kurtis Kraft 5
Jordan 4
McLaren-Ford 4
Honda 3

Constructor ∨
Race Wins ∨

**Constructors by points & race wins**

Field wells

| Group by | Size | Color |
|---|---|---|
| constructorname ⌄ | points (Sum) ⌄ | race_win (Sum) ⌄ |

Races Overview | **Constructors** ⌄ | Drivers | Details | +



**Constructors by points & race wins**

**Constructor details**