



**Compass Analytics Sports Analytics Hackathon User Guide** 

In collaboration with Databricks

Friday, October 13, 2023

# **Table of Contents**

| Context and Preamble Description | 3 |
|----------------------------------|---|
| File descriptions                | 3 |
| Data Description                 | 3 |
| Data Dictionary                  | 3 |
| Hackathon Goal                   | 3 |
| Guidance for different profiles  | 4 |
| Data Engineering                 | 4 |
| Data Analyst                     | 6 |
| Data Scientist                   | 6 |
| Submission Outputs               | 8 |
| Evaluation Criteria              | 8 |
| Ouick Links and Reference        | 9 |

# **Context and Preamble Description**

You are a beginner in an NBA fantasy league where the goal is to most effectively predict the plus/minus score line of games with the most accuracy. You are competing against NAB fanastics who spend night and day watching and researching basketball. You are lucky if you watch some of the games when the playoffs come on. You have heard of the power of data science and statistical analysis to predict results so to gain a competitive edge and increase your chances of winning the league you have decided to take a data-driven approach instead of using only intuition. You have found NBA play-by-play data since 1996 to 2023 aid you in your journey. Your fantasy league predicts the score margin at the start of each quarter. You believe you can gain an edge for last quarter predictions and make up the difference in basketball knowledge. Using your analytics skills you will overcome their basketball expertise, take home the league trophy and have bragging rights forever!

#### File descriptions

- Play\_by\_play\_YYYY-YY.parquet (YYYY-YY is the season eg. 2021-22)
  - Contains all the play-by-play data
- Example\_game\_data.xlsx
  - o Example of a single game from the dataset for reference

#### Data Description

- Each file contains play-by-play data for the given season. Some of the data is at the play-by-play level while some other data is at the game level
- Context for each step of the hackathon process is provided in the data dictionary file below. Please review the data dictionary extensively as it will be very useful

# Data Dictionary

• Data\_Dictionary\_Hackathon.xlsx

#### **Hackathon Goal**

The goal of the hackathon is to effectively predict the NBA Score margins for the 2022-2023 season using the historical data. The idea of the project is to get experience working with all aspects of an analysis from the point of view of a data engineer, a data analyst and a data scientist.

For the Data Engineer profile, the main goal is to take raw data stored as Parquet files in an S3 bucket and move it into Delta Live Tables within Unity Catalog in Databricks. This will give you an idea of the tools to migrate data from external locations, process the data using PySpark and generate clean data to be used for data analysis and data science.

For the Data Analyst profile, the main goal is to generate SQL queries to create datasets from the SQL warehouse from contextualized (gold-level) data. Using the data you will create a lakehouse dashboard on top of the warehouse to generate insights into player stats, team stats and season trends.

For the Data Scientist profile, the main goal is to perform exploratory data analysis on contextualized data, create resuasble datasets with feature engineering and perform various regression modelling techniques to predict the NBA Score margins.

Workspace URL: <a href="https://dbc-edd6d92b-cd2f.cloud.databricks.com">https://dbc-edd6d92b-cd2f.cloud.databricks.com</a>

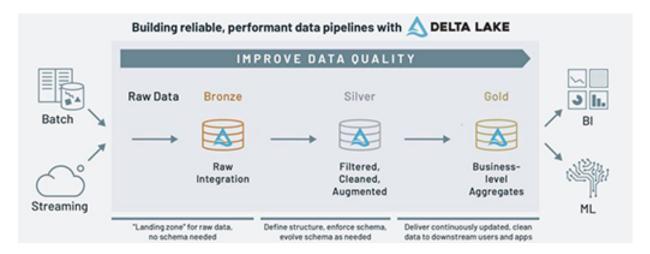
# **Guidance for different profiles**

#### **Data Engineering**

The goal of the data engineer profile is to take the raw parquet data files in an S3 bucket (s3://databricks-compass-hackathon-main-external-

location/Montreal\_Hackathon\_2023/NBA\_Games/) and clean the data using delta live tables in Databricks. The idea is to follow the Medallion Architecture

(<u>https://www.databricks.com/glossary/medallion-architecture</u>) which brings the data from raw data all the way to the gold layer.



The data follows the following steps:

- Bronze: Ingest raw data from source systems into delta tables with no schema required
- Silver: Perform data augmentation, ensure data quality and enforce schema of data
- Gold: Perform aggregations on data for downstream users such as in BI dashboards to gain insights

#### Data Engineering Tasks:

Create the following delta tables in your team's schema by replicating the tables and schema found in hackathon.example\_schema. The following list of tables should be generated:

- nba\_games\_bronze
- Overall game silver
- Play\_by\_play\_silver
- Play\_by\_event\_gold
- Play by event team gold
- Play by player gold
- Play\_by\_player\_per\_season\_gold
- Team stats season away gold
- Team stats season home gold

#### Bronze Table:

Use the montreal\_hackathon\_2023 external location (s3://databricks-compass-hackathon-main-external-location/Montreal\_Hackathon\_2023/NBA\_Games) to ingest the data from the

S3 bucket into databricks (<a href="https://docs.databricks.com/en/storage/amazon-s3.html">https://docs.databricks.com/en/storage/amazon-s3.html</a>). Use the dlt module to implement a delta live table pipeline using Python (<a href="https://docs.databricks.com/en/delta-live-tables/python-ref.html">https://docs.databricks.com/en/delta-live-tables/python-ref.html</a>)

The data stored is in parquet format and the external location is mentioned above. Include inferSchema=True and header="true. Add meta data information specifiying a comment relating to the actions performed in the query and the quality level of the data.

Silver Tables:

Overall game silver:

Separate out the game level data by selecting all of the columns defined in the data dictionary in the Overall Game data. Drop all duplicates using the corresponding PySpark command

Play\_by\_Play\_silver:

Select all of the play-by-play columns defined in the data dictionary as well as the game\_date, home\_team (concatenate the columns home\_team\_city and home\_team\_name together) and away\_team (concatenate the columns away\_team\_city and away\_team\_name together).

Perform the following data cleaning steps to move data to silver level:

- Cast all columns to the appropriate data type following the data dictionary
- Concatenate home\_description and visitor\_description to create play\_description
- Use regexp\_replace to replace "TIE" with "0" and cast to integer
- Apply a window function to fill in the score\_margin, score, away\_score, home\_score
  and team\_leading (Hint: use the following stack overflow article
  <a href="https://stackoverflow.com/questions/36019847/pyspark-forward-fill-with-last-observation-for-a-dataframe">https://stackoverflow.com/questions/36019847/pyspark-forward-fill-with-last-observation-for-a-dataframe</a>)
- Fill in missing values for the columns generated using one of 0, "0 0", or "Home"
- Order results by game id then event num in ascending order
- Drop any duplicates

Use a DLT expectation to ensure that score\_margin\_complete does not contain any empty rows

Gold Tables:

Play\_by\_\*:

Create the following 3 tables which start with Play\_by\_\* by grouping on the following columns and aggregating to get the count. (Use groupBy, agg with F.count and F.expr, and alias to rename). The expected columns are shown in the data dictionary

- Play by event gold: "season" and "event type"
- Play\_by\_event\_team\_gold: "season", "event\_type" and "home\_team"
- Play\_by\_player\_gold: "season", "event\_type" and "player1\_name"

Team\_Stats\_season\*:

Create the following 2 tables which start with Team\_Stats\_season by grouping on the following columns and getting the count and the average of the results. (Use groupBy, agg

with F.count and F.expr, F.avg and alias to rename). The expected columns are shown in the data dictionary.

- Team\_stats\_season\_home\_gold: "Home\_team\_name", "Season"
- Team\_stats\_season\_away\_gold: "Away\_team\_name", "Season"

#### Data Analyst

If you were not able to complete the data engineering steps during the first half of the hackathon please use the example\_schema data provided at hackathon.example\_schema to create your data sources below.

Dashboard: Descriptive Statistics

Utilize the gold tables created in the data engineering steps to gain insights into player stats, team stats and number of events.

SQL Manipulation in Queries Tab:

Create a table in your schema (eg. Hackathon.your\_team.event\_type\_lookup) called event\_type\_lookup from the data in the Event\_Type\_Lookup excel tab to update the event\_type ids with descriptions in english.

SQL Manipulation in Data Tab of Lakehouse Dashboard:

Create a SQL query to get a data source called Season\_Events\_Counts by joining the event\_type\_lookup with the play\_by\_event\_gold data. Perform the same steps to generate the table called Player\_Season\_Event\_Counts.

Create a SQL query to get team\_stats\_season\_home\_gold and team\_stats\_season\_away\_gold into a single table. Join the data on the team name and the season. Filter out team names of teams which are not official NBA teams using a where clause

To test out the queries you can use the SQL Editor tab to query your data sources directly.

Visualizations in the Canvas Tab:

Utilizing the generated datasets using SQL create visualizations to best capture the relationships in the data. Please include examples of all types of features such as visualization, text box and filters.

Be creative and present the data in the most effective way to gain insights quickly into teams, players or number of events.

Machine Learning Results (Bonus):

Visualize the results of specific games of interest and demonstrate the effectiveness of the predictions. Find a game of interest and visualize the actual score margin during the game and the predicted score margin of the 4<sup>th</sup> quarter

#### Data Scientist

#### Notebook 1: Data Exploration and Feature Store Creation

First step is to load the play\_by\_play\_silver data into the notebook. To perform the data cleaning steps we are going to use the pandas\_api for PySpark. First we want to perform some exploratory data analysis such as checking the data types, using the info command as well as checking the number of missing values. Now, remove all columns with a significant number of missing values, duplicate identifiers, uninformative/unusable columns and columns with duplicate information. This will remove a majority of the columns.

After, to get an understanding of your data set visualize the 20 games by taking the first 20 game ids from the dataset and use matplotlib to visualize

Next, we will be filtering for only plays during the regulation time to simplify the analysis as this should not impact our ability to calculate the score margin over the course of a game. Next, we want to keep only official NBA teams and convert any older team names to the most recent team name. Please refer to the data dictionary for the list of official team names and team conversions.

Check for distributions and correlations of numerical variables by selecting the integer and float columns and using the describe and corr methods.

Perform some feature engineering by creating some seasonality variables usng the game\_date and season columns. Extract the month and day from the game\_date by first converting the column to datetime format and getting the month and day. Next use the first half of the wraparound season (eg. 2020 if 2020-21) to get the year.

Create the column momentum\_40 which will take the difference between the current score margin and the score margin 40 plays ago. Perform this within each game by creating a temporary column with the values shifted 40 plays previously then take the difference between the current margin and the shifted margin. Drop the temporary column and all rows with null values in momentum 40.

Finally, convert any included categorical variables to dummy variables and save in feature store within your schema. The data set is very large so please filter for one of the home teams such as "Toronto\_Raptors" to speed up the analysis. Please include any other additional columns you believe will help in the analysis.

#### Notebook 2: Model Evaluation and Hyperparameter Tuning

In the second notebook, read in the feature store data from the created feature store and drop the unnecessary columns from the dataset such as the unique identifiers. Split the dataset into your training and testing datasets by setting all data from before 2022-2023 season or not in the 4<sup>th</sup> quarter to the training set. Only keep 4<sup>th</sup> quarter results as the testing set.

Use the Sci-kit Learn Library to implement a variety of regression modelling techniques such as Linear Regression and Random Forest. Please test out at least 3 different regression models which can be found here https://scikit-learn.org/stable/modules/classes.html.

To evaluate the models please use the crossvalidate function to implement 5-fold cross-validation on the training dataset and evaluate the results using r-squared, negative mean squared error and negative mean absolute error. For each model generated print out the evaluation metrics on both the training and validation sets during each iteration.

**Optional:** After finding the best model using the standard settings, implement hyperparameter tuning using the GridSearchCV method. Select at least two parameters to adjust in your model to find the optimal parameters. Store the best estimator for predictions.

Using the test set make predictions using the fully trained model and evaluate the results through performance metrics and creating visualizations comparing the actual 4<sup>th</sup> quarter results and the predicted results.

Save the actual results with predictions in a table within your schema.

# **Submission Outputs**

#### Data Engineer:

- Notebook containing DLT (Delta Live Table Pipeline) moving data from S3 bucket through the Medallion Architecture (Bronze (Raw Data) --> Silver (Clean Data) --> Gold (Aggregated))
- Delta Live Table Pipeline with Specified Tables stored in your team schema

#### Data Analyst:

- Lakehouse Dashboard with visualizations of Gold Level Data for Insights
- Creative Design and Good User Experience in Canvas Tab
- SQL Queries to generate data tables in Data Tab

#### Data Scientist:

- Notebook(s) containing exploratory data analysis with some visualizations, appropriate numerical and categorical feature engineering, implementation of regression models using scikit-learn with proper training sets and testing sets and visualizations of actual results with predictions using seaborn or matplotlib
- Optional: Use Databricks Feature Store, perform hyperparameter tuning and creating a quick dashboard for results

#### **Evaluation Criteria**

Submissions will be evaluated on the 3 main components: Data Engineering, Data Analyst and Data Scientist Roles

#### Data Engineering:

- Bronze table built (ingest data from S3) 1 point
- Silver table built (data cleaning) 5 points
- Gold tables built (1 point for every table; total of 5 tables to be built) up to 5 points

#### TOTAL DATA ENGINEERING POINTS: up to 11 points available

#### Data Analyst:

- Use of SQL queries to generate tables 3 points
- Descriptive/exploratory dashboard built 5 points
- Creative and easy-to-use dashboard filters incorporated 1 point
- Professional and uses design best practices 1 point
- Demonstrate ability to draw 3 insights 1 point

#### TOTAL DATA ANALYST POINTS: up to 11 points available

#### Data Scientist:

- Lowest MAE (Mean Absolute Error) by team's ranking
  - o 11 points for 1st
  - o 9 points for 2<sup>nd</sup>
  - o 8 points for 3<sup>rd</sup>
  - o 7 points for 4<sup>th</sup>
  - o 6 points for 5<sup>th</sup>
  - o 5 points for 6<sup>th</sup>
  - o 4 points for 7<sup>th</sup>

# TOTAL DATA SCIENTIST POINTS: up to 11 points available

#### **Bonus Points**

- Use of Feature Store 1 point
- Additional feature engineering 1 point
- Use of external data 1 point
- Results visualization in notebook 1 point

#### TOTAL BONUS POINTS: up to 4 points available

# Total available points: 37 points

#### **Quick Links and Reference**

- https://spark.apache.org/docs/latest/api/python/reference/index.html
- <a href="https://spark.apache.org/docs/latest/api/python/user\_guide/pandas\_on\_spark/index\_ntml">https://spark.apache.org/docs/latest/api/python/user\_guide/pandas\_on\_spark/index\_ntml</a>
- https://pandas.pydata.org/docs/
- https://scikit-learn.org/stable/modules/classes.html#
- https://docs.databricks.com/en/data-governance/unity-catalog/index.html
- <a href="https://docs.databricks.com/en/delta-live-tables/index.html">https://docs.databricks.com/en/delta-live-tables/index.html</a>
- <a href="https://docs.databricks.com/en/dashboards/lakeview.html">https://docs.databricks.com/en/dashboards/lakeview.html</a>
- https://scikit-learn.org/stable/modules/classes.html