# Individual Project Report - Rishabh Singh - G26592249

## Project Overview

The project focused on leveraging big data and analytics tools to gather, store, and visualize extensive baseball game data. Specifically, the data encompassed Major League Baseball (MLB) games from 1901 to 2023. My role in this project was multifaceted, involving setting up a data pipeline, data storage, and data visualization.

## Setting up Airflow and DAGs

- **Task:** Establish an Apache Airflow environment and create Directed Acyclic Graphs (DAGs).
- **Objective:** Automate the data extraction process from the MLB schedule API and handle subsequent data processing tasks.
- **Execution:** I successfully set up an Airflow environment, leveraging its scheduling and workflow management capabilities. Two DAGs were written:
    - **DAG 1:** Extracted data from the MLB schedule API. This process involved fetching game data spanning over a century (1901-2023).
    - **DAG 2:** Called the game_pk API using the game IDs (game_pk) extracted by DAG 1.

## Data Storage in S3

- **Task:** Store the retrieved data in an AWS S3 bucket.
- **Data Volume:** The initial data from the schedule API totaled approximately 8 GB. The data from the game_pk API was substantially larger, around 175 GB.
- **Structure and Format:**
    - The schedule API data was stored along with a CSV file containing all game IDs.
    - The game_pk data was organized in a structured folder hierarchy: year/month/day/game_pk.json.
- **Data Lake:** Both sets of data were integrated into our S3-based data lake, providing a centralized and scalable storage solution.

## Data Visualization in Google Colab

- **Task:** Develop data visualizations to represent the cleaned and processed MLB data.
- **Tools Used:** Google Colab, leveraging Python libraries for data analysis and visualization.
- **Process:** Utilized the cleaned data prepared by my teammate. This data was extracted from our S3 data lake and then analyzed and visualized in Google Colab. The visualizations focused on

providing insights into various aspects of the MLB data, such as game trends, player performance, and historical comparisons.

## Challenges Faced and Overcome

- **Data Volume Management:** Handling large volumes of data (especially the 175 GB from the game_pk API) was challenging. Efficient data processing and storage strategies were crucial.
- **Complex Data Extraction:** Extracting and structuring data from APIs spanning over a century required meticulous planning and execution.
- **Visualization of Complex Data Sets:** Converting large and complex data sets into comprehensible visualizations demanded a deep understanding of both the data and the visualization tools.
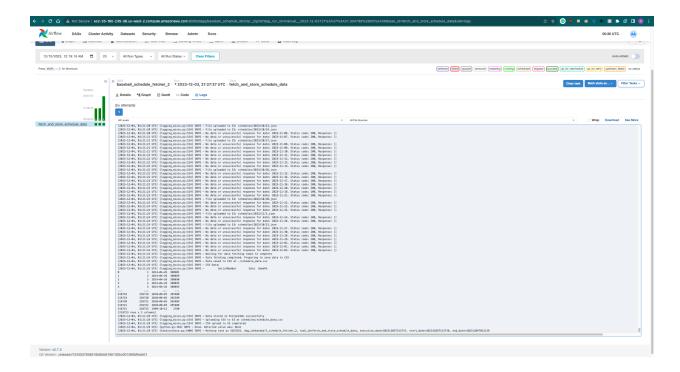
## Conclusion

This project not only demonstrated my ability to handle large-scale data using modern big data tools but also emphasized the importance of clear and insightful data visualization. Working with historical MLB data presented unique challenges and learning opportunities, especially in terms of data management and visualization techniques.

# Screenshots Attached

## DAG 1

```python
import os
import datetime
import requests
import boto3
import sqlalchemy
import json
import pandas as pd
from concurrent.futures import ThreadPoolExecutor, as_completed
from airflow import DAG
from airflow.operators.python_operator import PythonOperator

# DAG configuration
default_args = {
    'owner': 'airflow',
    'depends_on_past': False,
    'start_date': datetime.datetime.today(),
    'retries': 1,
    'retry_delay': datetime.timedelta(minutes=5),
}

dag = DAG('baseball_schedule_fetcher_2', default_args=default_args, schedule_interval=None)

def upload_to_s3(bucket_name, object_key, data):
    """
    Upload a file to an S3 bucket
    """
    s3_client = boto3.client('s3')
    try:
        s3_client.put_object(Bucket=bucket_name, Key=object_key, Body=json.dumps(data))
        print(f"File uploaded to S3: {object_key}")
    except Exception as e:
        print(f"Error uploading to S3: {e}")

def fetch_schedule_data(date):
    schedule_url = f"https://baseballsavant.mlb.com/schedule?date={date.strftime('%Y-%m-%d')}"
    schedule_response = requests.get(schedule_url)
    data = []

    if schedule_response.status_code == 200 and schedule_response.json():
        try:
            schedule_data = schedule_response.json()

            # Extracting gamePk values from the schedule data
            if 'schedule' in schedule_data and 'dates' in schedule_data['schedule']:
                for date_entry in schedule_data['schedule']['dates']:
                    for game in date_entry['games']:
                        game_pk = game.get('gamePk')
                        if game_pk:
                            data.append([date, game_pk])
            # Upload JSON response directly to S3
            s3_path = f"schedules/{date.year}/{date.month}/{date.day}.json"
            upload_to_s3('mlb-data-store', s3_path, schedule_data)
        except Exception as e:
            print(f"Error processing schedule data for date: {date}: {e}")
    else:
        # Log a message if the response is not successful or if there's no JSON data
        print(f"No data or unsuccessful response for date: {date}. Status code: {schedule_response.status_code}, Response: {schedule_response.text}")

    return data
```

```python
61  def fetch_and_store_data(ds, **kwargs):
62      print("Starting data fetch and store process")
63
64      start_date = datetime.date(1901, 1, 1)
65      end_date = datetime.date.today()
66      print(f"Fetching data from {start_date} to {end_date}")
67
68      dates = [start_date + datetime.timedelta(days=i) for i in range((end_date - start_date).days + 1)]
69
70      futures = []
71      with ThreadPoolExecutor(max_workers=20) as executor:
72          for date in dates:
73              futures.append(executor.submit(fetch_schedule_data, date))
74              print(f"Scheduled data fetch for date: {date}")
75
76      print("Waiting for data fetching tasks to complete")
77      results = []
78      for future in as_completed(futures):
79          results.extend(future.result())
80
81      print("Data fetching completed. Preparing to save data to CSV")
82      # Convert results to DataFrame and save as CSV
83      df = pd.DataFrame(results, columns=['Date', 'GamePk'])
84      df.insert(0, 'SerialNumber', range(1, 1 + len(df)))
85      csv_file_path = './schedule_data.csv'
86      df.to_csv(csv_file_path, index=False)
87      print(f"Data saved to CSV at {csv_file_path}")
88
89      # Read and print the CSV file
90      df = pd.read_csv(csv_file_path)
91      print("CSV Data:")
92      print(df)
93
94      # Store data in PostgreSQL
95      # Replace with your PostgreSQL connection details
96      postgres_connection = 'postgresql+psycopg2://airflow:airflow@localhost/airflow'
97      engine = sqlalchemy.create_engine(postgres_connection)
98      try:
99          df.to_sql('schedule_data', engine, index=False, if_exists='replace')
100         print("Data stored in PostgreSQL successfully")
101     except Exception as e:
102         print(f"Error storing data in PostgreSQL: {e}")
103
104     # Upload the CSV file to S3
105     s3_csv_path = f"schedules/schedule_data.csv"
106     print(f"Uploading CSV to S3 at {s3_csv_path}")
107     s3_client = boto3.client('s3')
108     s3_client.upload_file(csv_file_path, 'mlb-data-store', s3_csv_path)
109     print("CSV upload to S3 completed")
110
111 # Task definition
112 fetch_store_task = PythonOperator(
113     task_id='fetch_and_store_schedule_data',
114     provide_context=True,
115     python_callable=fetch_and_store_data,
116     dag=dag,
117 )
118
119 fetch_store_task
120
```

Airflow

DAGs   Cluster Activity   Datasets   Security   Browse   Admin   Docs

00:36 UTC

12/15/2023, 12:19:16 AM     25     All Run Types     All Run States     Clear Filters                     Auto-refresh

Press **shift** + **/** for Shortcuts

deferred | failed | queued | removed | restarting | running | scheduled | skipped | success | up_for_reschedule | up_for_retry | upstream_failed | no_status

DAG                                          Run                              Task
baseball_schedule_fetcher_2   ▶ 2023-12-03, 21:37:37 UTC   fetch_and_store_schedule_data

Clear task   Mark state as...   Filter Tasks

⚠ Details   ⚙ Graph   ⌗ Gantt   <> Code   ▤ Logs

(by attempts)

All Levels     All File Sources                     Wrap   Download   See More

```
[2023-12-04, 01:21:20 UTC] {logging_mixin.py:154} INFO - File uploaded to S3: schedules/2023/10/23.json
[2023-12-04, 01:21:20 UTC] {logging_mixin.py:154} INFO - File uploaded to S3: schedules/2023/10/24.json
[2023-12-04, 01:21:21 UTC] {logging_mixin.py:154} INFO - No data or unsuccessful response for date: 2023-11-06. Status code: 200, Response: []
[2023-12-04, 01:21:21 UTC] {logging_mixin.py:154} INFO - No data or unsuccessful response for date: 2023-11-07. Status code: 200, Response: []
[2023-12-04, 01:21:21 UTC] {logging_mixin.py:154} INFO - File uploaded to S3: schedules/2023/10/27.json
[2023-12-04, 01:21:21 UTC] {logging_mixin.py:154} INFO - No data or unsuccessful response for date: 2023-11-08. Status code: 200, Response: []
[2023-12-04, 01:21:21 UTC] {logging_mixin.py:154} INFO - No data or unsuccessful response for date: 2023-11-09. Status code: 200, Response: []
[2023-12-04, 01:21:22 UTC] {logging_mixin.py:154} INFO - No data or unsuccessful response for date: 2023-11-10. Status code: 200, Response: []
[2023-12-04, 01:21:22 UTC] {logging_mixin.py:154} INFO - No data or unsuccessful response for date: 2023-11-11. Status code: 200, Response: []
[2023-12-04, 01:21:22 UTC] {logging_mixin.py:154} INFO - No data or unsuccessful response for date: 2023-11-13. Status code: 200, Response: []
[2023-12-04, 01:21:22 UTC] {logging_mixin.py:154} INFO - No data or unsuccessful response for date: 2023-11-14. Status code: 200, Response: []
[2023-12-04, 01:21:22 UTC] {logging_mixin.py:154} INFO - File uploaded to S3: schedules/2023/10/28.json
[2023-12-04, 01:21:22 UTC] {logging_mixin.py:154} INFO - No data or unsuccessful response for date: 2023-11-15. Status code: 200, Response: []
[2023-12-04, 01:21:23 UTC] {logging_mixin.py:154} INFO - No data or unsuccessful response for date: 2023-11-16. Status code: 200, Response: []
[2023-12-04, 01:21:23 UTC] {logging_mixin.py:154} INFO - No data or unsuccessful response for date: 2023-11-17. Status code: 200, Response: []
[2023-12-04, 01:21:23 UTC] {logging_mixin.py:154} INFO - No data or unsuccessful response for date: 2023-11-18. Status code: 200, Response: []
[2023-12-04, 01:21:23 UTC] {logging_mixin.py:154} INFO - No data or unsuccessful response for date: 2023-11-20. Status code: 200, Response: []
[2023-12-04, 01:21:23 UTC] {logging_mixin.py:154} INFO - No data or unsuccessful response for date: 2023-11-19. Status code: 200, Response: []
[2023-12-04, 01:21:23 UTC] {logging_mixin.py:154} INFO - No data or unsuccessful response for date: 2023-11-21. Status code: 200, Response: []
[2023-12-04, 01:21:23 UTC] {logging_mixin.py:154} INFO - File uploaded to S3: schedules/2023/10/30.json
[2023-12-04, 01:21:23 UTC] {logging_mixin.py:154} INFO - No data or unsuccessful response for date: 2023-11-22. Status code: 200, Response: []
[2023-12-04, 01:21:24 UTC] {logging_mixin.py:154} INFO - No data or unsuccessful response for date: 2023-11-23. Status code: 200, Response: []
[2023-12-04, 01:21:24 UTC] {logging_mixin.py:154} INFO - No data or unsuccessful response for date: 2023-11-25. Status code: 200, Response: []
[2023-12-04, 01:21:24 UTC] {logging_mixin.py:154} INFO - File uploaded to S3: schedules/2023/11/1.json
[2023-12-04, 01:21:24 UTC] {logging_mixin.py:154} INFO - No data or unsuccessful response for date: 2023-11-24. Status code: 200, Response: []
[2023-12-04, 01:21:24 UTC] {logging_mixin.py:154} INFO - No data or unsuccessful response for date: 2023-11-26. Status code: 200, Response: []
[2023-12-04, 01:21:24 UTC] {logging_mixin.py:154} INFO - File uploaded to S3: schedules/2023/10/31.json
[2023-12-04, 01:21:24 UTC] {logging_mixin.py:154} INFO - No data or unsuccessful response for date: 2023-11-27. Status code: 200, Response: []
[2023-12-04, 01:21:24 UTC] {logging_mixin.py:154} INFO - No data or unsuccessful response for date: 2023-11-30. Status code: 200, Response: []
[2023-12-04, 01:21:24 UTC] {logging_mixin.py:154} INFO - No data or unsuccessful response for date: 2023-11-26. Status code: 200, Response: []
[2023-12-04, 01:21:24 UTC] {logging_mixin.py:154} INFO - No data or unsuccessful response for date: 2023-11-26. Status code: 200, Response: []
[2023-12-04, 01:21:24 UTC] {logging_mixin.py:154} INFO - No data or unsuccessful response for date: 2023-12-02. Status code: 200, Response: []
[2023-12-04, 01:21:24 UTC] {logging_mixin.py:154} INFO - No data or unsuccessful response for date: 2023-12-01. Status code: 200, Response: []
[2023-12-04, 01:21:24 UTC] {logging_mixin.py:154} INFO - No data or unsuccessful response for date: 2023-12-03. Status code: 200, Response: []
[2023-12-04, 01:21:24 UTC] {logging_mixin.py:154} INFO - Waiting for data fetching tasks to complete
[2023-12-04, 01:21:24 UTC] {logging_mixin.py:154} INFO - Data fetching completed. Preparing to save data to CSV
[2023-12-04, 01:21:25 UTC] {logging_mixin.py:154} INFO - Data saved to CSV at ./schedule_data.csv
[2023-12-04, 01:21:25 UTC] {logging_mixin.py:154} INFO - CSV Data:
0            1  2014-04-26  300001
1            2  2014-04-26  300009
2            3  2014-04-26  300010
3            4  2014-04-26  300095
4            5  2014-04-26  300003
...        ...         ...     ...
218718  218719  2010-08-03  265400
218719  218720  2010-08-03  265399
218720  218721  2010-08-03  265404
218721  218722  2010-08-03  265406
218722  218723  1999-10-12    2700
[218723 rows x 3 columns]
[2023-12-04, 01:21:28 UTC] {logging_mixin.py:154} INFO - Data stored in PostgreSQL successfully
[2023-12-04, 01:21:28 UTC] {logging_mixin.py:154} INFO - Uploading CSV to S3 at schedules/schedule_data.csv
[2023-12-04, 01:21:29 UTC] {logging_mixin.py:154} INFO - CSV upload to S3 completed
[2023-12-04, 01:21:29 UTC] {python.py:154} INFO - Done. Returned value was: None
[2023-12-04, 01:21:29 UTC] {taskinstance.py:1400} INFO - Marking task as SUCCESS. dag_id=baseball_schedule_fetcher_2, task_id=fetch_and_store_schedule_data, execution_date=20231203T213737, start_date=20231203T213738, end_date=20231204T012129
```
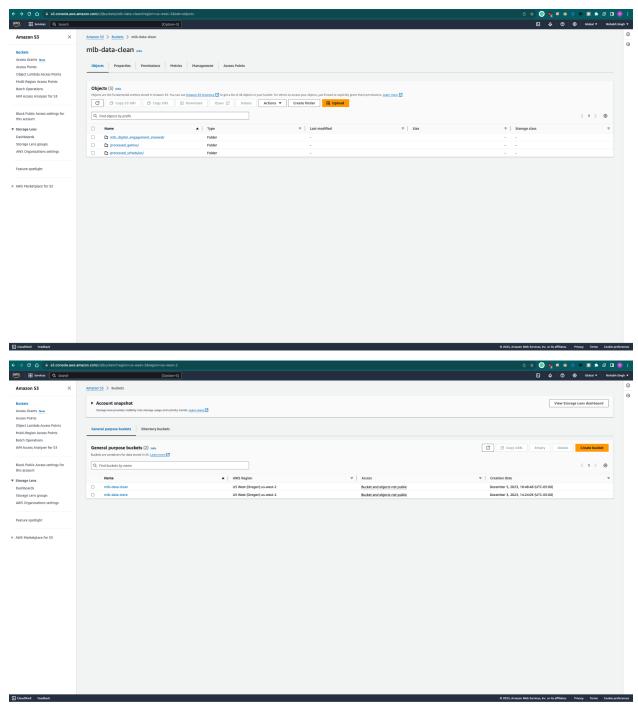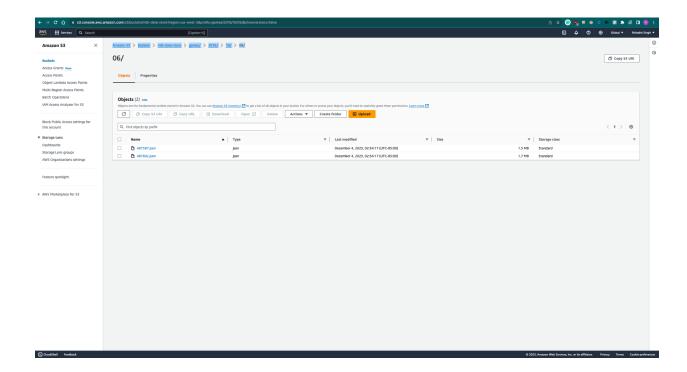
# DAG 2

```python
1  import os
2  import json
3  import requests
4  import pandas as pd
5  import boto3
6  from concurrent.futures import ThreadPoolExecutor, as_completed
7  from airflow import DAG
8  from airflow.operators.python_operator import PythonOperator
9  from datetime import datetime, timedelta
10
11 # DAG configuration
12 default_args = {
13     'owner': 'airflow',
14     'start_date': datetime(2023, 1, 1),
15     'retries': 1,
16     'retry_delay': timedelta(minutes=5)
17 }
18
19 dag = DAG('fetch_and_store_mlb_data',
20           default_args=default_args,
21           schedule_interval=None)
22
23 def upload_to_s3(bucket_name, object_key, data):
24     """
25     Upload data to an S3 bucket
26     """
27     s3_client = boto3.client('s3')
28     try:
29         s3_client.put_object(Bucket=bucket_name, Key=object_key, Body=json.dumps(data))
30         print(f"Data uploaded to S3: {bucket_name}/{object_key}")
31     except Exception as e:
32         print(f"Error uploading to S3: {e}")
33
```

```python
34  def fetch_and_save_game_data(date, game_pk, bucket):
35      print(f"Fetching data for game_pk {game_pk}...")
36      api_url = f"https://baseballsavant.mlb.com/gf?game_pk={game_pk}"
37      try:
38          response = requests.get(api_url)
39          if response.status_code == 200:
40              data = response.json()
41              if "error" in data and data["error"] == "Invalid Game PK.":
42                  return f"Invalid Game PK for game_pk {game_pk}. Skipping."
43
44              object_key = f"games/{date.year}/{date.month:02d}/{date.day:02d}/{game_pk}.json"
45              upload_to_s3(bucket, object_key, data)
46              return f"Data for game_pk {game_pk} uploaded to S3 at {object_key}"
47          else:
48              return f"Failed to fetch data for game_pk {game_pk}: Status code {response.status_code}"
49      except Exception as e:
50          return f"Exception for game_pk {game_pk}: {e}"
51
52  def fetch_game_data_and_store():
53      print("Loading DataFrame from CSV...")
54      df = pd.read_csv('./schedule_data.csv')
55      print(df.head())
56      print("DataFrame loaded. Processing data.")
57      bucket = 'mlb-data-store'  # Replace with your S3 bucket name
58
59      with ThreadPoolExecutor(max_workers=20) as executor:
60          futures = {executor.submit(fetch_and_save_game_data, pd.to_datetime(row['Date']), row['GamePk'], bucket): row for index, row in df.iterrows()}
61
62      for future in as_completed(futures):
63          result = future.result()
64          print(result)
65
66      print("Data fetching and storing process completed")
67
68  fetch_store_task = PythonOperator(
69      task_id='fetch_and_store_mlb_data_task',
70      python_callable=fetch_game_data_and_store,
71      dag=dag,
72  )
73
74  fetch_store_task
75
```

ec2-35-160-235-36.us-west-2.compute.amazonaws.com:8080/dags/fetch_and_store_mlb_data/grid?tab=logs&dag_run_id=manual__2023-12-04T01%3A42%3A16.191205%2B00%3A00&task_id=fetch_and_store_mlb_data_task

Airflow

DAGs    Cluster Activity    Datasets    Security ▾    Browse ▾    Admin ▾    Docs ▾

00:36 UTC

Press **Esc?** + **/** for Shortcuts

DAG
**fetch_and_store_mlb_data**  ▶ 2023-12-04, 01:42:16 UTC  /  fetch_and_store_mlb_data_task

Clear task    Mark state as... ▾    Filter Tasks ▾

⚠ Details    ᵇᵍ Graph    ☑ Gantt    <> Code    ☰ Logs

(by attempts)

[ 1 ]

All Levels ▾                                                     All File Sources ▾                          ☐ Wrap    Download    See More

fetch_and_store_mlb_data_task

[2023-12-04, 23:49:07 UTC] (logging_mixin.py:154) INFO - Data for game_pk 661773 uploaded to S3 at games/2022/08/28/661773.json
[2023-12-04, 23:49:07 UTC] (logging_mixin.py:154) INFO - Data for game_pk 134477 uploaded to S3 at games/1051/06/24/134477.json
[2023-12-04, 23:49:07 UTC] (logging_mixin.py:154) INFO - Data for game_pk 104180 uploaded to S3 at games/1983/04/24/104180.json
[2023-12-04, 23:49:07 UTC] (logging_mixin.py:154) INFO - Invalid Game PK for game_pk 41036. Skipping.
[2023-12-04, 23:49:07 UTC] (logging_mixin.py:154) INFO - Data for game_pk 134473 uploaded to S3 at games/1051/06/24/134473.json
...
[2023-12-04, 23:49:08 UTC] (python.py:194) INFO - Done. Returned value was: None
[2023-12-04, 23:49:08 UTC] (taskinstance.py:1400) INFO - Marking task as SUCCESS. dag_id=fetch_and_store_mlb_data, task_id=fetch_and_store_mlb_data_task, execution_date=20231204T014216, start_date=20231204T014210, end_date=20231204T234908

Version: **v2.7.3**
Git Version: .release:f1243537836516b8bb8156130bc001986bfbeb01

# S3

# Airflow