

Project Name: ETL pipeline Designing with Talend

Guided By

Sai Prakash

Team Members

Jagruti Airao

Jatan Sahu

Jigar Shekhat

Priyansh Agarwal

Purvang Maheria

Vipasha Vaghela

Vivek Soni

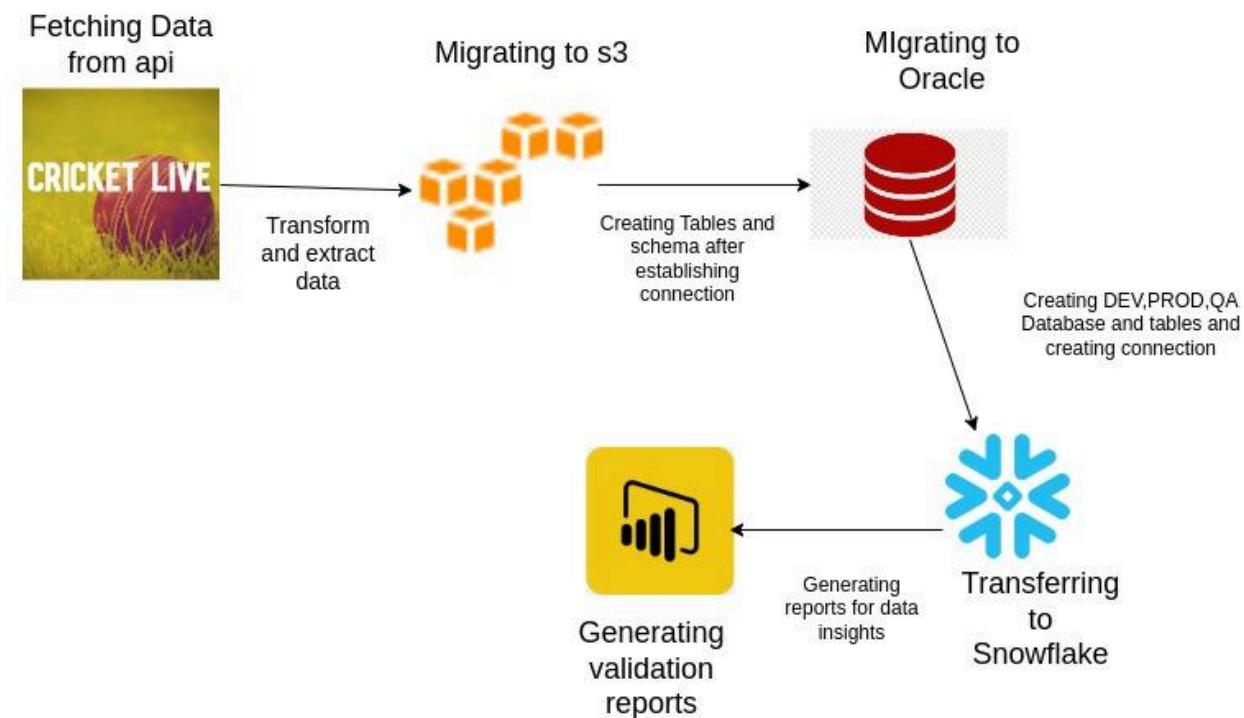
Table of Contents

Abstract.....	3
Project architecture	4
Methodology	5
Task 1: Data modeling and schemas	5
Task 2: Data fetching and S3 transfer	6
Task 3: S3 to oracle data transfer	15
Task 4: Oracle Database to Snowflake database	22
Task 5: Data Validation and BI Notebook.	30
i) Data Validation	30
ii) SnowFlake to PowerBI	35
Conclusion	41

Abstract:

This project aims to streamline the data management process through the implementation of various stages encompassing data modeling, acquisition, integration, migration, and validation. Initially, a data model is constructed utilizing the Star/Snowflake Schema, providing a foundation for structured data organization. Subsequently, cricket data is obtained from external sources via RAPID API and stored in Amazon S3 utilizing Python, with scheduled data acquisition ensuring timely updates. The acquired data is then transferred to Oracle utilizing Talend for further processing and analysis. Following this, a migration process is employed to transition data from Oracle to Snowflake, leveraging Talend's capabilities for seamless data migration. Finally, to ensure data integrity and reliability, a comprehensive data validation process is established, complemented by the generation of sample reports using Power BI, facilitating insightful data analysis and decision-making. Through the integration of these components, this project endeavors to optimize data management practices, enabling organizations to harness the full potential of their data assets.

Project Architecture:



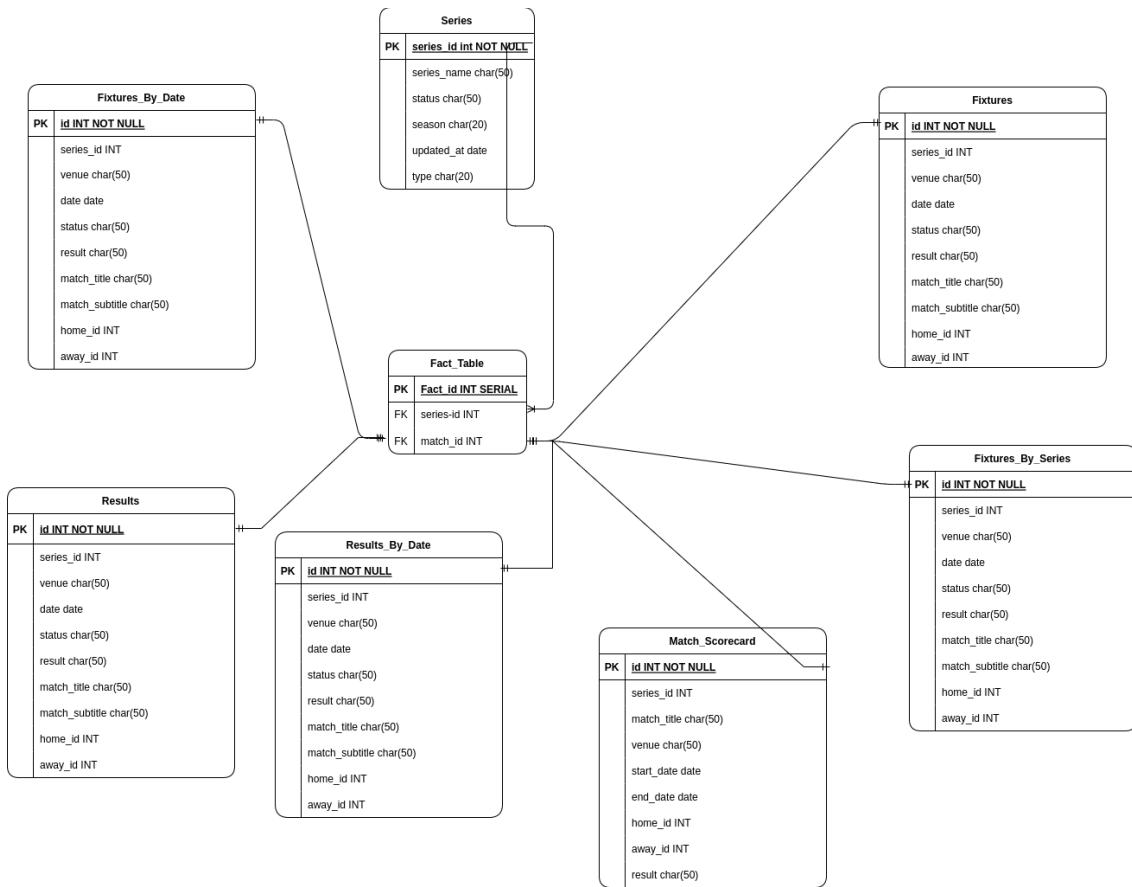
Methodology:

1. Prepare Data model (Star/ Snowflake Schema) with ER Diagram - Jatan and Priyansh

Star and snowflake schemas are both widely used in data warehousing to organize and structure relational databases for efficient querying and analysis. They have some similarities but differ in their structure and complexity.

Star schemas feature a denormalized structure where dimension tables are fully denormalized, resulting in simpler schema designs with fewer tables and straightforward queries. In contrast, snowflake schemas follow a normalized approach where dimension tables are partially or fully normalized into multiple related tables, leading to a more complex schema with additional joins required for querying. Consequently, while star schemas offer simplicity and faster query performance, snowflake schemas prioritize data integrity and storage efficiency at the cost of increased query complexity.

Star Schema :



2. Get Data from 3rd Party with API and send to s3 with Python and schedule it - Jatan and Priyansh

Getting data from API

For this we have used RAPID API as a third party source to gather the data using python.

[Click to view data source](#)

The screenshot shows the RapidAPI platform interface. At the top, there's a search bar and navigation links for API Hub, Organizations, Apps, My APIs, and a user profile. Below the header, the 'Cricket Live Data' API card is displayed, categorized under Sports, with a PREMIUM badge. It shows popularity (9.6 / 10), latency (715ms), service level (100%), and health check (N/A). The API version V1 (Current) is selected. The 'Endpoints' tab is active, showing the 'GET Series' endpoint. The description for this endpoint is 'Lists available cricket series' to query'. It requires a 'Personal Account' (Jatin Sahu) and a 'RapidAPI App' (default-application_8858372). The 'Request URL' is set to 'rapidapi.com'. Under 'Header Parameters', 'X-RapidAPI-Key' is defined as an ENUM with value '2fa7fc6606mshc9294d9833a48cep1e7a66jsn7b106a3105b3', and 'X-RapidAPI-Host' is defined as a STRING with value 'cricket-live-data.p.rapidapi.com'. On the right side, there are tabs for 'Code Snippets', 'Example Responses', and 'Results'. A Python code snippet is provided for the endpoint:

```

(url = "https://cricket-live-data.p.rapidapi.com/series"
headers = {
    "X-RapidAPI-Key": "2fa7fc6606mshc9294d9833a48cep1e7a66jsn7b106a3105b3",
    "X-RapidAPI-Host": "cricket-live-data.p.rapidapi.com"
}

response = requests.get(url, headers=headers)

print(response.json())

```

At the bottom, there are links for Home, Cricket Live Data, and various platform sections like About, Blog, Learn, Careers, Press, Contact, Terms, and Privacy. Language selection is set to English. Social media icons for Facebook, Twitter, and LinkedIn are also present.

As we are using free version of RAPIDAPI we have very limited access(250 requests only) so we have requested limited number of calls and ingested less amount of data.

For sending this data to S3 bucket we have written python code. As we have 7 tables and one which is venue table which is a nested file which contains venue id and the team name and team code. So we have created total 7 python files and 8 csv file, send all these structured (.csv) data into s3 bucket.

The data is in JSON format we have also done required preprocessing task for all the tables. **Python Code Snippets** [Click to view code](#)

1.Series

```

home > growlt250 > Courses > TALEND > Python Files for deploy > series_upload_to_s3.py > upload_to_s3
  1 import pandas as pd
  2 import requests
  3 import boto3
  4 from io import StringIO
  5 import datetime
  6
  7 # Global variables
  8 RAPIDAPI_KEY = "b72f452b32msheb263a2dd53ade5p151fa5jsn1016eed39cf"
  9 RAPIDAPI_HOST = "cricket-live-data.p.rapidapi.com"
10 AWS_ACCESS_KEY_ID = 'AKIARACMVC5URUSKCGX'
11 AWS_SECRET_ACCESS_KEY = 'crE7tDK1sczBdU2w6jM7Se/fx1FN5ywRRpeaKmc0'
12 S3_BUCKET_NAME = 'cricket-api-to-s3'
13 S3_UPLOADED_FILE_NAME = 'series'
14
15 # API endpoint function
16 def fetch_data_from_api(url, headers):
17     response = requests.get(url, headers=headers)
18     if response.status_code == 200:
19         return response.json()
20     else:
21         print("Failed to fetch data:", response.status_code)
22         return None
23
24 # Function to process data
25 def process_data(json_data):
26     all_series_data = []
27     for result in json_data['results']:
28         series_data = result['series']
29         series_type = result['type']
30         # Add 'type' to each series data
31         for data in series_data:
32             data['type'] = series_type
33         all_series_data.extend(series_data)
34     return pd.DataFrame(all_series_data)
35
36 # # Function to save DataFrame to CSV and upload to S3
37 # def upload_to_s3(df, bucket_name, s3_output_name):
38 #     csv_buffer = StringIO()
39 #     df.to_csv(csv_buffer, index=False)
40 #     s3 = boto3.client('s3', aws_access_key_id=AWS_ACCESS_KEY_ID, aws_secret_access_key=AWS_SECRET_ACCESS_KEY)
41 #     s3.put_object(Body=csv_buffer.getvalue(), Bucket=bucket_name, Key=s3_output_name)
42 #     print(f'Data has been uploaded to S3 bucket {bucket_name} with key {s3_output_name}')
43
44
45 def upload_to_s3(df, bucket_name, folder):
46     # Get current year and month
47     current_date = datetime.datetime.now()
48     year = current_date.year
49     month = current_date.month
50
51     # Construct S3 key
52     s3_output_name = f'{folder}/'
53     s3_output_name += f'{year}/'
54     s3_output_name += f'{month:02}/'
55     s3_output_name += f'{folder}_{year}_{month:02}_{current_date.day}.csv'
56
57     # Save DataFrame to CSV and upload to S3
58     csv_buffer = StringIO()
59     df.to_csv(csv_buffer, index=False)
60     s3 = boto3.client('s3', aws_access_key_id=AWS_ACCESS_KEY_ID, aws_secret_access_key=AWS_SECRET_ACCESS_KEY)
61     s3.put_object(Body=csv_buffer.getvalue(), Bucket=bucket_name, Key=s3_output_name)
62     print(f'Data has been uploaded to S3 bucket {bucket_name} with key {s3_output_name}')
63
64
65 # Main function
66 def main():
67     url = "https://cricket-live-data.p.rapidapi.com/series"
68     headers = {
69         "X-RapidAPI-Key": RAPIDAPI_KEY,
70         "X-RapidAPI-Host": RAPIDAPI_HOST
71     }
72     json_data = fetch_data_from_api(url, headers)
73     if json_data:
74         df_series = process_data(json_data)
75         cols = ['series.name', 'status']
76         for col in cols:
77             df_series[col] = df_series[col].str.replace(',', '')
78         upload_to_s3(df_series, S3_BUCKET_NAME, S3_UPLOADED_FILE_NAME)
79     else:
80         print("No data fetched from the API.")
81
82 if __name__ == "__main__":
83     main()

```

2. Fixtures

```

home > growlt256 > Courses > TALEND > Python files for deploy > fixtures_api_to_s3.py > ...
1 import pandas as pd
2 import requests
3 import boto3
4 from io import StringIO
5 import datetime
6
7 # Global variables
8 RAPIDAPI_KEY = "f381eaaa79msh51537dd6e014035p180695jsn67cf7d549118"
9 RAPIDAPI_HOST = "cricket-live-data.p.rapidapi.com"
10 AWS_ACCESS_KEY_ID = 'AKIARAACMVC5URUSKCGX'
11 AWS_SECRET_ACCESS_KEY = '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'
12 S3_BUCKET_NAME = 'cricket-api-to-s3'
13 S3_UPLOADED_FILE_NAME = 'fixtures'
14
15 def fetch_data_from_api(url):
16     headers = {
17         "X-RapidAPI-Key": RAPIDAPI_KEY,
18         "X-RapidAPI-Host": RAPIDAPI_HOST
19     }
20     response = requests.get(url, headers=headers)
21     if response.status_code == 200:
22         json_data = response.json()
23         df_date = []
24         for i in range(len(json_data['results'])):
25             entry = []
26             ide = json_data['results'][i]['id']
27             series_id = json_data['results'][i]['series_id']
28             venue = json_data['results'][i]['venue']
29             date_with_timestamp = json_data['results'][i]['date']
30             status = json_data['results'][i]['status']
31             result = json_data['results'][i]['result']
32             match_title = json_data['results'][i]['match_title']
33             match_subtitle = json_data['results'][i]['match_subtitle']
34             home_team_id = json_data['results'][i]['home']['id']
35             away_team_id = json_data['results'][i]['away']['id']
36             entry.append([ide, series_id, venue, date_with_timestamp, status, result, match_title, match_subtitle,
37                         home_team_id, away_team_id])
38             df_date.extend(entry)
39     return pd.DataFrame(df_date, columns = ['id', 'series_id', 'venue', 'date_with_timestamp', 'status', 'result',
40                                         'match_title', 'match_subtitle', 'home_team_id', 'away_team_id'])
41 else:
42     print("Failed to fetch data from:", url)
43     return None
44
45 # Function to save DataFrame to CSV and upload to S3
46 def upload_to_s3(df, bucket_name, folder):
47     current_date = datetime.datetime.now()
48     year = current_date.year
49     month = current_date.month
50
51     s3_output_name = f"{folder}/"
52     s3_output_name += f"{year}/"
53     s3_output_name += f"{month:02}/"
54     s3_output_name += f"{folder}_{year}_{month:02}_{current_date.day}.csv"
55
56     csv_buffer = StringIO()
57     df.to_csv(csv_buffer, index=False)
58     s3 = boto3.client('s3', aws_access_key_id=AWS_ACCESS_KEY_ID, aws_secret_access_key=AWS_SECRET_ACCESS_KEY)
59     s3.put_object(Body=csv_buffer.getvalue(), Bucket=bucket_name, Key=s3_output_name)
60     print(f"Data has been uploaded to S3 bucket {bucket_name} with key {s3_output_name}")
61
62     fixtures_df = fetch_data_from_api("https://cricket-live-data.p.rapidapi.com/fixtures")
63     cols = ['venue', 'result', 'match_title', 'match_subtitle']
64     for col in cols:
65         fixtures_df[col] = fixtures_df[col].str.replace(',', '')
66     if fixtures_df is not None:
67         upload_to_s3(fixtures_df, S3_BUCKET_NAME, S3_UPLOADED_FILE_NAME)
68
69
70

```

3.Fixtures_by_series:

```
home > growlt256 > Courses > TALEND > Python Files for deploy > Fixtures_by_series_api_to_s3.py > fetch_fixtures_by_series

1 import os
2
3 # Global variables
4 RAPIDAPI_KEY = "b72f452b32msheb263a2dd53ade5p151fa5jsn1016eed39cfcc"
5 RAPIDAPI_HOST = "cricket-live-data.p.rapidapi.com"
6 AWS_ACCESS_KEY_ID = "AKIARAAACMVC5URUSKCGX"
7 AWS_SECRET_ACCESS_KEY = "cRE7tOKlsczBdU2w6jM7Se/fx1FN5yRRpeaKmc0"
8 S3_BUCKET_NAME = 'cricket-api-to-s3'
9 S3_UPLOADED_FILE_NAME = 'fixtures_by_series'
10
11 base_url = "https://cricket-live-data.p.rapidapi.com/fixtures-by-series/"
12 headers = {
13     "X-RapidAPI-Key": RAPIDAPI_KEY,
14     "X-RapidAPI-Host": RAPIDAPI_HOST
15 }
16
17 def fetch_fixtures_by_series(date_list):
18     df_date = []
19     df_venue = []
20     for date in date_list:
21         url = base_url + str(date)
22         response = requests.get(url, headers=headers)
23         if response.status_code == 200:
24             json_data = response.json()
25             for i in range(len(json_data['results'])):
26                 entry = []
27                 entry_venue = []
28                 ide = json_data['results'][i]['id']
29                 series_id = json_data['results'][i]['series_id']
30                 venue = json_data['results'][i]['venue']
31                 date_with_timestamp = json_data['results'][i]['date']
32                 status = json_data['results'][i]['status']
33                 result = json_data['results'][i]['result']
34                 match_title = json_data['results'][i]['match_title']
35                 match_subtitle = json_data['results'][i]['match_subtitle']
36                 home_team_id = json_data['results'][i]['home']['id']
37                 away_team_id = json_data['results'][i]['away']['id']
38                 venue_id = json_data['results'][i]['home']['id']
39                 name = json_data['results'][i]['home']['name']
40                 code = json_data['results'][i]['home']['code']
41
42                 entry.append([ide, series_id, venue, date_with_timestamp, status, result, match_title, match_subtitle,
43                               home_team_id, away_team_id])
44                 entry_venue.append([venue_id, name, code])
45
46                 df_date.extend(entry)
47                 df_venue.extend(entry_venue)
48
49             return (pd.DataFrame(df_date, columns = ['id', 'series_id', 'venue', 'date_with_timestamp', 'status', 'result',
50                                   'match_title', 'match_subtitle', 'home_team_id', 'away_team_id']),
51                   pd.DataFrame(df_venue, columns = ['id', 'name', 'code']))
52
53
54 def upload_to_s3(df, bucket_name, S3_UPLOADED_FILE_NAME):
55     current_date = datetime.datetime.now()
56     year = current_date.year
57     month = current_date.month
58
59     s3_output_name = f"{S3_UPLOADED_FILE_NAME}/"
60     s3_output_name += f"(year)@"
61     s3_output_name += f"(month:02)@"
62     s3_output_name += f"(S3_UPLOADED_FILE_NAME)_{year}_{month:02}_{current_date.day}.csv"
63
64     # Save DataFrame to CSV and upload to S3
65     csv_buffer = StringIO()
66     df.to_csv(csv_buffer, index=False)
67     s3 = boto3.client('s3', aws_access_key_id=AWS_ACCESS_KEY_ID, aws_secret_access_key=AWS_SECRET_ACCESS_KEY)
68     s3.put_object(Body=csv_buffer.getvalue(), Bucket=bucket_name, Key=s3_output_name)
69     print(f"Data has been uploaded to S3 bucket {bucket_name} with key {s3_output_name}")
70
71
72
73 series_ids = [2002, 1430, 978, 812, 833, 608, 756, 425]
74
75 fixtures_df = fetch_fixtures_by_series(series_ids)[0]
76 venue_df = fetch_fixtures_by_series(series_ids)[1]
77
78 cols = ['venue', 'result', 'match_title', 'match_subtitle']
79 for col in cols:
80     fixtures_df[col] = fixtures_df[col].str.replace(',', '')
81
82 cols = ['name']
83 for col in cols:
84     venue_df[col] = venue_df[col].str.replace(',', '')
85
86 upload_to_s3(fixtures_df, S3_BUCKET_NAME, S3_UPLOADED_FILE_NAME)
87 upload_to_s3(venue_df, S3_BUCKET_NAME, 'Venue')
```

4.Fixtures_by_date

```

home > growlt256 > Courses > TALEND > Python files for deploy > fixtures_by_date_api_to_s3.py > ...
1 import pandas as pd
2 import requests
3 import boto3
4 from io import StringIO
5 import datetime
6
7 # Global variables
8 RAPIDAPI_KEY = "b72f452b32msheb263a2dd53ade5p15lfa5jsn10l6eed39cfc"
9 RAPIDAPI_HOST = "cricket-live-data.p.rapidapi.com"
10 AWS_ACCESS_KEY_ID = 'AKIARAACMV5URUSKCGX'
11 AWS_SECRET_ACCESS_KEY = 'crE7tDK1scz8dU2w6jM7Se/fx1FN5ywRRpeaKmc0'
12 S3_BUCKET_NAME = 'cricket-api-to-s3'
13 S3_FOLDER_NAME = 'fixtures_by_date'
14
15 # API endpoint base URL and headers
16 base_url = "https://cricket-live-data.p.rapidapi.com/fixtures-by-date/"
17 headers = {
18     "X-RapidAPI-Key": RAPIDAPI_KEY,
19     "X-RapidAPI-Host": RAPIDAPI_HOST
20 }
21
22 arr_date = [
23     '2023-04-02', '2023-04-03', '2023-04-04', '2023-04-05', '2024-03-22', '2024-03-23', '2024-03-23',
24     '2024-03-24', '2024-03-24', '2024-03-25', '2024-03-26', '2024-03-27', '2024-03-28', '2024-03-29',
25     '2024-03-30', '2024-03-31', '2024-03-31', '2024-04-01', '2024-04-02', '2024-04-03'
26 ]
27
28 def fetch_fixtures_by_date(arr_date):
29     df_date = []
30     for date in arr_date:
31         url = base_url + date
32         response = requests.get(url, headers=headers)
33         if response.status_code == 200:
34             json_data = response.json()
35             for i in range(len(json_data['results'])):
36                 entry = []
37                 ide = json_data['results'][i]['id']
38                 series_id = json_data['results'][i]['series_id']
39                 venue = json_data['results'][i]['venue']
40                 date_with_timestamp = json_data['results'][i]['date']
41                 status = json_data['results'][i]['status']
42                 result = json_data['results'][i]['result']
43                 match_title = json_data['results'][i]['match_title']
44                 match_subtitle = json_data['results'][i]['match_subtitle']
45                 home_team_id = json_data['results'][i]['home']['id']
46                 away_team_id = json_data['results'][i]['away']['id']
47                 entry.append(ide, series_id, venue, date_with_timestamp, status, result, match_title, match_subtitle,
48                             home_team_id, away_team_id)
49             df_date.extend(entry)
50     return pd.DataFrame(df_date, columns = [
51         'id', 'series_id', 'venue', 'date_with_timestamp', 'status', 'result',
52         'match_title', 'match_subtitle', 'home_team_id', 'away_team_id'])
53
54 # Function to save DataFrame to CSV and upload to S3
55 def upload_to_s3(df, bucket_name, folder, file_name):
56     current_date = datetime.datetime.now()
57     year = current_date.year
58     month = current_date.month
59
60     s3_output_name = f"{folder}/{year:02}/{month:02}/{file_name}_{year:02}_{month:02}_{current_date.day}.csv"
61
62     csv_buffer = StringIO()
63     df.to_csv(csv_buffer, index=False)
64     s3 = boto3.client('s3', aws_access_key_id=AWS_ACCESS_KEY_ID, aws_secret_access_key=AWS_SECRET_ACCESS_KEY)
65     s3.put_object(Body=csv_buffer.getvalue(), Bucket=bucket_name, Key=s3_output_name)
66     print(f'Data has been uploaded to S3 bucket {bucket_name} with key {s3_output_name}')
67
68 df_fixtures_by_date = fetch_fixtures_by_date(arr_date)
69
70 # upload_to_s3(df_fixtures_by_date, S3_BUCKET_NAME, S3_FOLDER_NAME, S3_FOLDER_NAME)
71 cols = ['venue', 'status', 'result', 'match_title', 'match_subtitle']
72 for col in cols:
73     df_fixtures_by_date[col] = df_fixtures_by_date[col].str.replace(',', '')
74 upload_to_s3(df_fixtures_by_date, S3_BUCKET_NAME, S3_FOLDER_NAME, S3_FOLDER_NAME)

```

5.Result:

```
home > growlt250 > Courses > TALEND > Python Files for deploy > results_api_to_s3.py > ...
1 import pandas as pd
2 import requests
3 import boto3
4 from io import StringIO
5 import datetime
6
7 # Global variables
8 RAPIDAPI_KEY = "f381eaaa79msh51537dd6e014035p180695jsn67cf7d549118"
9 RAPIDAPI_HOST = "cricket-live-data.p.rapidapi.com"
10 AWS_ACCESS_KEY_ID = 'AKIARAACMV5URUSKCGX'
11 AWS_SECRET_ACCESS_KEY = 'crE7tDK1sczBdU2w6jM7Se/fx1FN5ywRRpeaKmc0'
12 S3_BUCKET_NAME = 'cricket-api-to-s3'
13 S3_FOLDER_NAME = 'results'
14
15
16 def fetch_data_from_api(url):
17     headers = {
18         "X-RapidAPI-Key": RAPIDAPI_KEY,
19         "X-RapidAPI-Host": RAPIDAPI_HOST
20     }
21     response = requests.get(url, headers=headers)
22     if response.status_code == 200:
23         json_data = response.json()
24         df_date = []
25         for i in range(len(json_data['results'])):
26             entry = []
27             ide = json_data['results'][i]['id']
28             series_id = json_data['results'][i]['series_id']
29             venue = json_data['results'][i]['venue']
30             date_with_timestamp = json_data['results'][i]['date']
31             status = json_data['results'][i]['status']
32             result = json_data['results'][i]['result']
33             match_title = json_data['results'][i]['match_title']
34             match_subtitle = json_data['results'][i]['match_subtitle']
35             home_team_id = json_data['results'][i]['home']['id']
36             away_team_id = json_data['results'][i]['away']['id']
37             entry.append([ide, series_id, venue, date_with_timestamp, status, result, match_title, match_subtitle,
38                         home_team_id, away_team_id])
39             df_date.extend(entry)
40     return pd.DataFrame(df_date, columns = ['id', 'series_id', 'venue', 'date_with_timestamp', 'status', 'result',
41                                         'match_title', 'match_subtitle', 'home_team_id', 'away_team_id'])
42 else:
43     print("Failed to fetch data from:", url)
44     return None
45
46 # Function to save DataFrame to CSV and upload to S3
47 def upload_to_s3(df, bucket_name, folder, file_name):
48
49     current_date = datetime.datetime.now()
50     year = current_date.year
51     month = current_date.month
52
53     s3_output_name = f"{folder}/{year}/{month:02}/{file_name}_{year}_{month:02}_{current_date.day}.csv"
54
55
56     csv_buffer = StringIO()
57     df.to_csv(csv_buffer, index=False)
58
59     s3 = boto3.client('s3', aws_access_key_id=AWS_ACCESS_KEY_ID, aws_secret_access_key=AWS_SECRET_ACCESS_KEY)
60     s3.put_object(Body=csv_buffer.getvalue(), Bucket=bucket_name, Key=s3_output_name)
61
62     print(f"Data has been uploaded to S3 bucket {bucket_name} with key {s3_output_name}")
63
64
65 results_df = fetch_data_from_api("https://cricket-live-data.p.rapidapi.com/results")
66
67 cols = ['venue', 'result', 'match_title', 'match_subtitle']
68 for col in cols:
69     results_df[col] = results_df[col].str.replace(',', '')
70
71
72 # Upload to S3
73 if results_df is not None:
74     upload_to_s3(results_df, S3_BUCKET_NAME, S3_FOLDER_NAME, S3_FOLDER_NAME)
```

6.Result_by_date

```
home > growlt256 > Courses > TALEND > Python Files for deploy > results_by_date_api_to_s3.py > ...
1 import pandas as pd
2 import requests
3 import boto3
4 from io import StringIO
5 import datetime
6
7 # Global variables
8 RAPIDAPI_KEY = "f381eaaa79msf51537dd6e014035p180695jsn67cf7d549118"
9 RAPIDAPI_HOST = "cricket-live-data.p.rapidapi.com"
10 AWS_ACCESS_KEY_ID = "AKTARAACMVC5URUSKCGX"
11 AWS_SECRET_ACCESS_KEY = 'crE7tDK1sczBdU2w6jM75e/fx1FN5ywRRpeaKmc0'
12 S3_BUCKET_NAME = 'cricket-api-to-s3'
13 S3_UPLOADED_FILE_NAME = 'results-by-date'
14
15 # API endpoint base URL and headers
16 base_url = "https://cricket-live-data.p.rapidapi.com/results-by-date/"
17 headers = {
18     "X-RapidAPI-Key": RAPIDAPI_KEY,
19     "X-RapidAPI-Host": RAPIDAPI_HOST
20 }
21 def fetch_result_by_date(date_list):
22     df_date = []
23     for date in date_list:
24         url = base_url + date
25         response = requests.get(url, headers=headers)
26         if response.status_code == 200:
27             json_data = response.json()
28             for i in range(len(json_data['results'])):
29                 entry = []
30                 ide = json_data['results'][i]['id']
31                 series_id = json_data['results'][i]['series_id']
32                 venue = json_data['results'][i]['venue']
33                 date_with_timestamp = json_data['results'][i]['date']
34                 status = json_data['results'][i]['status']
35                 result = json_data['results'][i]['result']
36                 match_title = json_data['results'][i]['match_title']
37                 match_subtitle = json_data['results'][i]['match_subtitle']
38                 home_team_id = json_data['results'][i]['home']['id']
39                 away_team_id = json_data['results'][i]['away']['id']
40                 entry.append([ide, series_id, venue, date_with_timestamp, status, result, match_title, match_subtitle,
41                               home_team_id, away_team_id])
42             df_date.extend(entry)
43     return pd.DataFrame(df_date, columns = ['id', 'series_id', 'venue', 'date_with_timestamp', 'status', 'result',
44                                         'match_title', 'match_subtitle', 'home_team_id', 'away_team_id'])
45
46
47 # Function to save DataFrame to CSV and upload to S3
48 def upload_to_s3(df, bucket_name, folder):
49     # Get current year and month
50     current_date = datetime.datetime.now()
51     year = current_date.year
52     month = current_date.month
53
54     # Construct S3 key
55     s3_output_name = f"{folder}/"
56     s3_output_name += f"{year}/"
57     s3_output_name += f"(month:02)/"
58     s3_output_name += f"{folder}_{year}_{month:02}_{current_date.day}.csv"
59
60     # Save DataFrame to CSV and upload to S3
61     csv_buffer = StringIO()
62     df.to_csv(csv_buffer, index=False)
63     s3 = boto3.client('s3', aws_access_key_id=AWS_ACCESS_KEY_ID, aws_secret_access_key=AWS_SECRET_ACCESS_KEY)
64     s3.put_object(Body=csv_buffer.getvalue(), Bucket=bucket_name, Key=s3_output_name)
65     print(f'Data has been uploaded to S3 bucket {bucket_name} with key {s3_output_name}')
66
67 # List of dates
68 arr_dates = ['2024-03-19', '2024-03-20', '2024-03-21', '2024-03-22', '2024-03-23',
69             '2024-03-23', '2024-03-24', '2024-03-24', '2024-03-25', '2024-03-26',
70             '2024-03-27', '2024-03-28', '2024-03-29', '2024-03-30', '2024-03-31',
71             '2024-03-31', '2024-04-01']
72
73 # Fetch data by date
74 df_result_by_date = fetch_result_by_date(arr_dates)
75
76 cols = ['venue', 'result', 'match_title', 'match_subtitle']
77 for col in cols:
78     df_result_by_date[col] = df_result_by_date[col].str.replace(',', '')
79
80 # Upload to S3
81 upload_to_s3(df_result_by_date, S3_BUCKET_NAME, S3_UPLOADED_FILE_NAME)
82
```

7.Match_scorecard

```
home > growlt256 > Courses > TALEND > Python Files for deploy > ⚡ match_scorecard_api_to_s3.py > ...
1 import pandas as pd
2 import requests
3 import boto3
4 from io import StringIO
5 import datetime
6
7 # Global variables
8 RAPIDAPI_KEY = "f38leaaa79msh51537dd6e014035p180695jsn67cf7d549118"
9 RAPIDAPI_HOST = "cricket-live-data.p.rapidapi.com"
10 AWS_ACCESS_KEY_ID = 'AKIARAACMVC5URUSKCGX'
11 AWS_SECRET_ACCESS_KEY = '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'
12 S3_BUCKET_NAME = 'cricket-api-to-s3'
13 S3_UPLOADED_FILE_NAME = 'match_scorecard'
14
15 match_arr = [2844251, 2844253, 2844255, 2844257, 2844259, 2844261, 2844263, 2844265,
16           2844267, 2844269, 2844271, 2844273, 2844275, 2844277, 2844279,
17           2844281, 2844283, 2844285, 2844287, 2844289]
18
19 def fetch_match_data(match_ids):
20     entry = []
21     for num in match_ids:
22         url = f'https://cricket-live-data.p.rapidapi.com/match/{num}'
23         headers = {
24             "X-RapidAPI-Key": RAPIDAPI_KEY,
25             "X-RapidAPI-Host": RAPIDAPI_HOST
26         }
27         response = requests.get(url, headers=headers)
28         json_data = response.json()
29         ide = json_data['results'][0]['fixture'][0]['id']
30         series_id = json_data['results'][0]['fixture'][0]['series_id']
31         match_title = json_data['results'][0]['fixture'][0]['match_title']
32         venue = json_data['results'][0]['fixture'][0]['venue']
33         start_date = json_data['results'][0]['fixture'][0]['start_date']
34         end_date = json_data['results'][0]['fixture'][0]['end_date']
35         home = json_data['results'][0]['fixture'][0]['home'][0]['id']
36         away = json_data['results'][0]['fixture'][0]['away'][0]['id']
37         try:
38             result = json_data['results'][0]['live_details'][0]['match_summary'][0]['status']
39         except:
40             result = "Match not completed"
41         record = [ide, series_id, match_title, venue, start_date, end_date, home, away, result]
42         entry.append(record)
43     return pd.DataFrame(entry, columns=['id', 'series_id', 'match_title', 'venue', 'start_date',
44                                     'end_date', 'home_team_id', 'away_team_id', 'result'])
45
46 def upload_to_s3(df, bucket_name, folder):
47     # Get current year and month
48     current_date = datetime.datetime.now()
49     year = current_date.year
50     month = current_date.month
51
52     # Construct S3 key
53     s3_output_name = f'{folder}/'
54     s3_output_name += f'{year}/'
55     s3_output_name += f'{month:02}/'
56     s3_output_name += f'{folder}_{year}_{month:02}_{current_date.day}.csv'
57
58     # Save DataFrame to CSV and upload to S3
59     csv_buffer = StringIO()
60     df.to_csv(csv_buffer, index=False)
61     s3 = boto3.client('s3', aws_access_key_id=AWS_ACCESS_KEY_ID, aws_secret_access_key=AWS_SECRET_ACCESS_KEY)
62     s3.put_object(Body=csv_buffer.getvalue(), Bucket=bucket_name, Key=s3_output_name)
63     print(f'Data has been uploaded to S3 bucket {bucket_name} with key {s3_output_name}')
64
65 # Fetch match data
66 df_match = fetch_match_data(match_arr)
67
68 cols = ['match_title', 'venue', 'result']
69
70 for col in cols:
71     df_match[col] = df_match[col].str.replace(' ', '')
72
73 # Upload to S3
74 upload_to_s3(df_match, S3_BUCKET_NAME, S3_UPLOADED_FILE_NAME)
75
76
```

S3 Bucket interface

The screenshot shows the AWS S3 console. On the left, the navigation pane is visible with sections like Buckets, Storage Lens, and Feature spotlight. The main area displays the contents of the 'cricket-api-to-s3' bucket. The top navigation bar includes 'Amazon S3', 'Services', 'Search', and account information ('Stockholm', 'Vivek @ 0687-2450-0667'). Below the navigation bar, there are tabs for 'Objects', 'Properties', 'Permissions', 'Metrics', 'Management', and 'Access Points'. The 'Objects' tab is selected, showing a list of 9 objects. The table headers are 'Name', 'Type', 'Last modified', 'Size', and 'Storage class'. The objects listed are: data/, fixtures_by_date/, fixtures_by_series/, fixtures/, match_scorecard/, results-by-date/, results/, series/, and Venue/. Each object is a folder.

Problems we have faced:

1. Most of the time we faced api errors because of our free version.
2. While debugging code according to the requirement of upstream.
3. We have faced a problem while dealing with nested JSON files.

3. Send S3 data to Oracle with Talend - Vivek

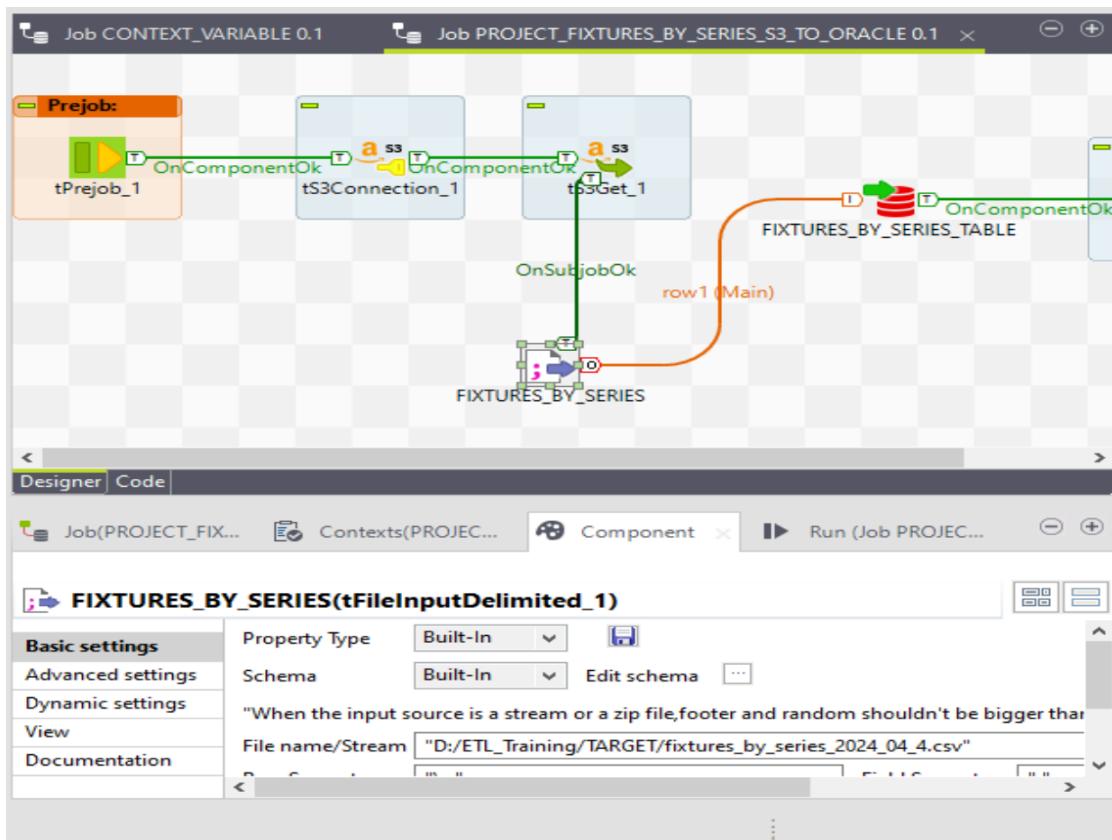
Here we have to prepare 7 individual jobs as we have 7 csv files in S3. So to process each file we have to create individual job for each file. Configuration in these jobs are same. So attaching only one jobs configuration screen shots. We have created metadata of all the csv files in metadata section to give schema information in our component.

1. **tS3Connection:** Start by connecting to the TeamSpeak 3 server using **ts3connection**, then using **ts3get** to retrieve information or perform actions like querying user data.

2. After processing data from the TeamSpeak server, using tfileinputdelimited to read CSV data, transform it as needed, and then using tdboutput to write the transformed data into a database table.
3. Once the data processing is complete and the connection is no longer needed, using ts3close to gracefully close the connection to the TeamSpeak 3 server.

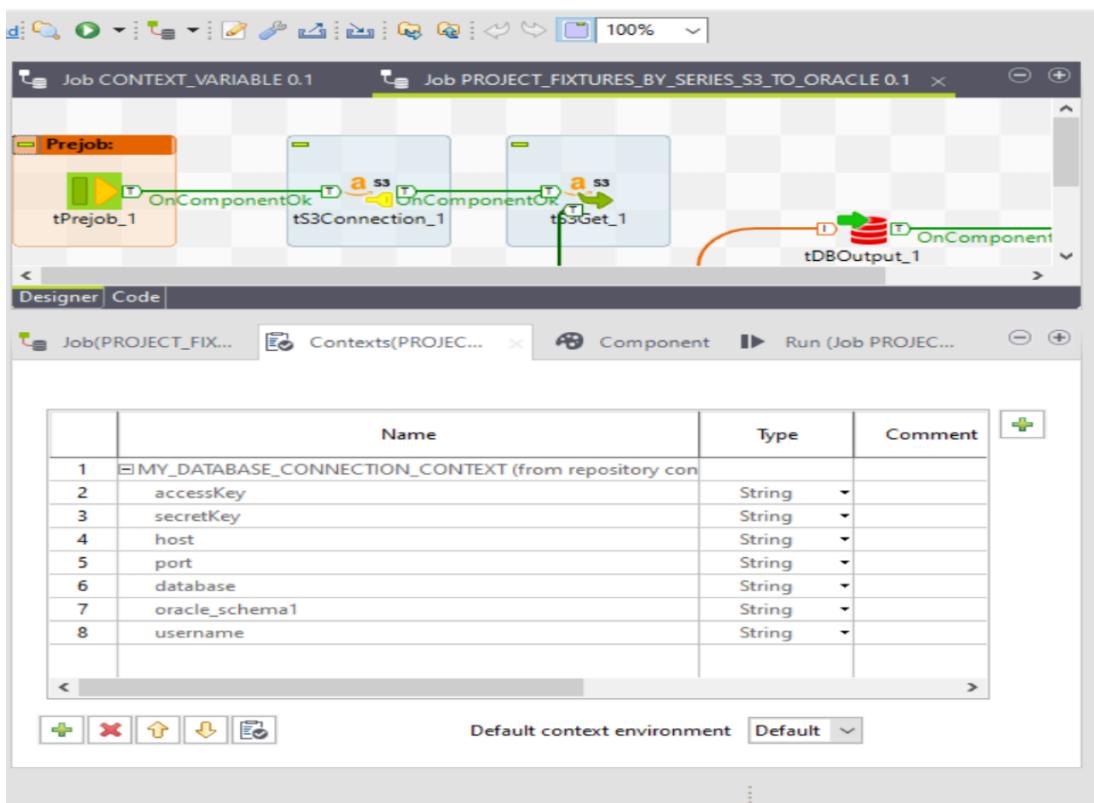
Here are the screenshots for fixtures by series data sets. All the connection and context parameter screenshots are attached.

1.

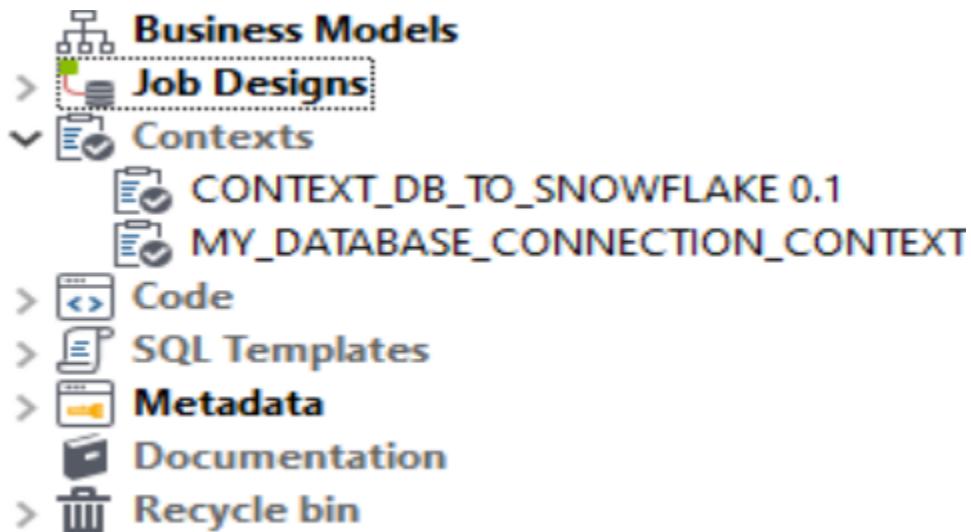


Setting up the context parameter for a particular job. Here database connection details are the same for all the jobs. But if in any scenario we want to change the database connection details then we can change it from the .csv file that we have created as a context parameter.

2.

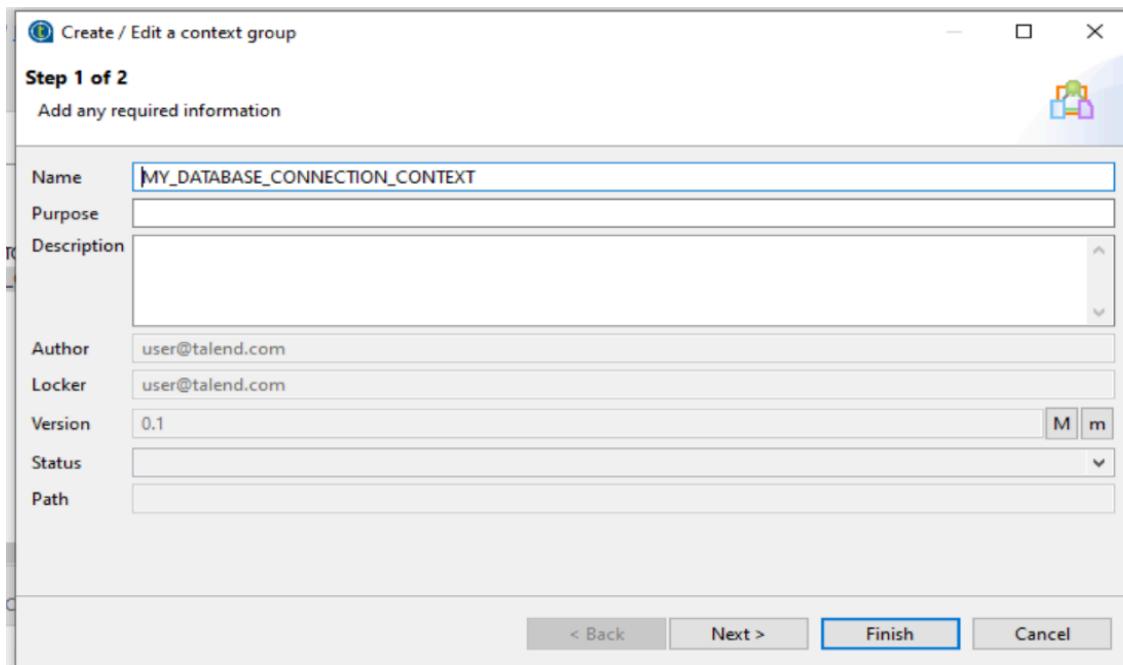


3.

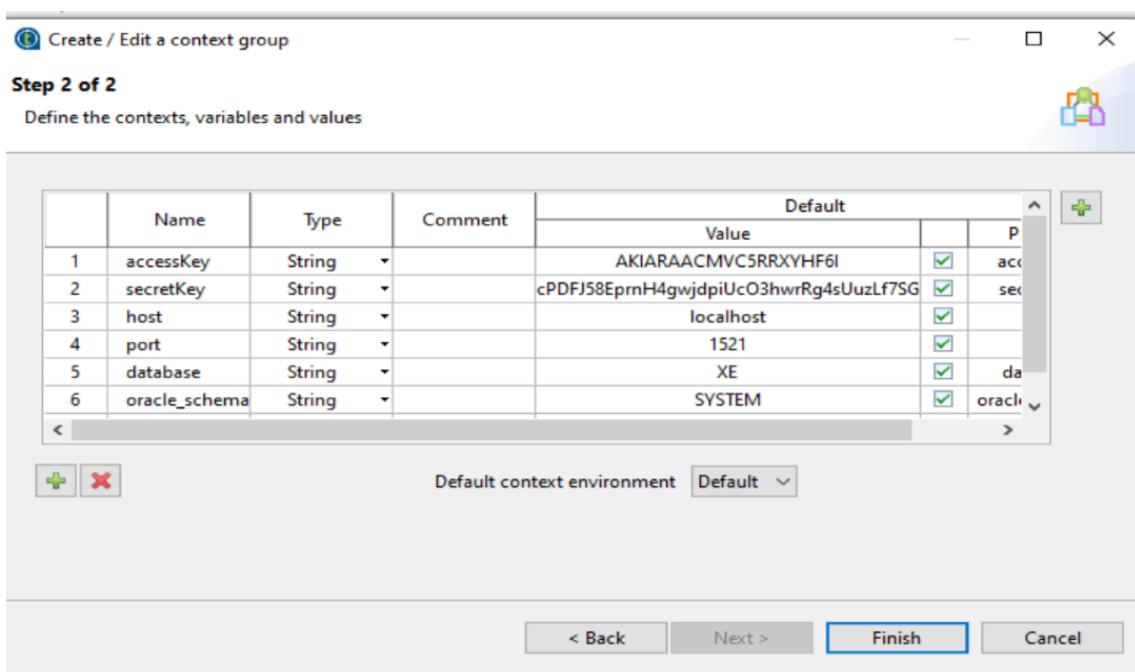


Parameter name in context file and in context group should be the same, otherwise it will give an error.

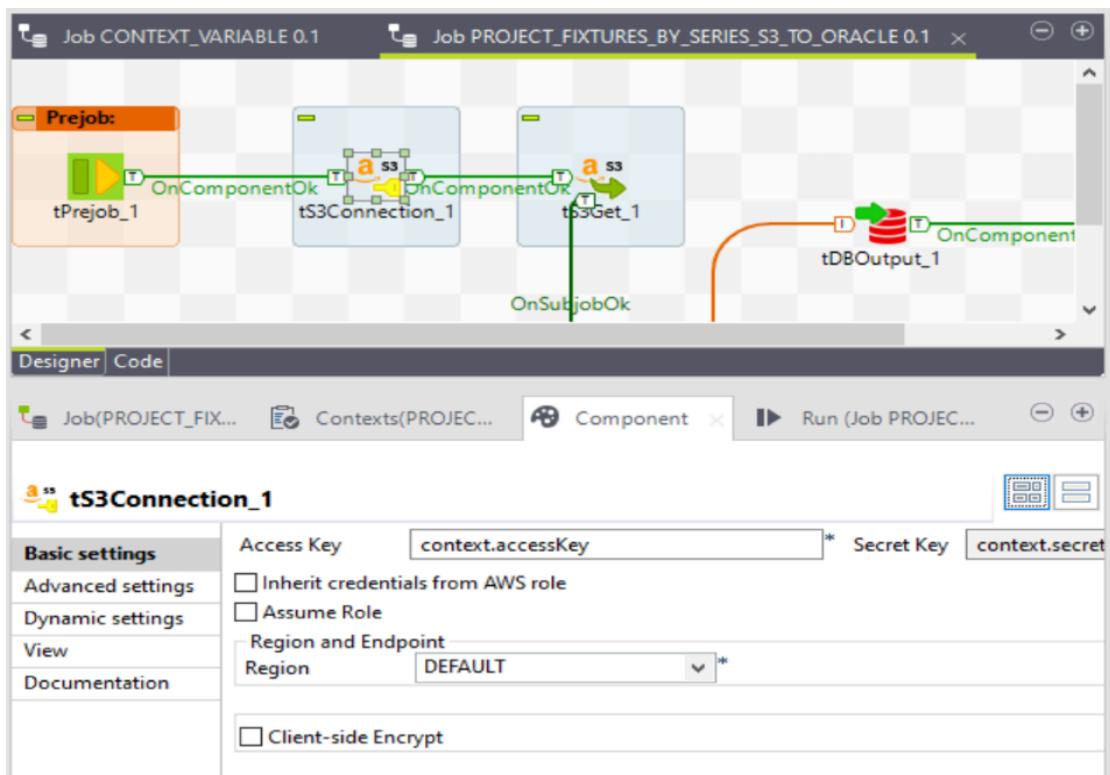
4.



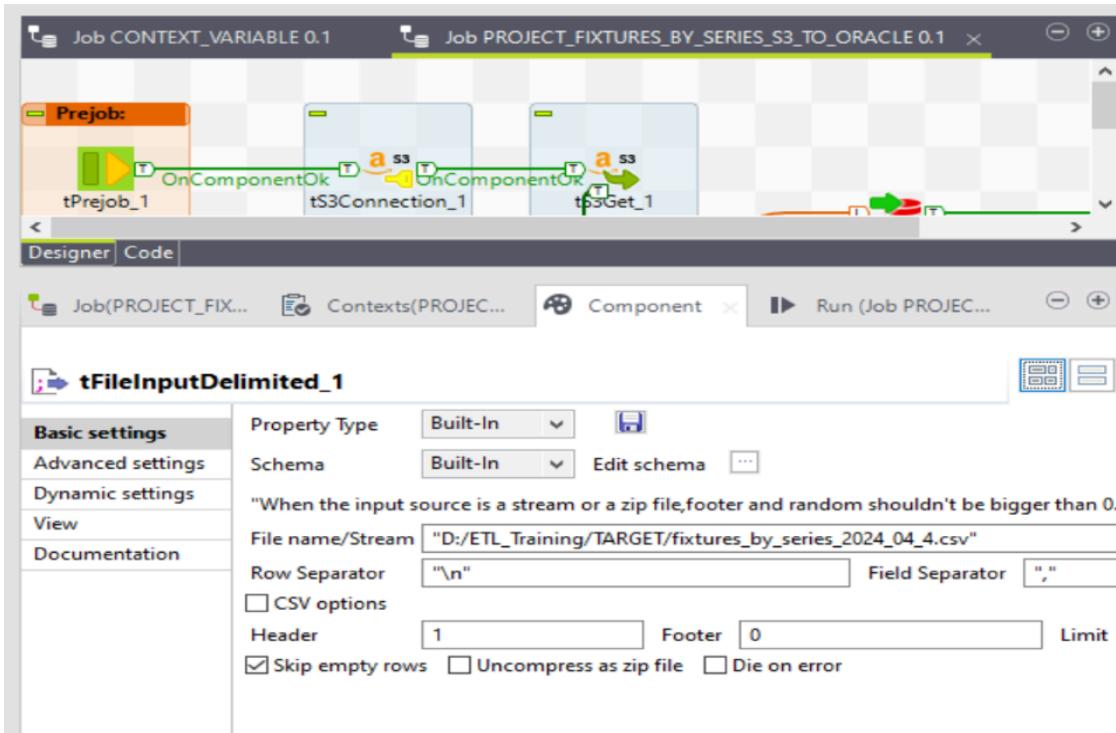
5.



6.



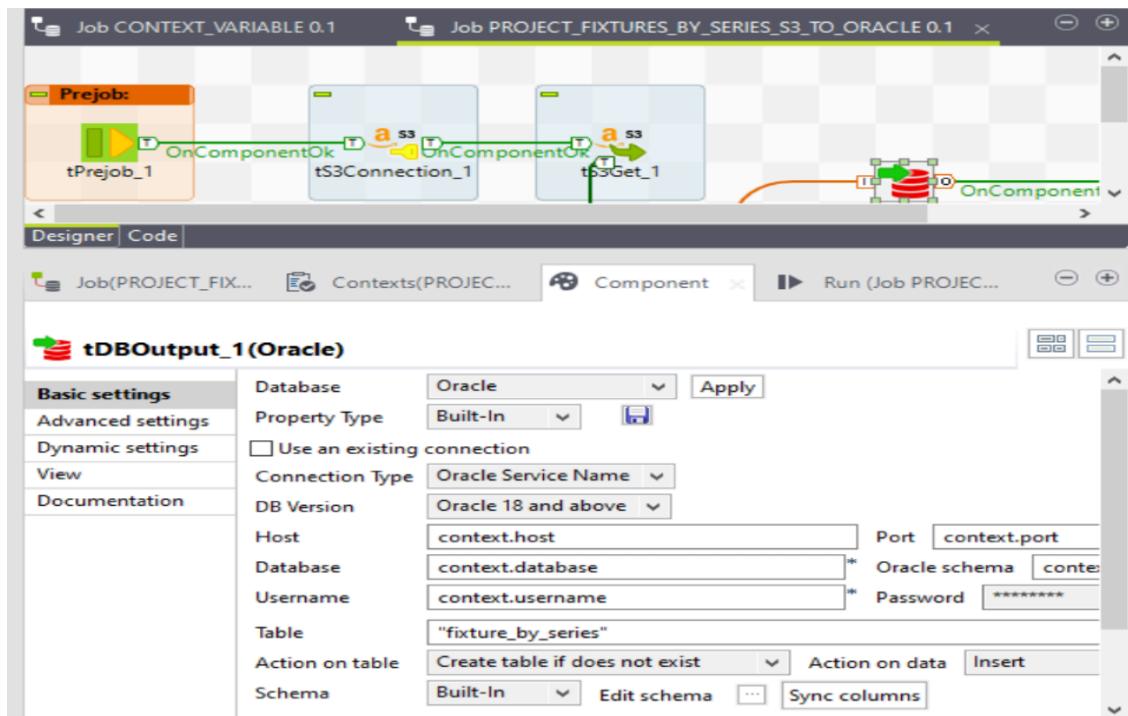
7.



8.



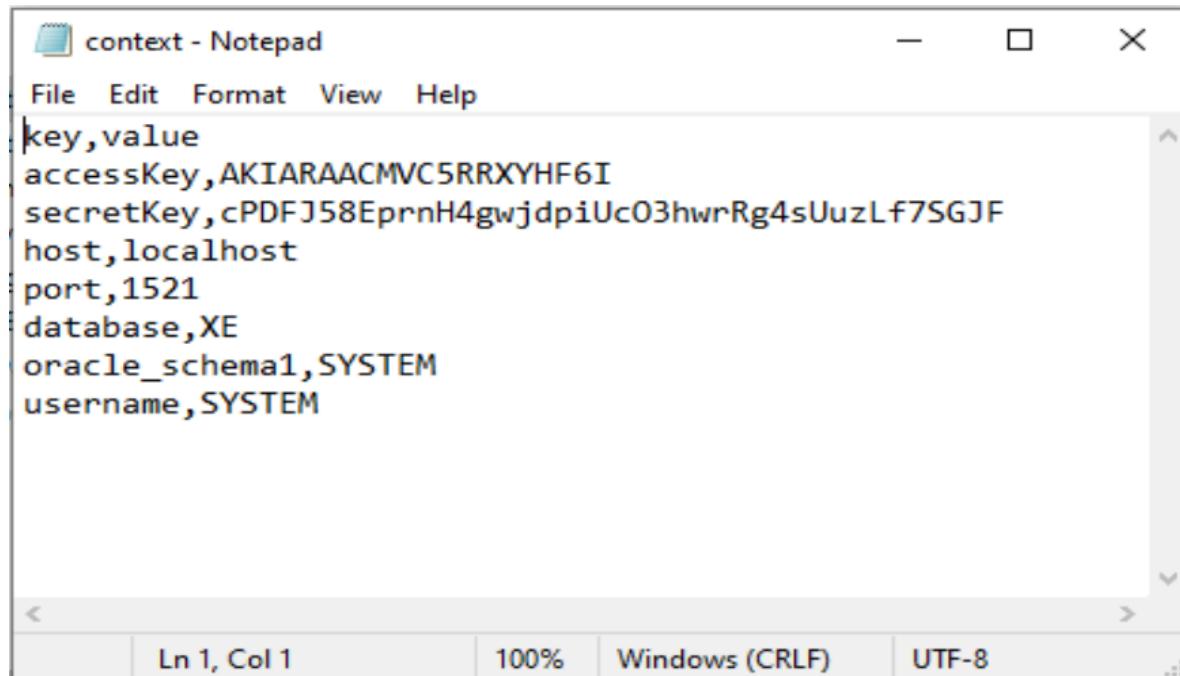
9.



10.

- > Employee 0.1
- > fixture_by_series 0.1
- > Fixtures_by_date 0.1
- > fixtures 0.1
- > Intern_ETL_stdents 0.1
- > Maths_Score 0.1
- > My_Placemente_Dataset 0.1
- > normlization_data 0.1
- > RESULT_BY_DATE 0.1
- > resultbydate 0.1
- > results 0.1
- > sampledenorm 0.1
- > series 0.1
- > venue 0.1

11.



The screenshot shows a Windows Notepad window with the title bar 'context - Notepad'. The menu bar includes File, Edit, Format, View, and Help. The main content area contains the following text:

```
key,value
accessKey,AKIARAACMVC5RRXYHF6I
secretKey,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
host,localhost
port,1521
database,XE
oracle_schema1,SYSTEM
username,SYSTEM
```

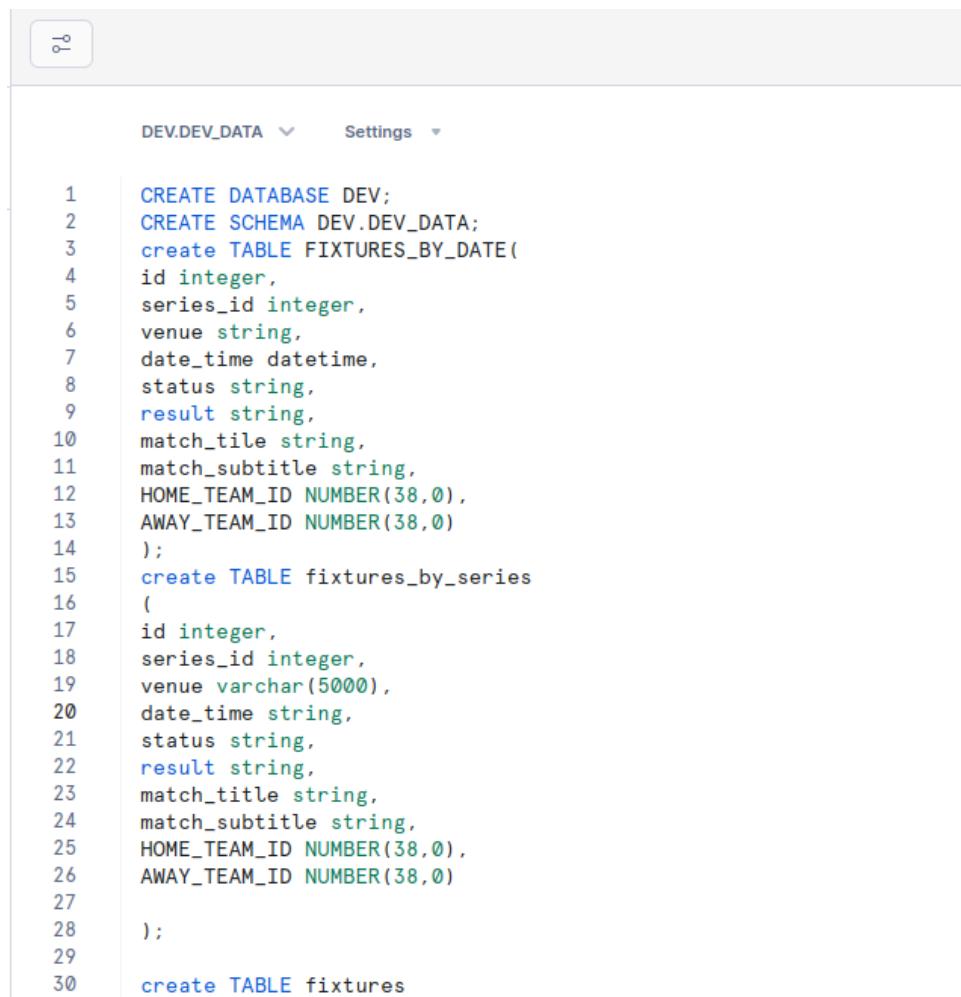
The bottom status bar shows 'Ln 1, Col 1', '100%', 'Windows (CRLF)', and 'UTF-8'.

This is the .csv file for setting up the connection details.

4. Migrate Oracle data to Snowflake with Talend - [Purvang,Jagruti]

Creating Schema and Database in snowflake:

DEV, PROD and QA:



The screenshot shows a database editor interface for the DEV.DEV_DATA schema. The top bar includes a refresh icon, a dropdown menu, and a 'Settings' button. The main area displays the following SQL code:

```
1 CREATE DATABASE DEV;
2 CREATE SCHEMA DEV.DEV_DATA;
3 create TABLE FIXTURES_BY_DATE(
4     id integer,
5     series_id integer,
6     venue string,
7     date_time datetime,
8     status string,
9     result string,
10    match_tile string,
11    match_subtitle string,
12    HOME_TEAM_ID NUMBER(38,0),
13    AWAY_TEAM_ID NUMBER(38,0)
14 );
15 create TABLE fixtures_by_series
16 (
17     id integer,
18     series_id integer,
19     venue varchar(5000),
20     date_time string,
21     status string,
22     result string,
23     match_title string,
24     match_subtitle string,
25     HOME_TEAM_ID NUMBER(38,0),
26     AWAY_TEAM_ID NUMBER(38,0)
27 );
28 );
29
30 create TABLE fixtures
```

The screenshot shows a database management interface with the following details:

- Top Bar:** A small icon in a rounded square.
- Database Selection:** QA.QA_DATA
- Settings:** A dropdown menu.

The main area contains the following SQL code:

```
1 CREATE DATABASE QA;
2 CREATE SCHEMA QA.QA_DATA;
3 create TABLE FIXTURES_BY_DATE(
4     id integer,
5     series_id integer,
6     venue string,
7     date_time datetime,
8     status string,
9     result string,
10    match_title string,
11    match_subtitle string
12 );
13 create TABLE fixtures_by_series
14 (
15     id integer,
16     series_id integer,
17     venue varchar(5000),
18     date_time datetime,
19     status string,
20     result string,
21     match_title string,
22     match_subtitle string
23 );
24
25 create TABLE fixtures
26 (
27     id integer,
```

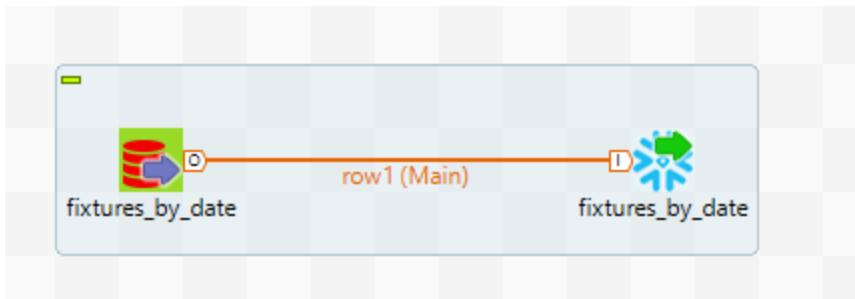
Below the code, there are two buttons:

- Results** (highlighted in blue)
- Chart**

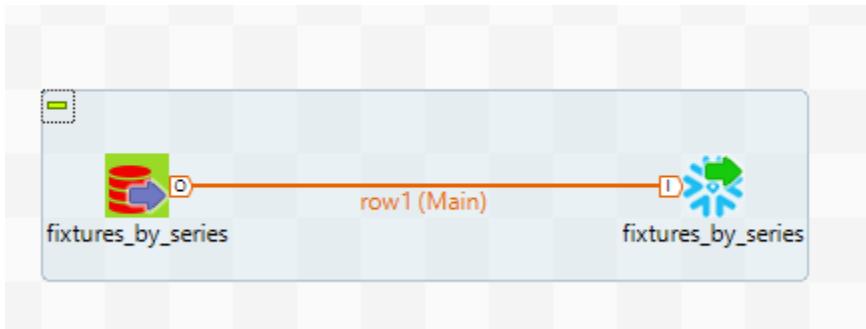
```
PROD.PROD_DATA ▾ Settings ▾

1   CREATE DATABASE PROD;
2   CREATE SCHEMA PROD.PROD_DATA;
3   create TABLE FIXTURES_BY_DATE(
4       id integer,
5       series_id integer,
6       venue string,
7       date_time datetime,
8       status string,
9       result string,
10      match_title string,
11      match_subtitle string
12  );
13  create TABLE fixtures_by_series
14  (
15      id integer,
16      series_id integer,
17      venue varchar(5000),
18      date_time datetime,
19      status string,
20      result string,
21      match_title string,
22      match_subtitle string
23  );
24
25  create TABLE fixtures
26  (
27      id integer,
28      series_id integer,
29      venue varchar(5000),
30      date_time datetime,
31      status string,
32      result string,
33      match_title string,
34      match_subtitle string
35  );
36
```

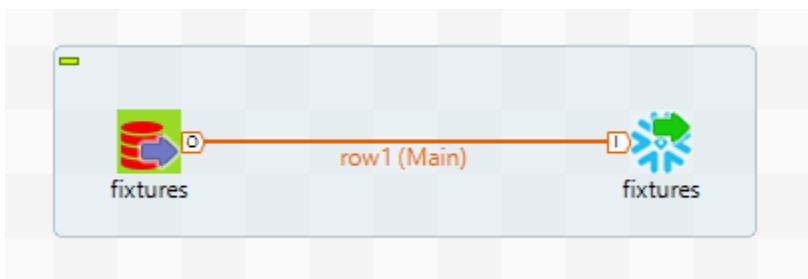
1. Fixtures_by_date:



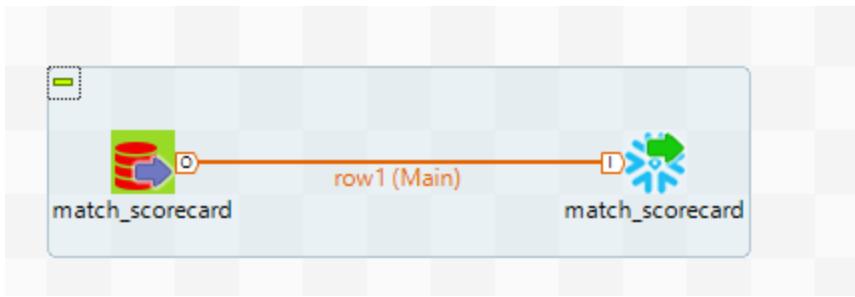
2. Fixtures_by_series:



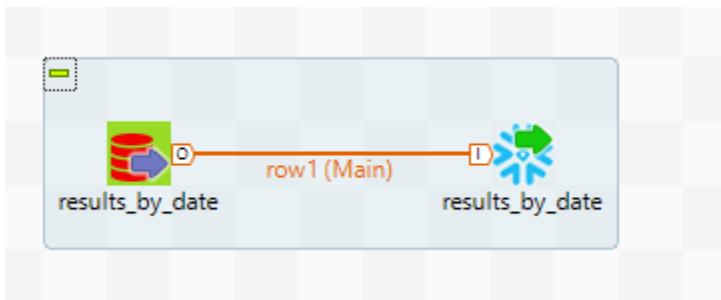
3. Results:



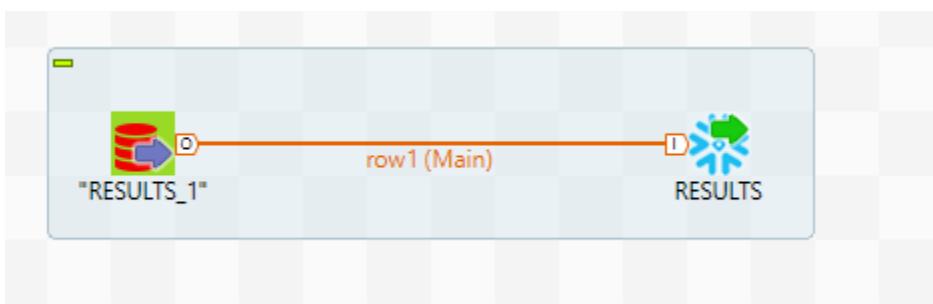
4. Match_scorecard:



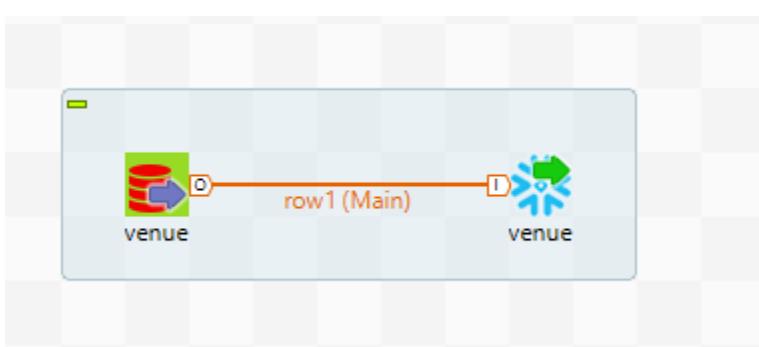
5. Results_by_date:



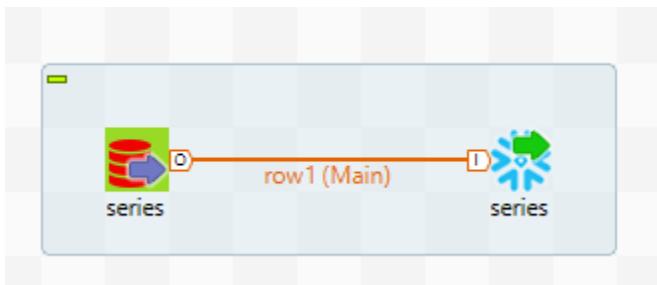
6. Results:



7. Venue:



8. Series:



Oracle Configuration:

Designer | Code

Component Run (Job venue)

series(tDBInput_1)(Oracle)

Basic settings

Database: Oracle | Apply

Property Type: Repository | DB (ORACLE_SERVICE_NAME):SCD_Final | ...

Use an existing connection

Connection Type: Oracle Service Name | ...

DB Version: Oracle 18 and above | ...

Host: "localhost" | Port: "1521" | ...

Database: "xe" | * Oracle schema: "" | ...

Username: "SYSTEM" | * Password: ***** | The variable attached to this parameter is: _SCHEMA_DB_

Schema: Built-In | Edit schema | ...

Table Name: "Venue_1" | ...

Query Type: Built-In | Guess Query | Guess schema

Query: "SELECT * FROM Venue_1"

Data source
This option only applies when deploying and running in the Talend Runtime

Activate Windows
Go to Settings to activate Windows.

Snowflake Configuration:

Designer | Code

Component | Run (Job venue)

series(tDBOutput_1)(Snowflake)

Basic settings

Database: Snowflake | Apply

Property Type: Built-In | Connection Component | Use this Component

Account: "wahnmfvqa98386"

Snowflake Region: AWS US West

Authentication Type: Basic

User Id: "PURVANG_DEV"

Password: *****

Warehouse: "COMPUTE_WH"

Schema: "DEV_DATA"

Database: "DEV"

Table: "VENUE"

Schema: Built-In | Edit schema | Sync columns

Table Action: CREATE_IF_NOT_EXISTS

Output Action: INSERT

Activate Windows
Go to Settings to activate Windows.

Snowflake:

Databases Worksheets

Pinned (0)

Fixtures_By_Series

Search objects

DEV

Tables

- FIXTURES
- FIXTURES_BY_DATE
- FIXTURES_BY_SERIES
- MATCH_SCORECARD
- RESULTS
- RESULTS_BY_DATE
- SERIES
- VENUE

```

6 COUNTRY VARCHAR(50)
)
8
9 -- SELECT * FROM emp2;
10 select * from match_scorecard
11 -----
12 -- SELECT * FROM results_by_dat
13 -----
14 SELECT * FROM fixtures_by_series;
15 -----
16 -----
17 SELECT * FROM RESULTS;
18 SELECT * FROM RESULTS WHERE ID=2834105;
19 SELECT COUNT(*) FROM RESULTS;
20 TRUNCATE TABLE RESULTS;
21
22 select get_ddl('table','results')
23 -- DROP TABLE RESULTS;
24
25 CREATE TABLE fixtures_by_series (
  id INT,
  name STRING,
  country STRING,
  code STRING
);
  
```

Results

ID	NAME	CODE
1	Chennai Super Kings	CSK
2	Punjab Kings	null
3	Kolkata Knight Riders	KKR
4	Rajasthan Royals	RAR
5	Gujarat Titans	null
6	Royal Challengers Bangalore	RCB
7	Chennai Super Kings	CSK
8	Sunrisers Hyderabad	SUH
9	Rajasthan Royals	RAR
10	Royal Challengers Bangalore	RCB
11	Lucknow Super Giants	null

Query Details

Query duration: 48ms

Rows: 313

Query ID: 01b385fd-0000-9d3d-...

ID: 101742

NAME: Delhi Capitals

NAME: Royal Challengers Bangalore

5. Prepare Data validation process and sample reports in Power BI for data prepared in Snowflake [Vipasha,Jagruti and Jigar]

The ETL process serves as the backbone of any data-driven decision-making process. It's the means through which raw data is transformed into meaningful insights. This makes data validation an integral part of ETL.

- **Maintain Data Quality:** Data validation helps in maintaining high data quality by identifying and rectifying errors, inconsistencies, and discrepancies in the data during the ETL process.
- **Ensure Accurate Reporting:** By validating data during the ETL process, you ensure that the final data used for reporting and analysis is accurate. This leads to reliable insights and informed decision-making.
- **Minimize Errors:** Through data validation, you can spot and rectify errors in the early stages of the ETL process, saving time and effort in handling them later.
- **Protect Data Integrity:** By ensuring that the transformation rules are correctly applied, and that the right data is loaded into the target system, data validation protects the integrity of your data.

Stages of Data Validation in ETL

Data validation in ETL isn't a one-time operation but rather a continuous process that spans across all stages of ETL. Let's take a closer look at how data validation is integrated at each step.

Extraction Stage: Validating the Source Data

The extraction phase is where data is collected from various source systems. It's the starting point of the ETL process and the first opportunity for data validation.

- **Data Completeness Check:** At this stage, you want to ensure that all required data has been extracted. This can be done by comparing record counts or using checksums from the source and the extracted data.
- **Data Accuracy Check:** Data pulled from source systems should match exactly what's in those systems. Basic field-level checks can be performed to verify accuracy.

- Initial Data Quality Check: If the quality of the source data is suspect, it may be helpful to perform some preliminary data quality checks to identify potential issues early.

Transformation Stage: Validating Transformation Rules

The transformation phase, where data is cleaned and transformed into a format suitable for the target system, is another crucial point for data validation.

- Validation of Transformation Rules: As data gets transformed, validating that the rules and logic applied are producing the expected results is crucial.
- Data Consistency Check: Ensuring that data remains consistent after transformation is key. For instance, the same type of data should not have different formats.
- Null Check: It's essential to check that mandatory fields are not left null after transformation.

Load Stage: Validating the Target Data

The final stage of ETL is loading the transformed data into the target system, usually a data warehouse. Here, validation ensures that the load operation was successful and accurate.

- Data Completeness Check: Similar to the extraction stage, you want to ensure that all data has been loaded into the target system.
- Data Integrity Check: Validating that relationships between data elements (like foreign key relationships) have been maintained during the load process is crucial.
- Reconciliation Check: This involves comparing the data in the source and target systems to ensure they match, confirming the ETL process was successful.

Techniques for Data Validation in ETL

Data validation in the ETL process encompasses a range of techniques designed to ensure data integrity, accuracy, and consistency. Here are some commonly utilized validation techniques:

Data Type Checks

Data type checks involve verifying that each data element is of the correct data type. This could include checking whether numeric fields contain numeric data, date fields contain valid dates, and so forth.

data quality issues even in the source system.

(i) Non-numerical type: Under this classification, we verify the accuracy of the non-numerical content. Examples are Emails, Pin codes, Phone in a valid format.

(ii) Domain analysis: In this type of test, we pick domains of data and validate for errors. There are three groupings for this:

- Based on Value: Here we create a list of values that can occur for a field (column in the table). Then validate if the column values are a subset of our list.
 - Example: Verify if the Gender column contains either M or F.
- Based on Range: Here we set minimum and maximum range for valid data values for a column, based on logical or business reasoning. We then validate if the column values fall within this range.
 - Example: 0 to 120 for Age.
- Reference File: Here the system uses an external validity file.
 - Example: Are Country codes valid, do they pick the right value from the reference file, are Country codes the same between QA and the Production environment? If the reference file had a country code updated, is it rightly updated in DB?
- Data Type Check: Example: Will Total Sales work correctly with Decimal (8, 16, or 20 bytes) or Double type?
- Data Length Check: Example: Will the data length for the Address field be sufficient with 500 chars? It might be a case where data migration is done as new geography is getting added to the company. The addresses of the new geography might have an exceedingly long format and sticking to the original length might error a use case.

- Index check: Example: Is there indexing done for the OrderId column in the target system? What if a merger of companies happened, requiring data migration and the Orders table grows 100 times in size in the target system?

Range Checks

Range checks validate that data values fall within acceptable ranges. For instance, if a data field is supposed to contain a person's age, a range check would confirm that the values fall within a plausible range, such as 0 to 120.

Constraint Checks

Constraint checks involve verifying that data adheres to predefined constraints. These could be unique constraints (e.g., every customer ID should be unique), primary key constraints (e.g., every row should have a unique identifier), or foreign key constraints (e.g., a reference to another table that must exist).

Consistency Checks

Consistency checks are used to ensure data values are consistent across datasets. For example, if you have two tables containing customer data, the customer names and IDs should be consistent across both.

Uniqueness Checks

Uniqueness checks are similar to constraint checks but specifically focus on ensuring that values in a certain field are unique where required. A typical example would be checking that each customer has a unique customer ID.

Null Data

In this type of test, we focus on the validity of null data and verification that the important column cannot be null.

- Example: Filter out all null data and validate if null is allowed.

- If there are important columns for business decisions, make sure nulls are not present.

Range Check

Data entity where ranges make business sense should be tested.

- Example: Order quantity per invoice cannot be more than 5K in the software category.
- Age should not be more than 120.

Referential Integrity Checks

Referential integrity checks involve validating that relationships between tables remain intact. This is particularly important in the loading stage of ETL, where maintaining relationships between data elements (like foreign keys) is crucial.

Using these techniques, you can help ensure that your data is clean, consistent, and reliable, regardless of its source. Remember, a single validation technique might not be enough. It's often a combination of these techniques that ensures thorough validation and helps maintain the quality and integrity of your data.

```

Data Validation.ipynb
File Edit View Insert Runtime Tools Help All changes saved
Comment Share Connect Colab AI

Table of contents + Code + Text
Mount Data
{x} File 1 = Fixtures
File 2 = Match Scorecard
File 3 = Series
File 4 = result by date
File 5 = result
+ Section

# Data Completeness Check
def data_completeness_check(source_data, extracted_data):
    return len(source_data) == len(extracted_data)

# Data Accuracy Check
def data_accuracy_check(source_data, extracted_data):
    return all(source_data[key] == extracted_data[key] for key in source_data)

# Initial Data Quality Check
def initial_data_quality_check(source_data):
    # Example: Checking for missing values in source data
    return all(value is not None for value in source_data.values())

# Validation of Transformation Rules
def validate_transformation_rules(transformed_data, expected_data):
    # Example: Simple comparison between transformed and expected data
    return transformed_data == expected_data

# Data Consistency Check
def data_consistency_check(transformed_data):
    # Example: Checking consistency in a list of values
    return len(set(transformed_data)) == 1

# Null Check
def null_check(transformed_data):
    # Example: Checking for null values in transformed data
    return all(value is not None for value in transformed_data.values())

```

I have validated and preprocessed the given 5 data tables: with fixing file format, flattening and other text preprocessing and transformed the data into clean CSV format from raw formats.

Jagruti has also performed on 2 files

Jagruti - Tried to extract winning team name from result. Also, since venue name was available in match_title. I removed that and created 2 more columns for teams playing.

[Data validation on fixture by date and series files](#)

[Data Validation Code Link](#)

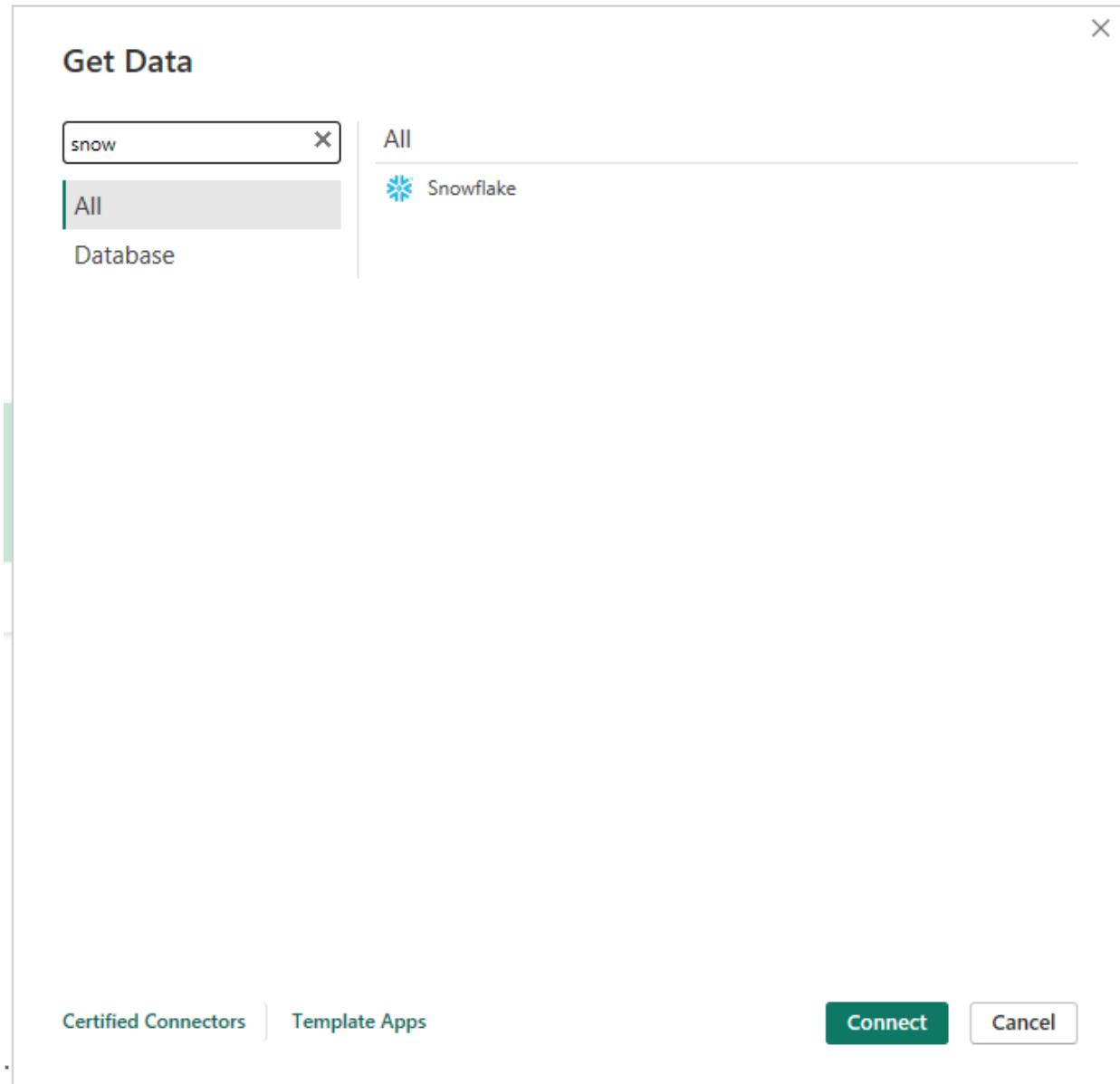
SnowFlake to PowerBI (Jigar Shekhat)

For analytic insight in powerBI, I am load data from DEV database of snowflake.

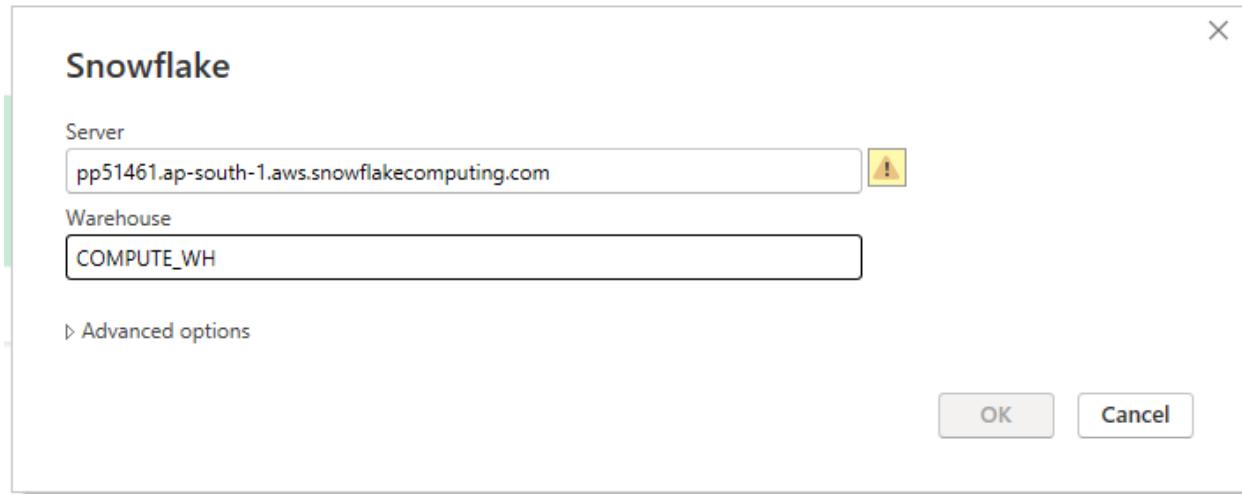
After loading Data, I prepared some sample report about that cricket data.

Connection with Snowflake :

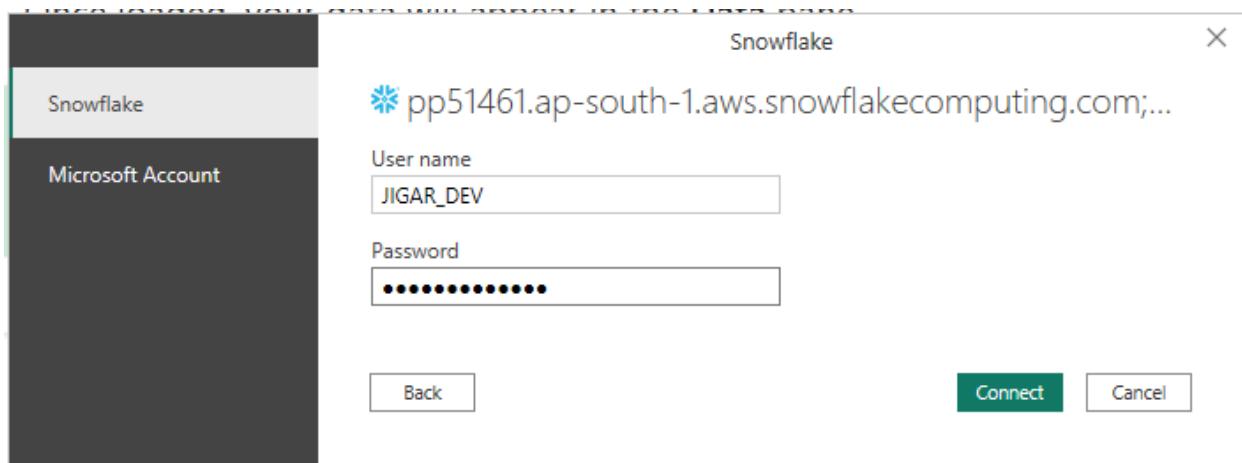
- 1. Select snowflake in get data tab.**



2. Enter Server and warehouse detail



3. After that, add your credentials of snowflake account



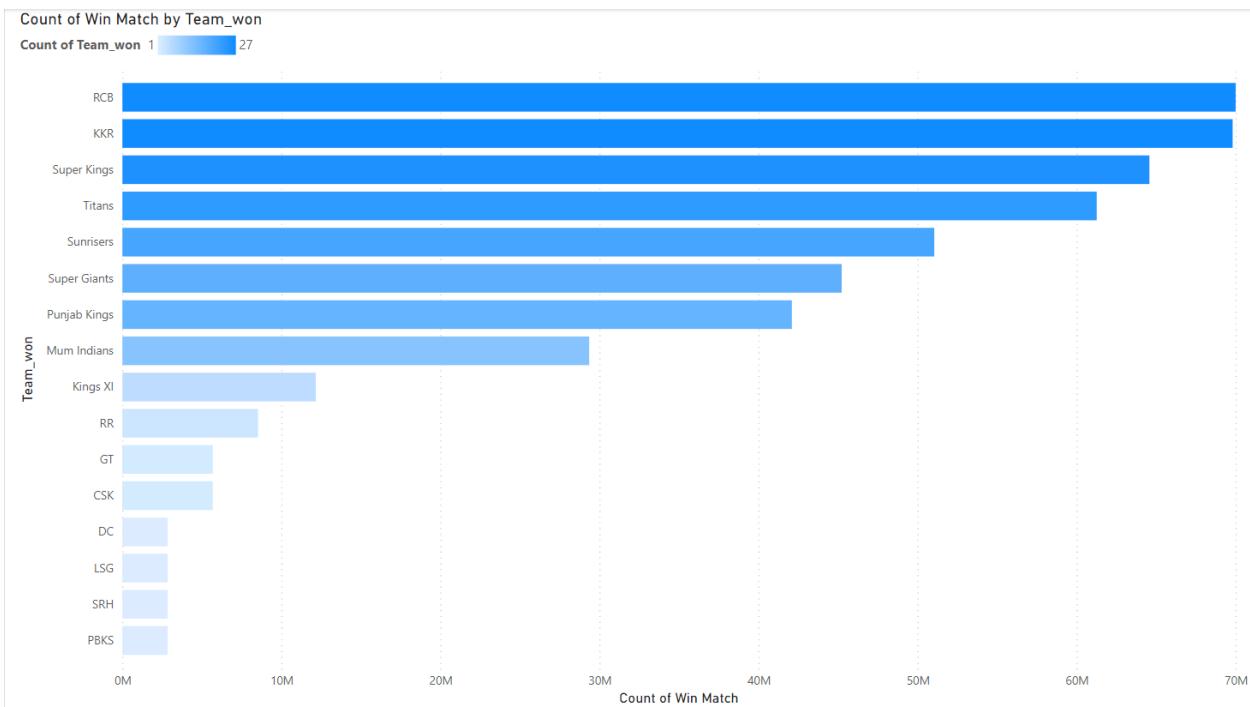
4. In this navigator, we can select tables which you have to load into our powerBI.

The screenshot shows the Power BI Navigator window. On the left, there's a tree view of database objects under 'pp51461.ap-south-1.aws.snowflakecomputing.com'. The 'SERIES' table is selected, indicated by a grey highlight. The main area displays the 'SERIES' table data with columns: SERIES_ID, SERIES_NAME, and STATUS. The table contains 34 rows of cricket series information. At the bottom, there are buttons for 'Select Related Tables', 'Load' (which is green), 'Transform Data', and 'Cancel'.

SERIES_ID	SERIES_NAME	STATUS
1202	India in Bangladesh Test Series	Result India 2-0 (2)
1239	Border-Gavaskar Trophy	Result India 2-1 (4)
1248	The Frank Worrell Trophy	Result Australia 2-0 (2)
1250	West Indies in Zimbabwe Test Series	Result West Indies 1-0 (2)
1406	Ireland in Bangladesh Test Match	Result Bangladesh 1-0 (1)
793	Pataudi Trophy	Result Drawn 2-2 (5)
865	New Zealand in England Test Series	Result England 3-0 (3)
866	Basil D'Oliveira Trophy	Result England 2-1 (3)
988	Warne-Muralitharan Trophy	Result Drawn 1-1 (2)
993	Sri Lanka in Bangladesh Test Series	Result Sri Lanka 1-0 (2)
1016	Bangladesh in West Indies Test Series	Result West Indies 2-0 (2)
1032	Pakistan in Sri Lanka Test Series	Result Drawn 1-1 (2)
873	New Zealand in India Test Series	Result India 1-0 (2)
874	Pakistan in Bangladesh Test Series	Result Pakistan 2-0 (2)
875	Afghanistan in Australia Test Match	Upcoming
876	The Ashes	Result Australia 4-0 (5)
877	Freedom trophy	Result South Africa 2-1 (3)
878	Sri Lanka in India Test Series	Result India 2-0 (2)
879	The Wisden Trophy	Result West Indies 1-0 (3)
880	Bangladesh in South Africa Test Series	Result South Africa 2-0 (2)
934	West Indies in Sri Lanka Test Series	Result Sri Lanka 2-0 (2)
935	Australia in Pakistan Test Series	Result Australia 1-0 (3)

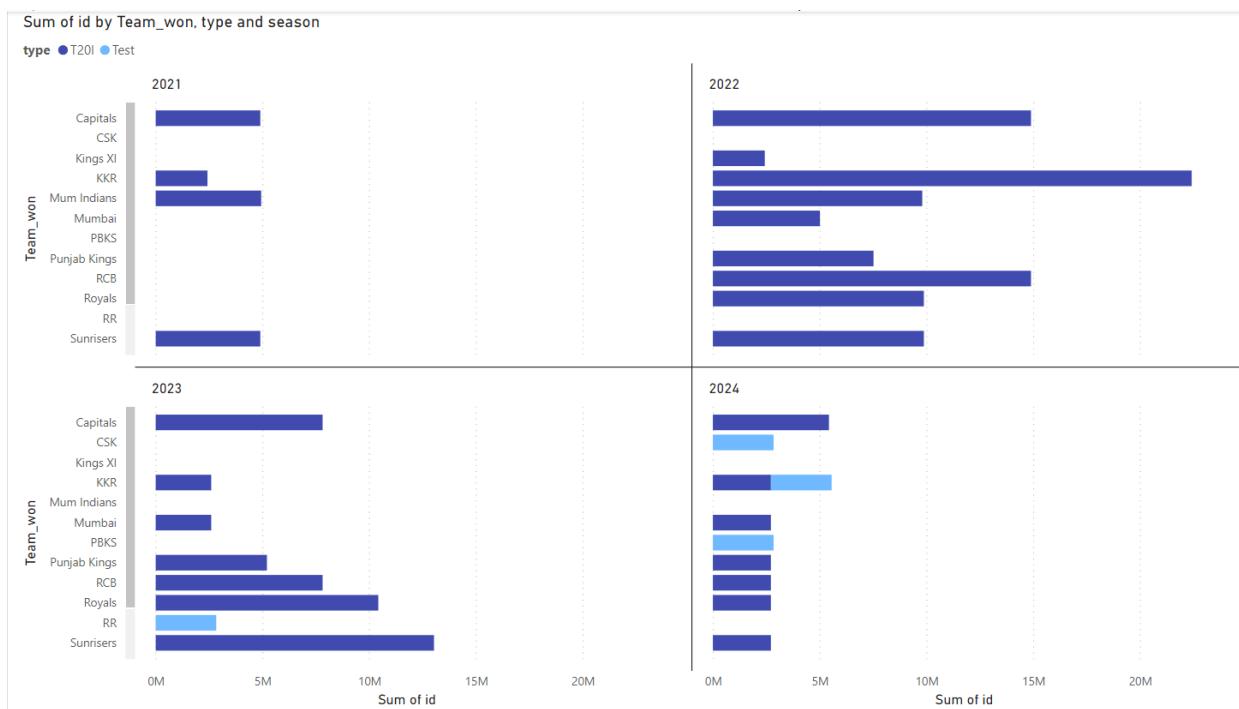
Report Screenshots :

1.



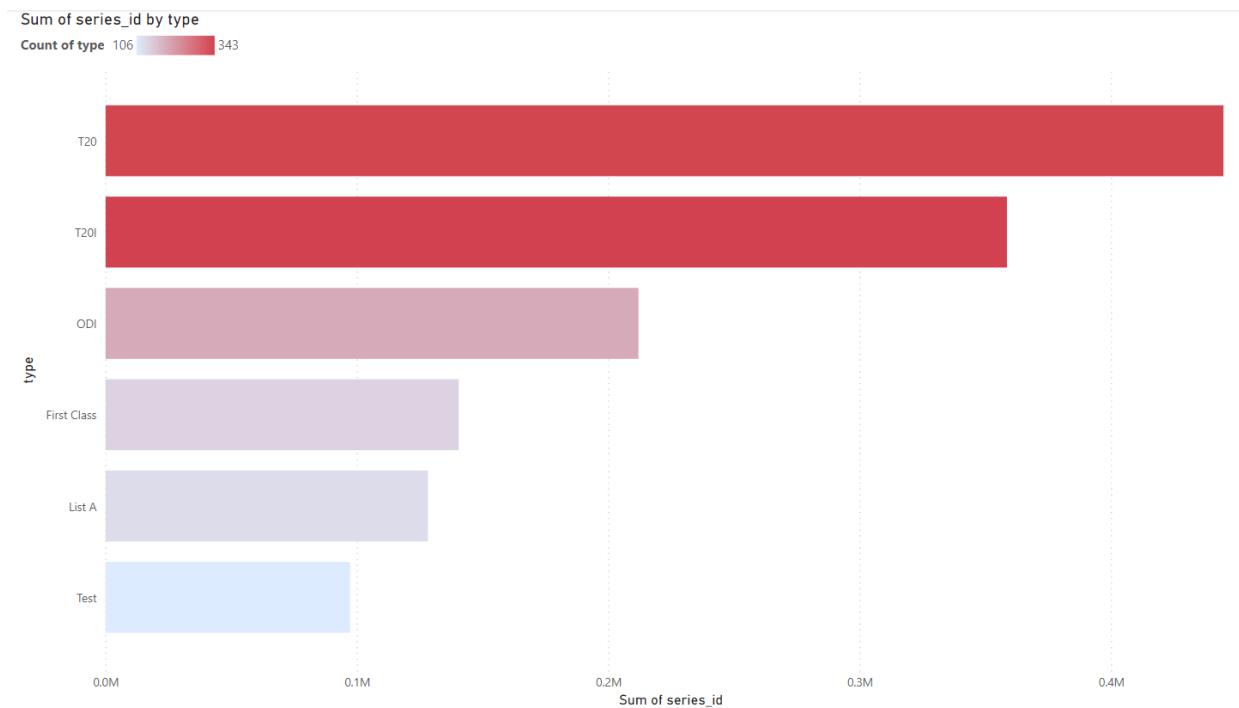
By this graph, we can represent which team has more wins than other teams.

2.



By this chart, we can show the winning team of each match by years (2021,2022,2023 and 2024).

3.



We can see that there is more data about T20, T20I and ODI series than other series in cricket api by this graph.

Conclusion:

This project emphasizes the significance of efficient data management and integration in today's competitive landscape. It entails developing a robust data model utilizing the Star/Snowflake Schema, acquiring Cricket data from third-party sources through RAPID APIs, storing it in Amazon S3 with Python, transferring it to Oracle, and migrating it to Snowflake using Talend. Additionally, the project involves implementing data validation and reporting in Power BI.

Efficient data management and integration are pivotal for organizational success. This project demonstrates the establishment of a strong data foundation and seamless processes, spanning acquisition, migration, and validation. By leveraging modern technologies such as Talend and Power BI, organizations can fully exploit their data assets, gaining valuable insights and maintaining a competitive edge in the data-driven landscape of today.