

Project Title

Data Engineering Capstone Project

Project Summary

-describe your project at a high-level-

The project follows the following steps:

- Step 1: Scope the Project and Gather Data
- Step 2: Explore and Assess the Data
- Step 3: Define the Data Model
- Step 4: Run ETL to Model the Data
- Step 5: Complete Project Write Up

```
In [26]: import sql
        %load_ext sql

The sql extension is already loaded. To reload it, use:
        %reload_ext sql

In [27]: # Do all imports and installs here
import pandas as pd
import psycopg2
from sqlalchemy import create_engine
from sqlalchemy.pooling import QueuePool
import os

In [28]: agmarknet_insert

Out[28]: 'INSERT INTO agmarknet (state,district ,market , commodity ,variety, arrival_date ,min_price ,max_price ,modal_price ) VALUES (%s, %s, %s, %s, %s, %s, %s, %s)'

Step 1: Scope the Project and Gather Data

Scope

This projects aims to collect and analyse the AGMARKNET Data. AGMARKNET is one Gov API that publish commodity price various mandi in INDIA. There is History data stored Manually hitting the API on python. Aim of the projects is to store all the history data into postgres database and update this data base on daily basis using scheduled AWS lambda on daily basis. Next connect the database to AWS quick sight and analyse the data base on daily basis
```

AGMARKNET DATA

Data come from hitting the GOVT API :

API URL: <https://api.data.gov.in/resource/c9ef84268-4588-465a-4308-a864a43d0070?api-key=579b464db66ec23bd6d98089d169143fc81ac74bc7ad727abc27056a8a&format=csv&limit=10000>

```
In [49]: #read in the data here
import pandas as pd
df=pd.read_csv('data.csv')

In [38]: pd.options.display.max_columns = None
df.head()

Out[38]:
   state district  market commodity variety arrival_date min_price max_price modal_price
0  Andhra Pradesh Chittoor Mulakalacheruvu Tomato Local 06/07/20 0.00 2000.0 3100.0 2100.0
1  Gujarat Amreli Dammagar Bhdnd(Ladies Finger) Bhdnd 06/07/20 0.00 900.0 1100.0 1000.0
2  Gujarat Amreli Dammagar Birinjai Other 06/07/20 0.00 900.0 1100.0 1000.0
3  Gujarat Amreli Dammagar Cabbage Cabbage 06/07/20 0.00 600.0 800.0 700.0
4  Gujarat Amreli Dammagar Coriander(Laaves) Coriander 06/07/20 0.00 4900.0 5100.0 5000.0

In [34]: df.info()

Out[51]:
<class 'pandas.core.frame.DataFrame'>
Int64Index: 903383 entries, 0 to 913928
Data columns (total 11 columns):
#   Column              Non-Null Count  Dtype
--  --
0   timestamp           903383 non-null  int64
1   time                903383 non-null  object
2   state               903383 non-null  object
3   district            903383 non-null  object
4   market              903383 non-null  object
5   commodity           903383 non-null  object
6   variety             903383 non-null  object
7   arrival_date        903383 non-null  object
8   min_x8028_Price     903383 non-null  float64
9   max_price           903383 non-null  float64
10  modal_price         903383 non-null  float64
dtypes: float64(3), int64(1), object(7)
memory usage: 82.7+ MB

Step 2: Explore and Assess the Data

Explore the Data

Identify data quality issues, like missing values, duplicate data, etc.
```

```
In [52]: df.describe()

Out[52]:
      timestamp      min_price      max_price      modal_price
count  9.023830e+05  903383.000000  903383.000000  903383.000000
mean    1.616027e+09    2948.979908    3448.637897    3220.634627
std     3.660346e+09    3899.397100    6216.771599    4981.900246
min     1.610201e+09     0.000000     0.000000     0.000000
25%    1.612979e+09    1000.000000    1250.000000    1175.000000
50%    1.615979e+09    1800.000000    2250.000000    2050.000000
75%    1.618562e+09    3600.000000    4325.000000    4000.000000
max     1.625075e+09   150000.000000  480000.000000  350000.000000
```

Daily scrit

```
In [34]: import requests
import json
import pandas as pd
import urllib3
from sqlalchemy import create_engine
from urllib3.poolmanager import PoolManager

url = "https://api.data.gov.in/resource/c9ef84268-4588-465a-4308-a864a43d0070?api-key=579b464db66ec23bd6d98089d169143fc81ac74bc7ad727abc27056a8a&format=csv&limit=10000"

payload={}
headers = {
    'accept': 'application/xml',
    'Cookie': 'BIGIPServerapi.data.gov.in=13+Hghu2GEPnrt6C6bshduWRH56yPaT0/So18x2yuzwQTb269vGBvG/gm86Fybv1nT2RE85VvyQ==; TS01a12685=9161d6cf399291b55df766b82117226b;'
}

response = requests.request("GET", url, headers=headers, data=payload)
x=response.text
listrx.split("\n")
columns=list(t[0].split(","))

State=[]
District=[]
Market=[]
Commodity=[]
Variety=[]
Arrival_Date=[]
Min_x8028_Price=[]
Max_x8028_Price=[]
Modal_x8028_Price=[]

for i in listrx[1:-1]:
    State.append(i.split(',')[0])
    District.append(i.split(',')[1])
    Market.append(i.split(',')[2])
    Commodity.append(i.split(',')[3])
    Variety.append(i.split(',')[4])
    Arrival_Date.append(i.split(',')[5])
    Min_x8028_Price.append(i.split(',')[8])
    Max_x8028_Price.append(i.split(',')[9])
    Modal_x8028_Price.append(i.split(',')[10])

df=pd.DataFrame({'State':State,'District':District,'Market':Market,'Commodity':Commodity,'Variety':Variety,'Arrival_Date':Arrival_Date,'Min_Price':Min_x8028_Price,'Max_x8028_Price':Max_x8028_Price,'Modal_Price':Modal_x8028_Price})
df.head(20)
```

	State	District	Market	Commodity	Variety	Arrival_Date	Min_Price	Max_Price	Modal_Price
0	Andhra Pradesh	Chittoor	Mulakalacheruvu	Tomato	Local	06/07/20 0.00	2000.0	3100.0	2100.0
1	Andhra Pradesh	Kurnool	Alilagadda	Jowar(Sorghum)	'Jowar' (White)	26/12/2021	4500	5200	5000
2	Andhra Pradesh	Kurnool	Alilagadda	Paddy(Dhan)(Common)	Sona	26/12/2021	1860	2180	1950
3	Bihar	"Easi Champaran"	Chakia	"Biraj(Pearl Millet)Cumbu"	Budl	26/12/2021	1100	1400	1200
4	Bihar	"Easi Champaran"	Chakia	Birinjai	"Arakshela Matiguln"	26/12/2021	1500	2000	1600
5	Bihar	"Easi Champaran"	Chakia	Cauliflower	"African Sarsori"	26/12/2021	2500	3000	2600
6	Bihar	"Easi Champaran"	Chakia	Onion	"Lit Sort"	26/12/2021	2500	3200	2600
7	Bihar	Gaya	Gaya	Cabbage	Cabbage	26/12/2021	1500	2500	2000
8	Bihar	Gaya	Gaya	Cauliflower	"African Sarsori"	26/12/2021	1500	2500	2000
9	Chattisgarh	Dantewada	Odum	Paddy(Dhan)(Common)	"Paddy Medium"	26/12/2021	1800	1900	1400
10	Chattisgarh	Kanker	Chesma	Paddy(Dhan)(Common)	Other	26/12/2021	1500	1500	1500
11	Chattisgarh	Kanker	Lakhanpuri	Paddy(Dhan)(Common)	Other	26/12/2021	1500	1500	1500
12	Chattisgarh	Kanker	Narhapur	Paddy(Dhan)(Common)	Other	26/12/2021	1500	1500	1500
13	Chattisgarh	Mahesamund	Bagahna	Paddy(Dhan)(Common)	MTU-1001	26/12/2021	1380	1380	1380
14	Chattisgarh	Mahesamund	Bagahna	Paddy(Dhan)(Common)	"Swarna Masuri (Newy"	26/12/2021	1400	1400	1400
15	Chattisgarh	Narayanpur	Narayanpur	Maze	Medium	26/12/2021	1300	1400	1350
16	Chattisgarh	Narayanpur	Narayanpur	Paddy(Dhan)(Common)	"Paddy Coarse"	26/12/2021	1200	1300	1250
17	Chattisgarh	Sukma	Dorampal	Paddy(Dhan)(Common)	"Paddy Coarse"	26/12/2021	1200	1250	1220
18	Chattisgarh	Sukma	Kuknar	Paddy(Dhan)(Common)	"Paddy Coarse"	26/12/2021	1200	1250	1220
19	Chattisgarh	Sukma	Kuknar	Paddy(Dhan)(Common)	"Paddy Coarse"	26/12/2021	1200	1250	1220

Cleaning Steps

Document steps necessary to clean the data

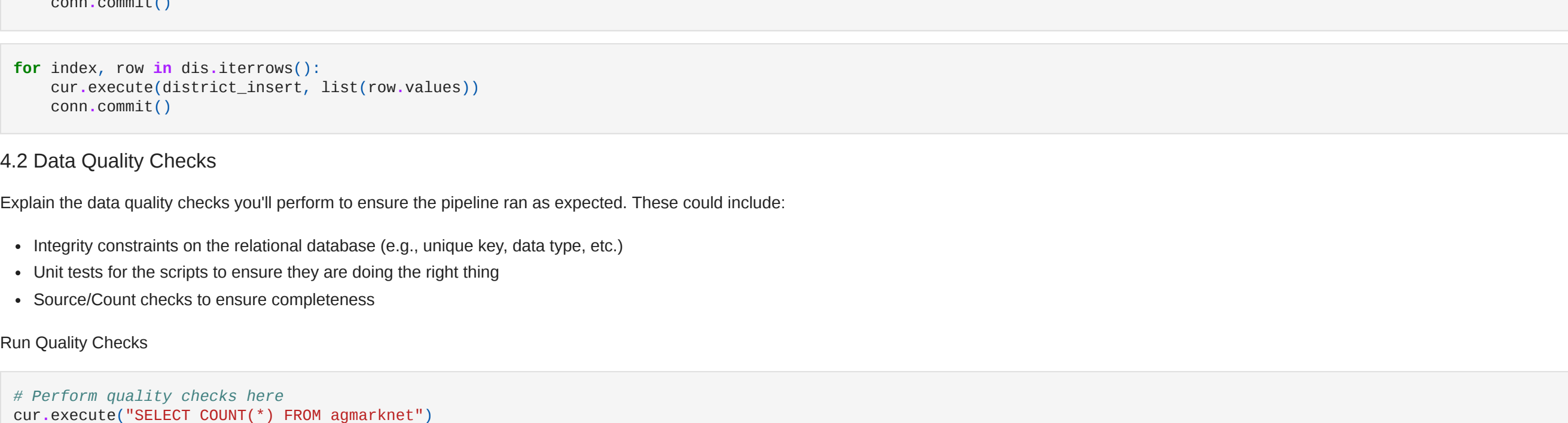
```
In [35]: df.head()

Out[35]:
   state district  market commodity variety arrival_date min_price max_price modal_price
0  Andhra Pradesh Chittoor Mulakalacheruvu Tomato Local 06/07/20 0.00 2000.0 3100.0 2100.0
1  Gujarat Amreli Dammagar Bhdnd(Ladies Finger) Bhdnd 06/07/20 0.00 900.0 1100.0 1000.0
2  Gujarat Amreli Dammagar Birinjai Other 06/07/20 0.00 900.0 1100.0 1000.0
3  Gujarat Amreli Dammagar Cabbage Cabbage 06/07/20 0.00 600.0 800.0 700.0
4  Gujarat Amreli Dammagar Coriander(Laaves) Coriander 06/07/20 0.00 4900.0 5100.0 5000.0
```

Step 3: Define the Data Model

3.1 Conceptual Data Model

Map out the conceptual data model and explain why you chose that model



	district	state
0	Agra	Uttar Pradesh
1	Ahmedabad	Gujarat
2	Almer	Rajasthan
3	Alappuzha	Kerala
4	Aligarh	Uttar Pradesh
...		
354	Wayanad	Kerala
355	West District	Tripura
356	West Garo Hills	Meghalaya
357	Yamuna Nagar	Haryana
358	Yapthula	Punjab
359 rows × 2 columns		

Step 4: Run Pipelines to Model the Data

4.1 Create the data model

Build the data pipelines to create the data model.

```
In [38]: # After running create tables.py, insert the data into the database
conn = psycopg2.connect("host=127.0.0.1 dbname=postgres user=postgres password=postgres")
cur = conn.cursor()

In [39]: for index, row in df.head(100).iterrows():
        cur.execute(agmarknet_insert, list(row.values))
        conn.commit()

In [40]: for i in df['market'].unique():
        cur.execute(markets_insert, [i])
        conn.commit()

In [41]: for i in df['state'].unique():
        cur.execute(state_insert, [i])
        conn.commit()

In [42]: for index, row in dis.iterrows():
        cur.execute(district_insert, list(row.values))
        conn.commit()
```

4.2 Data Quality Checks

Explain the data quality checks you'll perform to ensure the pipeline ran as expected. These could include:

- Integrity constraints on the relational database (e.g. unique key, data type, etc.)
 - Unit tests for the scripts to ensure they are doing the right thing
 - Source/Checks checks to ensure completeness
- Run Quality Checks

```
In [43]: # Perform quality checks here
cur.execute("SELECT COUNT(*) FROM agmarknet")
conn.commit()
if cur.rowcount < 1:
    print("No data found in table agmarknet")

cur.execute("SELECT COUNT(*) FROM markets")
conn.commit()
if cur.rowcount < 1:
    print("No data found in table markets")

cur.execute("SELECT COUNT(*) FROM state")
conn.commit()
if cur.rowcount < 1:
    print("No data found in table state")

cur.execute("SELECT COUNT(*) FROM district")
conn.commit()
if cur.rowcount < 1:
    print("No data found in table district")

In [42]: import sql
        %load_ext sql

DB_ENDPOINT = "127.0.0.1"
DB = 'postgres'
DB_USER = 'postgres'
DB_PASSWORD = 'postgres'
DB_PORT = '5432'

# postgres://username:password@host:port/database
conn_string = "postgres://{}:{}@{}:{}/{}?format=DB_USER,DB_PASSWORD,DB_ENDPOINT,DB_PORT,DB"

print(conn_string)
%sql <conn_string>

The sql extension is already loaded. To reload it, use:
        %reload_ext sql

postgres://postgres:postgres@127.0.0.1:5432/postgres

Data Check

In [45]: %time
        %sql
SELECT * FROM agmarknet limit 5

5 rows affected.
CPU times: user 2.77 ms, sys: 1.46 ms, total: 4.23 ms
Wall time: 3.3 ms

Out[45]:
   state district  market commodity variety arrival_date min_price max_price modal_price
0  Andhra Pradesh Chittoor Mulakalacheruvu Tomato Local 06/07/20 0.00 2000.0 3100.0 2100.0
1  Andhra Pradesh Chittoor Mulakalacheruvu Tomato Local 06/07/20 0.00 2000.0 3100.0 2100.0
2  Gujarat Amreli Dammagar Bhdnd(Ladies Finger) Bhdnd 06/07/20 0.00 900.0 1100.0 1000.0
3  Gujarat Amreli Dammagar Cabbage Cabbage 06/07/20 0.00 600.0 800.0 700.0

In [46]: %time
        %sql
SELECT * FROM markets limit 5

5 rows affected.
CPU times: user 3.77 ms, sys: 2.11 ms, total: 5.88 ms
Wall time: 4.49 ms

Out[46]:
market_id  market_name
0  Mulakalacheruvu
1  Dammagar
2  Rajpatti
3  Patnawa
4  Bhajpur

In [47]: %time
        %sql
SELECT * FROM district limit 5

5 rows affected.
CPU times: user 4.89 ms, sys: 2.49 ms, total: 7.38 ms
Wall time: 5.63 ms

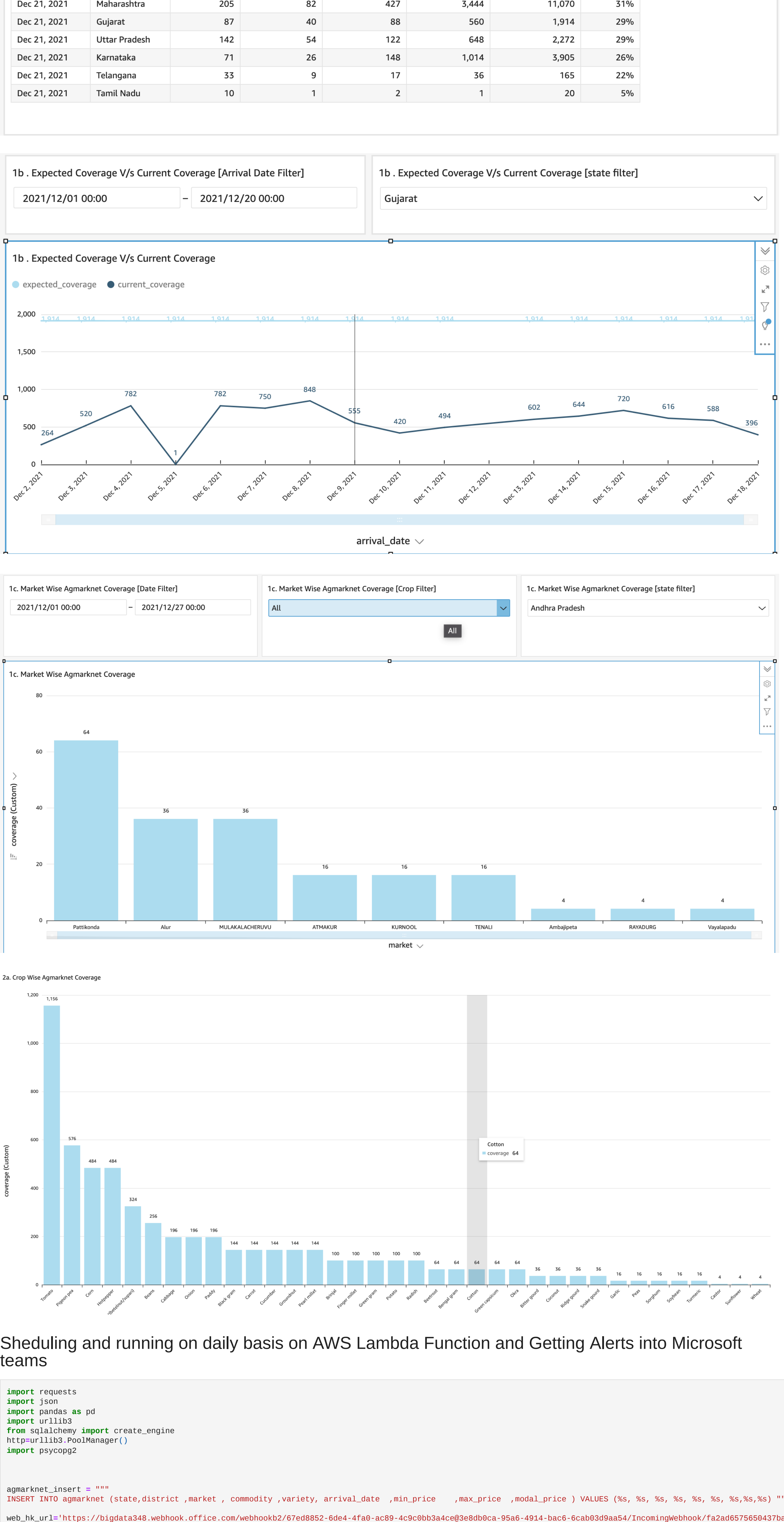
Out[47]:
district_id  district_name  state_name
0  1  Agra  Uttar Pradesh
1  2  Ahmedabad  Gujarat
2  3  Almer  Rajasthan
3  4  Alappuzha  Kerala
4  5  Aligarh  Uttar Pradesh

In [48]: %time
        %sql
SELECT * FROM state limit 5

5 rows affected.
CPU times: user 2.88 ms, sys: 1.5 ms, total: 4.38 ms
Wall time: 3.11 ms

Out[48]:
state_id  state_name
0  1  Andhra Pradesh
1  2  Gujarat
2  3  Haryana
3  4  Himachal Pradesh
4  5  Karnataka
```

Analysing on the AWS Quicksight



Scheduling and running on daily basis on AWS Lambda Function and Getting Alerts into Microsoft teams

```
In [64]: import requests
import json
import pandas as pd
import urllib3
from sqlalchemy import create_engine
from urllib3.poolmanager import PoolManager
import psycopg2

agmarknet_insert = """
INSERT INTO agmarknet (state,district ,market , commodity ,variety, arrival_date ,min_price ,max_price ,modal_price ) VALUES (%s, %s, %s, %s, %s, %s, %s, %s)
"""

web_hk_url="https://bigdata4b.webhook.office.com/webhookb2/67ed8852-6de4-4fa0-ac89-4c9cb0b3a4ce@9ebdb0ca-95a6-4914-bac6-6a0ab098a541/incomingwebhookb2/fad2a0575659437baf"

msg = {
    "context": "https://schema.org/extensions",
    "@type": "MessageCard",
    "themeColor": "#444477",
    "title": "Agmarknet Data Append",
    "text": "Data successfully updated ",
    "sections": [
        {
            "facts": [
                {
                    "name": "No of Records ",
                    "value": nor
                },
                {
                    "name": "Upload database",
                    "value": db
                },
                {
                    "name": "Upload table",
                    "value": table
                },
                {
                    "name": "Movement Type",
                    "value": movement_type
                }
            ]
        }
    ]
}

encoded_msg = json.dumps(msg, encode='utf-8')
resp = http.request("POST", web_hk_url, body=encoded_msg)

def lambda_handler(event, context):
    url = "https://api.data.gov.in/resource/c9ef84268-4588-465a-4308-a864a43d0070?api-key=579b464db66ec23bd6d98089d169143fc81ac74bc7ad727abc27056a8a&format=csv&limit=10000"

    payload={}
    headers = {
        'accept': 'application/xml',
        'Cookie': 'BIGIPServerapi.data.gov.in=13+Hghu2GEPnrt6C6bshduWRH56yPaT0/So18x2yuzwQTb269vGBvG/gm86Fybv1nT2RE85VvyQ==; TS01a12685=9161d6cf399291b55df766b82117226b;'
    }

    response = requests.request("GET", url, headers=headers, data=payload)
    x=response.text
    listrx.split("\n")
    columns=list(t[0].split(","))

    State=[]
    District=[]
    Market=[]
    Commodity=[]
    Variety=[]
    Arrival_Date=[]
    Min_x8028_Price=[]
    Max_x8028_Price=[]
    Modal_x8028_Price=[]

    for i in listrx[1:-1]:
        State.append(i.split(',')[0])
        District.append(i.split(',')[1])
        Market.append(i.split(',')[2])
        Commodity.append(i.split(',')[3])
        Variety.append(i.split(',')[4])
        Arrival_Date.append(i.split(',')[5])
        Min_x8028_Price.append(i.split(',')[8])
        Max_x8028_Price.append(i.split(',')[9])
        Modal_x8028_Price.append(i.split(',')[10])

    df=pd.DataFrame({'State':State,'District':District,'Market':Market,'Commodity':Commodity,'Variety':Variety,'Arrival_Date':Arrival_Date,'Min_Price':Min_x8028_Price,'Max_x8028_Price':Max_x8028_Price,'Modal_Price':Modal_x8028_Price})

    conn = psycopg2.connect("host=127.0.0.1 dbname=postgres user=postgres password=postgres")
    cur = conn.cursor()
    for index, row in df.iterrows():
        cur.execute(agmarknet_insert, list(row.values))
        conn.commit()

    nor=len(df)
    db="postgres"
    table="agmarknet"
    movement_type="append"
    msg[nor]=nor
    msg[db]=db
    msg[table]=table
    msg[movement_type]=movement_type

In [65]: lambda_handler('a','b')
```

Sample ALERT

 Agmarknet 1:18 am

Agmarknet Data Append

Data successfully updated

No of Records	1455
Upload database	postgres
Upload table	agmarknet
Movement Type	append

Reply

```
In [ ]:
```