



Data Visualization for Food Price in DKI Jakarta

by Mikael Dewabrata

Based on Info Pangan website

Problem Definition

DKI Jakarta has a information portal about food price in Jakarta covering several traditional market. While the data is there, it is not presented well. This project is to 'remake' the local government data presentation to be more clear and reachable. Using this project, I am trying to recreate the data pipeline and present it with a more proper data visualization.



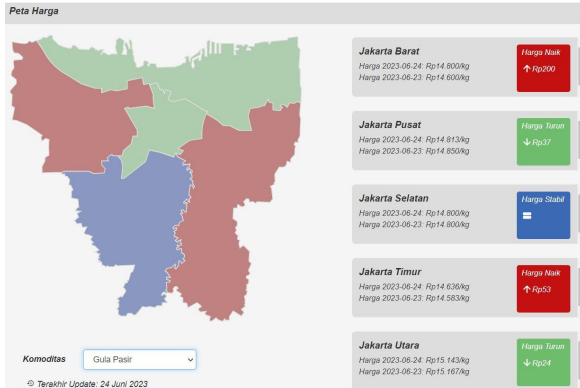
“Info Pangan DKI Jakarta is an excellent website, providing data for commodities in every traditional market in DKI Jakarta. However, the data is not well presented and quite slow.”

— Decision Making Questions



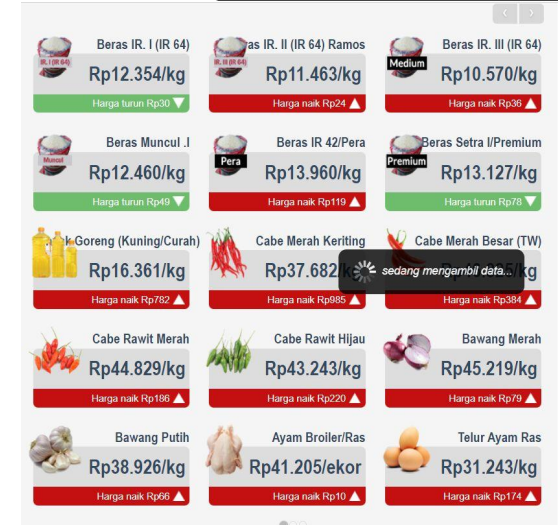
Data Selection

For this project, the data is taken from Informasi Pangan Jakarta (infopangan.jakarta.go.id). The information portal has complete pricing data from several traditional market in Jakarta.



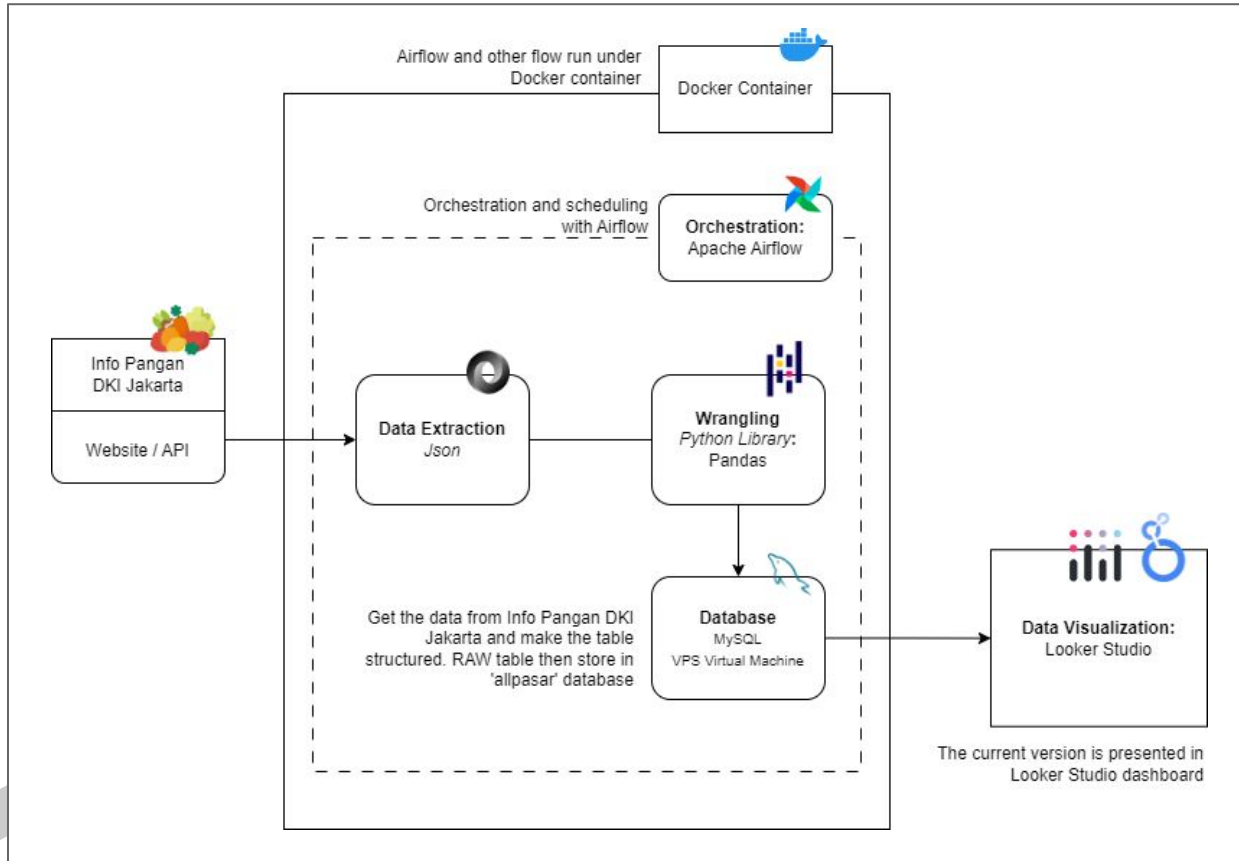
Data Problem

Though the data is updated regularly, the data is presented in very unstructured way. It's hard to read and not even using basic data visualization best practices.



* taken from infopangan.jakarta.go.id

Data Pipeline



This system is run on Virtual Machine in a cloud service VPS.

The documentation can be seen in GitHub, including files and coding example. Data Visualization itself is going to be presented in Looker Studio, shared publicly.

```
28 return transformed_date
29
30 Existing dataframe on list of IDs
31 df_list = pd.DataFrame({'idlist': [7,8,36,37,38,39,41,3,4,21,22,23,24,25,26,27,40,9,10,28,
32
33 _list = pd.DataFrame({'idlist': [7,8,36,37,38,39,41,3,4,21,22,23,24,25,26,27,40,9,10,28,29
34 ar = '2023'
35
36
37 f process_urls(df_list, year):
38     dfs = [] # List to store the separate dataframes
39     mnth = datetime.datetime.now().month
40
41     for index, row in tqdm(df_list.iterrows()):
42         id_param = row['idlist']
43         df_long_list = [] # List to store the individual market ID dataframes
44
45         url = f'https://infopangan.jakarta.go.id/api/price/series_by_location/public&types=
46
47         response = requests.get(url)
48         json_data = response.json()
49
50         data = json_data['data']
51
52         if len(data) > 0:
53             df_1 = pd.json_normalize(data)
54
55             columns_re = ['name', 'unit', 'id', 'low', 'high', 'average', 'location_name']
56             df_2 = df_1.drop(columns=columns_re)
57
58             columns_ren = {
59                 'commodity_id': 'CommodityID',
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100
```

Index	MarketID	mmodity	Date	Price
11	7	2	20230706	12500
12	7	20	20230706	5000
13	7	21	20230706	19000
14	7	22	20230706	15000
15	7	23	20230706	10000
16	7	24	20230706	10000
17	7	25	20230706	25000
18	7	26	20230706	150000
19	7	27	20230706	130000
20	7	28	20230706	50000
21	7	29	20230706	43000

Data Collection

To get the data, no need for scraping since the website already provide the endpoint. Easily this can be done by getting the data with JSON schema.

The more challenging process is to wrangle and clean the schema and make it a nicer table.

infopangan.jakarta.go.id/api/price/series_by_location?public=1&type=market&lid=41&m=6&y=2023

Index	Market ID	Market Name	Commodity	Date	Price
1	41	Pasar Jembatan Lima	Beras IR. II (IR 64) Ramos	20230607	13000
2	41	Pasar Jembatan Lima	Beras IR. II (IR 64) Ramos	20230621	35000

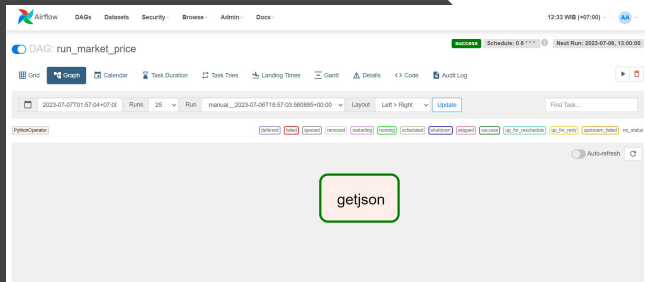
```
{"commodity_id":"2","name":"Beras IR. II (IR 64) Ramos","unit":"Rp","id":"2","low":"13000","high":"130000","average":"17720","location_name":"Pasar Jembatan Lima","series":{"1":"130000","2":"13000","3":"13000","4":"13000","5":"13000","6":"13000","7":"13000","8":"13000","9":"13000","10":"13000","11":"13000","12":"13000","13":"13000","16":"13000","17":"13000","18":"13000","19":"13000","20":"14000","21":"13000","22":"13000","24":"13000","25":"13000","26":"13000","27":"13000","28":"13000","29":"13000"}}
```

```
C:\Users\MIKAEL\Documents\Building DE Project Final\Others\dag_psear.py
yt_code_lay x get_api_DK.py x get_api_DK_All_giving.py x ariston_toped.py x get_API_All.py x get_API_latest.py x dag_psear.py x get_latest_API.py x
1 #!/usr/bin/env python3
2 # -*- coding: utf-8 -*-
3 """
4 Created on Wed Jul 5 16:35:46 2023
5
6 @author: mikael
7 """
8
9 from airflow import DAG
10 from airflow.operators.python_operator import PythonOperator
11 from datetime import datetime, timedelta
12 from pythonfile.run.get_latest_API import getjson
13
14 default_args = {
15     'owner': 'mikael',
16     'start_date': datetime(2023, 6, 26),
17     'retries': 1,
18     'retry_delay': timedelta(minutes=10),
19 }
20
21 dag = DAG(
22     'run_market_price',
23     default_args=default_args,
24     schedule_interval='0 6 * * *', # Run daily at 1 PM (13:00)
25 )
26
27 with dag:
28     get_json = PythonOperator(task_id='getjson',
29                             python_callable=getjson
30 )
```

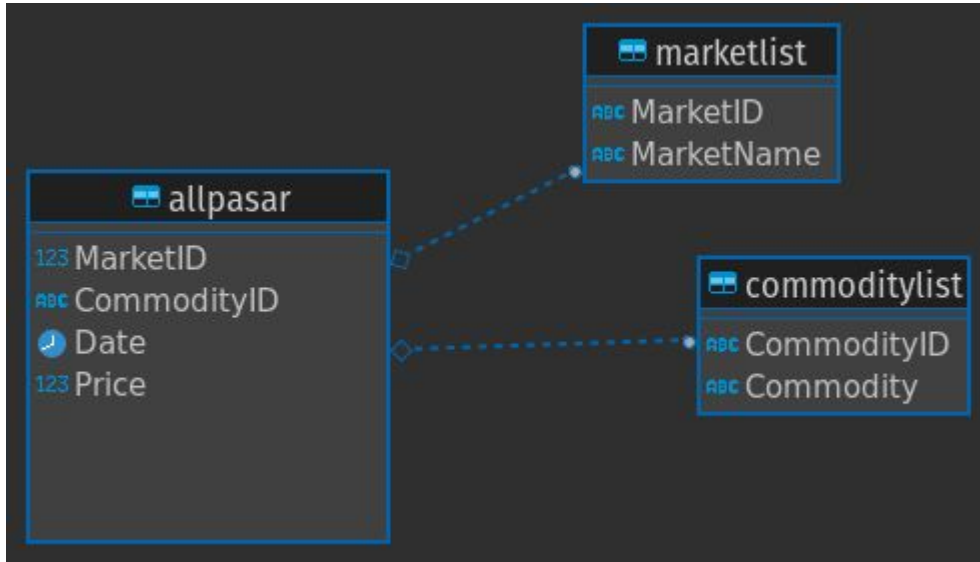
Data Orchestration

Using Airflow, the data then update daily so the dashboard can get the latest price for each commodity for each traditional market.

The Airflow is set under Docker container in VM environment using Linux Ubuntu.



ERD - Data Pangan DKI



The ERD is pretty simple, only consist of few columns. The main table has Market ID, Commodity ID, Date, and Price. To see further details on this information, we can expand to other table consist of Market Name and Commodity.

Properties Data ER Diagram

allpasar Enter a SQL expression to filter results (use Ctrl+Space)

	MarketID	CommodityID	Date	Price
1	20	14	2023-01-01	110,000
2	20	5	2023-01-01	10,000
3	20	15	2023-01-01	150,000
4	20	16	2023-01-01	150,000
5	20	17	2023-01-01	140,000
6	20	13	2023-01-01	60,000
7	20	18	2023-01-01	3,000
8	20	19	2023-01-01	21,000
9	20	6	2023-01-01	9,500
10	20	20	2023-01-01	14,000
11	20	21	2023-01-01	35,000
12	20	12	2023-01-01	50,000
13	20	22	2023-01-01	28,000
14	20	23	2023-01-01	38,000
15	20	7	2023-01-01	9,000
16	20	1	2023-01-01	38,000
17	20	24	2023-01-01	25,000
18	20	25	2023-01-01	9,500
19	20	11	2023-01-01	40,000
20	20	26	2023-01-01	15,000

Properties Data ER Diagram

marketlist Enter a SQL expression to filter results (use Ctrl+Space)

	MarketID	MarketName
1	7	Pasar Grogol
2	8	Pasar Glodok
3	36	Pasar Tomang Barat
4	37	Pasar Cengkareng
5	38	Pasar Kalideres
6	39	Pasar Pos Pengumben
7	41	Pasar Jembatan Lima
8	3	Pasar Senen Blok III - VI
9	4	Pasar Jembatan Merah
10	21	Pasar Tanah Abang Blok A-G

Properties Data ER Diagram

commoditylist Enter a SQL expression to filter results (use Ctrl+Space)

	CommodityID	Komoditas
1	1	Ayam Broiler/Ras
2	2	Bawang Merah
3	3	Bawang Putih
4	4	Beras IR 42/Pera
5	5	Beras IR. I (IR 64)
6	6	Beras IR. II (IR 64) Ramos
7	7	Beras IR. III (IR 64)
8	8	Beras Muncul .I
9	9	Beras Setra I/Premium
10	10	Cabe Merah Besar (TW)
11	11	Cabe Merah Keriting

Data Pipeline Summary Understanding

Once the data extraction already defined, cleaning and wrangling and storing to MySQL need to be done. Using data orchestration, we can make the dashboard presenting the current data.

Summary:

- Creating an extraction flow from infopangan.jakarta.go.id
- Store raw the HTML or table to MySQL in raw database
- Connect to MySQL , retrieve raw for past pricing data and clean using Pandas to store in new database
- Create a scheduled scraping script to take recent data daily and store to new database (clean) database
- Data orchestration for worker and scheduler will be dockerized and using Airflow
- The clean data will be taken as data source for Looker Studio

Few notes on the pipeline

Python:

- Creating script to get the data from Info Pangan DKI Jakarta and then store it to MySQL database
- Clean and wrangling the data for visualization to make a better data visualization
- Using Dash Plotly app for visualization (soon)
- Library: requests, sqlalchemy, pandas

Airflow:

- Setup and dockerize Airflow in Virtual Machine in VPS
- Prepare dependencies for MySQL connector
- Create DAG to run scheduling to extract the data daily every 1 PM WIB
- Create error report or duplication in DAG (soon)

Docker:

- Install Docker in Virtual Machine
- Prepare the YAML for Airflow and Image for required library
- Compose Airflow container, run in background. The scheduler and worker are both running daily



THANK YOU