

END-TO-END FRAUD DETECTION WITH MLOPS

BY: NAUFAL FAIZ N.

Tujuan Proyek

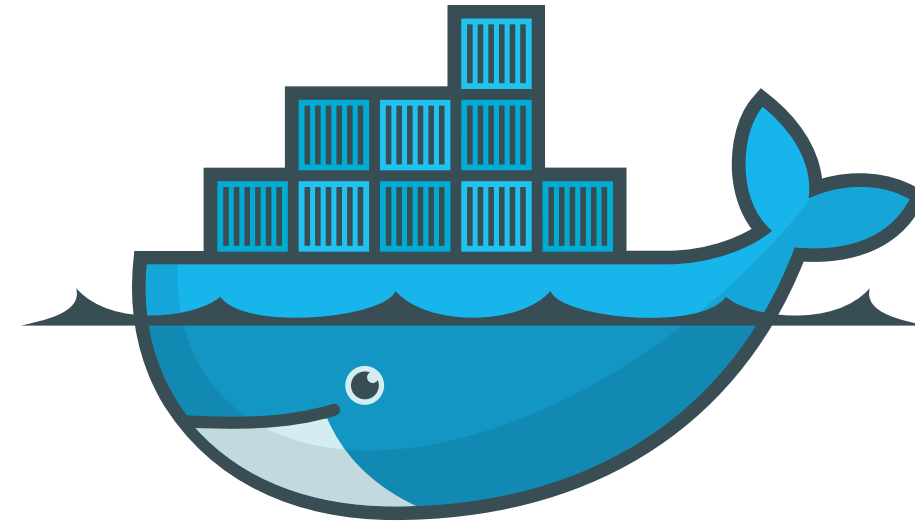
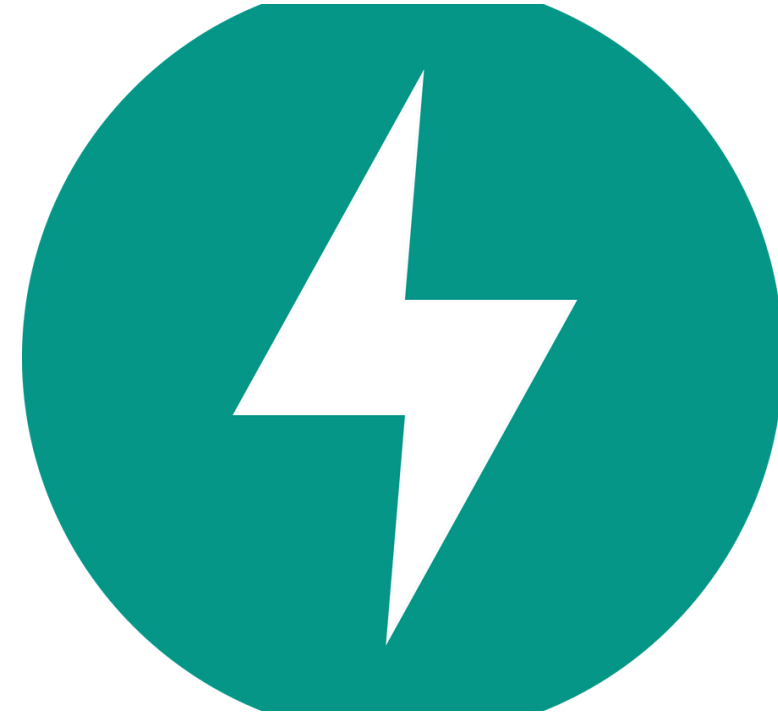
- Membangun sistem pendeteksi kecurangan (fraud detection) yang efisien dengan menggunakan teknik machine learning.
- Implementasi MLOps untuk pipeline yang mencakup pengolahan data, model training, dan deployment dengan integrasi Docker.

The Ideas

Proyek ini bertujuan untuk membangun sistem yang dapat mendeteksi kecurangan atau penipuan dengan menggunakan data teks dari dataset SMS Spam. Dataset yang digunakan dapat diakses di [id-nlp-resource](https://www.kaggle.com/datasets/alexisbcook/sms-spam-collection-dataset).



Liceria & Co.



**Tools yang
Digunakan**



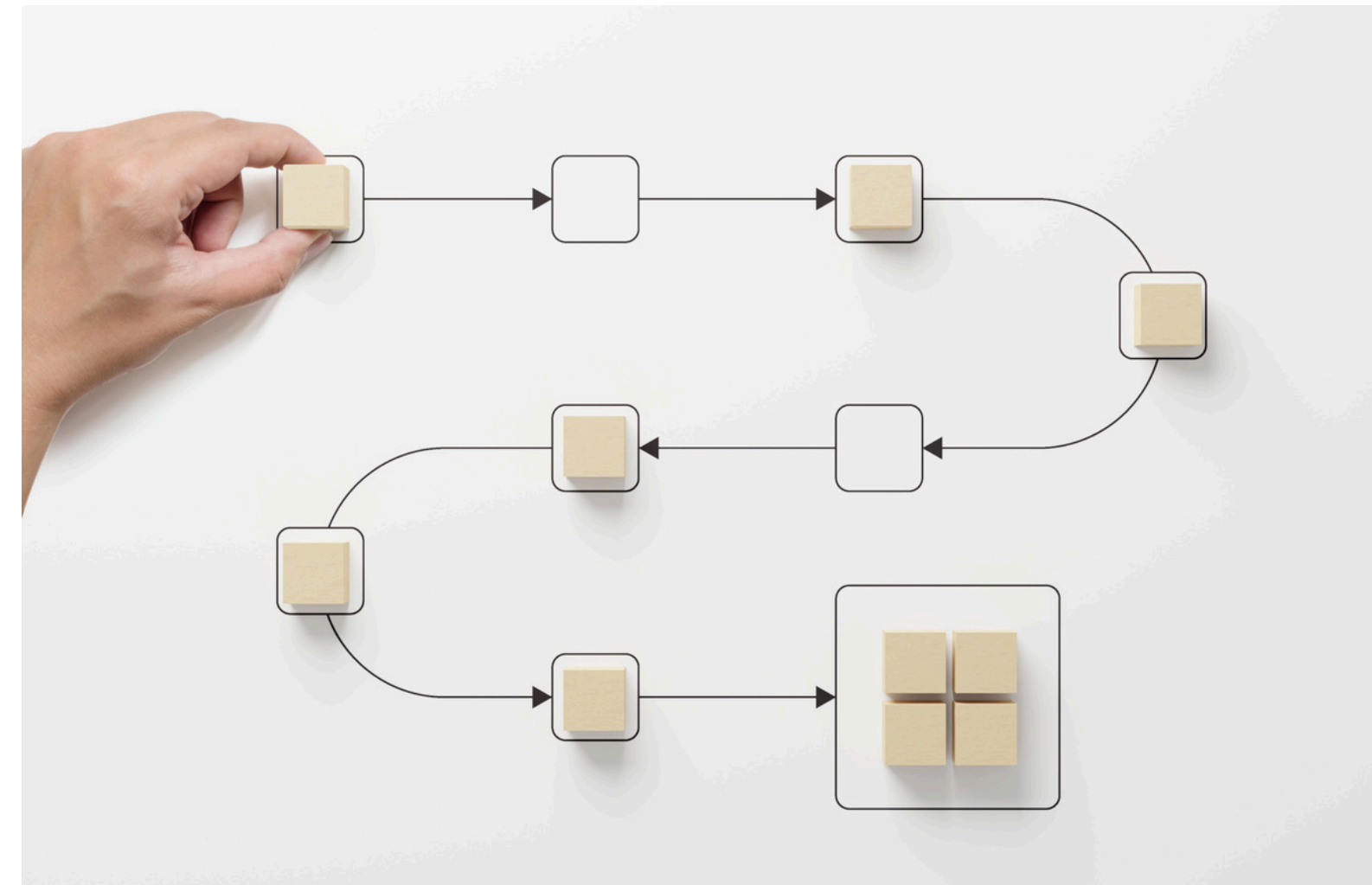
PostgreSQL

SQLAlchemy



Workflow Proyek

1. Data Collection (BigQuery): Data transaksi yang digunakan untuk deteksi kecurangan diambil dari BigQuery.
2. ETL Process (PostgreSQL): Data diproses dan dipindahkan ke PostgreSQL untuk disimpan dalam bentuk data mart.
3. Data Preprocessing: Melakukan pembersihan teks dan ekstraksi fitur menggunakan TF-IDF.
4. Model Training: Menggunakan Random Forest untuk membangun model deteksi kecurangan.
5. Model Deployment: Model dilatih dan dideploy menggunakan FastAPI di dalam Docker container.



🔍 query_transaction

▶ RUN

📄 SAVE QUERY (CLASSIC) ▼

👤 SHARE ▼

```
1 # show data
2 SELECT *
3 FROM `fraud-detection-project-449112.fraud_data.transactions`
4 LIMIT 10
5
6 # cek jumlah fraud dan Tidak fraud
7 SELECT label, COUNT(*) as count
8 FROM `fraud-detection-project-449112.fraud_data.transactions`
9 GROUP BY label
10
11 # cek nilai null
12 SELECT COUNT(*)
13 FROM `fraud-detection-project-449112.fraud_data.transactions`
14 WHERE Label IS NULL
15
16 # seleksi data duplikat
17 SELECT DISTINCT *
18 FROM `fraud-detection-project-449112.fraud_data.transactions`
19
20 # query untuk ekstraksi data
21 SELECT
22 | Teks,
23 | label
24 FROM `fraud-detection-project-449112.fraud_data.transactions`
```



```
1 from google.cloud import bigquery
2 from sqlalchemy import create_engine
3 import pandas as pd
4
5 # set up BigQuery client
6 client = bigquery.Client.from_service_account_json('C:\\Users\\NAUFAL FAIZ\\Documents\\Fraud ETL\\data\\fraud-detection-project-449112-d58ffa4a2213.json')
7
8 # query untuk menarik data
9 query = """
10 SELECT
11     Teks,
12     label
13 FROM `fraud-detection-project-449112.fraud_data.transactions`;
14 """
15 df = client.query(query).to_dataframe()
16
17 # menampilkan beberapa baris pertama
18 print(df.head())
19
20 # set up kredensial PostgreSQL Anda
21 db_url = 'postgresql://postgres:postgres123@localhost:5432/fraud_warehouse'
22
23 # buat koneksi ke PostgreSQL
24 engine = create_engine(db_url)
25
26 # save dataframe ke PostgreSQL
27 df.to_sql('transactions_mart', engine, if_exists='replace', index=False)
28
29 print("Data berhasil disimpan ke PostgreSQL")
```

Query

Query History

1 ▼ **SELECT** *

2 **FROM** transactions_mart

Data Output

Messages

Notifications

≡+

📄

▼

📋

▼

🗑️

🗄️

⬇️

📈

SQL

	Teks text	label bigint
1	Jika anda bermasalah dgn CC/KT@, stres dgn bunga, pelunasan disc s/d 75%...	1
2	Lelah byr min payment? Kami Solusinya, bantu secara LEGAL penutupan CC/...	1
3	Bisa Dgn BPKB rate 0.99% 3 hari CAIR . hub AYU 081584650877 (WA).mhn m...	1
4	"ROXI CELL" Hanya dengan Rp.100rb Anda bisa jadi agen pulsa elektrik ke se...	1
5	3 RAMADHAN Selamat Anda Pemenang Rp.100jt. PIN CODE 7Y7R8K9Z Info: ...	1
6	Anda brminat cash&kredit mtor scond brg istmwa tipe&merk apa sj.dsini mny...	1
7	ANDA MAU MENANG TOGEL 100% Tembus pasang shio 7 tunggal angka 07....	1
8	Anda mempunyai 3 pesan suara dari 083139195872; untuk mendengarkan pe...	1
9	Anda pmain togel sring kalah.mau bagi hasil. MBAH JOKO bisah mbntu anda ...	1
10	ANEKA SHOP Berbagi Promo Type Blackberry DAKOTA 2,4 jt, TOURCH 1,7 jt, ...	1
11	Anyonghaseyo! Ngaku KPOPers Sejati?? Ayo lengkapi koleksi Video Kpopmu ...	1
12	Artha Cell diskon Brg elektronik BB Dakota 3 juta BB Onyx3 2,8 jt Camera DLL ...	1
13	Ass, mengenai rumah dan tanah yg kemarin, sy sdah survei, sy mrasa cocok, ...	1
14	Ass, Sy Randy. Mengenai mobil yg sdh sy liat kondisi'y , kebetulan kami bermi...	1
15	Ass,maaf sy WAHYU.mengenai Mobil bpk yg mau di jual sy sudah lihat dan co...	1
16	Ass. Maaf sy by hj.suri yang kemarin survey rumah anda. Mengenai harga bis...	1
17	Ass.sv SUGIYONO va tempo hari mngeai mobil TOYOTANYA.sv brminat.coco...	1

Total rows: 904

Query complete 00:00:00.128

```
1 import pandas as pd
2 from sqlalchemy import create_engine
3 import string
4 import re
5 from Sastrawi.StopWordRemover.StopWordRemoverFactory import StopWordRemoverFactory
6 from Sastrawi.Stemmer.StemmerFactory import StemmerFactory
7 import nltk
8 nltk.download('punkt')
9 nltk.download('punkt_tab')
10 from nltk.tokenize import word_tokenize
11 from sklearn.feature_extraction.text import TfidfVectorizer
12
13 # buat koneksi ke PostgreSQL
14 db_url = 'postgresql://postgres:postgres123@localhost:5432/fraud_warehouse'
15 engine = create_engine(db_url)
16
17 # ambil data dari PostgreSQL
18 query = "SELECT * FROM transactions_mart"
19 df = pd.read_sql(query, engine)
20
21 print(df.head())
22
23 # fungsi cleaning
24 def cleaning(text):
25     # stopwords
26     stop_factory = StopWordRemoverFactory().get_stop_words()
27     # stemmer
28     stem_factory = StemmerFactory()
29     stemmer = stem_factory.create_stemmer()
30
31     text = text.lower() # ubah teks menjadi lower case
32     text = text.strip(' ') # menghapus spasi di awal dan di akhir
33     text = re.sub(r'\d+', '', text) # menghapus angka
34     text = text.translate(str.maketrans('', '', string.punctuation)) # menghapus punctuation
35     text = re.sub(r'\b[a-zA-Z]\b', '', text) # menghapus kata yang hanya terdiri dari satu huruf
36     text = re.sub(r'\s+', ' ', text) # menghapus spasi berlebih
37     text = word_tokenize(text) # tokenisasi
38     text = [word for word in text if word not in stop_factory] # menghapus stopwords
39     text = [stemmer.stem(word) for word in text] # stemming
40     text = ' '.join(text)
41     return text
42
43 # fungsi TF-IDF
44 def tfidf(df):
45     tfidf = TfidfVectorizer(max_features=1000)
46     X = tfidf.fit_transform(df['Teks_bersih']).toarray()
47     y = df['label']
48
49     return X, y, tfidf
```



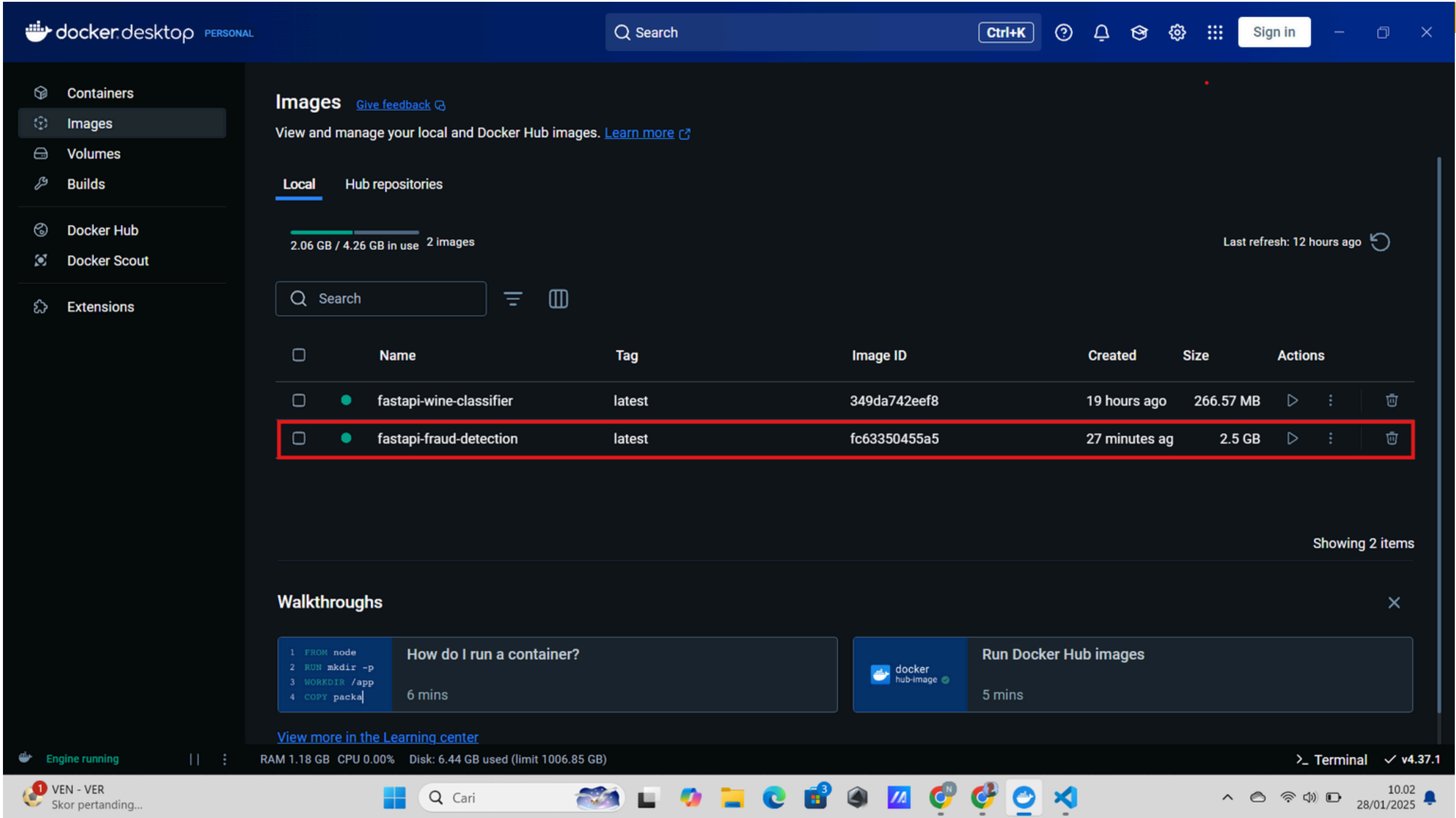
```
1 from preprocessing.preprocessing import *
2 from sklearn.model_selection import train_test_split
3 from sklearn.ensemble import RandomForestClassifier
4 from sklearn.metrics import classification_report
5 from sqlalchemy import create_engine
6 import pandas as pd
7 import pickle
8
9 # buat koneksi ke PostgreSQL
10 db_url = 'postgresql://postgres:postgres123@localhost:5432/fraud_warehouse'
11 engine = create_engine(db_url)
12
13 # ambil data dari PostgreSQL
14 query = "SELECT * FROM transactions_mart"
15 df = pd.read_sql(query, engine)
16
17 # cleaning dataframe
18 df['Teks_bersih'] = df['Teks'].apply(cleaning)
19 # ekstraksi fitur
20 X, y, tfidf = tfidf(df)
21
22 # save dataset hasil cleaning ke PostgreSQL dan csv
23 df.to_sql('transactions_mart_clean', engine, if_exists='replace', index=False)
24 print("Data berhasil disimpan ke PostgreSQL")
25
26 df.to_csv('transactions_mart_clean.csv', index=False)
27 print("Data berhasil disimpan ke CSV")
28
29 # train test split
30 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
31
32 # inisialisasi model
33 model = RandomForestClassifier()
34 # training model
35 model.fit(X_train, y_train)
36
37 # evaluasi model
38 y_pred = model.predict(X_test)
39
40 print(classification_report(y_test, y_pred))
41
42 # save model
43 with open('model.pkl', 'wb') as f:
44     pickle.dump(model, f)
45 print("Model berhasil disimpan")
46
47 # save tfidf vectorizer
48 with open('tfidf_vectorizer.pkl', 'wb') as f:
49     pickle.dump(tfidf, f)
50 print("TF-IDF Vectorizer berhasil disimpan")
```

```
1  from fastapi import FastAPI
2  from pydantic import BaseModel
3  import pickle
4  import pandas as pd
5  from sklearn.feature_extraction.text import TfidfVectorizer
6
7  # load model
8  with open('C:\\Users\\NAUFAL FAIZ\\Documents\\Fraud ETL\\model\\model.pkl', 'rb') as f:
9      model = pickle.load(f)
10
11 # load tfidf vectorizer
12 with open('C:\\Users\\NAUFAL FAIZ\\Documents\\Fraud ETL\\model\\tfidf_vectorizer.pkl', 'rb') as f:
13     tfidf = pickle.load(f)
14
15 # inisialisasi fastapi
16 app = FastAPI()
17
18 # pydantic untuk input data
19 class TransactionRequest(BaseModel):
20     text: str
21
22 # endpoint untuk prediksi fraud
23 @app.post('/predict')
24 def predict(transaction: TransactionRequest):
25     X_new = tfidf.transform([transaction.text])
26
27     prediction = model.predict(X_new)
28
29     result = "Fraud" if prediction[0] == 1 else "Not Fraud"
30     return {"Text": transaction.text, "Prediction": result}
31
32 # menjalankan aplikasi fastapi
33 if __name__ == '__main__':
34     import uvicorn
35     uvicorn.run(app, host="0.0.0.0", port=8000)
```

```
1  # menggunakan image Python yang ringan sebagai base image
2  FROM python:3.8-slim
3
4  # install dependencies yang diperlukan
5  RUN apt-get update && apt-get install -y \
6      build-essential \
7      libpq-dev \
8      postgresql-server-dev-all
9
10 # set working directory di dalam container
11 WORKDIR /app
12
13 # salin requirements.txt ke dalam container
14 COPY requirements.txt .
15
16 # salin model ke dalam container
17 COPY ./model /app/model
18
19 # install dependensi yang ada di requirements.txt
20 RUN pip install -r requirements.txt
21
22 # salin seluruh kode aplikasi ke dalam container
23 COPY . .
24
25 # Tentukan perintah untuk menjalankan aplikasi FastAPI dengan Uvicorn
26 CMD ["uvicorn", "main:app", "--host", "0.0.0.0", "--port", "8000"]
```



```
1  pandas
2  scikit-learn
3  nltk
4  sqlalchemy
5  psycopg2
6  psycopg2-binary
7  fastapi
8  uvicorn
```



THANK YOU