

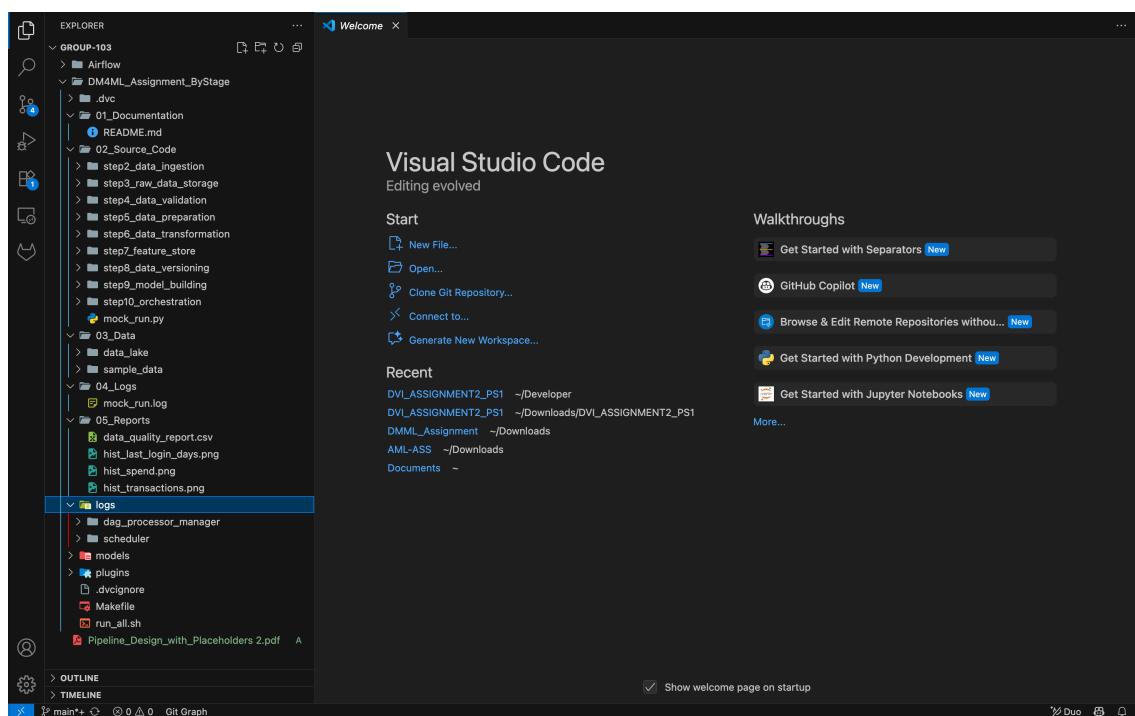
DMML Video Walkthrough

Hello Sir,

We were unable to upload the video along with the assignment as the portal has a file size limit, and the video exceeds that limit. Therefore, I am sharing the video link here for your reference:

<https://youtu.be/fW-mXFeO420>

Additionally, I have attached the screenshot of the code and pipeline.



data_ingestion.py

```

1 import os
2 import pandas as pd
3 import requests
4 import logging
5 from datetime import datetime
6 from pathlib import Path
7
8 # ----- CONFIG -----
9 RAW_DATA_DIR = Path("data/raw")
10 LOG_FILE = "logs/ingestion.log"
11 TRANSACTIONS_CSV = "sample_data/transactions.csv"
12 WEB_API_URL = "https://jsonplaceholder.typicode.com/posts" # Mock API
13
14 # Create necessary directories
15 RAW_DATA_DIR.mkdir(parents=True, exist_ok=True)
16 Path("logs").mkdir(parents=True, exist_ok=True)
17
18 # ----- LOGGING -----
19 logging.basicConfig(
20     filename=LOG_FILE,
21     level=logging.INFO,
22     format="%asctime)s (%(levelname)s %(message)s"
23 )
24
25 def log_message(message):
26     print(message)
27     logging.info(message)
28
29 # ----- INGESTION FUNCTIONS -----
30 def ingest_transactions():
31     """Fetch customer transactions from CSV."""
32     try:
33         df = pd.read_csv(TRANSACTIONS_CSV)
34         timestamp = datetime.now().strftime("%Y%m%d_%H%M%S")
35         output_path = RAW_DATA_DIR / f"transactions_{timestamp}.csv"
36         df.to_csv(output_path, index=False)
37         log_message(f"Transactions data ingested: {output_path}")
38     except Exception as e:
39         logging.error(f"Failed to ingest transactions: {e}")
40
41 def ingest_web_logs():
42     """Fetch web logs from mock API."""
43     try:
44         response = requests.get(WEB_API_URL, timeout=10)
45         response.raise_for_status()
46         data = response.json()
47         df = pd.DataFrame(data)
    
```

storage_raw_upload.py

```

1 import re
2 import shutil
3 import logging
4 from pathlib import Path
5 from datetime import datetime
6
7 # ----- CONFIG -----
8 RAW_DROP_DIR = Path("data/raw") # where Step 2 wrote raw files
9 DATA_LAKE_RAW = Path("data_lake/raw") # lake root (local)
10 LOG_FILE = "logs/storage_upload.log"
11
12 # Create dirs
13 DATA_LAKE_RAW.mkdir(parents=True, exist_ok=True)
14 Path("logs").mkdir(parents=True, exist_ok=True)
15
16 # ----- LOGGING -----
17 logging.basicConfig(
18     filename=LOG_FILE,
19     level=logging.INFO,
20     format="%asctime)s (%(levelname)s %(message)s"
21 )
22
23 def log(msg):
24     print(msg); logging.info(msg)
25
26 # ----- HELPERS -----
27 FILENAME_PATTERNS = {
28     "transactions": re.compile(r"^\w{8}_\d{6}\.csv$"),
29     "web_logs": re.compile(r"\w{8}_\d{6}\.csv$")
30 }
31
32 def detect_source_and_runid(file_name: str):
33     for source, pattern in FILENAME_PATTERNS.items():
34         m = pattern.match(file_name)
35         if m:
36             run_id = m.group(1) # YYYYMMDD_HHMMSS from Step 2
37             return source, run_id
38     return None, None
39
40 def run():
41     files = [p for p in RAW_DROP_DIR.glob("*.csv")]
42     if not files:
43         log("No raw files found to upload.")
44         return
45
46     for fpath in files:
47         source, run_id = detect_source_and_runid(fpath.name)
        if not source:
    
```

```
import pandas as pd
import logging
from pathlib import Path
from datetime import datetime

# ----- CONFIG -----
DATA LAKE_RAW = Path("data_lake/raw")
LOG_FILE = "logs/data_validation.log"
REPORT_FILE = Path("reports/data_quality_report.csv")

# Expected schema definitions for validation
EXPECTED_SCHEMAS = {
    "transactions": {
        "customer_id": "int64",
        "transaction_id": "object",
        "transaction_date": "object",
        "amount": "float64",
        "payment_method": "object",
        "transaction_status": "object"
    },
    "web_logs": {
        "user_id": "int64",
        "id": "int64",
        "title": "object",
        "body": "object"
    }
}

# Create necessary directories
Path("logs").mkdir(parents=True, exist_ok=True)
REPORT_FILE.parent.mkdir(parents=True, exist_ok=True)

# Logging setup
logging.basicConfig(
    filename=LOG_FILE,
    level=logging.INFO,
    format="%(asctime)s %(levelname)s %(message)s"
)

def log(msg):
    print(msg)
    logging.info(msg)

def validate_file(file_path: Path, source: str):
    df = pd.read_csv(file_path)
    issues = []

    # ... (rest of the validate_file function)
```

Pylance has detected type annotations in your code and recommends enabling type checking. Would you like to change this setting?

Source: Pylance

Yes No

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
from pathlib import Path
import logging
import os
import pandas as pd

# ----- CONFIG -----
DATA LAKE_RAW = Path("data_lake/raw/source=transactions") # We'll use transactions as example
CLEAN_DATA_DIR = Path("data/clean")
LOG_FILE = "logs/data_preparation.log"

# Create dirs
CLEAN_DATA_DIR.mkdir(parents=True, exist_ok=True)
Path("logs").mkdir(parents=True, exist_ok=True)
Path("reports").mkdir(parents=True, exist_ok=True)

# Logging
logging.basicConfig(
    filename=LOG_FILE,
    level=logging.INFO,
    format="%(asctime)s %(levelname)s %(message)s"
)

def log(msg):
    print(msg)
    logging.info(msg)

def load_latest_transactions():
    data_dir = "/opt/airflow/03_Data/data_sample_data"
    files = [f for f in os.listdir(data_dir) if f.endswith(".csv")]
    if not files:
        raise FileNotFoundError("No transaction files found.")

    latest_file = max(
        [os.path.join(data_dir, f) for f in files],
        key=os.path.getctime
    )
    print(f"Loading file: {latest_file}")
    return pd.read_csv(latest_file)

def clean_data(df):
    # Handle missing values
    for col in df.columns:
        if df[col].dtype == "object":
            df[col] = df[col].fillna("Unknown")
        else:
```

Pylance has detected type annotations in your code and recommends enabling type checking. Would you like to change this setting?

Source: Pylance

Yes No

data_transformation.py

```

1 import pandas as pd
2 import sqlite3
3 from pathlib import Path
4 import logging
5 from datetime import datetime
6
7 # ----- CONFIG -----
8 CLEAN_DATA_FILE = Path("data/clean/transactions_clean.csv")
9 DB_FILE = Path("data/feature_store.db")
10 LOG_FILE = "logs/data_transformation.log"
11
12 # Create dirs
13 Path("logs").mkdir(parents=True, exist_ok=True)
14 Path("data").mkdir(parents=True, exist_ok=True)
15
16 # Logging
17 logging.basicConfig(
18     filename=LOG_FILE,
19     level=logging.INFO,
20     format="%(asctime)s %(levelname)s %(message)s"
21 )
22 def log(msg):
23     print(msg)
24     logging.info(msg)
25
26 def load_clean_data():
27     if not CLEAN_DATA_FILE.exists():
28         raise FileNotFoundError("Clean data file not found. Run Step 5 first.")
29     return pd.read_csv(CLEAN_DATA_FILE)
30
31 def engineer_features(df):
32     df["transaction_date"] = pd.to_datetime(df["transaction_date"], errors="coerce")
33
34     agg_df = df.groupby("customer_id").agg(
35         total_spend=("amount", "sum"),
36         transaction_count=("transaction_id", "count"),
37         avg_transaction_value=("amount", "mean"),
38         first_transaction=("transaction_date", "min"),
39         last_transaction=("transaction_date", "max")
40     ).reset_index()
41
42     now = datetime.now()
43     agg_df["tenure_days"] = (now - agg_df["first_transaction"]).dt.days
44     agg_df["recency_days"] = (now - agg_df["last_transaction"]).dt.days
45
46     agg_df.drop(columns=["first_transaction", "last_transaction"], inplace=True)
47     return agg_df

```

feature_store.py

```

1 import sqlite3
2 from datetime import datetime
3 import pandas as pd
4 from pathlib import Path
5 import logging
6
7 # ----- CONFIG -----
8 DB_FILE = Path("data/feature_store.db")
9 LOG_FILE = "logs/feature_store.log"
10
11 Path("logs").mkdir(parents=True, exist_ok=True)
12 logging.basicConfig(
13     filename=LOG_FILE,
14     level=logging.INFO,
15     format="%(asctime)s %(levelname)s %(message)s"
16 )
17 def log(msg):
18     print(msg)
19     logging.info(msg)
20
21 def init_metadata():
22     conn = sqlite3.connect(DB_FILE)
23     cursor = conn.cursor()
24     cursor.execute("""
25         CREATE TABLE IF NOT EXISTS feature_metadata (
26             feature_name TEXT PRIMARY KEY,
27             description TEXT,
28             source TEXT,
29             version INTEGER,
30             created_at TEXT
31         )
32     """)
33     conn.commit()
34     conn.close()
35     log("Feature metadata table initialized.")
36
37 def register_features(metadata):
38     conn = sqlite3.connect(DB_FILE)
39     cursor = conn.cursor()
40
41     for feature_name, details in metadata.items():
42         cursor.execute("""
43             INSERT INTO feature_metadata (feature_name, description, source, version, created_at)
44             VALUES (?, ?, ?, ?, ?)
45             ON CONFLICT(feature_name) DO UPDATE SET
46                 description=excluded.description,
47                 source=excluded.source.

```

```
model_building.py

1 import sqlite3
2 import pandas as pd
3 import numpy as np
4 import logging
5 from pathlib import Path
6 from datetime import datetime
7 from sklearn.model_selection import train_test_split
8 from sklearn.linear_model import LogisticRegression
9 from sklearn.ensemble import RandomForestClassifier
10 from sklearn.metrics import (
11     accuracy_score, precision_score, recall_score, f1_score, roc_auc_score
12 )
13 import joblib
14 import mlflow
15 import mlflow.sklearn
16 from mlflow.models import infer_signature
17
18 # ----- CONFIG -----
19 DB_FILE = Path("data/feature_store.db")
20 MODEL_DIR = Path("models")
21 LOG_FILE = "logs/model_building.log"
22 EXPERIMENT_NAME = "churn_prediction"
23
24 MODEL_DIR.mkdir(parents=True, exist_ok=True)
25 Path("logs").mkdir(parents=True, exist_ok=True)
26
27 logging.basicConfig(
28     filename=LOG_FILE,
29     level=logging.INFO,
30     format="%(asctime)s %(levelname)s %(message)s"
31 )
32 def log(msg):
33     print(msg)
34     logging.info(msg)
35
36 def load_features():
37     conn = sqlite3.connect(DB_FILE)
38     df = pd.read_sql("SELECT * FROM customer_features", conn)
39     conn.close()
40     return df
41
42 def simulate_labels(df):
43     np.random.seed(42)
44     df["churn?"] = np.random.choice([0, 1], size=len(df), p=[0.7, 0.3])
45     return df
46
47 def evaluate(model, X, test, y):
48     pass
```

The screenshot shows a code editor with a Python file named `orchestrate_prefect.py` open. The file contains a series of Prefect tasks defined using the `@task` and `@flow` decorators. The tasks include data ingestion, raw data storage, validation, preparation, transformation, feature store tasks, and a main flow that trains a model. The code editor interface includes a sidebar with project navigation, a top bar with file operations, and a status bar at the bottom.

```
orchestrate_prefect.py
DM4ML_Accreditation_ByStage > 02_Source_Code > step10_orchestration > orchestrate_prefect.py > ...
1  from prefect import flow, task
2  import subprocess
3
4  def run(cmd: str):
5      print(f"RUN: {cmd}")
6      subprocess.run(cmd, shell=True, check=True)
7
8  @task
9  def ingest(): run("python data_ingestion.py")
10
11 @task
12  def raw_store(): run("python storage_raw_upload.py")
13
14 @task
15  def validate(): run("python data_validation.py")
16
17 @task
18  def prepare(): run("python data_preparation.py")
19
20 @task
21  def transform(): run("python data_transformation.py")
22
23 @task
24  def feature_store_task(): run("python feature_store.py")
25
26 @task
27  def train(): run("python model_building.py")
28
29 @flow(name="churn_end_to_end_pipeline")
30  def main(*args, **kwargs):
31      (function) def raw_store() -> None
32      raw_store()
33      validate()
34      prepare()
35      transform()
36      feature_store_task()
37      train()
38
39  if __name__ == "__main__":
40      main_flow()
```

GROUP-103

- ✓ Airflow
 - ✓ airflow-docker
 - > dags
 - > logs
 - > plugins
 - ! .env
 - ✓ docker-compose.yaml
 - Dockerfile
 - requirements.txt

✓ DM4ML_Assignment_ByStage

✓ Pipeline_Design_with_Placeholders 2.pdf

docker-compose.yaml

```

1  docker > airflow-docker > docker-compose.yaml > {} services
2    services:
3      postgres:
4        image: postgres:13
5        environment:
6          POSTGRES_USER: airflow
7          POSTGRES_PASSWORD: airflow
8          POSTGRES_DB: airflow
9        volumes:
10          - postgres_data:/var/lib/postgresql/data
11
12      redis:
13        image: redis:latest
14
15      airflow-webserver:
16        build:
17        image: airflow-custom:latest
18        restart: always
19        depends_on:
20          - postgres
21          - redis
22          - airflow-scheduler
23        environment:
24          - AIRFLOW_CORE_EXECUTOR=CeleryExecutor
25          - AIRFLOW_DATABASE_SQLALCHEMY_CONN=postgresql+psycopg2://airflow@postgres/airflow
26          - AIRFLOW_CELERY_BROKER_URL=redis://redis:6379/0
27          - AIRFLOW_CELERY_RESULT_BACKEND=db+postgresql://airflow:airflow@postgres/airflow
28          - AIRFLOW__CORE__FERNET_KEY=${AIRFLOW__CORE__FERNET_KEY}
29        volumes:
30          - ./dags:/opt/airflow/dags
31          - ./logs:/opt/airflow/logs
32          - ./plugins:/opt/airflow/plugins
33          - ../../DM4ML_Assignment_ByStage/02_Source_Code:/opt/airflow/02_Source_Code
34          - ../../DM4ML_Assignment_ByStage/03_Data:/opt/airflow/03_Data
35          - ../../DM4ML_Assignment_ByStage:/opt/airflow/DM4ML_Assignment_ByStage
36        ports:
37          - "8080:8080"
38        command: webserver
39
40      airflow-scheduler:
41        build: .
42        image: airflow-custom:latest
43        restart: always
44        depends_on:
45          - postgres
46        environment:

```

✓ main*+ ↻ ⌂ 0 ⌂ 0 Git Graph

Duo Ln 1, Col 1 Spaces:2 UTF-8 LF ⌂ Compose ⌂ Prettier ⌂

GROUP-103

- ✓ Airflow
 - ✓ airflow-docker
 - > dags
 - > logs
 - > plugins
 - ! .env
 - ✓ docker-compose.yaml
 - Dockerfile
 - requirements.txt

✓ DM4ML_Assignment_ByStage

✓ Pipeline_Design_with_Placeholders 2.pdf

Pipeline_Design_with_Placeholders 2.pdf

8 of 10

Automatic Zoom

Airflow UI – Grid View after successful run

Start Pipeline

Get Data

Process Data

Feature Selection

Train Model

Evaluate Model

End Pipeline

Pipeline Summary

Total Tasks: 7

File Explorer

Outline

Timeline

Git Graph

Duo

