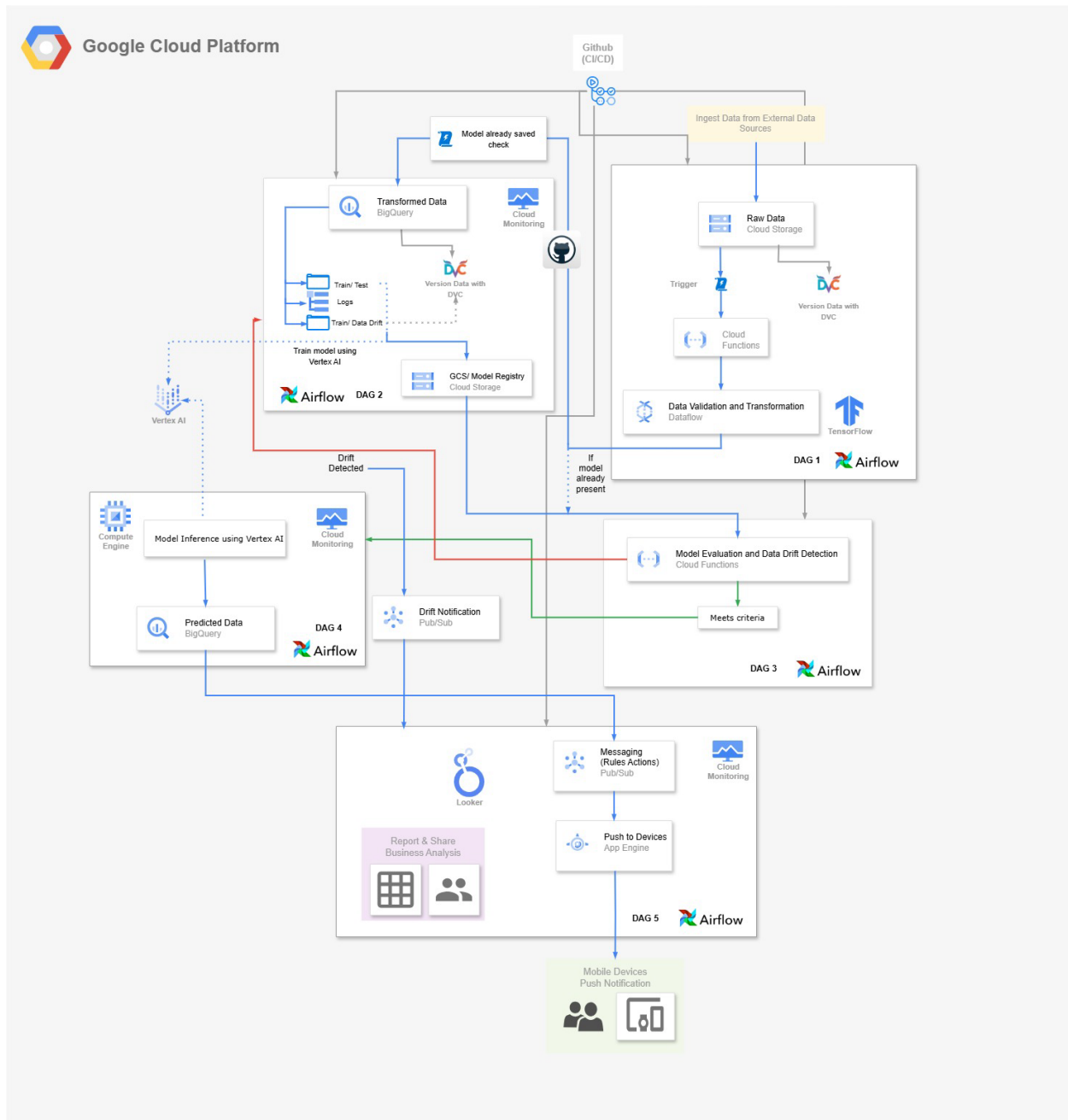


Comprehensive Data Pipeline Documentation for Telecom Churn Prediction and Device Upgrade Analysis



Data Preparation: Split the raw data locally into train.csv and test.csv for training and validation, respectively. holdout.csv is divided into smaller batches (e.g., holdout_batch_1.csv) to enable batch processing and drift detection.

```
bash
gs://data-source-telecom-customers/
├── raw_data/
│   ├── train/
│   │   └── train.csv # Raw training data
│   ├── test/
│   │   └── test.csv # Raw test data
│   └── validation_data/
│       ├── holdout_batch_1.csv # Holdout batch 1
│       ├── holdout_batch_2.csv # Holdout batch 2
│       ├── holdout_batch_3.csv # Holdout batch 3
│       └── holdout_batch_4.csv # Holdout batch 4
└── preprocessed_data/
    ├── train/
    │   └── train_preprocessed.csv # Preprocessed training data
    ├── test/
    │   └── test_preprocessed.csv # Preprocessed test data
    └── validation_data/
        ├── holdout_batch_1_preprocessed.csv # Preprocessed holdout batch 1
        ├── holdout_batch_2_preprocessed.csv # Preprocessed holdout batch 2
        ├── holdout_batch_3_preprocessed.csv # Preprocessed holdout batch 3
        └── holdout_batch_4_preprocessed.csv # Preprocessed holdout batch 4
```

The screenshot shows the Google Cloud console interface for an ML Ops Project. The top navigation bar includes the Google Cloud logo, project name "MLOpsProject", a search bar, and utility icons. A banner at the top indicates a free trial status with \$300.00 credit remaining. The left sidebar contains navigation links for Overview, Buckets, Monitoring, Settings, Marketplace, and Release Notes. The main content area displays the "Buckets" page, which includes a message about soft delete settings, action links like "CREATE" and "REFRESH", and a table listing buckets with columns for Name, Created, Location type, Location, Default storage class, and Last modified.

Bucket details

GO TO PATH REFRESH LEARN

< OBJECTS CONFIGURATION PERMISSIONS PROTECTION LIFECYCLE OBSERVABILITY INVENTORY REPORTS >

Folder browser

- data-source-telecom-customers
 - data/
 - preprocessed_data/
 - raw_data/
 - test/
 - train/
 - validation_data/
 - inference_logs/

Buckets > data-source-telecom-customers > data > raw_data

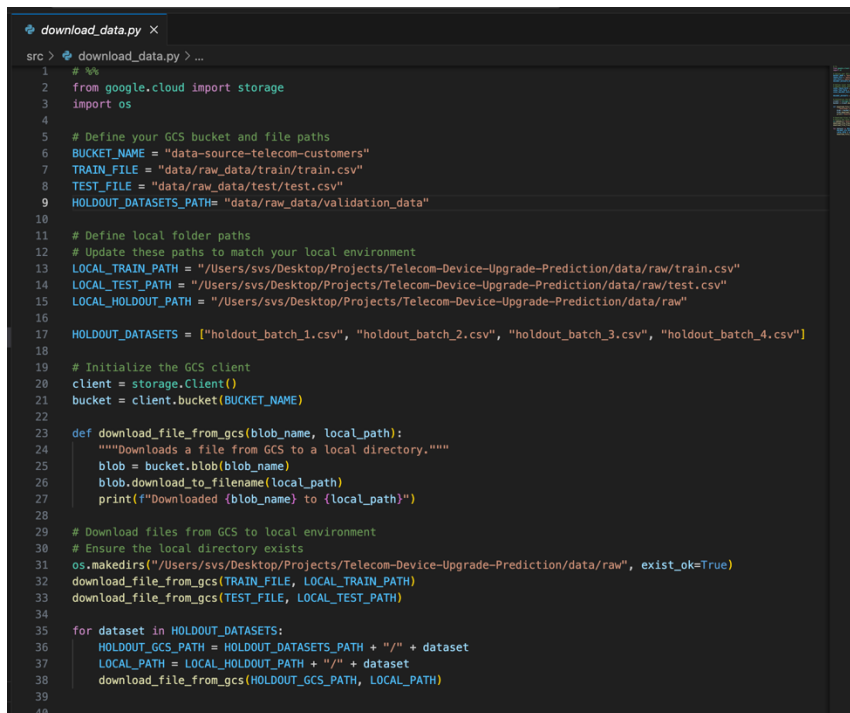
CREATE FOLDER UPLOAD TRANSFER DATA OTHER SERVICES

Filter by name prefix only Filter objects and folders Show Live objects only

<input type="checkbox"/>	Name	Size	Type	Created	Storage class
<input type="checkbox"/>	test/	—	Folder	—	—
<input type="checkbox"/>	train/	—	Folder	—	—
<input type="checkbox"/>	validation_data/	—	Folder	—	—

2. Download Data from GCS to Local Machine

After storing data in the GCS bucket, it's downloaded to the local environment for preprocessing. This step uses a script that fetches the raw train.csv, test.csv, and holdout data files directly from the GCS bucket to your local machine, readying the data for further cleaning and transformation.

A screenshot of a code editor window titled 'download_data.py'. The code is written in Python and is designed to download data from a Google Cloud Storage (GCS) bucket to a local environment. The script includes comments and code for defining GCS bucket and file paths, local folder paths, initializing the GCS client, and downloading files. The code is as follows:

```
1 # %%
2 from google.cloud import storage
3 import os
4
5 # Define your GCS bucket and file paths
6 BUCKET_NAME = "data-source-telecom-customers"
7 TRAIN_FILE = "data/raw_data/train/train.csv"
8 TEST_FILE = "data/raw_data/test/test.csv"
9 HOLDOUT_DATASETS_PATH = "data/raw_data/validation_data"
10
11 # Define local folder paths
12 # Update these paths to match your local environment
13 LOCAL_TRAIN_PATH = "/Users/svs/Desktop/Projects/Telecom-Device-Upgrade-Prediction/data/raw/train.csv"
14 LOCAL_TEST_PATH = "/Users/svs/Desktop/Projects/Telecom-Device-Upgrade-Prediction/data/raw/test.csv"
15 LOCAL_HOLDOUT_PATH = "/Users/svs/Desktop/Projects/Telecom-Device-Upgrade-Prediction/data/raw"
16
17 HOLDOUT_DATASETS = ["holdout_batch_1.csv", "holdout_batch_2.csv", "holdout_batch_3.csv", "holdout_batch_4.csv"]
18
19 # Initialize the GCS client
20 client = storage.Client()
21 bucket = client.bucket(BUCKET_NAME)
22
23 def download_file_from_gcs(blob_name, local_path):
24     """Downloads a file from GCS to a local directory."""
25     blob = bucket.blob(blob_name)
26     blob.download_to_filename(local_path)
27     print(f"Downloaded {blob_name} to {local_path}")
28
29 # Download files from GCS to local environment
30 # Ensure the local directory exists
31 os.makedirs("/Users/svs/Desktop/Projects/Telecom-Device-Upgrade-Prediction/data/raw", exist_ok=True)
32 download_file_from_gcs(TRAIN_FILE, LOCAL_TRAIN_PATH)
33 download_file_from_gcs(TEST_FILE, LOCAL_TEST_PATH)
34
35 for dataset in HOLDOUT_DATASETS:
36     HOLDOUT_GCS_PATH = HOLDOUT_DATASETS_PATH + "/" + dataset
37     LOCAL_PATH = LOCAL_HOLDOUT_PATH + "/" + dataset
38     download_file_from_gcs(HOLDOUT_GCS_PATH, LOCAL_PATH)
39
40
```

Note: This step is optional if you're working with our repository, as it already includes the raw train.csv, test.csv, and holdout files from the GCS bucket. You can proceed directly to preprocessing if the files are already present locally.

3. Data Loader

The data loader script (src/data_loader.py) is designed to retrieve datasets for processing. The script first checks if the raw data files (e.g., train.csv, test.csv) are available locally; if they are not found, they are automatically pulled from version control using DVC. Once available, the data is loaded into a panda DataFrame, preparing it for preprocessing in subsequent steps.

```
src > data_loader.py
4 import os
5 import sys
6 from config import RAW_DATA_PATH, PROCESSED_DATA_PATH
7
8 def load_data(file_path=RAW_DATA_PATH): #adjust path as needed
9     """
10     Load the dataset from the specified file path, using DVC if the file is
11
12     Args:
13         file_path (str): Path to the dataset file.
14
15     Returns:
16         pd.DataFrame: Loaded dataset as a DataFrame, or None if an error occurs.
17     """
18     if not os.path.exists(file_path):
19         print("Data file not found locally, pulling data using DVC...")
20         os.system(f'dvc pull {file_path}')
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
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
```

ModuleNotFoundError: No module named 'src'

PS C:\Users\A V NITHYA\MLopsProject\Telecom-Device-Upgrade-Prediction> python src\data_loader.py

Data loaded successfully

	CustomerID	Churn	MonthlyRevenue	...	PrizmCode	Occupation	MaritalStatus
0	3302766	No	47.73	...	Other	Other	Unknown
1	3039522	No	75.86	...	Town	Other	Unknown
2	3131142	No	28.75	...	Other	Other	No
3	3331902	No	10.00	...	Town	Professional	Yes
4	3045986	No	42.49	...	Other	Other	No

This approach ensures that the latest version of the dataset is always accessible, even if it is not stored locally, enabling a seamless transition to further stages in the pipeline.

4. Pre-processing Steps

```
[5 rows x 58 columns]
PS C:\Users\A V NITHYA\MLopsProject\Telecom-Device-Upgrade-Prediction> python src\preprocessing.py
Data preprocessing completed
Churn MonthlyRevenue MonthlyMinutes ... PrizmCode Occupation MaritalStatus
0 0 47.73 933.0 ... 0 3 1
1 0 75.86 655.0 ... 3 3 1
2 0 28.75 110.0 ... 0 3 0
3 0 10.00 256.0 ... 3 4 2
4 0 42.49 16.0 ... 0 3 0

[5 rows x 57 columns]
PS C:\Users\A V NITHYA\MLopsProject\Telecom-Device-Upgrade-Prediction>
```

In this step, preprocessing steps were performed to prepare the raw dataset for model training. The process was managed through a dedicated script, preprocessing.py, which includes several functions to handle data cleaning, encoding, and transformation:

- **Removal of Identifiers:** Columns with unique identifiers, such as customerID, were removed to protect sensitive information and retain only meaningful features for model training.
- **Handling Missing Values:** Missing values were imputed according to the category of the record (Churn or Non-Churn). Numerical columns were filled with median values, while categorical columns used the mode, ensuring that no essential information was lost.
- **Encoding Categorical Variables:**
 - Binary Variables: Used one-hot encoding and encoded as integers (e.g., "Yes" mapped to 1, "No" to 0).
 - Ordinal Variables: Mapped using predefined dictionaries (e.g., CreditRating, IncomeGroup) to retain their natural order.
 - Multi-class Variables: Label encoding was applied to make these variables suitable for model input.
- **Scaling of Numerical Columns:** Scaling was assessed but found unnecessary due to the use of tree-based models, which are not sensitive to feature scaling.

- **Testing Script: test_preprocessing.py:** To ensure each function in preprocessing.py worked as intended, a separate testing script, test_preprocessing.py, was developed using the unittest framework. Tests were created for each function, and specific issues, such as IntCastingNaNError, were addressed by adding .fillna() handling. This approach helped ensure compatibility and error-free execution across various data scenarios, reinforcing the robustness of the preprocessing pipeline.

5. Feature engineering

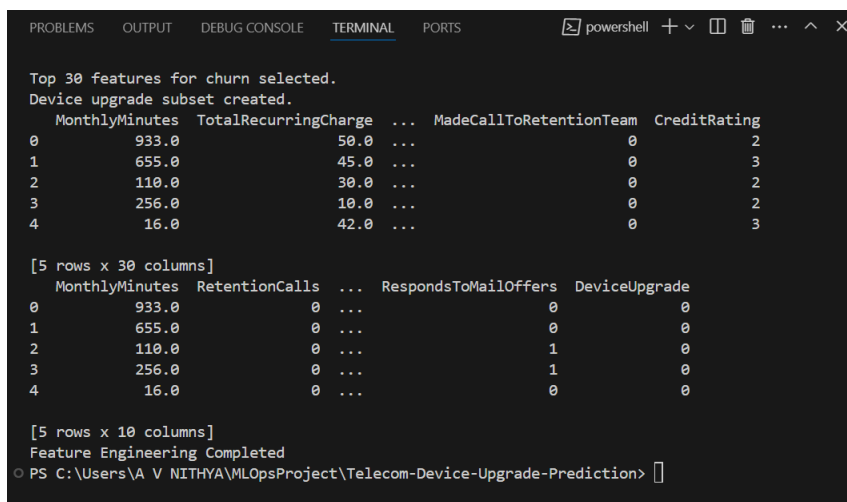
Extensive feature engineering was conducted to enhance model performance for predicting customer churn and identifying potential device upgrade candidates.

Feature Selection for Churn Prediction: Through experimentation using the SelectKBest method, it was determined that the optimal number of features (k) lies between 25 and 30 for customer churn prediction. Custom functions were developed to dynamically determine the best k within this range, selecting the most relevant features for the churn prediction model.

Creating a Device Upgrade Feature: In addition to churn prediction, an objective was to identify customers likely to upgrade their devices. Since the dataset lacked a direct indicator for device upgrades, a custom DeviceUpgrade feature was created using the following criteria:

- MonthlyMinutes > 3000
- RetentionCalls > 2
- RetentionOffersAccepted > 0
- HandsetWebCapable == 0
- HandsetRefurbished == 1
- CurrentEquipmentDays > 340
- CreditRating > 5
- MadeCallToRetentionTeam == 1
- RespondsToMailOffers == 1

These thresholds were defined based on domain-specific research, expert knowledge, and literature review, ensuring that the criteria were grounded in relevant insights.



```

PROBLEMS OUTPUT DEBUG CONSOLE TERMINAL PORTS powershell + - [] ... ^ x

Top 30 features for churn selected.
Device upgrade subset created.
  MonthlyMinutes  TotalRecurringCharge  ...  MadeCallToRetentionTeam  CreditRating
0           933.0             50.0  ...              0              2
1          655.0             45.0  ...              0              3
2          110.0             30.0  ...              0              2
3          256.0             10.0  ...              0              2
4           16.0             42.0  ...              0              3

[5 rows x 30 columns]
  MonthlyMinutes  RetentionCalls  ...  RespondsToMailOffers  DeviceUpgrade
0           933.0              0  ...              0              0
1          655.0              0  ...              0              0
2          110.0              0  ...              1              0
3          256.0              0  ...              1              0
4           16.0              0  ...              0              0

[5 rows x 10 columns]
Feature Engineering Completed
PS C:\Users\A V NITHYA\MLOpsProject\Telecom-Device-Upgrade-Prediction>

```

6. Testing and Debugging

Unit Testing:

- Comprehensive tests were conducted using unittest and pytest to validate the functionality of the preprocessing.py and feature_engineering.py scripts.
- Errors encountered during testing, such as IntCastingNaNError, were resolved by implementing robust handling of missing values before typecasting.

Debugging:

- Issues with .astype() errors, particularly when converting columns with NaN values to integers, were addressed by adding appropriate checks and handling methods.
- Warnings from libraries were managed to ensure that all functions remain compatible with future library updates, reinforcing the reliability and stability of the codebase.

This rigorous testing and debugging process ensured that both preprocessing and feature engineering workflows were robust, accurate, and compatible with current and future data scenarios.

7. Airflow DAG

1. Configuring Airflow

1.1 Initialize the Airflow Database

```
airflow db init
```

1.2 Start the Airflow Web Server and Scheduler

- Start the web server to access the Airflow UI:
`airflow webserver --port 8080`
- In a separate terminal, start the scheduler:
`airflow scheduler`

1.3 Accessing the Airflow UI

- Open a browser and go to `http://localhost:8080`.
- Log in with default credentials or create an admin user if needed.

2. DAG Creation: telecom_pipeline

The telecom_pipeline DAG automates data loading, preprocessing, and feature engineering for telecom data to predict device upgrades. The DAG includes three primary tasks:

2.1 Setting Up the DAG

- Place the DAG script (telecom_pipeline.py) in the `$AIRFLOW_HOME/dags` directory.
- The DAG consists of three tasks:

Task 1: Load Data - Loads raw data using `load_data` from `data_loader.py` and saves it temporarily.

Task 2: Preprocess Data - Preprocesses the loaded data using `preprocess_data`, saving the output for further analysis.

Task 3: Feature Engineering - Performs feature engineering to select key features for churn and device upgrade prediction.

2.2 Task Overview

- **load_data_task:** Loads raw data and saves it to `/tmp/loaded_data.csv`.

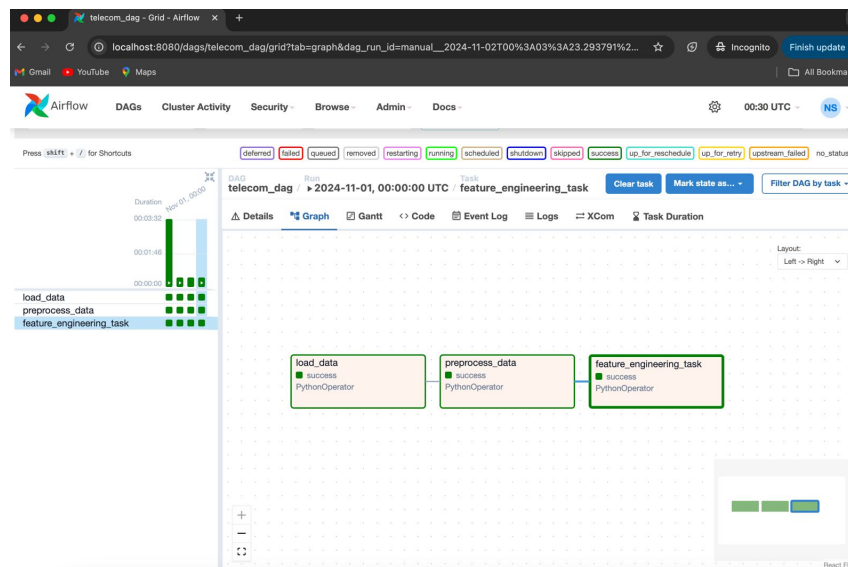
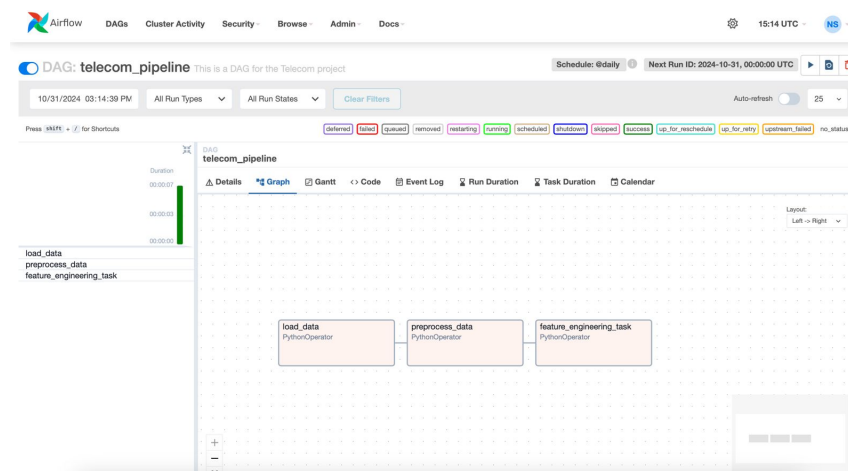
- **preprocess_data_task:** Reads, preprocesses, and saves data to /tmp/preprocessed_data.csv.
- **feature_engineering_task:** Reads preprocessed data, applies feature selection for churn prediction and device upgrade analysis, and saves results for further steps.

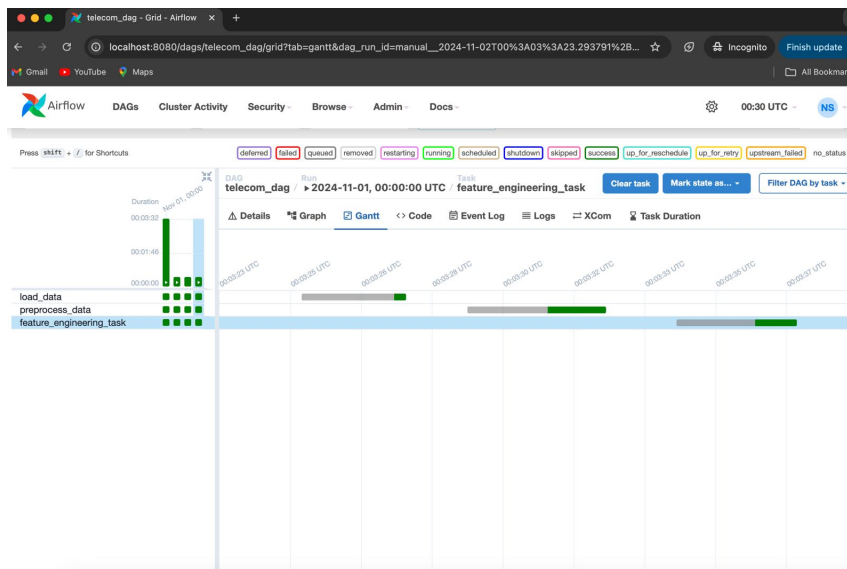
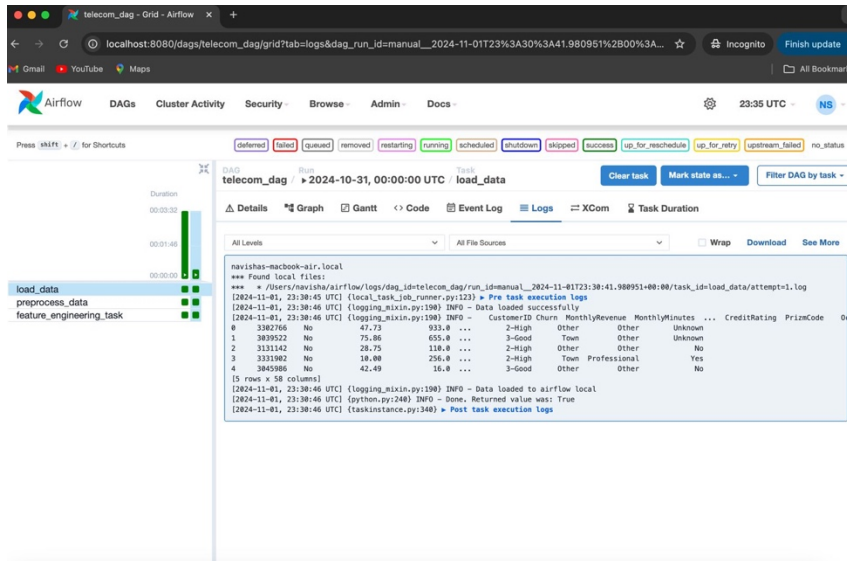
2.3 Running the DAG

- After defining the DAG, access the Airflow UI (<http://localhost:8080>) and manually trigger the DAG to confirm functionality.

3. Summary

The Airflow DAG (telecom_pipeline) automates the workflow for loading, preprocessing, and engineering features, forming a foundation for more complex telecom data processing tasks. This setup facilitates reliable, repeatable data processing in Airflow.





8. Schema Generation

Schema generation and validation are set up using TensorFlow Data Validation (TFDV), ML Metadata (MLMD), and Airflow for automation.

- **Dependencies:** Install required libraries (tensorflow, tensorflow-data-validation, apache-airflow, ml-metadata) listed in requirements.txt.
- **Schema Generation:** generate_schema.py uses TFDV to create a schema from the dataset, saved as schema.pbtxt.
- **Data Validation:** validate_data.py checks new data against the schema, identifying anomalies (e.g., missing values or data type issues).
- **Logging with MLMD:** Schema versions and validation runs are logged in MLMD for tracking changes and ensuring consistency.
- **Automation with Airflow:** An Airflow DAG schedules daily schema generation and validation, automating data quality checks.
- **Alerts:** Notifications are added to alert the team to validation failures.

```

1  feature {
2      name: "CustomerID"
3      type: INT
4      presence {
5          min_fraction: 1.0
6          min_count: 1
7      }
8      shape {
9          dim {
10             size: 1
11         }
12     }
13 }
14 feature {
15     name: "Churn"
16     type: BYTES
17     domain: "Churn"
18     presence {
19         min_fraction: 1.0
20         min_count: 1
21     }
22     shape {
23         dim {
24             size: 1
25         }
26     }
27 }
28 feature {
29     name: "MonthlyRevenue"
30     value_count {
31         min: 1
32         max: 1

```

data/schema/schema.pbtxt

9. Data slicing

Data slicing is performed on the test dataset to create targeted subsets, enabling detailed evaluation of model performance and potential bias mitigation. This process involves segmenting the data based on specific variables that are relevant to our analysis, specifically Churn, IncomeGroup, and CreditRating.

- **Churn-Based Slicing:** The dataset is first divided into two subsets based on the Churn column, creating separate DataFrames for customers who have churned (Churn = Yes) and those who have not (Churn = No). This allows for targeted analysis of model performance for each group, ensuring fair representation.
- **Income Group and Credit Rating Slicing:** Further slicing is performed on IncomeGroup and CreditRating. DataFrames are generated for each unique value in these columns, allowing for group-based analysis. This helps assess how the model performs across different income levels and credit ratings, identifying potential biases in predictions related to these variables.
- **Equipment Days Range:** A new column, `equipment_days_range`, is created by binning `CurrentEquipmentDays` into intervals (e.g., 0-100, 100-200, etc.). These bins provide an additional segmentation based on how long customers have had their current equipment, giving insight into performance variations for customers at different stages of their device lifecycle.

These slices create distinct, manageable views of the dataset, facilitating a detailed evaluation of model performance and helping identify any biases across different customer segments.

```
data_slicing.py M X download_data.py data_loader.py M main.py
src > data_slicing.py
1 # Code to slice test data for later use (performance evaluation and mitigation)
2 import pandas as pd
3 from config import TEST_DATA_PATH
4
5 # Slice the data by churn category (0 or 1)
6 def slice_by_churn(data):
7     churn_0 = data[data['Churn'] == 'No']
8     churn_1 = data[data['Churn'] == 'Yes']
9     return churn_0, churn_1
10
11 # Slice the data by a specific column and return a dictionary of DataFrames
12 def slice_by_column(data, column_name):
13     grouped_data = data.groupby(column_name)
14     return {group: subset for group, subset in grouped_data}
15
16 # Create a new column for binning based on 'CurrentEquipmentDays' and slice
```