

COURSE PROJECT REPORT

CREDIT-CARD FRAUD DETECTION PIPELINE

Vikrant Deshpande

INFO-I535 Management, Access, & Use of Big Data

1. Introduction (brief overview of the problem or topic, 1-2 paragraphs)

In order to prevent customers from being charged for items they did not buy, credit card companies must be able to accurately identify fraudulent credit card transactions. An excellent article by Stripe^[1] highlights that companies lose more loyal customers (approximately 33%) due to blocked transactions. On the other hand, fraudulent transactions going unnoticed could lead to potential lawsuits, operational dispute fees, etc.

2. Background (why this problem or topic was selected, why it is interesting / important to address, 1-2 paragraphs)

A subtle line needs to be drawn and my main focus here is to reduce blocked-transactions of legitimate customers. We can design the Null and Alternative Hypothesis as "A transaction is legitimate until proven otherwise":

H_0 : *The transaction is legitimate*

H_1 : *The transaction is fraudulent*

A Type-II error would be a False-Negative: expected H_0 but wrongly predicted H_1 . This is what we want to minimize to deliver tangible business-value.

My aim is to provide stakeholders with an automated cloud-framework that enables distinction between fraud and legitimacy, by leveraging Machine Learning. Simultaneously, the framework should allow for downstream analysis like dashboarding through effective data-management.

As such, I've created and deployed an end-to-end Machine Learning pipeline on GCP using Airflow as the pipeline orchestrator, and Cloud-Functions to monitor source-file changes on Google-Cloud-Storage.

The dataset is openly available as part of a Data Science challenge^[7]. Here's some metadata that I've interpreted:

Column	Data-Type	Description
accountNumber	Integer	Account-Number of customer who made the transaction
customerId	Integer	
creditLimit	Integer	Current credit-limit of customer who made the transaction
availableMoney	Float	Available credit for customer
transactionDateTime	String	Timestamp of transaction
transactionAmount	Float	Amount for transaction
merchantName	String	Name of merchant (POS)
acqCountry	String	Country where credit card was acquired
merchantCountryCode	String	Country code for merchant (POS)
posEntryMode	Float	2-digit code for info on how card was entered (https://help.globalmerchantportal.com/en/code-glossary/pos-entry-modes)
posConditionCode	Float	2-digit code to identifies mode of initiation of the transactions.
merchantCategoryCode	String	General category (fast food, entertainment, etc.)
currentExpDate	String	Card expiry date
accountOpenDate	String	Account opened (Start) date
dateOfLastAddressChange	String	Last updated date of address
cardCVV	Integer	Expected CVV available on credit-card
enteredCVV	Integer	Entered CVV at POS
cardLast4Digits	Integer	Last 4 digits of Credit Card
transactionType	String	Purchase, Reversal, Address Verification, etc.
currentBalance	Float	CreditLimit - AvailableMoney
merchantCity	Float	City for POS
merchantState	Float	State for POS
merchantZip	Float	Zipcode for POS
cardPresent	Boolean	True/False for card presented in-person
posOnPremises	Float	POS on premise of store
isFraud	Boolean	Target-Label for transaction being marked as Fraud

3. Methodology (steps taken to address the problem, including the technological setup 1-2 pgs)

As with any major Machine Learning project, I started with exploring the data and creating modularized functions. I released a package on PyPi [2] wherein I modularized my code for the entire training and prediction phase. Eventually, this containerized package was used in Airflow DAGs within a `PythonVirtualenvOperator()`, on a Cloud Composer environment.

Overall, here's a visual picture of what my containerized Airflow DAGs and package accomplish:

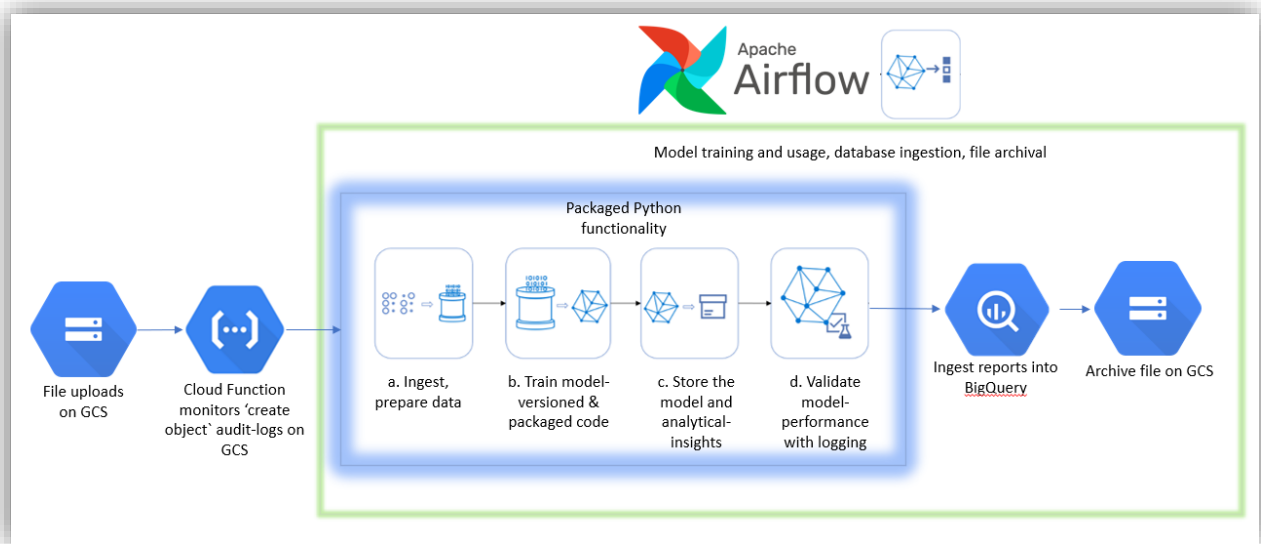


Fig.1.1. Overview of Airflow DAGs and containerized Python package

I'm maintaining my storage objects in the "i535-course-project-bucket" GCS Storage bucket. There are 2 Airflow DAGs, one for the model-training phase, and one for using the registered model to make predictions. The training-phase DAG gets triggered when a "Train_transactions.csv" file is uploaded onto Google-Cloud-Storage in the specified bucket. Similarly, for the prediction DAG but for "Test_transactions.csv".

2 Google Cloud Functions in Python: "gcs-trainfile-watcher-trigger-function", and "gcs-testfile-watcher-trigger-function", perform the source-file-upload monitoring and eventual Airflow DAG triggering activity.

The following 4 points should help summarize the package quickly:

a. Ingest and prepare data:

A series of data-cleansing, standardization, scaling, and feature-engineering steps take place:

1. **read_src_file_as_df()**- The Airflow task reads the "Train/Test_transactions.csv" file from GCS, and converts it into a dataframe.

2. **add_time_dependent_features()**- For the input datetime columns, engineer new date and time features. Also drop out the original column.
3. **levenshtein_distance()**- Rather than just binary non-equality, engineer a score of mismatch for enteredCVV vs cardCVV: a single digit wrong might be just a minor blunder.
4. **drop_unary_columns()**- Filter out features with ≤ 1 unique values.
5. **get_reversals_report()**- Create a report of reversal transactions.
6. **get_multiswipe_transactions()**- Create a report of multiswipe transactions.
7. **convert_boolean_to_int()**- Convert Boolean-columns in the data to integer encodings.
8. **encode_categorical_cols()**- Apply Ordinal-Encoding to each of categorical-columns and keep track of the transformations, ie- store the sklearn models as pickle objects on GCS for use during prediction.
9. **scaledown_numerical_cols()**- Apply MinMax-Scaling to each of numerical-columns and keep track of the transformations, ie- store the sklearn models as pickle objects on GCS for use during prediction.
10. **drop_irrelevant_columns()**- Filters out the non-numerical columns from data.

b. Train model- versioned and packaged code:

1. **encode_categorical_cols()**, **scaledown_numerical_cols()**- Train Ordinal-Encoder and MinMax-Scaler sklearn models for standardizing/transforming data.
2. **train_random_forest_classifier()**- Apply grid-search-cv for a random-forest-classifier.

c. Store the model and analytical insights:

1. **encode_categorical_cols()**, **scaledown_numerical_cols()**, and **train_random_forest_classifier()**- Store the respective models as pickle objects on GCS for use during batch-prediction.
2. The **model.predict()** outputs are stored in GCS for the train and test phase as "Random_Forest_Validation_Set_predictions.csv" and "Random_Forest_Test_Set_predictions.csv" respectively.
3. The analytical csv-reports from **get_reversals_report()** and **get_multiswipe_transactions()**, "Overall_Report_Reversals.csv", "Overall_Report_Multiswipes.csv" are later ingested into BigQuery as a Dataset.

d. Validate model-performance with logging:

The package supports efficient logging within a containerized environment (configurable verbosity of logs), with appropriate try-catch blocks and a **print_neat_metrics()** to evaluate model performance after each training-phase execution.

The next Airflow task is for ingestion of the csv reports currently on GCS, into BigQuery. I've also added in a logic to capture source date and last updated date changes, to mimic a data-warehouse SCD-Type 2.

Finally, the source CSV files are archived into a folder on the same bucket's "/archive/" folder. Their file-names are appended with the date when the archival happened ie- the date of being processed.

4. Results (1-3 pgs)

A working demonstration can be found on the wikis of my Github Repository: [here](#). Nevertheless, the following screenshots should showcase the algorithmic flow of tasks.

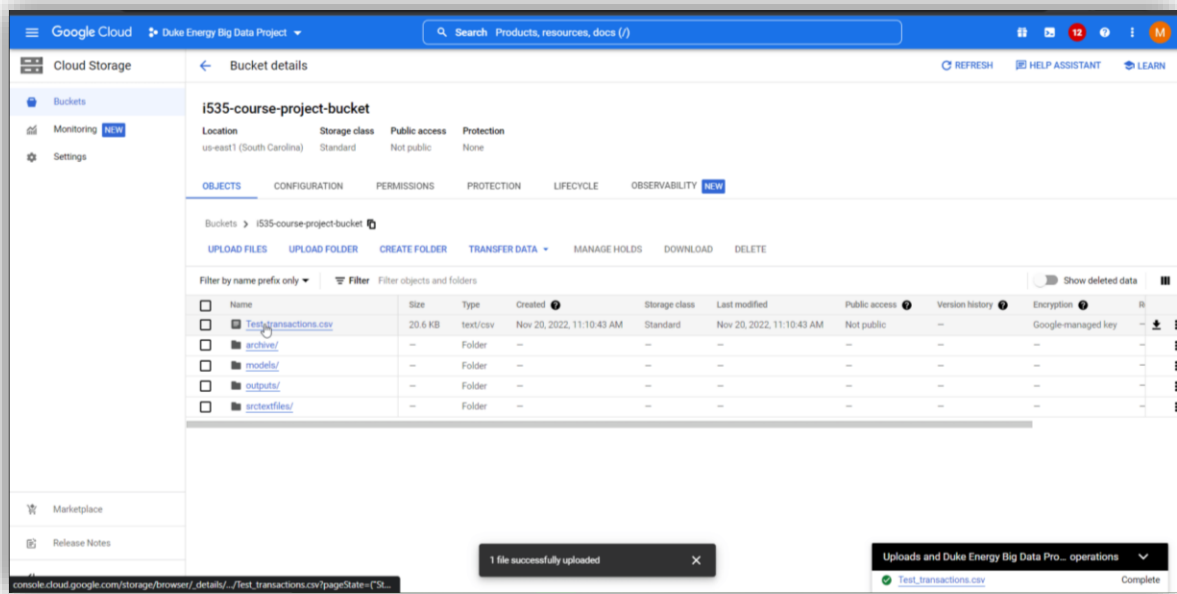


Fig.4.1. File upload onto GCS

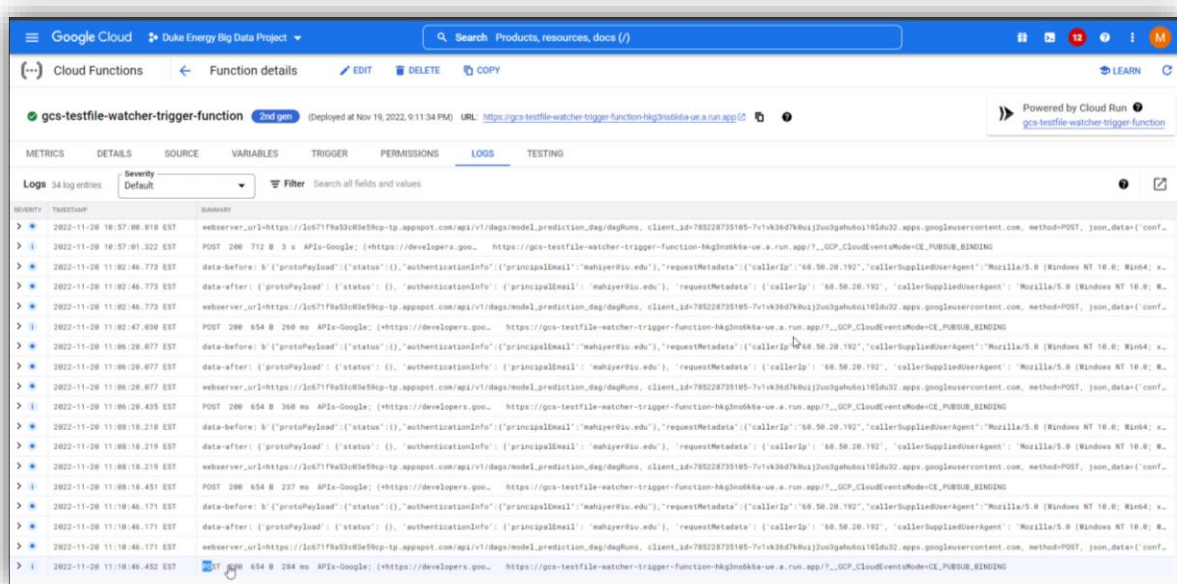


Fig.4.2. Cloud Function makes POST-request to Cloud Composer (Airflow) after file-upload

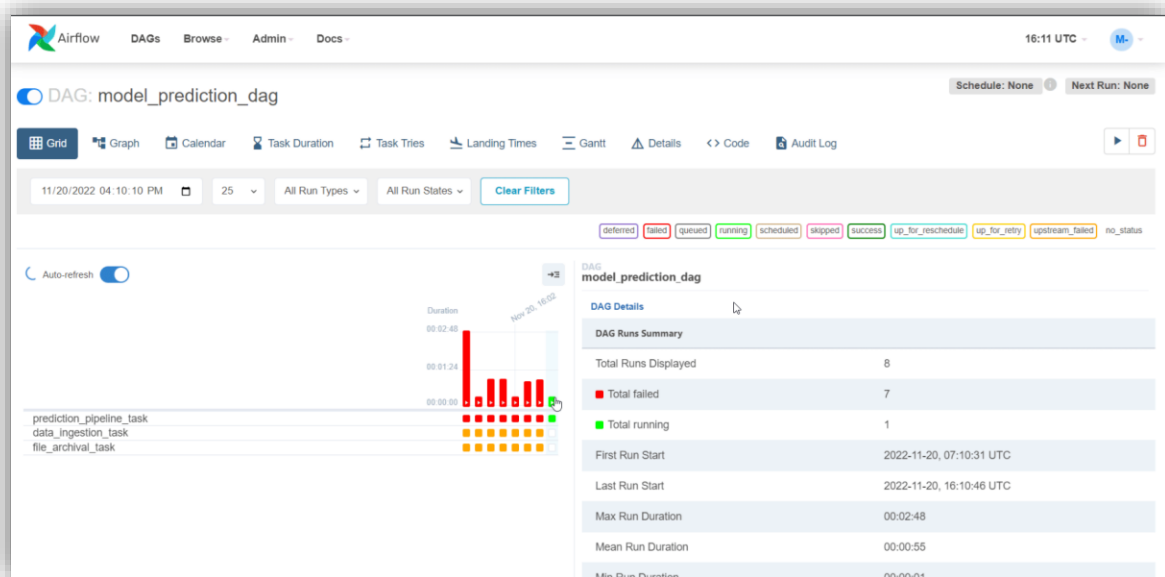


Fig.4.3. New Airflow run triggered due to POST-request from Cloud Function

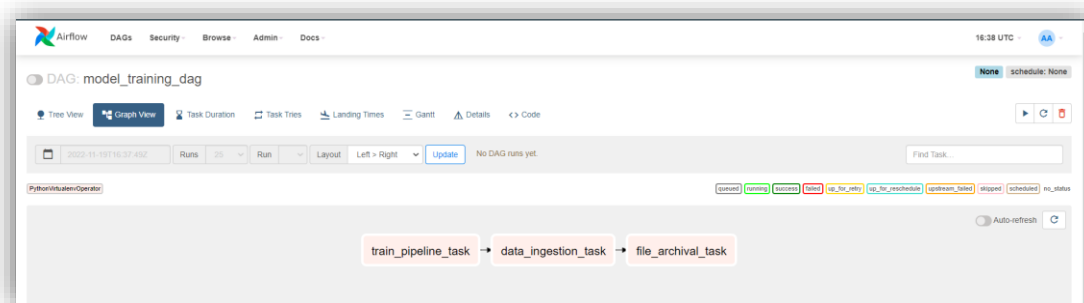


Fig. 4.4. Airflow DAG to train and register a model, ingest data into BigQuery, and archive files on a GCS bucket

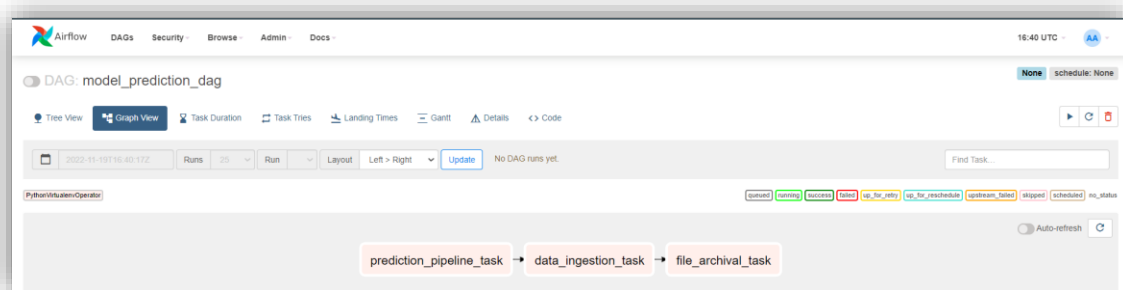


Fig.4.5. Airflow DAG to use registered-model for batch prediction on test-set, ingest data into BigQuery, and archive files on a GCS bucket

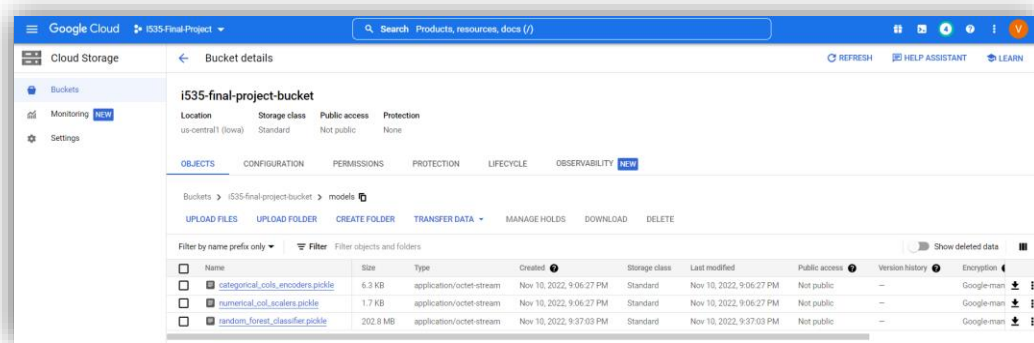


Fig.4.5. Models are cached in GCS bucket's `/models/` directory

Ideally, the model objects would be cached into a Model Registry like Vertex-AI on GCP. However, in the interest of time, I chose to simply cache and retrieve these pickled objects from GCS itself.

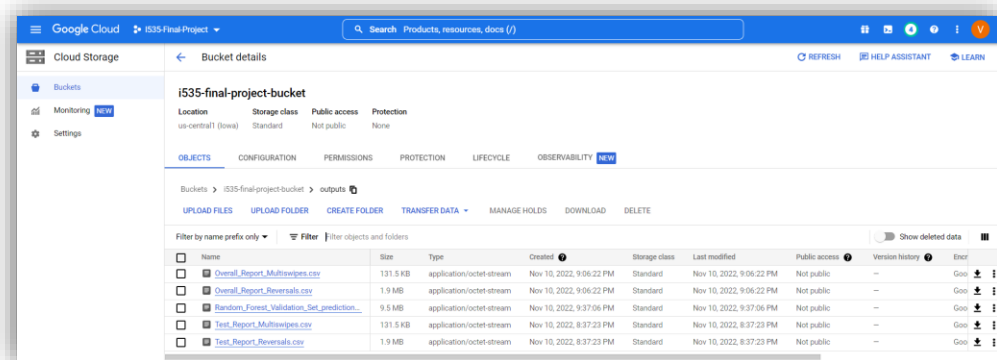


Fig.4.6. Analytical reports captured and stored within GCS bucket's `/outputs/` directory. The next Airflow Task loads these into BigQuery tables.

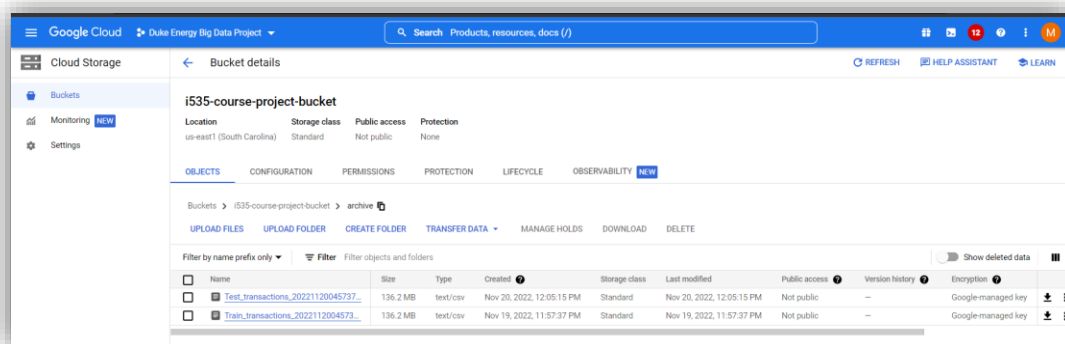


Fig.4.6. Final Airflow task archives source files in GCS bucket again

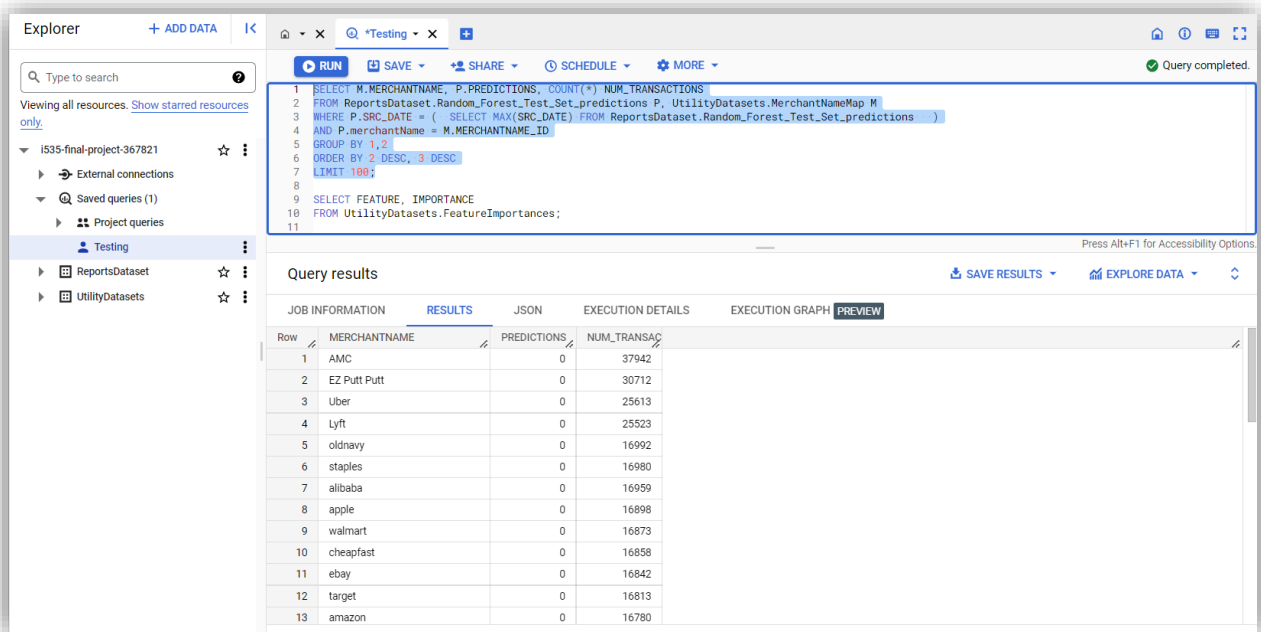


Fig.4.7. Data Management allows for downstream analytics and dashboarding

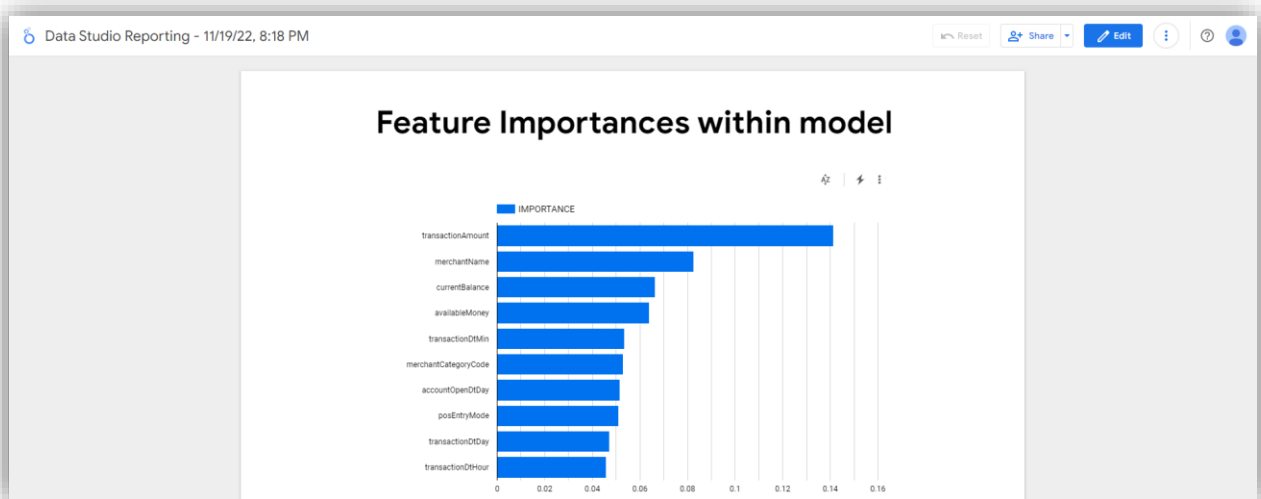


Fig.4.8. Data Studio dashboard of latest model performance ingested via pipeline, visible [here](#)

5. Discussion (an interpretation of the results + a discussion of how you employed the technologies / skills from this course + any barriers or failures you encountered, 1-2 pgs)

The steps to deploy the Airflow DAGs and GCS-monitoring functions may sound simple, but I ran into multiple challenges while trying to move my codebase from independent components, into the cloud. While deploying my code within Airflow, I realized that the `PythonVirtualenvOperator()` has some limitations in terms of not being able to read anything outside of its own local containerized scope. Hence it was necessary to create my own package^[2], publish it on the official PyPi server, and then use the modular components within the Airflow tasks- a great learning experience!

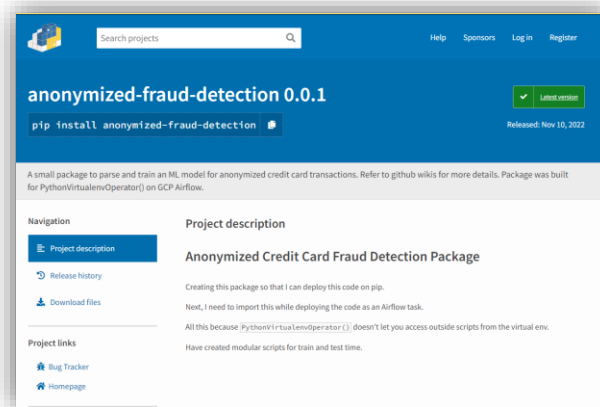


Fig.1.1. Published own custom package- necessary for use within PythonVirtualenvOperator()

Similarly while creating the cloud functions to monitor source-file changes, I ran into a bunch of issues due to lack of up-to-date documentation of GCP EventArc Triggers. Mainly, there are scarce documented examples^[3] on how to configure Path patterns for monitoring specific file-creation events on a bucket. I managed to figure out the right configuration after 3 days of rigorous experimentation- a great learning experience!

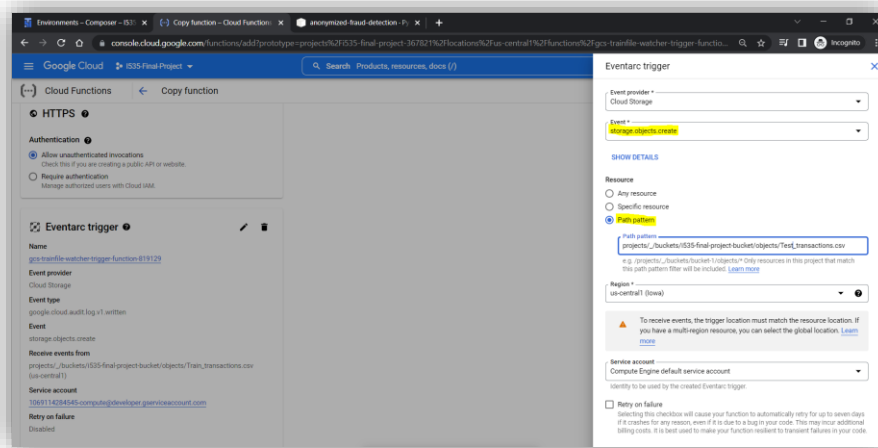


Fig.1.1. Cloud function configuration to monitor specific files on bucket

Other issues I encountered were:

1. Unable to setup connection with GCS bucket while reading and writing files.
2. Unable to store model pickle files on GCS bucket.
3. Unable to read the 580mb source text file due to limitations of compute resources (for some reason I'm unable to allocate a "Large" Cloud Composer environment).
4. Cloud Function trigger errors, due to function-"context" received in bytes rather than dictionary.
5. Cloud Function unable to communicate with Composer v1.

Eventually, I managed to fix these errors (except #3) and got my pipeline up and running!

The workaround I used for #3 was to convert the 580mb TXT file into a 360mb .csv file, and things worked fine. I didn't personally do any compression, the CSV format was just inherently more efficient in storage than the TXT format.

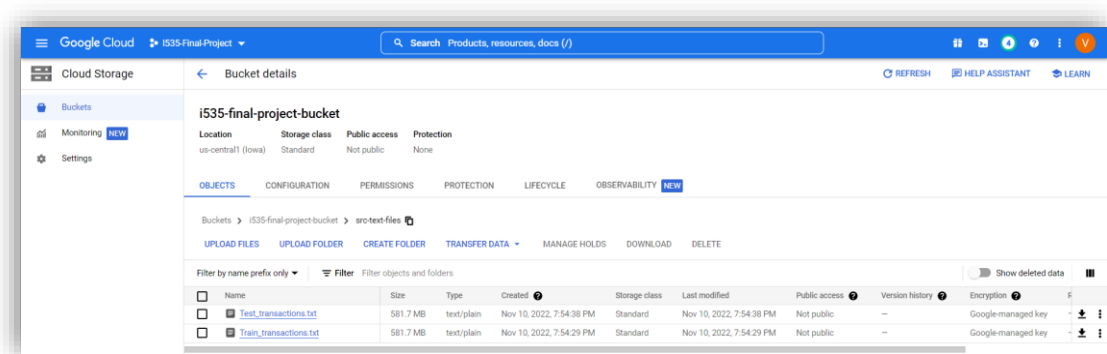


Fig.1.1. Ideally, .txt should be ingested but due to IAM access constraints, had to resort to CSV files

The workaround for #5 was frustrating because GKE services were going down around 15th November to 22nd November (autoscaling is a feature for Composer v2). Hence, I had to resort to Composer v1, but the services were still unable to communicate. After hours of searching, the solution was to make sure the Cloud Composer v1 environment has this variable configured^[5]:

``api.composer_auth_user_registration_role = Op``

This basically ensures that anyone can make REST-requests to the Airflow server API.

Additionally, I also had to come up with a data management strategy, hence I decided to implement a file-archival step, as well as using a timestamp to indicate creation-date in the database table. An advanced implementation would use last-updated-date if we're not looking to completely overwrite each row (SCD Type 2 table). For now, the data is ingested into the table as a fresh new load.

The screenshot shows the Google Cloud BigQuery console. The query editor at the top contains the following SQL:

```
1 SELECT * FROM `1535-final-project-367821.ReportsDataset.Overall_Report_Reversals`
2 WHERE SRC_DATE = (SELECT MAX(SRC_DATE) FROM `1535-final-project-367821.ReportsDataset.Overall_Report_Reversals`)
```

The query results are displayed in a table with the following columns: Row, customerid, transactionDate, transactionAmount, merchantName, transactionType, cardpy, isFraud, nextTransAmount, nextTransDate, and nextMerchantName. The results show 8 rows of transaction data.

Row	customerid	transactionDate	transactionAmount	merchantName	transactionType	cardpy	isFraud	nextTransAmount	nextTransDate	nextMerchantName
1	687511674	2016-04-29	3.0	Lyft	PURCHASE	fal..	fal..	3.0	2016-04-29	Lyft
2	321944399	2016-08-10	16.0	Lyft	PURCHASE	fal..	fal..	16.0	2016-09-01	Lyft
3	926332998	2016-05-16	17.0	Lyft	PURCHASE	fal..	fal..	17.0	2016-05-16	Lyft
4	133105237	2016-03-02	353.25	Lyft	PURCHASE	fal..	fal..	353.25	2016-03-20	Lyft
5	529994423	2016-09-22	676.0	Lyft	PURCHASE	fal..	fal..	676.0	2016-10-06	Lyft
6	564665867	2016-04-23	209.25	Lyft	PURCHASE	fal..	fal..	209.25	2016-04-23	Lyft
7	507456073	2016-04-03	241.5	Lyft	PURCHASE	fal..	fal..	241.5	2016-04-25	Lyft
8	442799740	2016-07-27	519.5	Lyft	PURCHASE	fal..	fal..	519.5	2016-08-10	Lyft

Fig.1.1. BigQuery tables uses `src_date` to mark active records

6. Conclusion (1 paragraph)

In conclusion, the project incorporated data management, data cleaning/wrangling, reporting, software engineering, containerization, deployment, and most importantly, using Github to maintain README tutorials, Wikis, documentation, and versioning of code. Versioning of models/data is made possible on the Google Cloud Storage bucket itself, and historical-datasets can be accessed on BigQuery using the Source-date field that distinguishes between active/inactive records.

Additionally, I'd enabled my entire Airflow pipeline to run on a Google Kubernetes Engine cluster using the Cloud Composer v2 service. This service supports autoscaling, load balancing, and distributed processing for Airflow tasks. However, in the last week of the course, this service was intermittently down, so I had to resort to Cloud Composer v1 (without autoscaling).

All in all, a really great and enjoyable learning experience!

7. References

1. <https://stripe.com/guides/primer-on-machine-learning-for-fraud-protection>
2. <https://pypi.org/project/anonymized-fraud-detection/>
3. <https://medium.com/google-cloud/applying-a-path-pattern-when-filtering-in-eventarc-f06b937b4c34>
4. <https://cloud.google.com/composer/docs/how-to/using/triggering-with-gcf>
5. <https://cloud.google.com/composer/docs/access-airflow-api>
6. https://github.com/vikrantdeshpande09876/Credit_Card_Fraud_Detection/
7. <https://github.com/CapitalOneRecruiting/DS>