

# Comprehensive Documentation:

## Machine Learning Model Pipeline

### 1. Overview

This section introduces the project's objective: building an end-to-end pipeline for machine learning using GCP. The pipeline incorporates preprocessing, automated workflows with Cloud Composer, and model training and inferences with Vertex AI.

#### Objective:

- Transition from local data preprocessing to cloud-based orchestration.
- Develop and evaluate machine learning models using both AutoML and custom built ML models.
- Ensure reproducibility, scalability, and fairness in model development.

### 2. Data Loading

#### 2.1. Implementing Airflow DAG using Cloud Composer

- The transition from a local processing environment to the cloud was facilitated using Google Cloud Composer. After preprocessing the data locally in Visual Studio Code, it was essential to ensure scalability and automation, prompting the move to Google Cloud Platform (GCP). The data was uploaded to Google Cloud Storage (GCS), which served as a centralized repository for storing raw and processed data. The use of Cloud Composer enabled the orchestration of complex workflows through Airflow's Directed Acyclic Graphs (DAGs). These DAGs automated the execution of preprocessing modules, ensuring seamless integration with the cloud infrastructure. Additionally, all dependencies and Python libraries required for preprocessing were mirrored in the Composer environment, maintaining consistency with the local setup. This transition allowed the preprocessing pipeline to handle larger datasets with greater efficiency while ensuring that every step was reproducible and scalable for future tasks.

### 3. Cloud Composer Workflow

Google Cloud Composer was employed to manage the data processing workflow by automating tasks and ensuring seamless integration of preprocessing modules developed locally. By leveraging Composer's orchestration capabilities, complex workflows were structured into easily manageable steps, ensuring reproducibility and scalability. This section outlines the setup of

Cloud Composer, the details of DAG 1 for data processing, and an overview of challenges encountered along with their resolutions.

### 3.1. Setting Up Cloud Composer

- The primary purpose of setting up Cloud Composer was to orchestrate complex workflows using Airflow's Directed Acyclic Graphs (DAGs). This setup provided a robust, managed environment for automating the data processing pipeline. To ensure compatibility, the Composer environment was configured with all necessary google cloud services, such as google-cloud-aiplatform, google-cloud-storage
- . Additionally, integration with Google Cloud Storage (GCS) allowed seamless data ingestion and output handling, ensuring that the environment was ready for executing DAGs designed to process and transform data efficiently.

### 3.2. DAG 1: Data Processing Workflow

The first DAG (DAG 1) implemented a fully automated data processing pipeline. It orchestrated the flow of tasks required to clean and preprocess raw data stored in Google Cloud Storage (GCS).

#### Pipeline Steps:

1. **Load Raw Data:** The workflow began by ingesting raw data files directly from GCS buckets into memory. This step ensured compatibility with the processing pipeline and validated the initial state of the data.
2. **Integrate Local Modules:** Custom Python modules previously developed locally were integrated into the DAG. These modules executed the preprocessing steps, including:
  - Feature engineering to derive additional attributes.
  - Cleaning operations to handle missing values, duplicates, and inconsistencies.
3. **Save Processed Data:** The cleaned and transformed datasets were written back to GCS in a structured format, ensuring they were ready for subsequent machine learning tasks.

#### Outputs:

The DAG produced fully cleaned and preprocessed datasets stored in GCS. These outputs maintained consistency and quality, making them suitable for downstream operations like training machine learning models in Vertex AI.

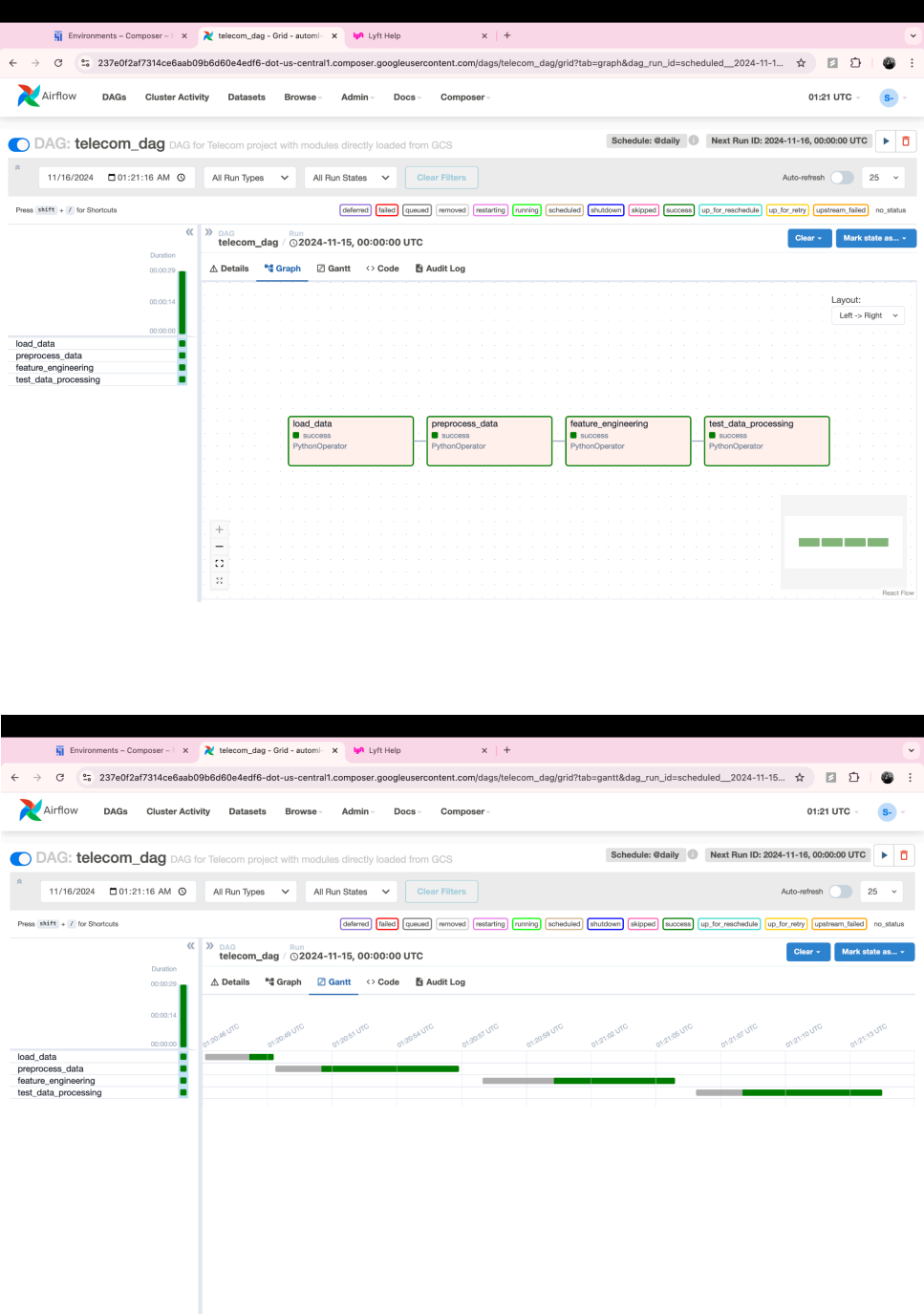


Fig: Data pipeline on cloud composer

## 4. Model Training on Vertex AI

The model training process on Vertex AI involved multiple approaches to leverage AutoML capabilities and custom model development, each offering distinct benefits and limitations. Vertex AI provided a streamlined environment for training, evaluating, and deploying machine learning models with minimal manual intervention while maintaining flexibility for custom workflows. Below is a detailed explanation of the approaches used.

### 4.1. Approach 1: AutoML via Colab Notebooks

To gain greater control over the AutoML process, a second approach involved writing Python scripts in Google Colab to interact with the Vertex AI APIs. These scripts automated tasks such as uploading data, initiating training jobs, and retrieving evaluation metrics. Long-running training tasks (typically lasting over two hours) were scheduled on GCP to optimize resource usage and manage workflow efficiently.

#### Benefits:

- This approach provided more control over data preprocessing, model configuration, and evaluation compared to the UI-based AutoML.
- Customizability allowed fine-tuning of AutoML parameters to better suit the specific problem.

#### Challenges:

- Scheduling long-running tasks on GCP required careful resource planning to avoid interruptions or inefficiencies.
- Triggering longrunning notebooks is not possible using cloud functions.
- The extended runtime for training necessitated optimized infrastructure and monitoring.

### 4.2. Approach 2: Custom Model Development

For greater flexibility, a custom model was developed using XGBoost. This approach involved manually training the model on the preprocessed dataset, optimizing hyperparameters through trial and error, and evaluating performance. The entire process required coding expertise to design and implement the workflow.

#### Workflow:

- The dataset was loaded from Google Cloud Storage (GCS), and preprocessing steps were applied programmatically to ensure data consistency and quality.
- XGBoost was used as the model framework, and hyperparameters were tuned manually to achieve the best possible results.
- Once trained, the model was Dockerized to create a portable, consistent runtime environment for execution.

- The Docker image was built and pushed to Google Container Registry (GCR) to facilitate easy access.
- In Google Colab, the Docker image was pulled from GCR and run to execute the custom-built machine learning code for XGBoost, ensuring scalability and reproducibility across different environments.

Drawbacks:

- The model's accuracy was limited to ~65%, which was suboptimal compared to AutoML outputs.
- There was no automated generation of performance metrics, requiring manual implementation of validation and visualization tasks.
- Dockerization and deployment added complexity to the workflow, increasing the time and resources needed for completion.

### 4.3. Approach 3: Airflow-Integrated AutoML

This approach combined the automation of AutoML with the workflow orchestration capabilities of Airflow. A new DAG was designed to handle the entire AutoML process seamlessly within Cloud Composer. DAG operators were programmed to check for existing models and initiate training only if necessary. This logic optimized resource usage by reusing pre-trained models whenever available, avoiding redundant computations.

#### Pipeline Steps:

1. **Model Check:** The DAG first checked whether a pre-trained model existed in the GCP Model Registry.
2. **Training Execution:** If no model was found, AutoML training was triggered using Vertex AI APIs using `CreateAutoMLTabularTrainingJobOperator` class
3. **Model Evaluation:** After training, the DAG integrated model evaluation metrics to select the best-performing model for deployment.

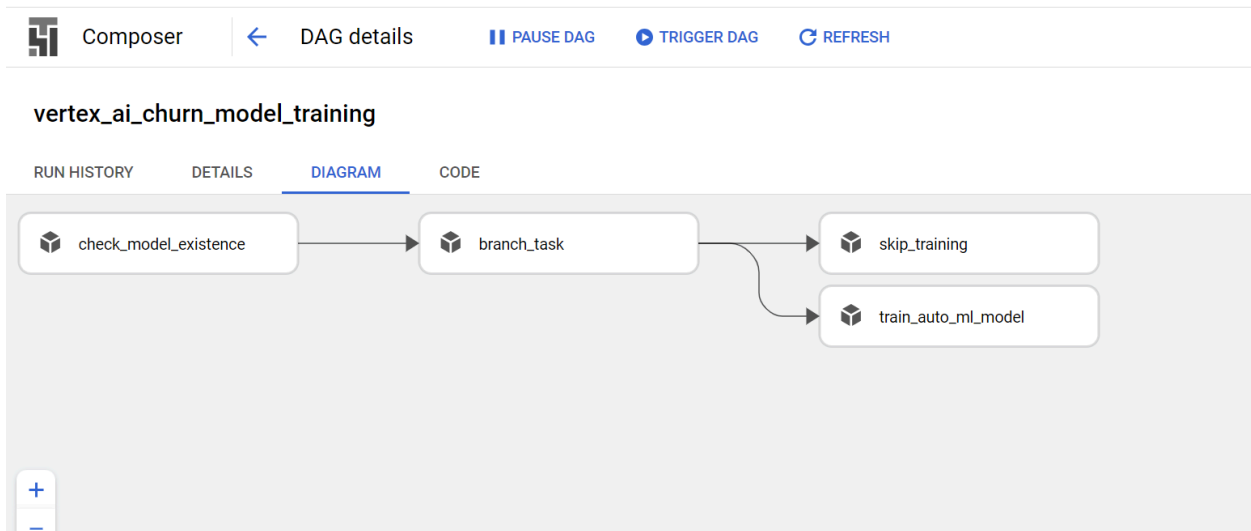
#### Advantages:

- This approach fully automated the training workflow, reducing manual intervention and improving efficiency.
- The integration of AutoML with Airflowcomposer ensured a seamless connection between preprocessing, training, and evaluation tasks.

### 4.4 Conditional Model Utilization Workflow

The workflow is designed to first check the Model Registry for an existing trained model before proceeding with any training steps. If a trained model is found in the registry, the pipeline skips the training process and directly uses the pre-existing model for predictions or further tasks. This

approach saves time and computational resources by leveraging previously validated models. However, if no trained model is available in the registry, the pipeline automatically initiates the training process. It begins by splitting the available data into training and testing subsets, ensuring an appropriate validation strategy. The model is then trained using the training data, evaluated on the test set, and registered in the Model Registry for future use, maintaining a seamless and efficient model lifecycle.



#### 4.5 Notification and alerts

Once the model is trained and batch predictions are completed, an email notification is automatically sent to confirm the successful completion of the process.

The screenshot shows the 'Evaluate' page for a model named 'churn\_model' in Version 1. The page has tabs for 'EVALUATE', 'DEPLOY & TEST', 'BATCH PREDICT', 'VERSION DETAILS', and 'LINEAGE'. Below the tabs, there is a description of pointwise evaluation and buttons for 'CREATE EVALUATION' and 'REFRESH'. A table titled 'Pointwise evaluations' displays the following data:

<input type="checkbox"/>	Name	Status	Evaluation dataset	Created	PR AUC	PROC AUC	
<input type="checkbox"/>	untitled_775753918437012...	Succeeded	—	Nov 11, 2024, 7:31:35 PM	0.777	0.787	⋮

EVALUATE

DEPLOY & TEST

BATCH PREDICT

VERSION DETAILS

LINEAGE

untitled\_77---4370120910

COMPARE

CREATE EVALUATION

Threshold target

←

Evaluation details

Filter

Confidence threshold 0.5

All labels

0.777

0.833

0.461

All labels

PR AUC0.777

ROC AUC0.787

Log loss0.557

Micro-average F10.7250704

Macro-average F10.5060168

Micro-average precision72.5%

Micro-average recall72.5%

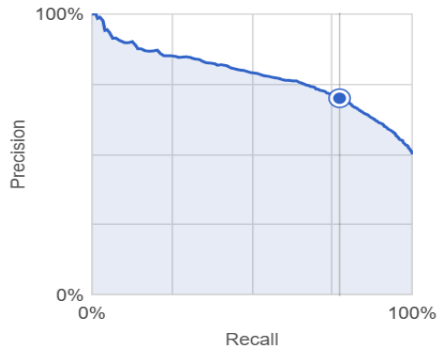
Total items0

Training items0

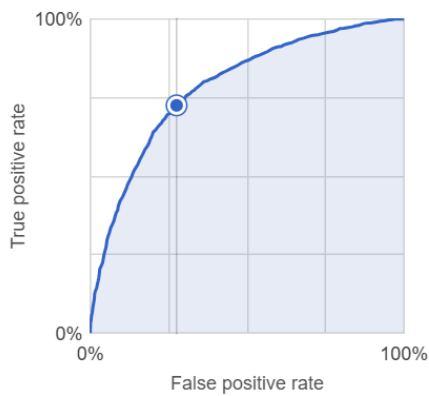
Validation items0

Test items0

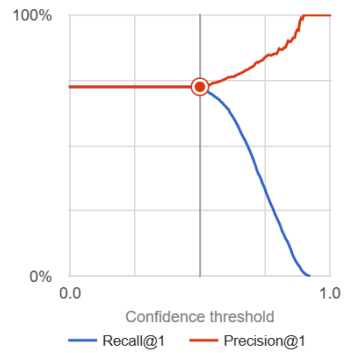
Precision-recall curve ?



ROC curve ?



### Precision-recall by threshold ?



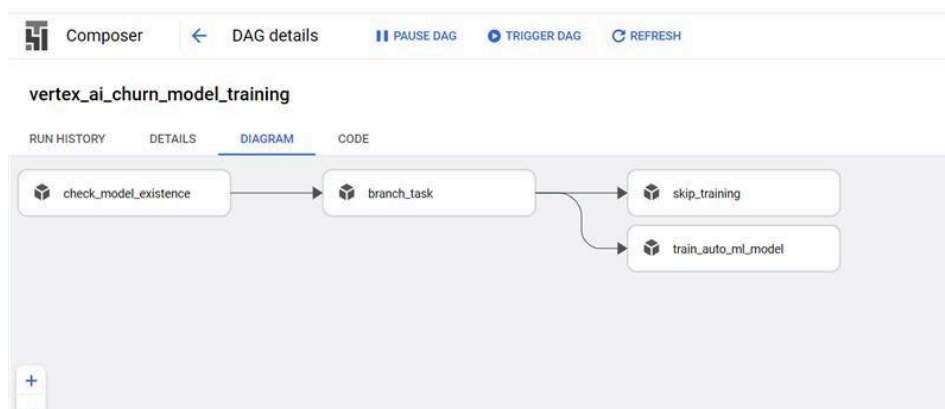
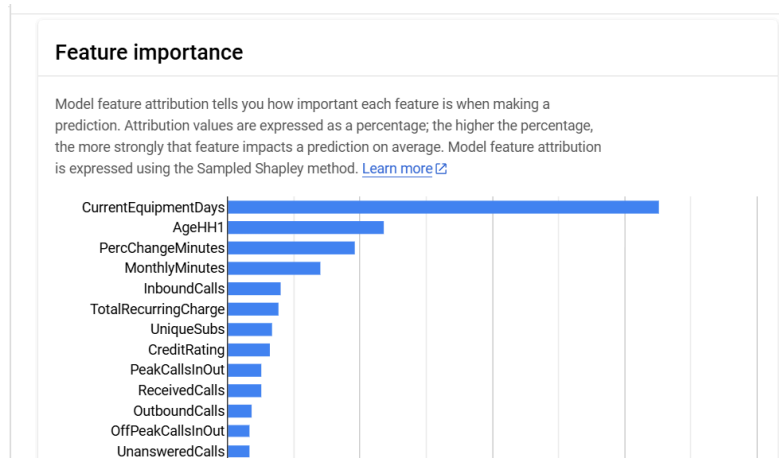
### Confusion matrix

☒ Item counts

A confusion matrix shows how the model classified each label in the evaluation dataset. The blue, bold cells indicate a correct prediction. A data item is moved to the dropped column if it does not meet the confidence threshold for any label.

True label	Predicted label		
	0	1	Dropped
0	<b>98%</b>	2%	0%
1	90%	<b>10%</b>	0%





## 5. Model Validation and Bias Detection

### 5.1. Validation

- For the validation of the dataset, AutoML handled the entire process seamlessly, ensuring a robust and efficient workflow. It performed dataset validation to check for inconsistencies, missing values, and other issues before proceeding with model training. Additionally, AutoML conducted hyperparameter tuning to optimize the model's performance by systematically adjusting parameters to achieve the best results. Once training was complete, AutoML evaluated multiple candidate models and automatically selected the best-performing model based on metrics such as precision, recall, F1 score, and ROC AUC. This comprehensive process eliminated the need for manual intervention while ensuring that the most accurate and reliable model was chosen for deployment.

The screenshot shows the BigQuery Studio interface. On the left is a navigation sidebar with 'Analysis' selected, containing links to BigQuery Studio, Data transfers, Scheduled queries, Analytics Hub, Dataform, Partner Center, and Orchestration. The main area is titled 'Viewing resources.' and shows a list of resources including 'Big\_query\_batch\_prediction' and 'predictions\_2024\_11\_13T14\_57\_56\_356Z\_511'. The 'predictions\_2024\_11\_13T14\_57\_56\_356Z\_511' resource is selected, and its 'SUMMARY' tab is active. The table preview shows the following data:

Row	MadeCallToRet...	e.CreditRating	explanatio...	outputDisplayNa...	predicted_Churn.classes	predi... scores
1	0.0	-0.00152288...	0		0	0.75395160...
2	0.0	0.03825228...	0		1	0.24604834...
3	0.0	0.06752718...	0		0	0.87575483...
4	0.0	0.01306454...	0		1	0.12424510...
5	0.0	0.0	0		0	0.86832690...
6	0.0	0.0	0		1	0.13167309...
7	0.0	0.09213446...	0		0	0.72782963...

## 5.2 Inferences after model validation

After the model was validated, a test set was used to generate model inferences in batches. These batch inferences were then efficiently processed and stored in BigQuery, Google Cloud's data warehouse, for further analysis and integration with downstream systems. By leveraging BigQuery, the predictions could be queried and analyzed at scale, ensuring that the results were both accessible and actionable for subsequent workflows.

The screenshot shows the BigQuery Studio interface. On the left is a navigation sidebar with 'Analysis' selected, containing links to BigQuery Studio, Data transfers, Scheduled queries, Analytics Hub, Dataform, Partner Center, and Release Notes. The main area is titled 'Explorer' and shows a list of resources including 'Big\_query\_batch\_prediction' and 'predictions\_2024\_11\_13T14\_57\_56\_356Z\_511'. The 'predictions\_2024\_11\_13T14\_57\_56\_356Z\_511' resource is selected, and its 'SUMMARY' tab is active. The table preview shows the following data:

Row	explanatio...	outputDisplayNa...	predicted_Churn.classes	predi... scores
1			0	0.75395160...
2			1	0.24604834...
3			0	0.87575483...
4			1	0.12424510...
5			0	0.86832690...
6			1	0.13167309...
7			0	0.72782963...

The screenshot shows the BigQuery Studio interface. On the left is a navigation sidebar with 'Analysis' selected, containing links to BigQuery Studio, Data transfers, Scheduled queries, Analytics Hub, Dataform, Partner Center, and Release Notes. The main area is titled 'Explorer' and shows a list of resources including 'Big\_query\_batch\_prediction' and 'predictions\_2024\_11\_13T14\_57\_56\_356Z\_511'. The 'predictions\_2024\_11\_13T14\_57\_56\_356Z\_511' resource is selected, and its 'SUMMARY' tab is active. The table preview shows the following data:

Row	UniqueSubs	predicted_Churn.classes	predi... scores
1	1	0	0.59071540...
2	1	1	0.40928459...
3	1	0	0.32709494...
4	1	1	0.67290508...
5	3	0	0.38333764...
6	3	1	0.61666226...
7	1	0	0.88043582...



represented distinct ranges or categories within the column, allowing us to assess whether the model's predictions varied significantly between them. The results of these bias detection analyses, including predictions and inferences for each slice, were stored in a **BigQuery** table. This enabled efficient querying and deeper exploration of potential disparities, ensuring that the model's outputs were fair and consistent across the defined groups.

Row	UniqueSubs	predicted_Churn.classes	predi... scores
1	1	0	0.51342076...
		1	0.48657914...
2	2	0	0.79082363...
		1	0.20917634...
3	1	0	0.66153705...
		1	0.33846282...
4	1	0	0.70162242...
		1	0.29837751...
5	4	0	0.15637695...
		1	0.84362304...
6	1	0	0.87552112...

## 6. Experiment Tracking and Hyperparameter Tuning

- Vertex AI efficiently managed experiment tracking and streamlined the entire machine learning workflow. It maintained comprehensive logs of hyperparameter configurations, model performance metrics, and results visualizations, such as confusion matrices and feature importance graphs, ensuring that every experiment was traceable and reproducible. Additionally, Vertex AI facilitated robust model versioning by automatically saving and organizing trained models in a centralized model registry. This registry provided easy access to past versions, enabling comparisons and ensuring reproducibility across iterations. Furthermore, Vertex AI handled artifact management by maintaining an artifacts registry, which included datasets, trained models, and other resources essential for experiment tracking and deployment. This cohesive ecosystem ensured a seamless and well-documented development process.

Model Registry

+

CREATE

📄

IMPORT

↻

REFRESH

📖

LEARN

Models are built from your datasets or unmanaged data sources. There are many different types of machine learning models available on Vertex AI, depending on your use case and level of experience with machine learning. [Learn more](#)

Region

us-central1 (Iowa)

?

≡

Filter

Enter a property name

?

☰

<input type="checkbox"/>	Name	Default version	Deployment status	Description	Type	Source	Updated ↓	
<input type="checkbox"/>	<a href="#">churn_model_1</a>	1	—	—	<div>📄</div> <div>Tabular</div>	AutoML training	Nov 14, 2024, 4:08:24 PM	⋮
<input type="checkbox"/>	<a href="#">churn_model</a>	1	—	—	<div>📄</div> <div>Tabular</div>	AutoML training	Nov 11, 2024, 7:31:45 PM	⋮

👤

Artifact Registry

🔔

⬅

Digests for xgboost-train

🗑

DELETE

📖

SETUP INSTRUCTIONS

↻

REFRESH

☑

Hide OCI alternative artifacts

≡

Filter

Enter property name or value

?

☰

<input type="checkbox"/>	Name	Description	Tags ?	Created	Updated ↓	
<input type="checkbox"/>	<a href="#">c4037c259e21</a>		latest	2 days ago	2 days ago	⋮
<input type="checkbox"/>	<a href="#">04c0546e0dc1</a>			2 days ago	2 days ago	⋮
<input type="checkbox"/>	<a href="#">dc3510df350e</a>			2 days ago	2 days ago	⋮
<input type="checkbox"/>	<a href="#">13f72d016dd7</a>			4 days ago	4 days ago	⋮
<input type="checkbox"/>	<a href="#">128adf03a481</a>			4 days ago	4 days ago	⋮
<input type="checkbox"/>	<a href="#">32b5e06e6944</a>			4 days ago	4 days ago	⋮
<input type="checkbox"/>	<a href="#">66cc01cf9e3e</a>			4 days ago	4 days ago	⋮
<input type="checkbox"/>	<a href="#">d99c79624a7b</a>			4 days ago	4 days ago	⋮
<input type="checkbox"/>	<a href="#">9569aee4df31</a>			4 days ago	4 days ago	⋮
<input type="checkbox"/>	<a href="#">5b5a5e62f90e</a>			4 days ago	4 days ago	⋮
<input type="checkbox"/>	<a href="#">70adebd8c658</a>			4 days ago	4 days ago	⋮

☰

Repositories

⚙

Settings

📖

Release Notes

<1

## Conclusion

This project successfully demonstrates the development of a fully cloud-native machine learning pipeline, eliminating the need for any local dependencies. By leveraging Google Cloud Platform (GCP) services such as Cloud Composer, Vertex AI, BigQuery, and Google Cloud Storage, all components of the pipeline—from data preprocessing to model training and deployment—were executed seamlessly on the cloud. The transition from local workflows to GCP ensured scalability, efficiency, and reproducibility, enabling the handling of large datasets and complex processes without relying on local environments.

The use of Cloud Composer to orchestrate Airflow workflows streamlined the data preprocessing and model training tasks, integrating seamlessly with GCP's managed services. Vertex AI facilitated robust model development, including AutoML and custom models, while

maintaining centralized experiment tracking, model versioning, and artifact management. The pipeline was designed to operate entirely within the cloud ecosystem, ensuring that all data, models, and results were processed, stored, and accessed through GCP services.

This cloud-centric approach highlights the scalability and flexibility of modern cloud-based workflows, showcasing the ability to manage end-to-end machine learning processes without any reliance on local systems. By fully utilizing GCP's ecosystem, this project establishes a robust, reproducible, and future-proof framework for deploying advanced machine learning models in a cloud-native environment.