

Building ML End-to-End

Google Cloud Platform

Presenter: Federico Arroyo

Federico Arroyo



- 6 years in data science consulting
- Data Science Manager @ Alvarez & Marsal
- Technical on-the-ground expertise
 - Cloud-based ML (AWS, Azure, GCP)
 - Time-Series Modelling
 - Natural Language Processing
 - MLOps implementation



Professional Experience

Companies



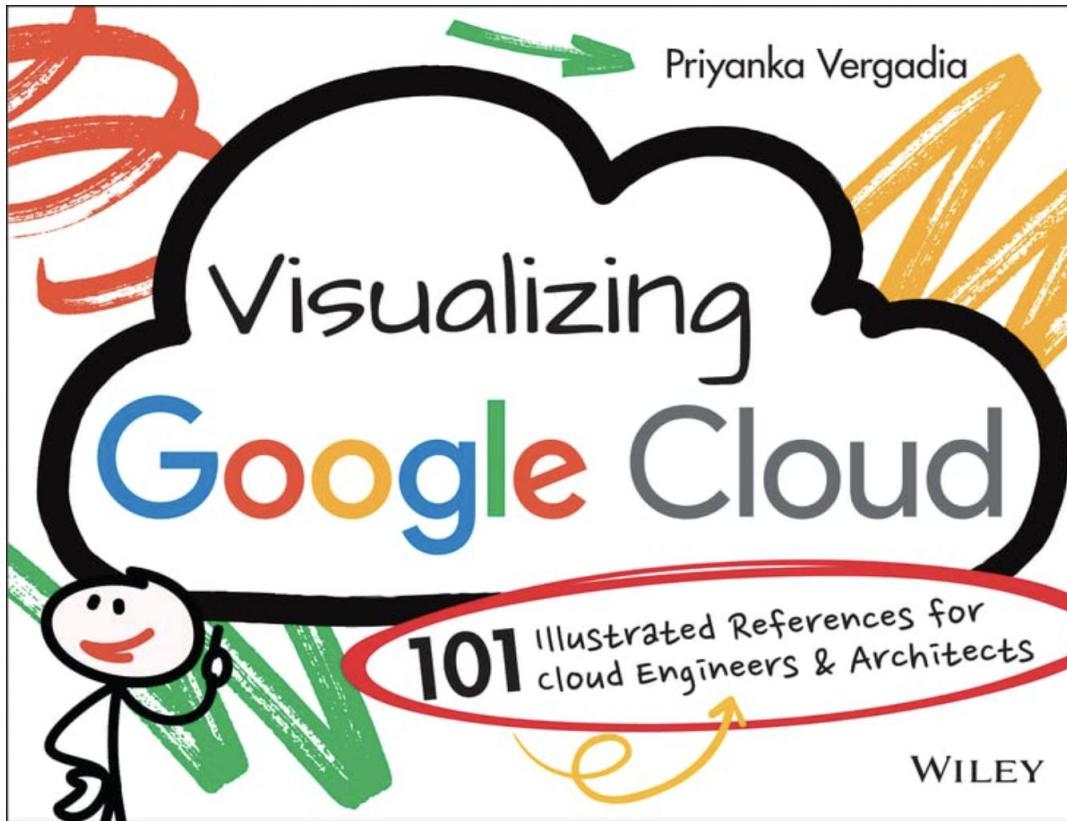
Commercial Clients



Learning Objectives for Today

1. Learn about Cloud History and its practical use cases for Machine Learning and AI applications
2. Understand the Machine Learning Lifecycle as it pertains to google cloud platform
3. Apply understanding in a live demo detecting fraud for a large financial institution using Machine Learning

Want to learn more?

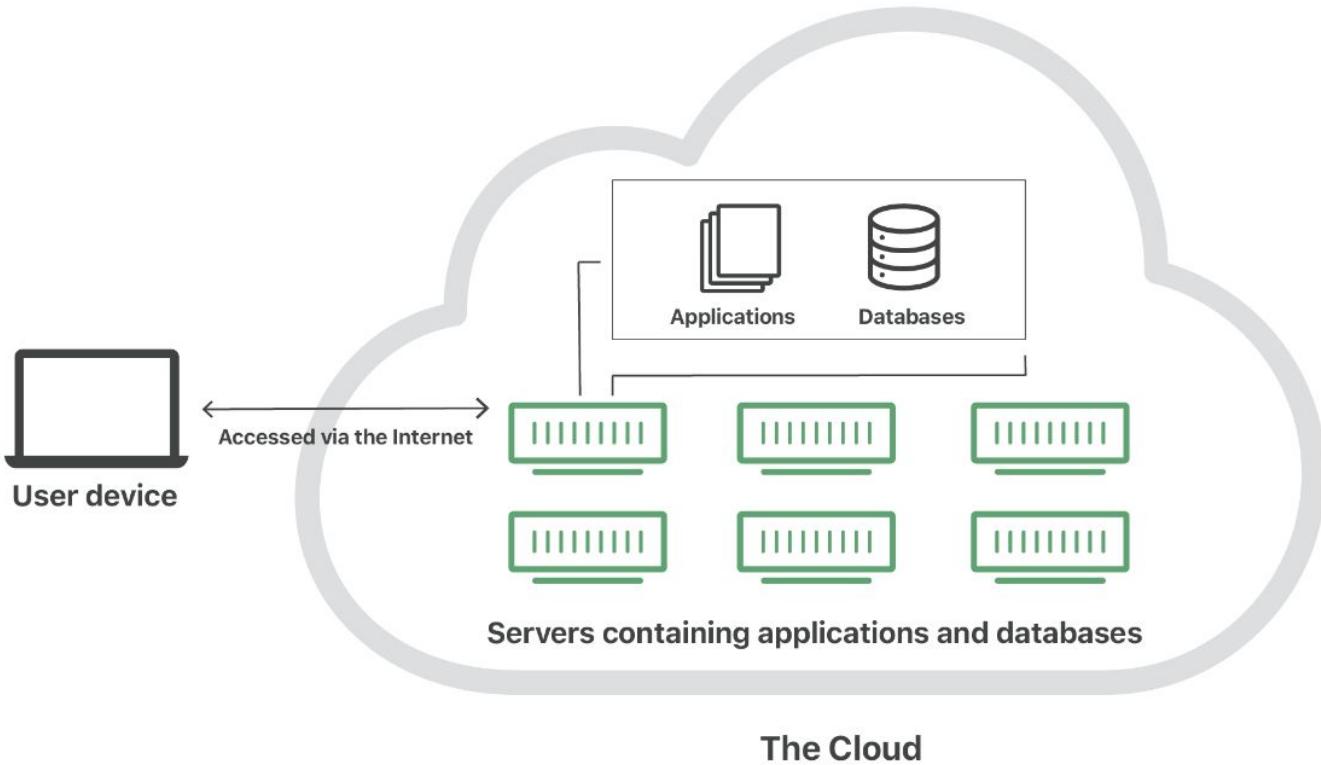




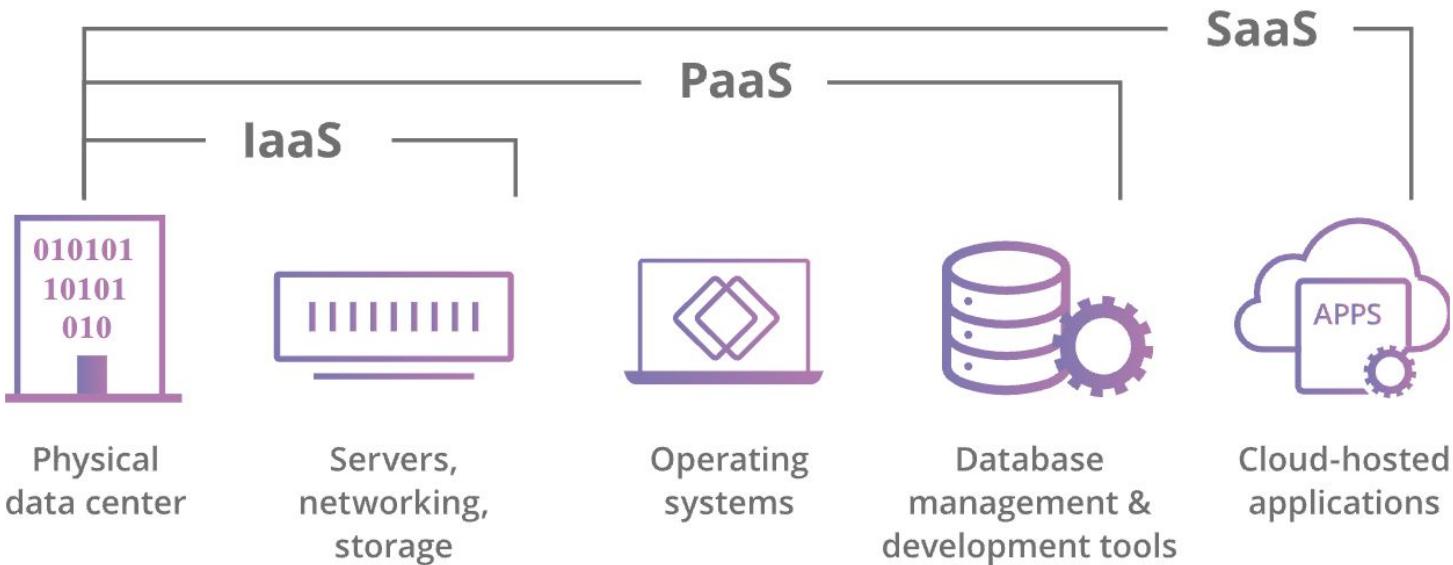
My boss told me to put my files in the cloud.

TheLifeOfDanger.com

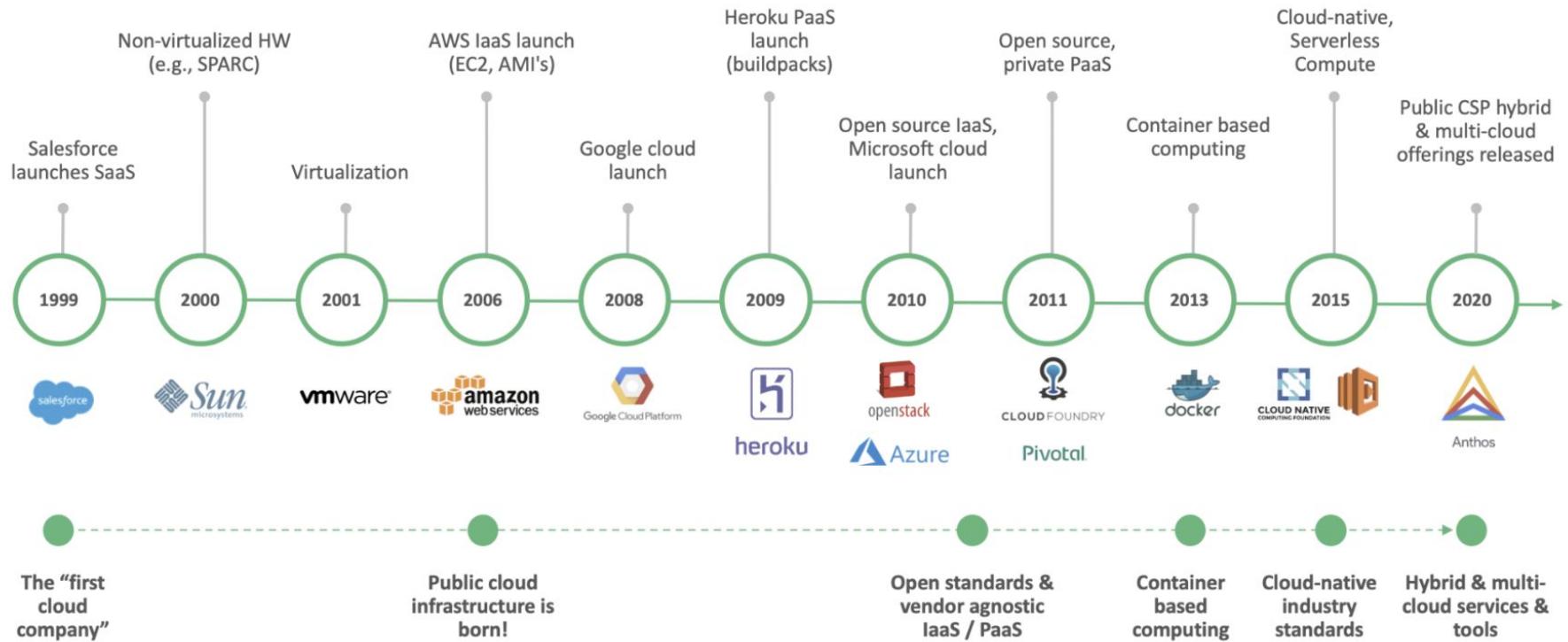
What is the cloud?



What are the main service models?



History



Source: Cloud-Native Computing Foundation: <https://www.cncf.io/blog/2017/05/15/developing-cloud-native-applications/> | Cloud Native and Container Technology Landscape

Service comparison

[LINK](#)

SERVICE TYPE	DESCRIPTION	aws	Azure	Google Cloud
STORAGE	Object storage	For storing any files you regularly use	Simple Storage Service (S3)	Blob Storage
	Archive storage	Low cost (but slower) storage for rarely used files	S3 Glacier Instant, Glacier Flexible, Glacier	Blob Cool/Cold/Archive tiers
	File storage	For storing files needing hierarchical organization	Elastic File System (EFS), FSx	Avers vEXT, Files
	Block storage	For storing groups of related files	Elastic Block Storage	Disk Storage
	Hybrid storage	Move files between on-prem & cloud	Storage Gateway	StorageSimple, Migrate
	Edge/offline storage	Move offline data to the cloud	Snowball	Data Box
	Backup	Prevent data loss	Backup	Backup and Disaster Recovery
DATABASE	Relational DB management	Standard SQL DB (PostgreSQL, MySQL, SQL Server, etc.)	Relational Database Service (RDS), Aurora	SQL, SQL Database
	NoSQL: Key-value	Redis-like DBs for semi-structured data	DynamoDB	Cosmos DB, Table storage
	NoSQL: Document	MongoDB/CouchDB-like DBs for hierarchical JSON data	DocumentDB	Cosmos DB
	NoSQL: Column store	Cassandra/HBase-like DBs for structured hierarchical data	Keyspaces	Cosmos DB
	NoSQL: Graph	Neo4j-like DBs for connected data	Neptune	N/A
	Caching	Redis/Memcached-like memory for calculation	ElastiCache	Cache for Redis, HPC Cache
	Time Series DB	DB tuned for time series data	Timestamp	Time Series Insights
	Blockchain	Dogecoin, etc.	Managed Blockchain	Blockchain Service, Blockchain Workbench, Confidential Ledger

Questions about Cloud?



Let's jump into GCP!



Google Cloud Platform

How to get started?

90-day, \$300 Free Trial: New Google Cloud and Google Maps Platform users can take advantage of a 90-day trial period that includes \$300 in free Cloud Billing credits to explore and evaluate Google Cloud and Google Maps Platform products and services. You can use these credits toward one or a combination of products.

What do I get?

If you upgrade before the trial is over: Any remaining, unexpired Free Trial credits remain in your Cloud Billing account. You can continue to use the resources you created during the Free Trial without interruption.

For resources you use in excess of what's covered by any remaining credit, your form of payment on file is charged (credit card or bank account).

If you upgrade within 30 days of the end of the trial: Your resources are marked for deletion, but you might be able to recover them. [Learn more about data deletion on Google Cloud](#).

If you upgrade more than 30 days after the end of the trial, your Free Trial resources are lost.

How to get started?

90-day, \$300 Free Trial: New Google Cloud and Google Maps Platform users can take advantage of a 90-day trial period that includes \$300 in free Cloud Billing credits to explore and evaluate Google Cloud and Google Maps Platform products and services. You can use these credits toward one or a combination of products.

\$300



90 days

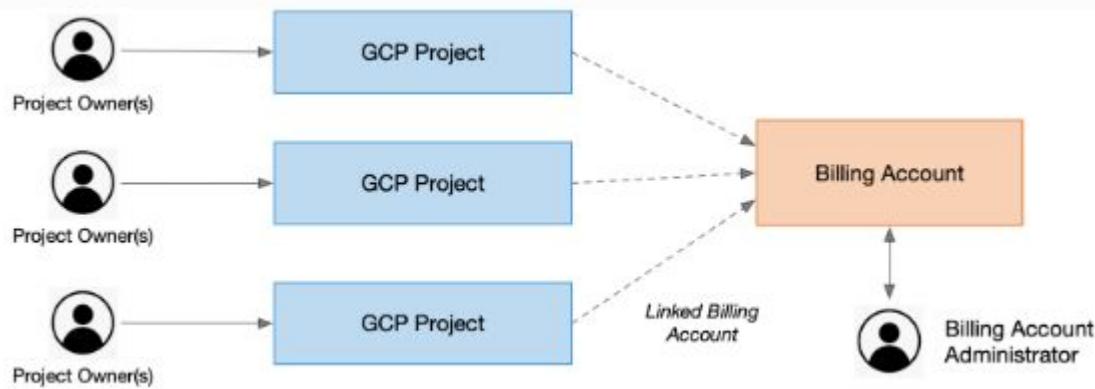
If you upgrade before the trial is over: Any remaining, unexpired Free Trial credits remain in your Cloud Billing account. You can continue to use the resources you created during the Free Trial without interruption.

For resources you use in excess of what's covered by any remaining credit, your form of payment on file is charged (credit card or bank account).

If you upgrade within 30 days of the end of the trial: Your resources are marked for deletion, but you might be able to recover them. [Learn more about data deletion on Google Cloud.](#)

If you upgrade more than 30 days after the end of the trial, your Free Trial resources are lost.

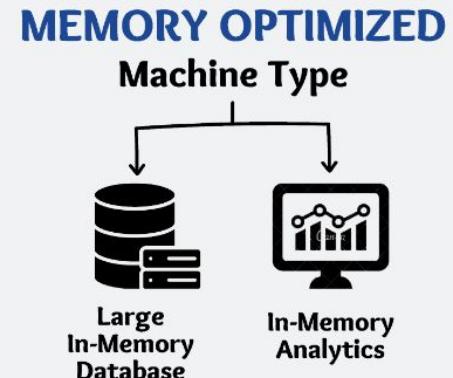
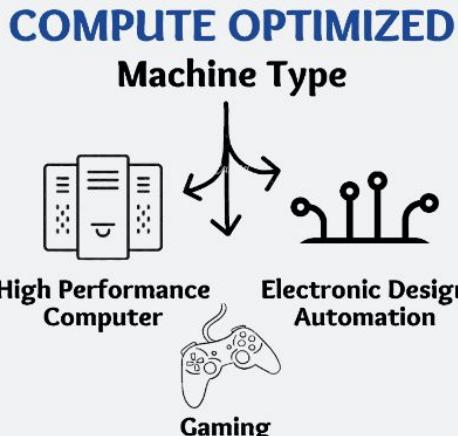
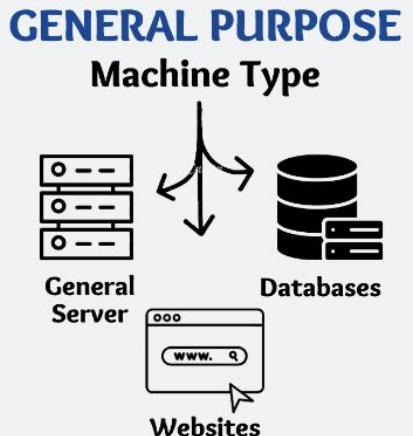
Billing



Compute

CATEGORY	PURPOSE
Standard	Balanced between processing power and memory. Fits most common application needs
High-Memory	Emphasis is put on memory over processing power for tasks that need accessible non-disk storage quickly
High-CPU	Higher CPU usage for high-intensity applications that require processing over memory
Shared-core	A single virtual CPU, backed by a physical CPU, that can run for a period of time. These machines are not for use cases that require an ongoing server or significant power. The micro shared-core machine also provides bursting capability when the virtual CPU requires more power than the single physical core. Bursting is for a short, intermittent period based on need.

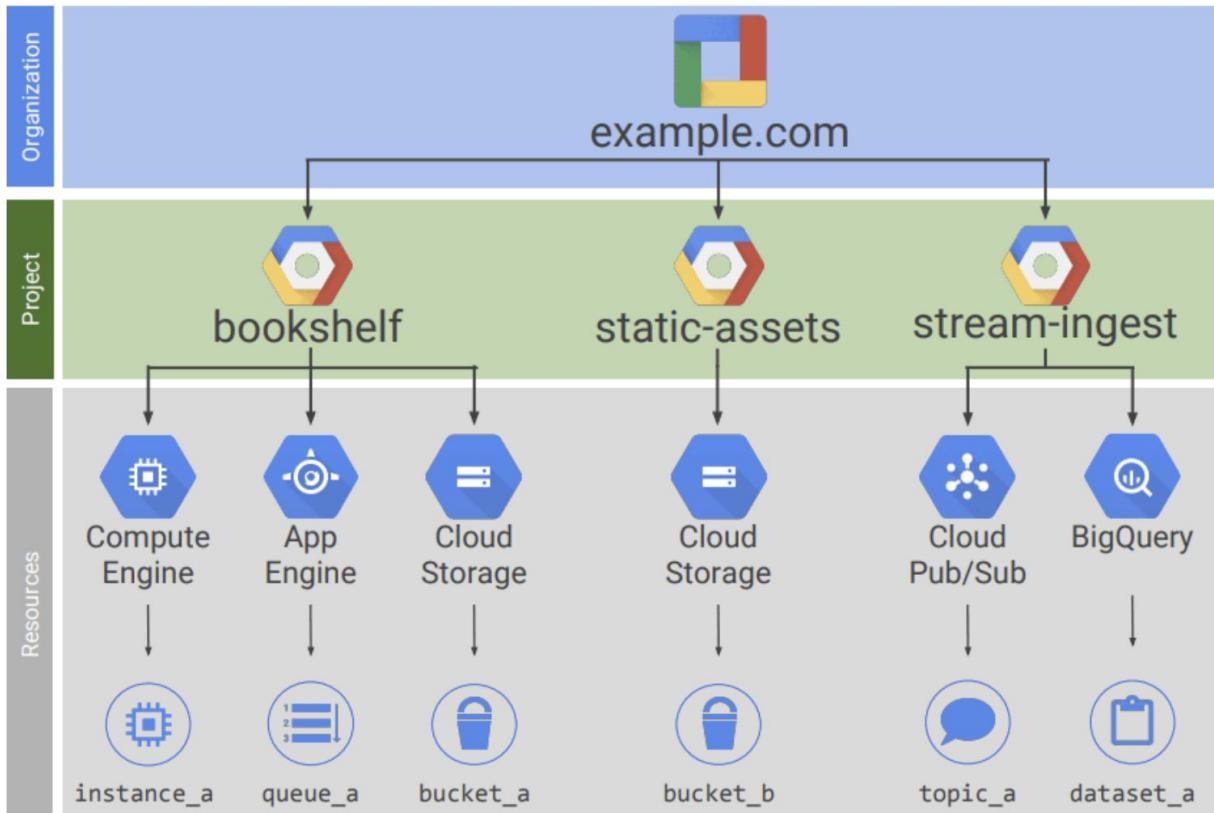
THERE ARE **3** MACHINE TYPE FAMILY



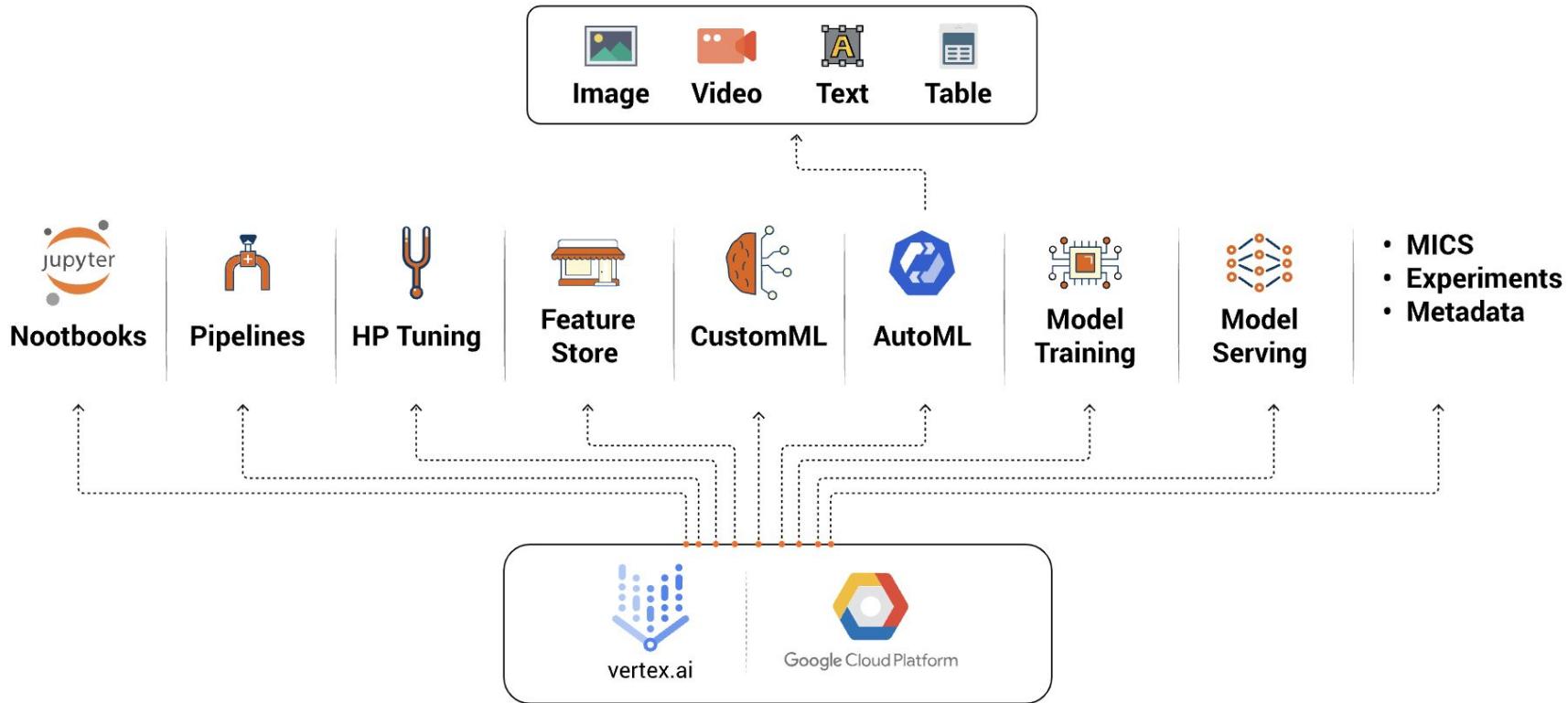
Storage

Cloud Datastore		Bigtable	Cloud Storage	Cloud SQL	Cloud Spanner	BigQuery
Type	NoSQL document	NoSQL Wide Column	Blob Storage	Relational SQL - OLTP	Relational SQL - OLTP	Relational SQL - OLAP
Transactions	Yes	Single-Row	No	Yes	Yes	No
Complex Queries	No	No	No	Yes	Yes	Yes
Capacity	Terabytes	Petabytes	Petabytes	Upto ~10TB	Petabytes	Petabytes
Unit Size	1MB/entity	~10 MB/cell ~100 MB/row	5TB per object	Depends on DB Engine	10,240 MiB/row	19 MB/row
Best for	Getting started, App Engine application	Flat data, heavy read/write events, analytical data	Structured or unstructured binary or object data	Web frameworks, existing applications	Large-scale database applications	Interactive querying, offline analytics

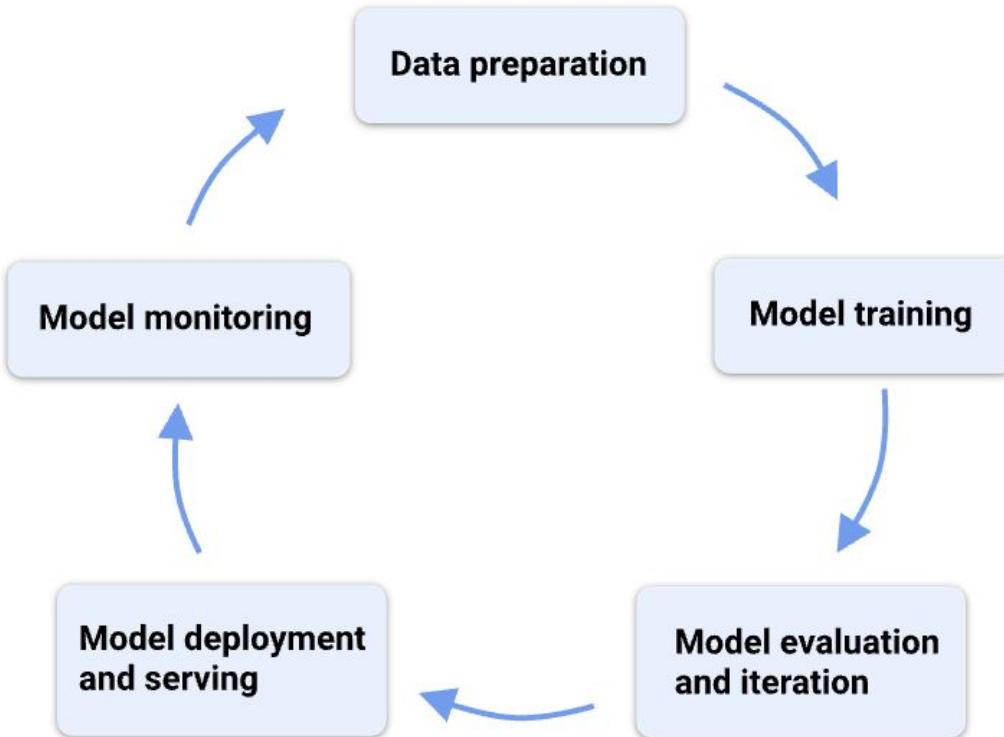
GCP Hierarchy



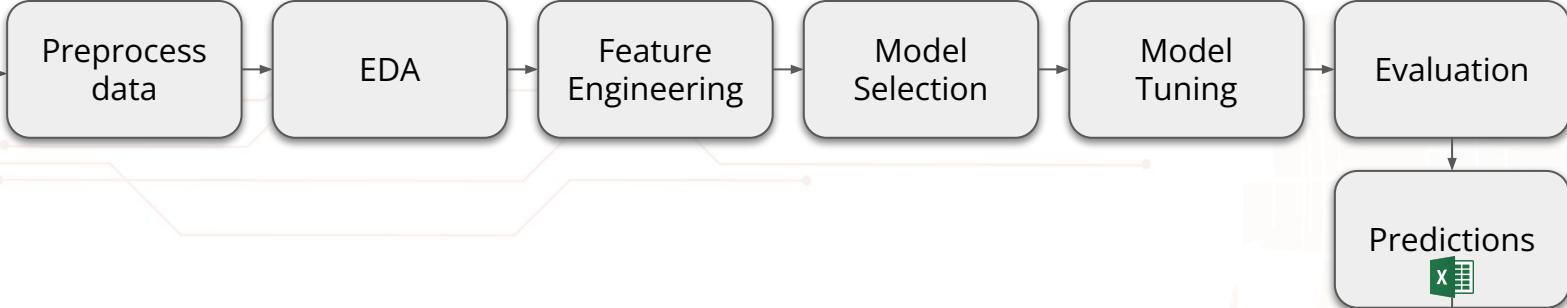
Vertex AI



Machine learning workflow



Machine Learning Lifecycle



The very beginning...



Caveat

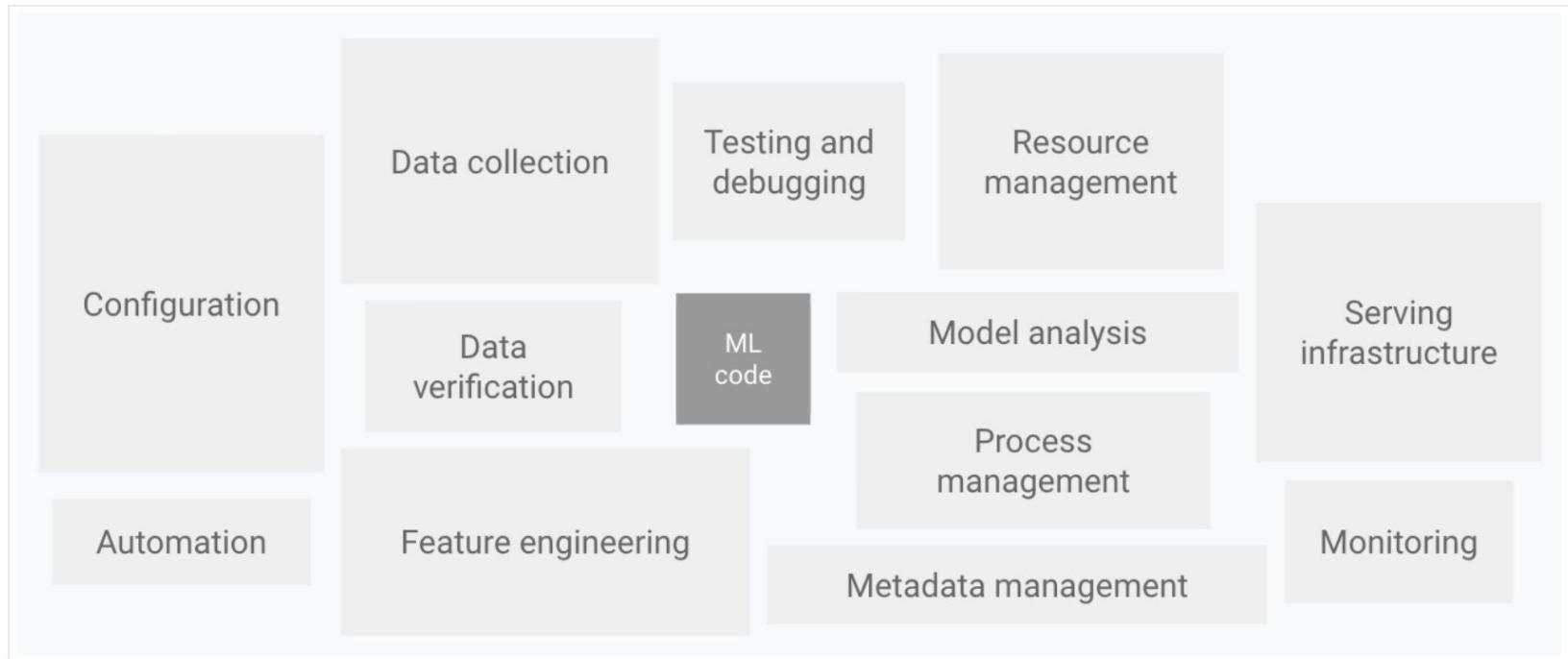
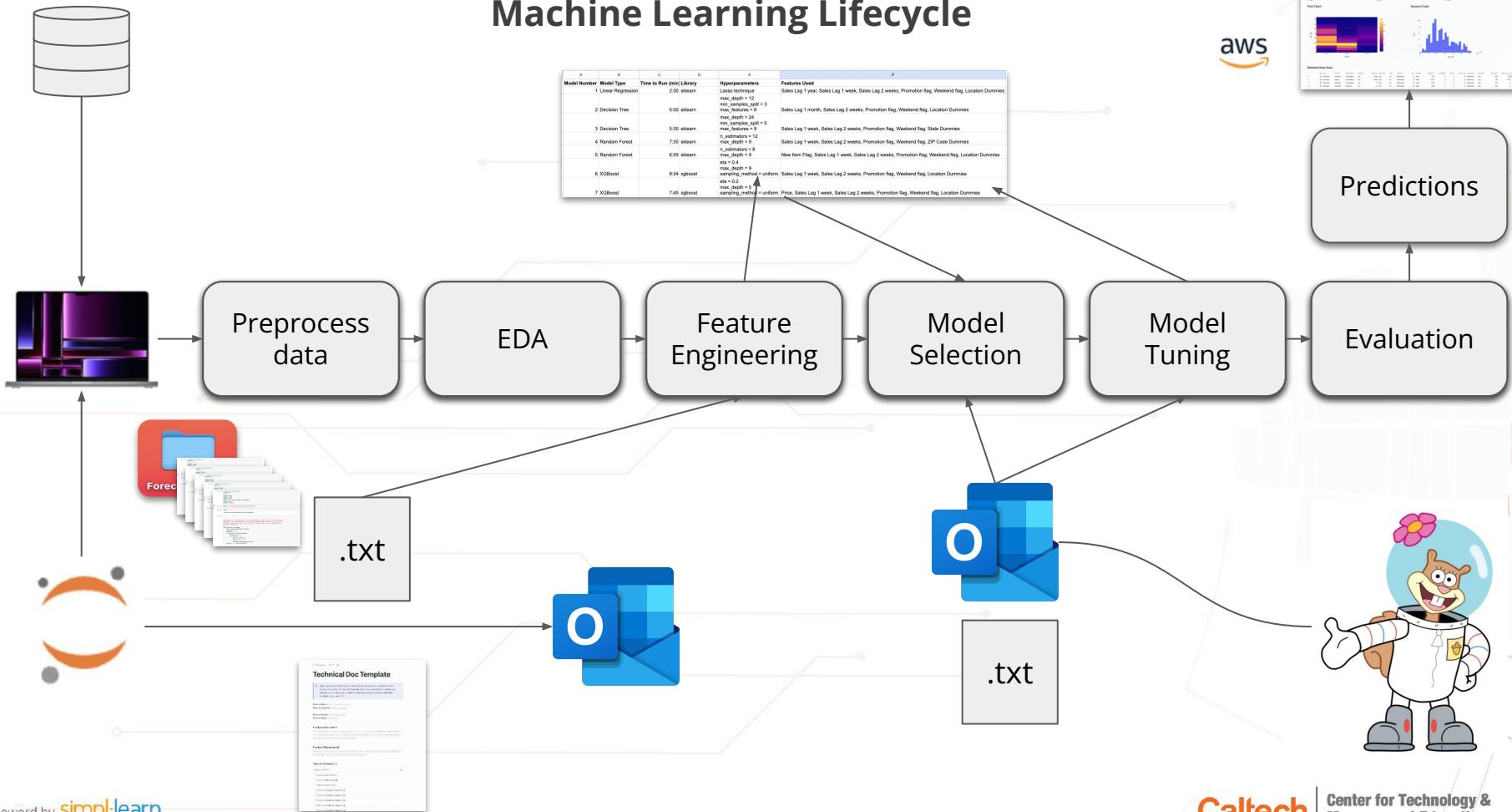


Figure 1. Elements for ML systems. Adapted from [Hidden Technical Debt in Machine Learning Systems](#).

Machine Learning Lifecycle

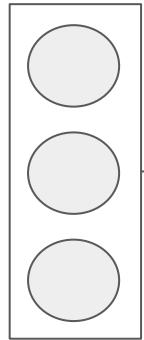


ML lifecycle that only focuses on modelling



MLOps Pipeline

Data Sources



ETL

Modelling Lifecycle

Model

CICD Stage: Build, Test, Package, Deploy Pipelines

Container Registry

Automated Pipelines

Feature Store

Data Engineering

ML Model Engineering

Model Registry

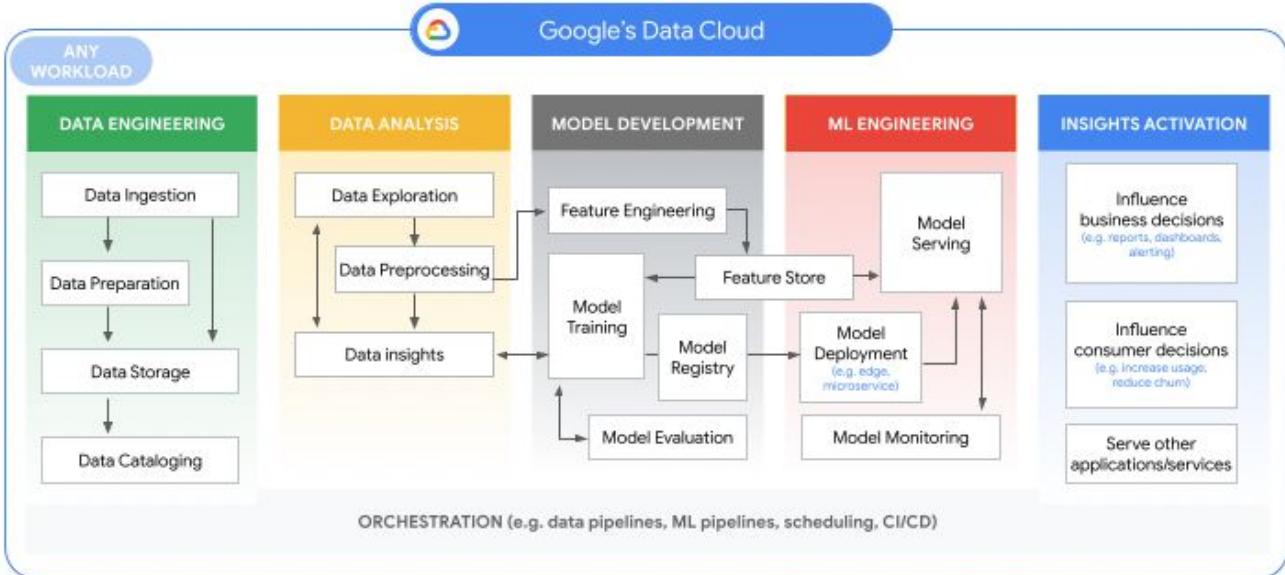
CD Stage:
ML Model Serving

ML Prediction Service

Trigger

Performance Monitoring

Data to AI on Google Cloud



Google Cloud

LINK

Want to learn more?



StatMike

@statmike-channel · 1.8K subscribers · 12 videos

Hi, I'm Mike 🙌 >

github.com/statmike and 1 more link

Subscribed

Home Videos Playlists Community

For You

Vertex AI Environment Setup 24:22

Environment Setup - Vertex AI for ML Operations [notebook 00]
14K views · 1 year ago

BigQuery Machine Learning Vertex AI End-To-End with SQL 54:07

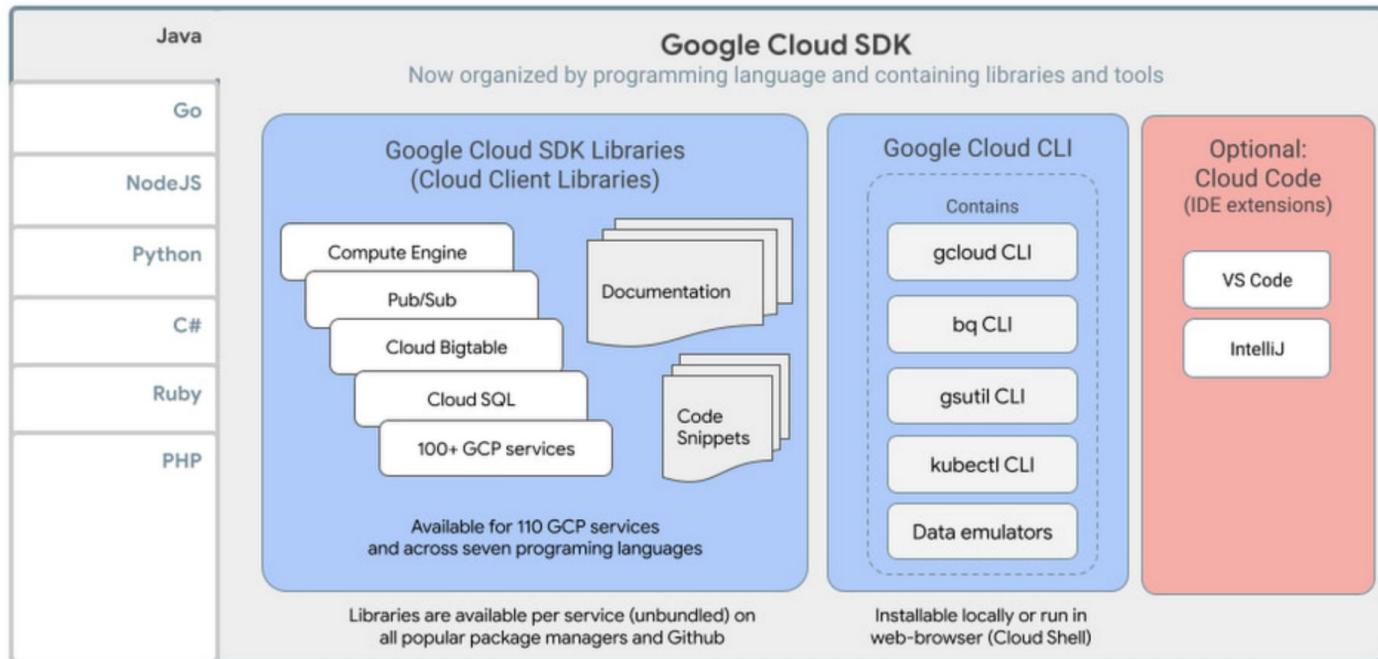
End-To-End: ML with SQL in BigQuery (BQML) [notebook 03a]
3.5K views · 1 year ago

Vertex AI End-To-End: Pipeline Orchestration 57:49

End-To-End: Pipeline Orchestration (KFP) - AutoML in Vertex AI for ML Operations [notebook 02c]
6.6K views · 1 year ago

Cloud SDK

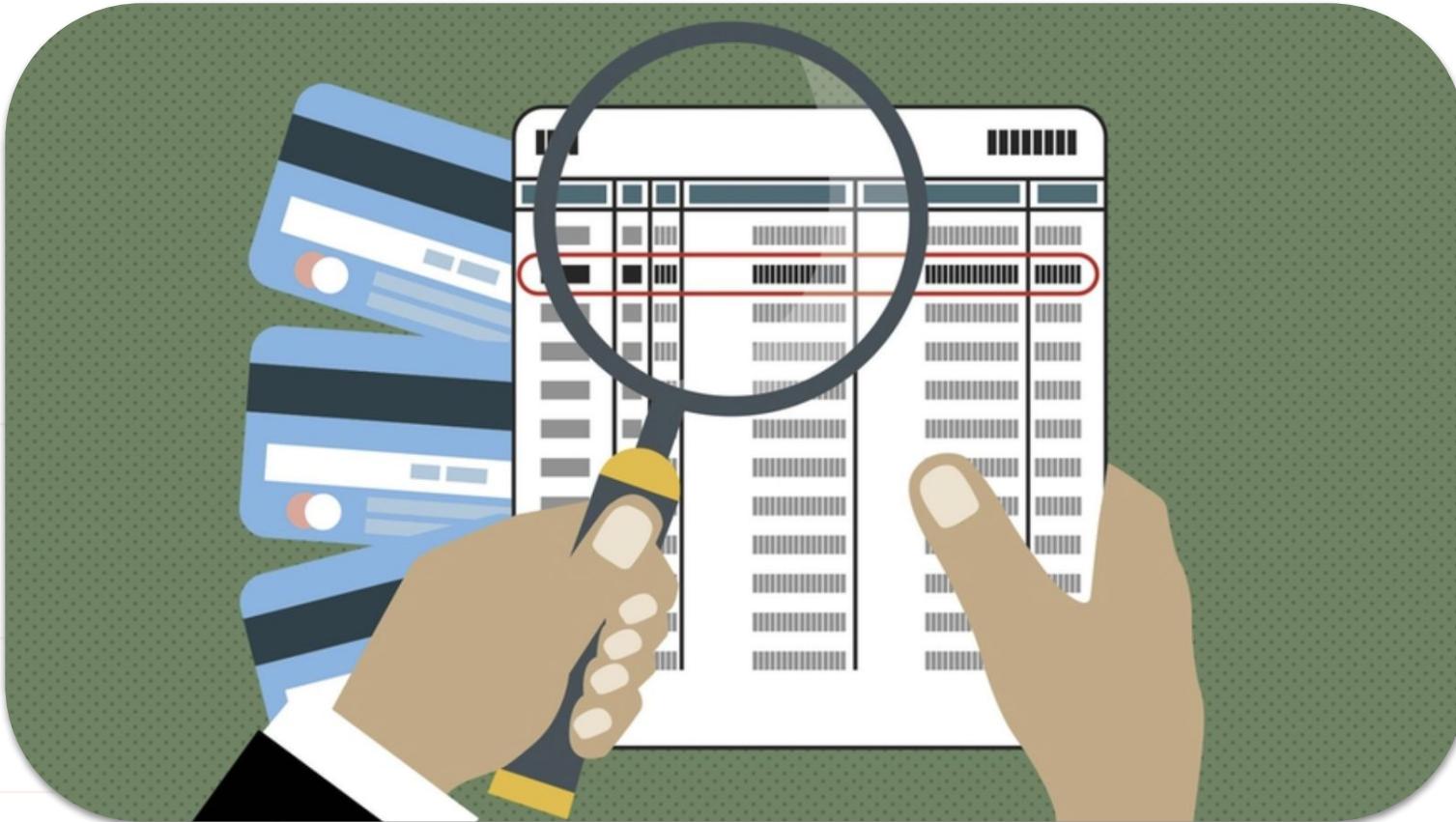
[LINK](#) (how to)



Google API Client Libraries

Now de-emphasized for Google Cloud but will remain available and continually updated to support legacy projects.

Today's example



Spinning up Workbench (notebook)

Need a [guide](#)?

← Deploy to notebook

To deploy this file, create a notebook with the required environment.

Notebook name *
tensorflow-2-11-20231030-175631

Must start with a letter followed by up to 62 lowercase letters, numbers, or hyphens (-) and cannot end with a hyphen

Region *
us-central1 (Iowa) ▾ ?

Zone *
us-central1-a ▾ ?

Notebook properties 

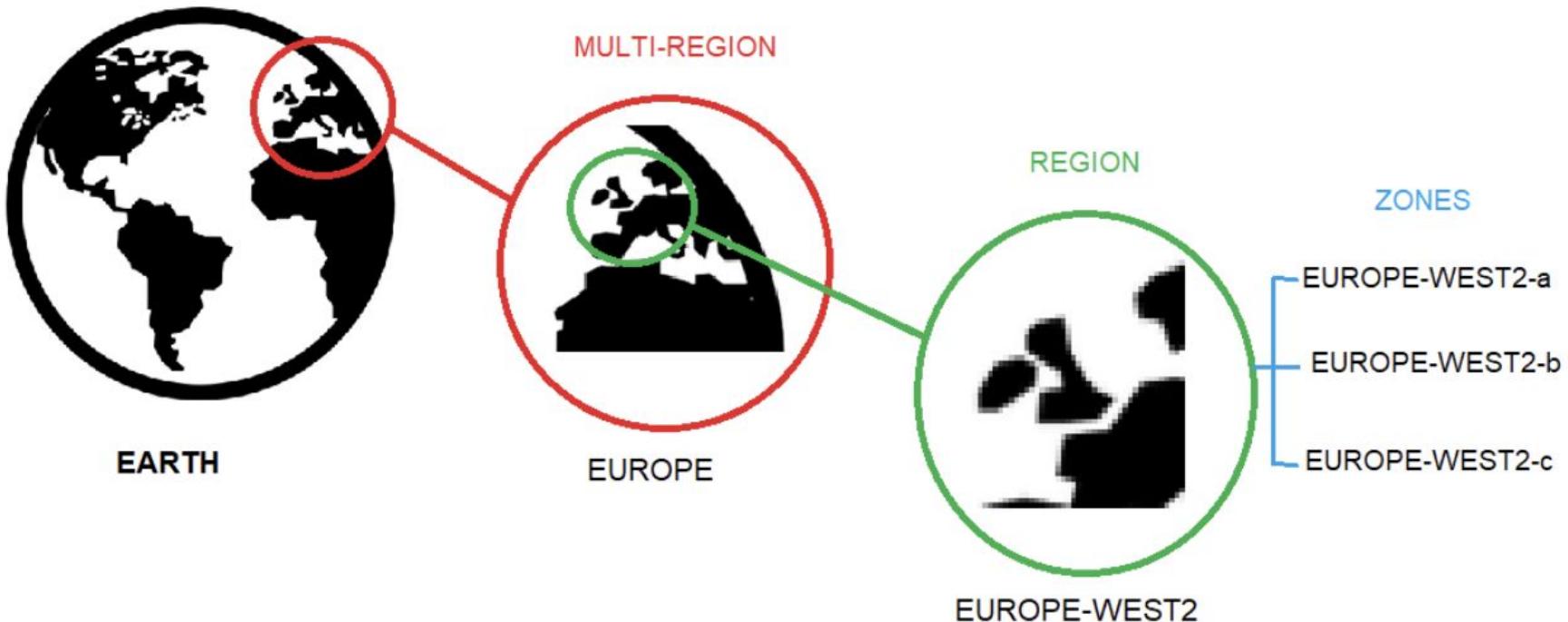
Operating System	Debian 11
Environment 	TensorFlow Enterprise 2.11 (Intel® MKL-DNN/MKL)
Machine type	4 vCPUs, 16 GB RAM
Boot disk	100 GB Standard persistent disk
Data disk	100 GB Standard persistent disk
Subnetwork	default(10.128.0.0/20) ▾
Permission	Compute Engine default service account
Estimated cost 	\$117.60 monthly, \$0.16 hourly

ADVANCED OPTIONS

CANCEL

CREATE

Regions / Zones

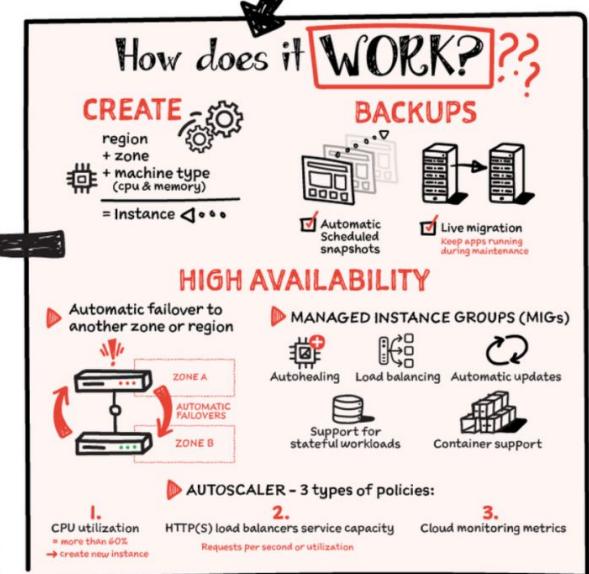
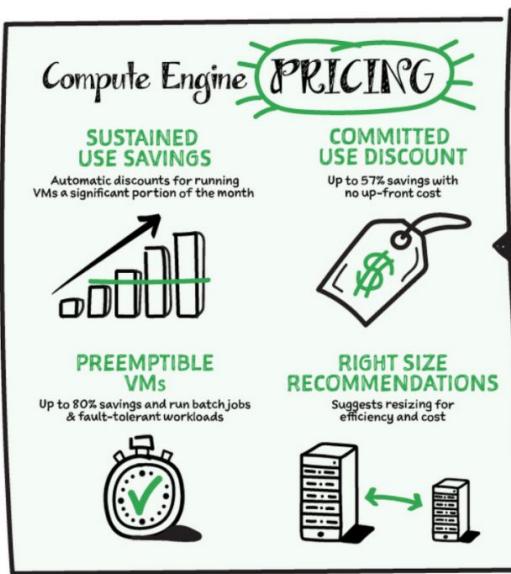
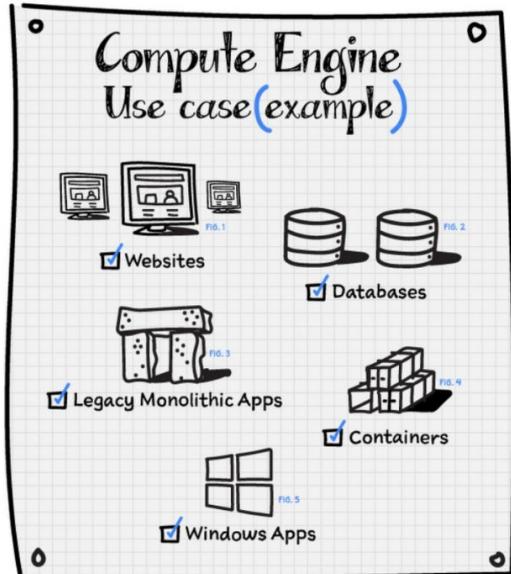
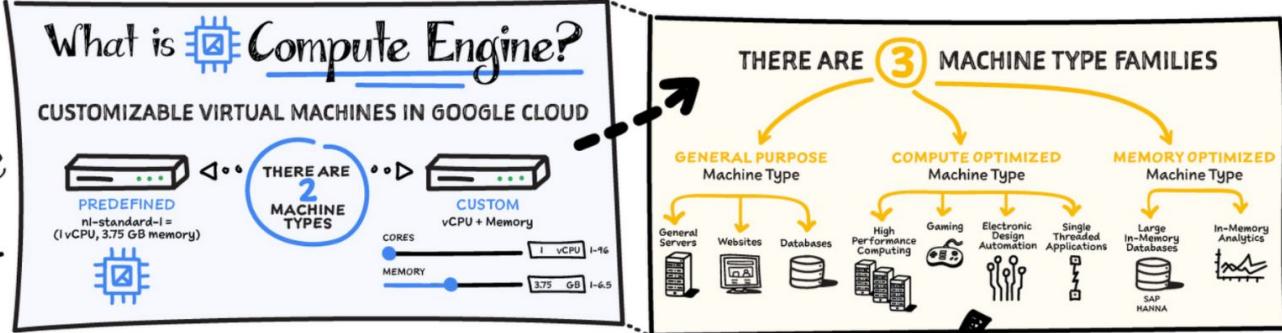


Compute Engine

#GCPsketchnotes



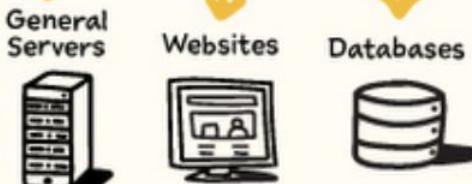
@PVERGADIA @THECLOUDGIRL.DEV



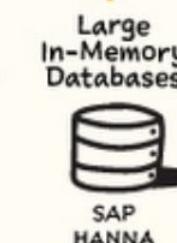
Machine Families

THERE ARE **3** MACHINE TYPE FAMILIES

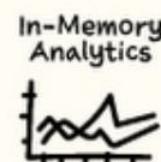
GENERAL PURPOSE
Machine Type



COMPUTE OPTIMIZED
Machine Type



MEMORY OPTIMIZED
Machine Type



Machine Types

		Workload type				
		General-purpose workloads		Compute-optimized	Memory-optimized	Accelerator-optimized
E2	N2, N2D, N1	C3, C3D	Tau T2D, Tau T2A	H3, C2, C2D	M3, M2, M1	A2, G2
Day-to-day computing at a lower cost	Balanced price/performance across a wide range of machine types	Consistently high performance for a variety of workloads	Best per-core performance/cost for scale-out workloads	Ultra high performance for compute-intensive workloads	Highest memory to compute ratios for memory-intensive workloads	Optimized for accelerated high performance computing workloads

INSTANCES

EXECUTIONS

SCHEDULES

View:

INSTANCES

USER-MANAGED NOTEBOOKS

MANAGED NOTEBOOKS

Notebooks have JupyterLab 3 pre-installed and are configured with GPU-enabled machine learning frameworks. [Learn more](#)

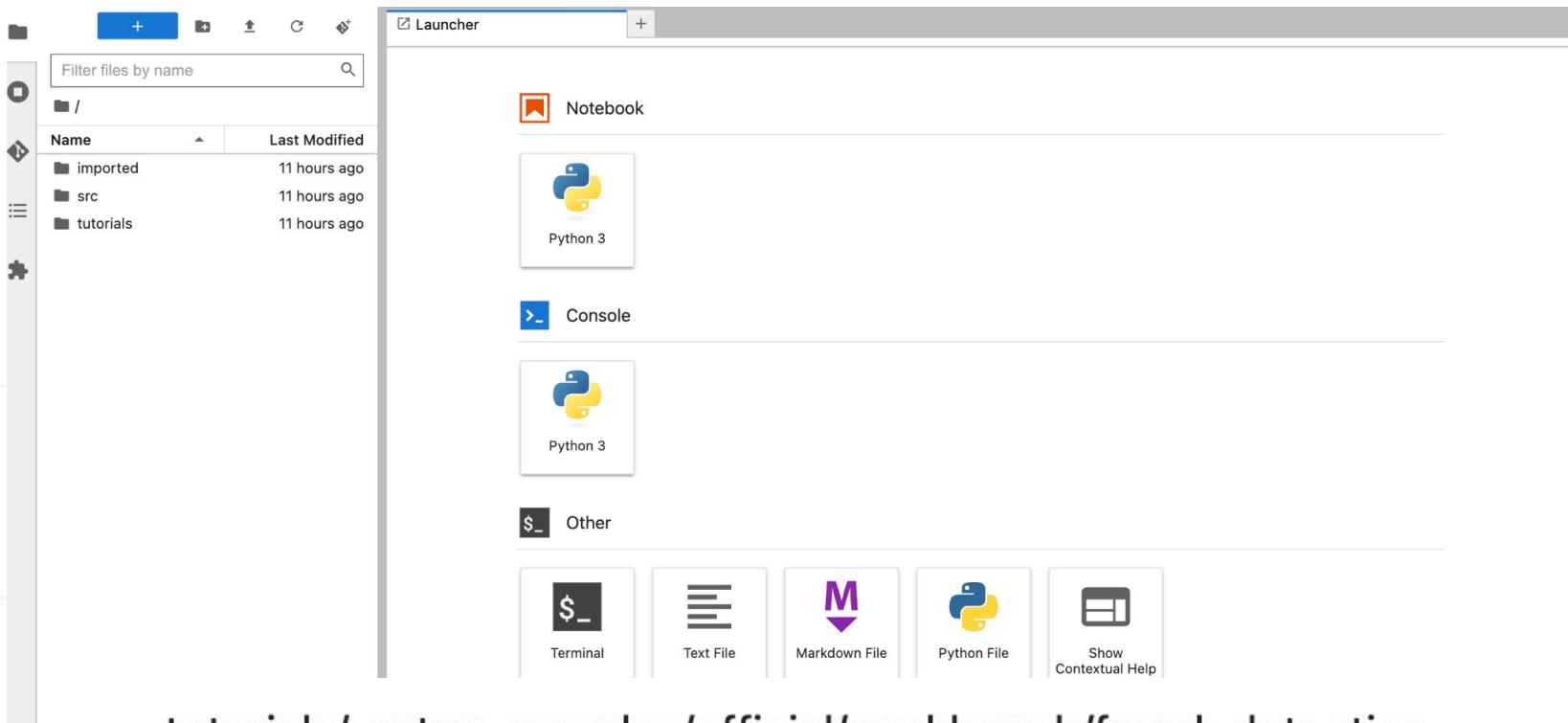
 Filter

<input type="checkbox"/>	<input checked="" type="radio"/>	Notebook name 	Zone	Auto upgrade	Environment	Machii
<input type="checkbox"/>		tensorflow-2-11-20231030-175631	OPEN JUPYTERLAB	us-central1-a	TensorFlow:2.11	Efficie Instan vCPUs RAM

Wait for instance to spin up (Water Break)



Jupyterlab



tutorials/vertex_samples/official/workbench/fraud_detection

Installation

Install the following packages required to execute this notebook.

```
! pip3 install --upgrade --quiet google-cloud-aiplatform \
               witwidget \
               fsspec \
               gcsfs
! pip3 install --quiet scikit-learn==1.2 \
               protobuf==3.20.1
```

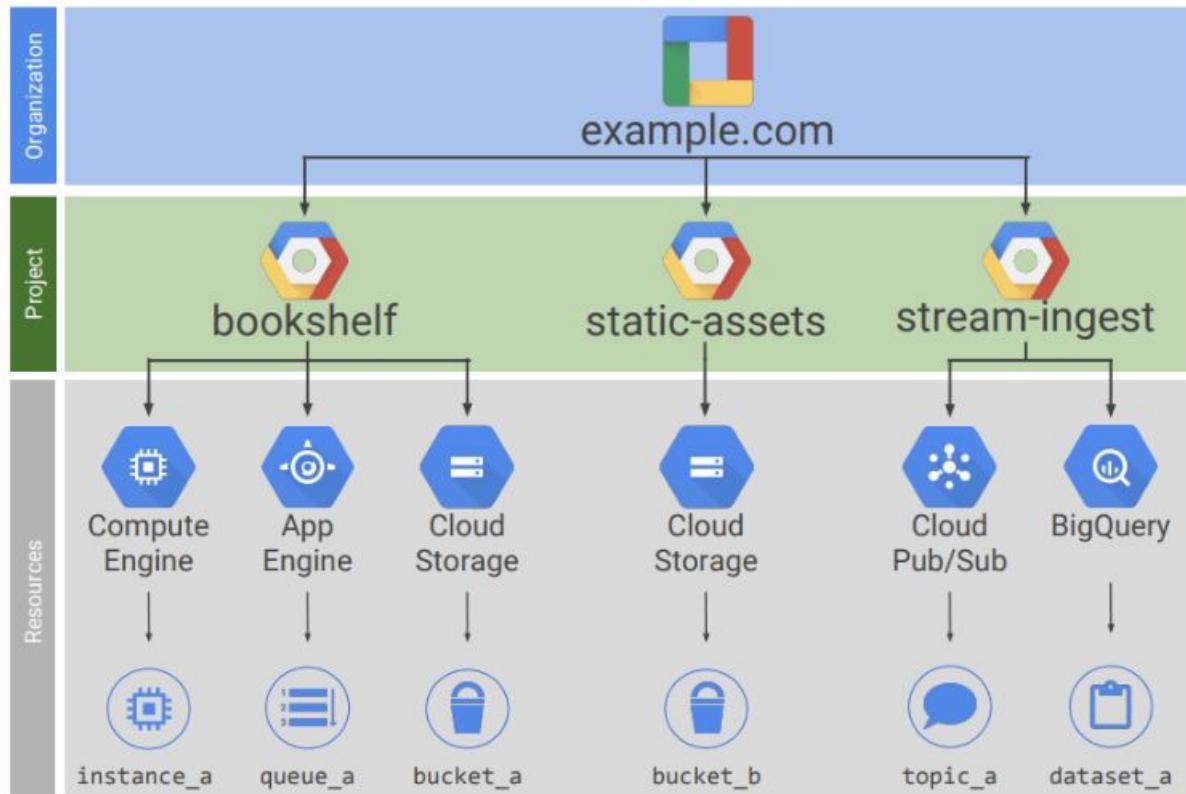
Before you begin

Set up your Google Cloud project

The following steps are required, regardless of your notebook environment.

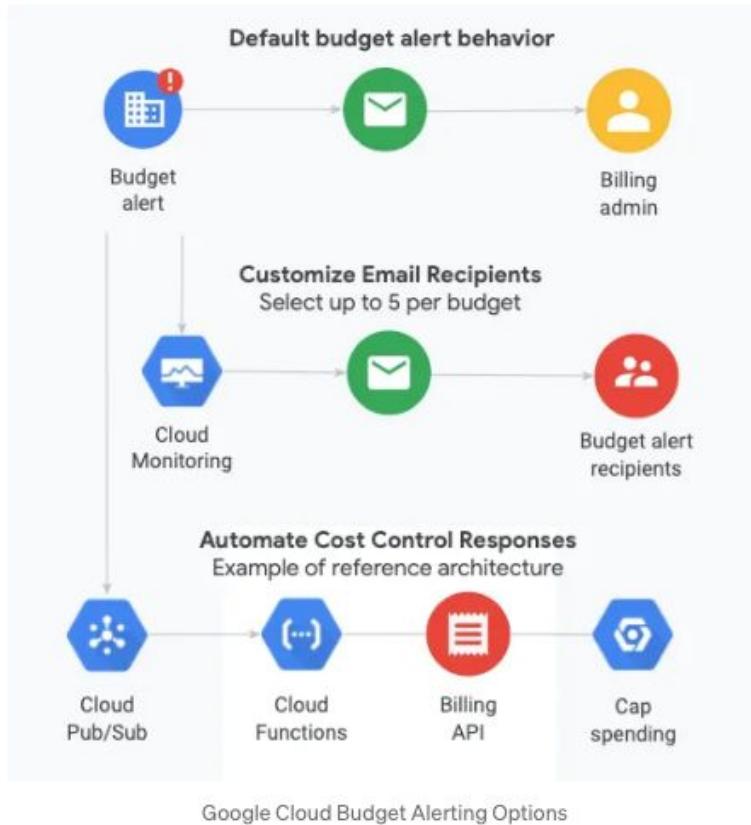
1. [Select or create a Google Cloud project](#). When you first create an account, you get a \$300 free credit towards your compute/storage costs.
2. [Make sure that billing is enabled for your project](#).
3. [Enable the Vertex AI API](#).
4. If you are running this notebook locally, you need to install the [Cloud SDK](#).

Creating a Google Project



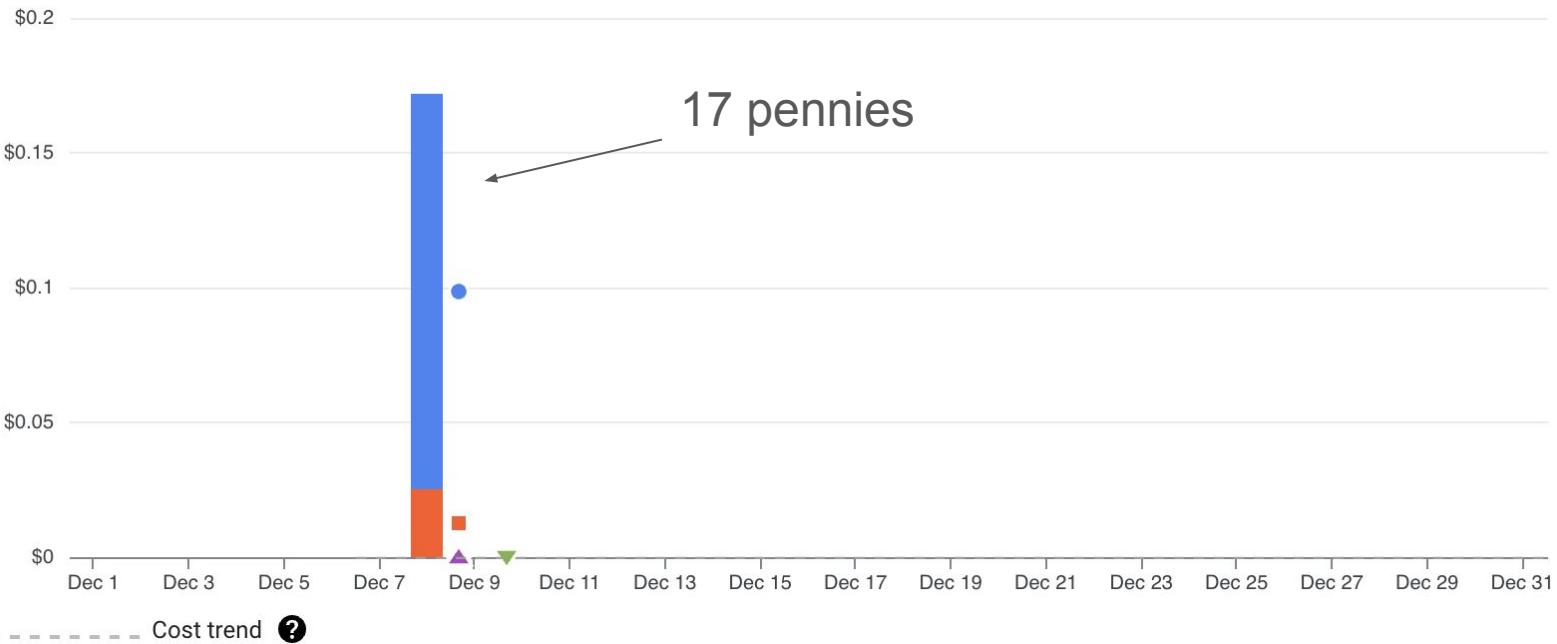
Enable Billing

Create billing alerts [tutorial](#)

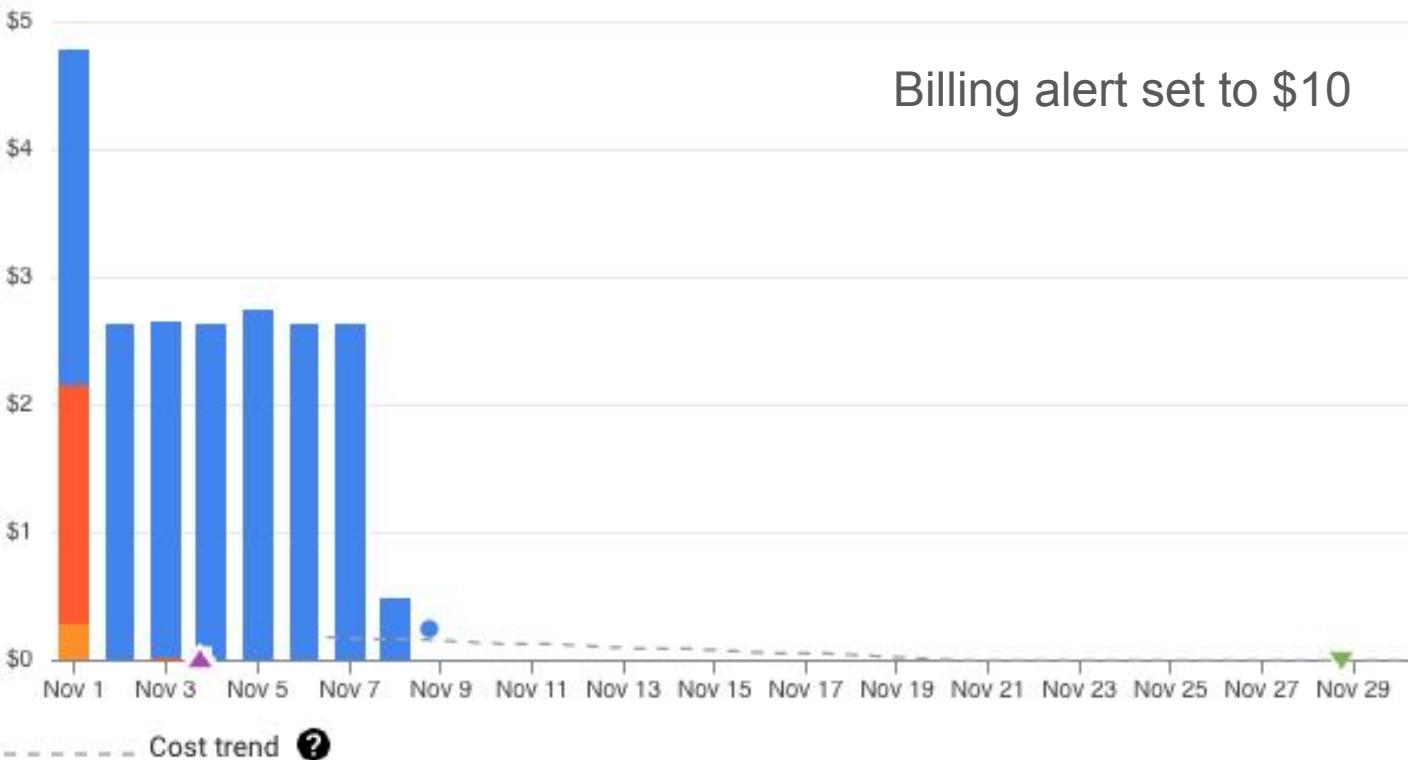


I do mine @ ~\$10

Cost of running this lab



Cost if you don't turn things off



Enable APIs

Get started with Vertex AI

Vertex AI empowers machine learning developers, data scientists, and data engineers to take their projects from ideation to deployment, quickly and cost-effectively. [Learn more](#)

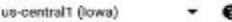
Try an interactive tutorial to learn how to train, evaluate, and deploy a Vertex AI AutoML or custom-trained model

[VIEW TUTORIALS](#)

[ENABLE VERTEX AI API](#)

Region

us-central1 (Iowa)



Prepare your training data

Collect and prepare your data, then import it into a dataset to train a model

[+ CREATE DATASET](#)

Train your model

Train a best-in-class machine learning model with your dataset. Use Google's AutoML, or bring your own code.

[+ TRAIN NEW MODEL](#)

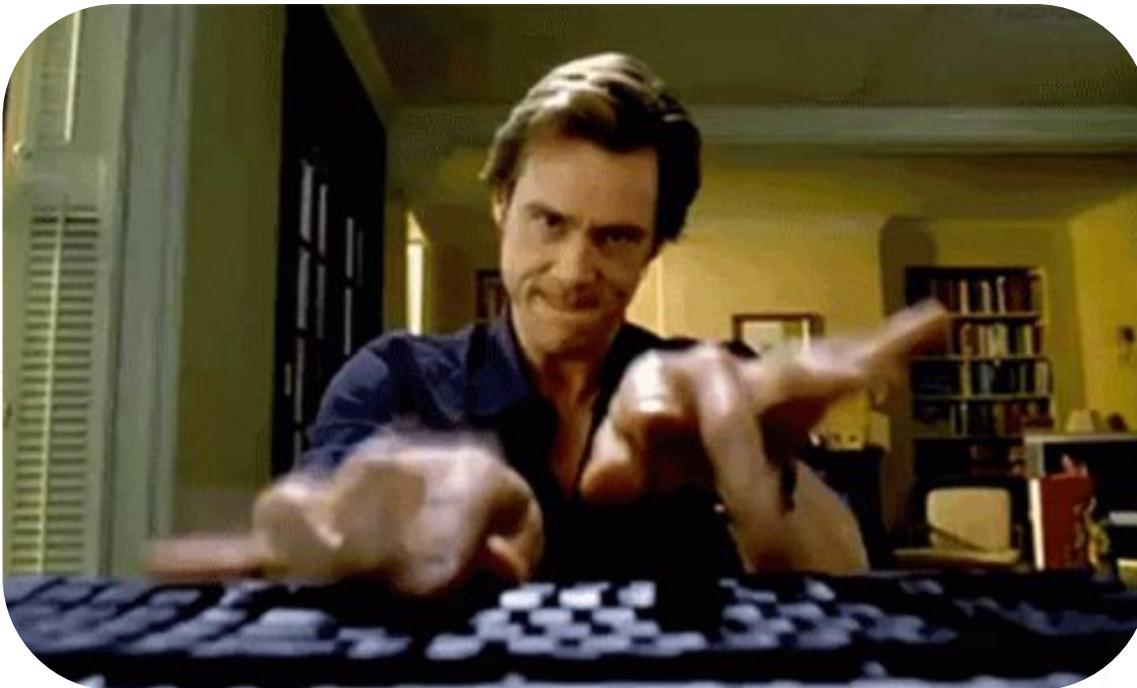
Get predictions

After you train a model, you can use it to get predictions, either online as an endpoint or through batch requests

[+ CREATE BATCH PREDICTION](#)



Back to Vertex AI



Setting the Project ID

Set your project ID

If you don't know your project ID, try the following:

- Run `gcloud config list`
- Run `gcloud projects list`
- See the support page: [Locate the project ID](#)

```
PROJECT_ID = "[your-project-id]" # @param {type:"string"}  
  
# set the project id  
! gcloud config set project $PROJECT_ID
```

Project ID = caltech-class

Create buckets

Create a Cloud Storage bucket

Create a storage bucket to store intermediate artifacts such as datasets.

```
[6]: BUCKET_URI = f"gs://your-bucket-name-{PROJECT_ID}-unique" # @param {type:"string"}
```

Only if your bucket doesn't already exist: Run the following cell to create your Cloud Storage bucket.

```
[7]: ! gsutil mb -l {REGION} -p {PROJECT_ID} {BUCKET_URI}
```

```
Creating gs://your-bucket-name-caltech-class-unique/...
```

Check on new bucket

Cloud Storage

Buckets + CREATE ⌂ REFRESH

Buckets

Monitoring

Settings

Transfer New

Power near real-time analytics and replication with event-driven transfers

You can now capture changes faster at your Google Cloud Storage and Amazon S3 sources via event-driven transfers, enabling you to act on your data in near real time. To get started, create a transfer job with a Pub/Sub or AWS SQS-based event stream configured to send event notifications when objects are created or updated.

[CREATE TRANSFER JOB](#) [LEARN MORE](#)

Analytics New

Preview the new Cloud Storage monitoring dashboard

Check out the new Cloud Storage monitoring dashboards! Powered by Cloud Operations, these dashboards for each project.

[TRY NOW](#)

Filter buckets

<input type="checkbox"/> Name ↑	Created	Location type	Location	Default storage class <small>?</small>	Last modified
your-bucket-name-caltech-class-unique	Oct 31, 2023, 5:17:29 AM	Region	us-central1	Standard	Oct 31, 2023, 5:17:29

Check on new bucket

Cloud Storage Bucket details

Buckets Monitoring Settings

your-bucket-name-caltech-class-unique

Location	Storage class	Public access	Protection
us-central1 (Iowa)	Standard	Subject to object ACLs	None

OBJECTS CONFIGURATION PERMISSIONS PROTECTION LIFECYCLE OBSERVABILITY INVENTORY REPORTS

Buckets > your-bucket-name-caltech-class-unique

UPLOAD FILES UPLOAD FOLDER CREATE FOLDER TRANSFER DATA MANAGE HOLDS DOWNLOAD DELETE

Filter by name prefix only Filter objects and folders

Name	Size	Type	Created	Storage class	Last modified	Public access	Version history	Encryption	Retention
No rows to display									

Import libraries

```
[*]:  
import os  
import pickle  
import sys  
import warnings  
  
import matplotlib.pyplot as plt  
import numpy as np  
import pandas as pd  
from google.cloud import aiplatform, storage  
from IPython.display import display  
from sklearn.ensemble import RandomForestClassifier  
from sklearn.metrics import (average_precision_score, classification_report,  
                             confusion_matrix, f1_score)  
from sklearn.model_selection import train_test_split  
from witwidget.notebook.visualization import WitConfigBuilder, WitWidget  
  
warnings.filterwarnings("ignore")
```

Initialize Vertex AI SDK for Python

Initialize the Vertex AI SDK for Python for your project.

```
[9]: aiplatform.init(project=PROJECT_ID, location=REGION, staging_bucket=BUCKET_URI)
```

Load dataset

Load the dataset from the public csv file path using Pandas.

```
<]: # set the dataset path  
DATASET_SOURCE_PATH = "gs://cloud-samples-data/vertex-ai/managed_notebooks/fraud_detection/fraud_detection_data.csv"  
# read the csv data using pandas  
df = pd.read_csv(DATASET_SOURCE_PATH)
```

Analyze the dataset

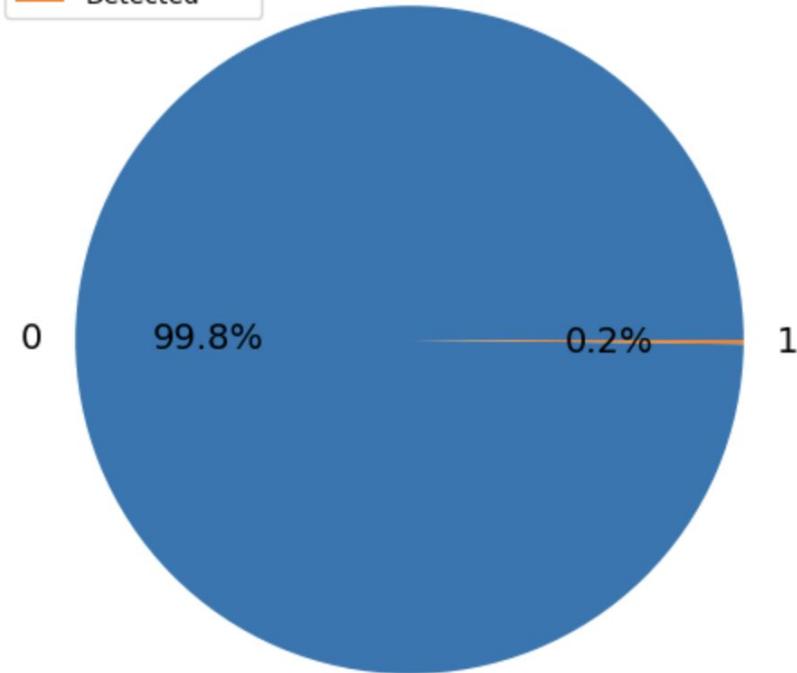
Take a quick look at the dataset and the number of rows:

```
# print the shape of dataframe  
print("shape : ", df.shape)  
# display the dataframe  
df.head()
```

Imbalanced data

% of fraud transaction detected

Not Detected
Detected



Prepare data for modeling

To prepare the dataset for training, a few columns need to be dropped that contain unique data ('nameOrig','nameDest','isFlaggedFraud'). The categorical field "type" which describes the type of transaction and is important for fraud detection needs to be hot encoded.

```
: # drop the unnecessary fields
df.drop(["nameOrig", "nameDest", "isFlaggedFraud"], axis=1, inplace=True)
# encode the "type" field
X = pd.concat([df.drop("type", axis=1), pd.get_dummies(df["type"])], axis=1)
X.head()
```

Remove the outcome variable from the training data.

```
: # copy the target data
y = X[["isFraud"]]
# remove the target field from the features
X = X.drop(["isFraud"], axis=1)
```

Train-test strategy

Split the data and assign 70% for training and 30% for testing.



For splitting, you specify the following parameters to Sklearn's `train_test_split` method:

- `*arrays` : The feature array(X) and the target array(y).
- `test_size` : Percentage(float) or number(integer) of test samples.
- `random_state` : Controls the shuffling applied to the data before applying the split. Pass an int for reproducible output across multiple function calls.
- `stratify` : If none, no stratified sampling is performed.

As the data is imbalanced, you use stratified sampling while splitting. Learn more about [stratified sampling and other parameters for train-test-splitting](#).

Model selection

Fit a Random Forest model



Fit a simple Random Forest classifier on the preprocessed training dataset.

Note: Setting `n_jobs` to -1 while defining the `RandomForestClassifier` object allows it to parallelize the training process using all processors.

Learn more about [Random Forest algorithm](#) and Sklearn's `RandomForestClassifier`.

```
# create a randomforestclassifier object
forest = RandomForestClassifier(n_jobs=-1, verbose=1)
# fit the model on the data
forest.fit(X_train, y_train)
```

This will take about 7 min

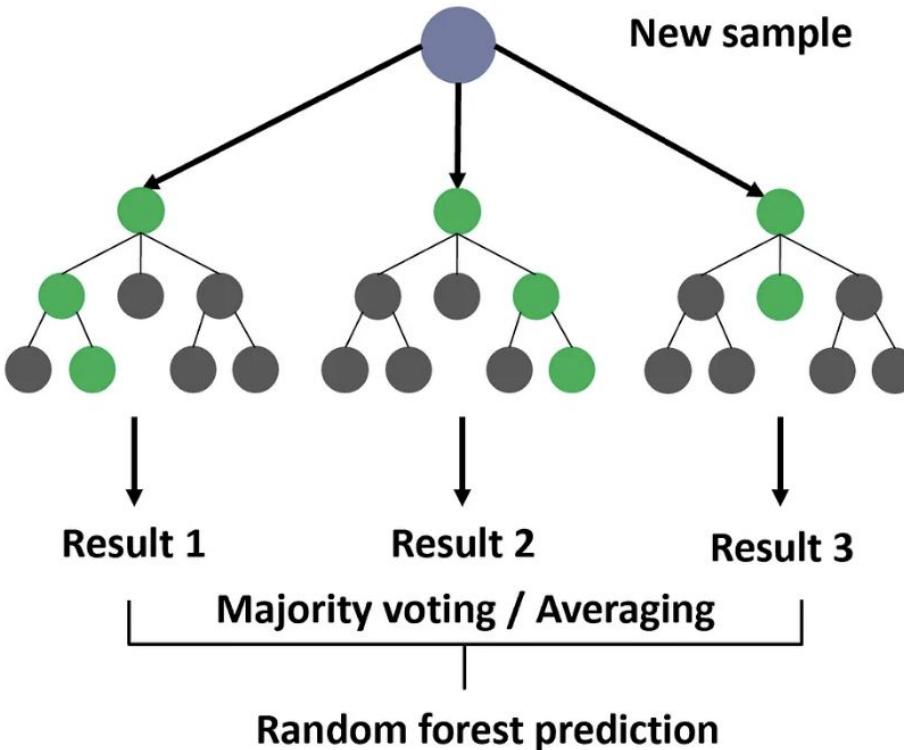
Model training

```
# create a randomforestclassifier object
forest = RandomForestClassifier(n_jobs=-1, verbose=1)
# fit the model on the data
forest.fit(X_train, y_train)
```

```
[Parallel(n_jobs=-1)]: Using backend ThreadingBackend with 4 concurrent workers.
[Parallel(n_jobs=-1)]: Done 42 tasks      | elapsed:  3.2min
[Parallel(n_jobs=-1)]: Done 100 out of 100 | elapsed:  7.2min finished
```

```
▼      RandomForestClassifier
RandomForestClassifier(n_jobs=-1, verbose=1)
```

Random Forest Overview (while training)



Analyze results

AP : 0.9376763239087629

F1 - score : 0.8665307047360074

Confusion_matrix :

```
[[1906285      37]
 [   552     1912]]
```

classification_report

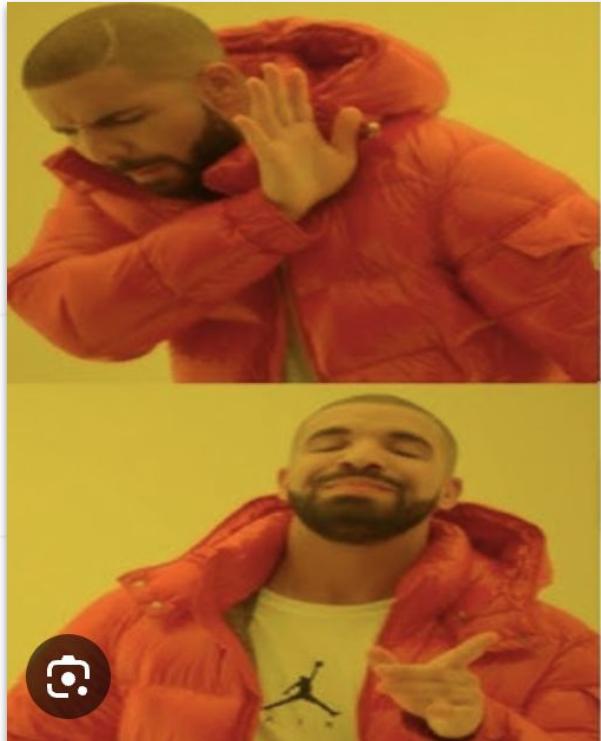
	precision	recall	f1-score	support
0	1.00	1.00	1.00	1906322
1	0.98	0.78	0.87	2464
accuracy			1.00	1908786
macro avg	0.99	0.89	0.93	1908786
weighted avg	1.00	1.00	1.00	1908786

What do we use for classification?

Accuracy?

$$\text{Accuracy} = \frac{\text{Number of Correct Predictions}}{\text{Total Number of Predictions}} \times 100$$

Going beyond accuracy



Accuracy

Precision & Recall F_1 Score

Classification Evaluation Metrics

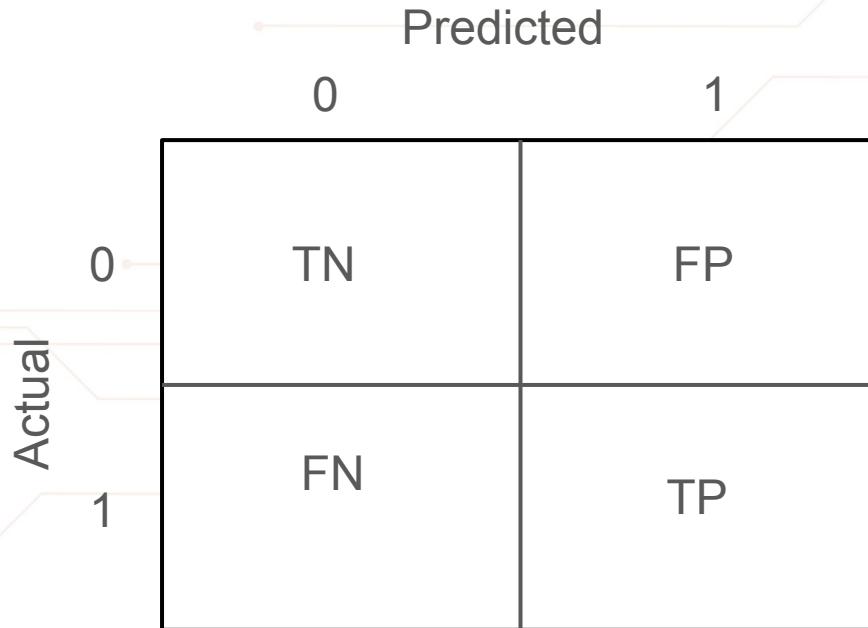
Metric	Formula	Evaluation focus
Accuracy	$ACC = \frac{TP + TN}{TP + TN + FP + FN}$	Overall effectiveness of a classifier
Precision	$PRC = \frac{TP}{TP + FP}$	Class agreement of the data labels with the positive labels given by the classifier
Sensitivity	$SNS = \frac{TP}{TP + FN}$	Effectiveness of a classifier to identify positive labels. Also called true positive rate (TPR)
Specificity	$SPC = \frac{TN}{TN + FP}$	How effectively a classifier identifies negative labels. Also called true negative rate (TNR)
F_1 score	$F_1 = 2 \frac{PRC \cdot SNS}{PRC + SNS}$	Combination of precision (PRC) and sensitivity (SNS) in a single metric
Geometric mean	$GM = \sqrt{SNS \cdot SPC}$	Combination of sensitivity (SNS) and specificity (SPC) in a single metric
Area under (ROC) curve	$AUC = \int_0^1 SNS \cdot dSPC$	Combined metric based on the receiver operating characteristic (ROC) space (Powers, 2011)

Source: https://www.researchgate.net/figure/Classification-performance-metrics-based-on-the-confusion-matrix_tbl3_324952663

Simple, right?



The **Confusion Matrix**: Where it all begins



Example: Deadly Disease classification model

		Predicted: NO	Predicted: YES
n=165	Actual: NO	50	10
	Actual: YES	5	100

Accuracy

$$\text{Accuracy} = \frac{TP+TN}{TP+TN+FP+FN}$$

Overall, how often is the classifier correct?

$$(TP+TN)/\text{total} = (100+50)/165 = 0.91$$

91% accurate

n=165	Predicted: NO	Predicted: YES
Actual: NO	TN = 50	FP = 10
Actual: YES	FN = 5	TP = 100



Let's celebrate !!

getclippy.co

Not so fast buddy

- What if 91% accuracy is a terrible model?
- What if we are inaccurately classifying a large group of people?
- How much does misclassification hurt us/the organization/business/project?
- Do we care more about misclassifying False Positives?
- What about False Negatives?

The Accuracy Paradox

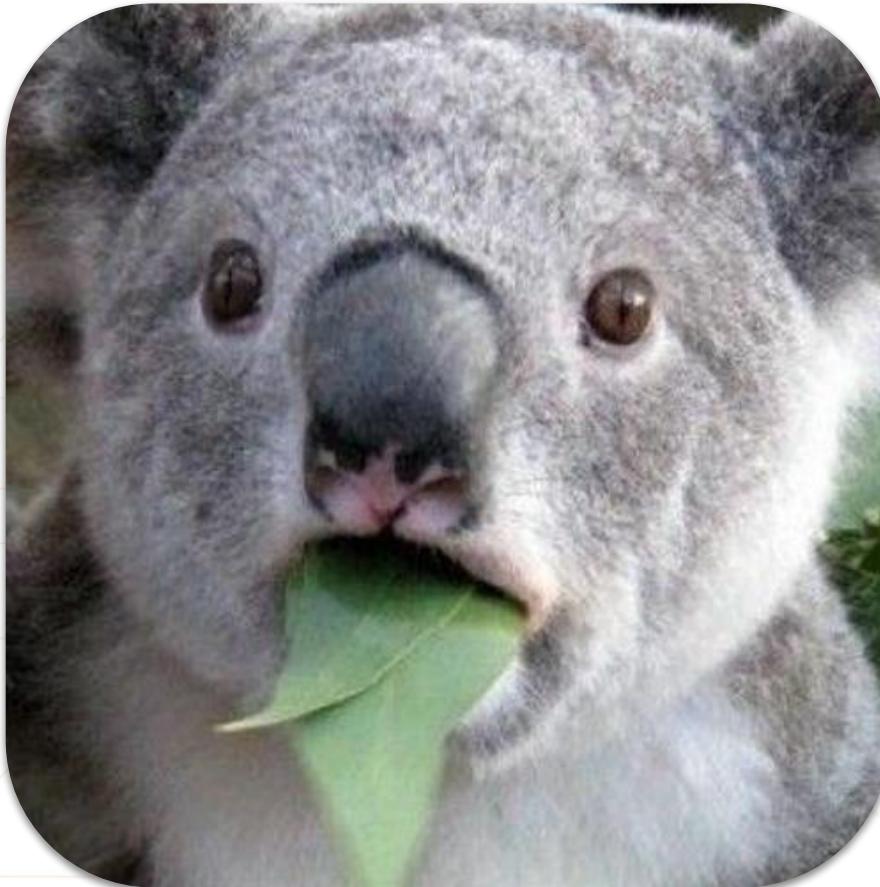
Example [\[edit\]](#)

For example, a city of 1 million people has ten terrorists. A profiling system results in the following [confusion matrix](#):

Predicted class \ Actual class	Fail	Pass	Sum
Fail	10	0	10
Pass	990	999000	999990
Sum	1000	999000	1000000

Even though the accuracy is $\frac{10 + 999000}{1000000} \approx 99.9\%$, 990 out of the 1000 positive predictions are incorrect. The precision of $\frac{10}{10 + 990} = 1\%$ reveals its poor performance. As the classes are so unbalanced, a better metric is the F1 score = $\frac{2 \times 0.01 \times 1}{0.01 + 1} \approx 2\%$ (the recall being $\frac{10 + 0}{10} = 1$).

Wait... what?



Precision and Recall

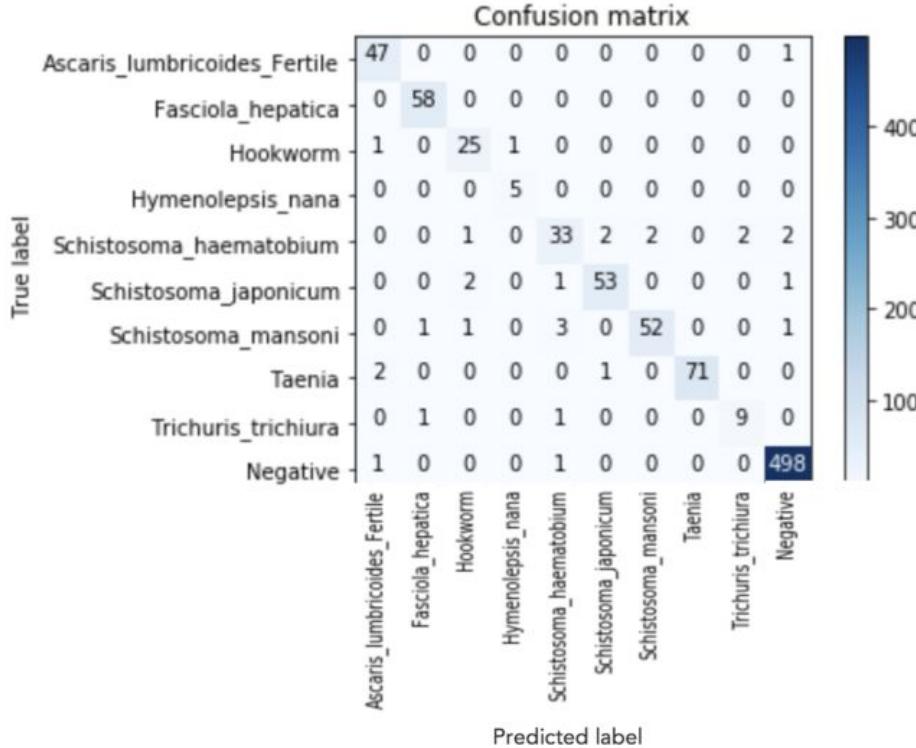
- **Recall:** When it's actually yes, how often does it predict yes?
 - $TP/\text{actual yes} = 100/105 = 0.95$
- **Precision:** When it predicts yes, how often is it correct?
 - $TP/\text{predicted yes} = 100/110 = 0.91$

$$\text{Recall} = \frac{TP}{TP+FN}$$

$$\text{Precision} = \frac{TP}{TP+FP}$$

n=165	Predicted: NO	Predicted: YES
Actual: NO	TN = 50	FP = 10
Actual: YES	FN = 5	TP = 100

Multiclass Confusion Matrix



Holistically, what are we talking about?

$$\text{Recall} = \frac{TP}{TP+FN}$$

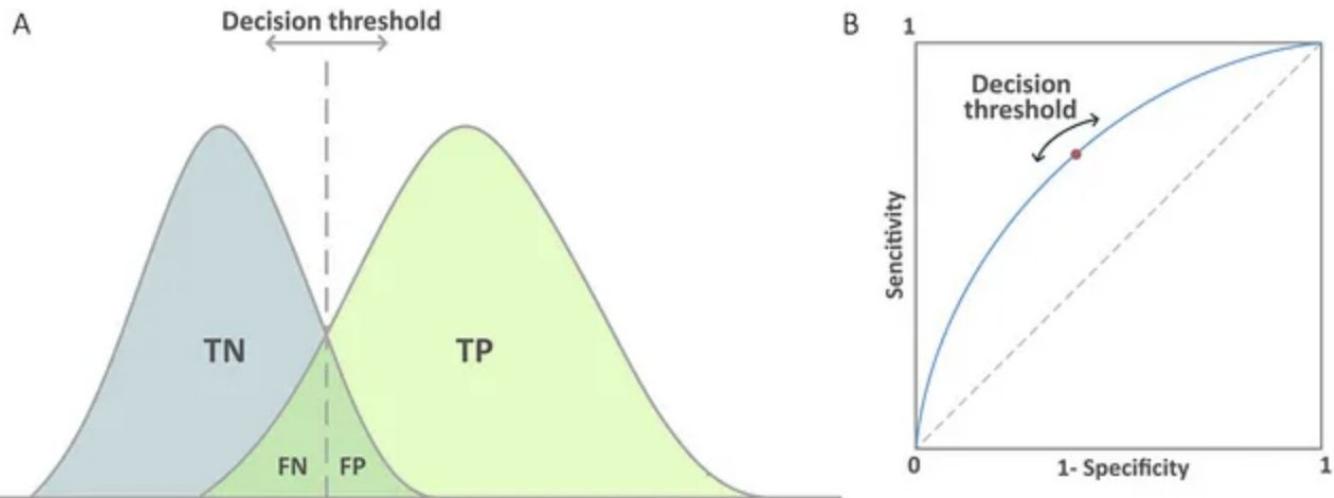
$$\text{Precision} = \frac{TP}{TP+FP}$$

What about some harmony?



$$F1 \text{ Score} = \frac{2 \times \text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$

Probability thresholds



The goal of this post is to outline how to move the decision threshold to the left in Figure A, reducing false negatives and maximizing sensitivity.

Source: <https://towardsdatascience.com/fine-tuning-a-classifier-in-scikit-learn-66e048c21e65>

Let's check a cool site out

Dive in!



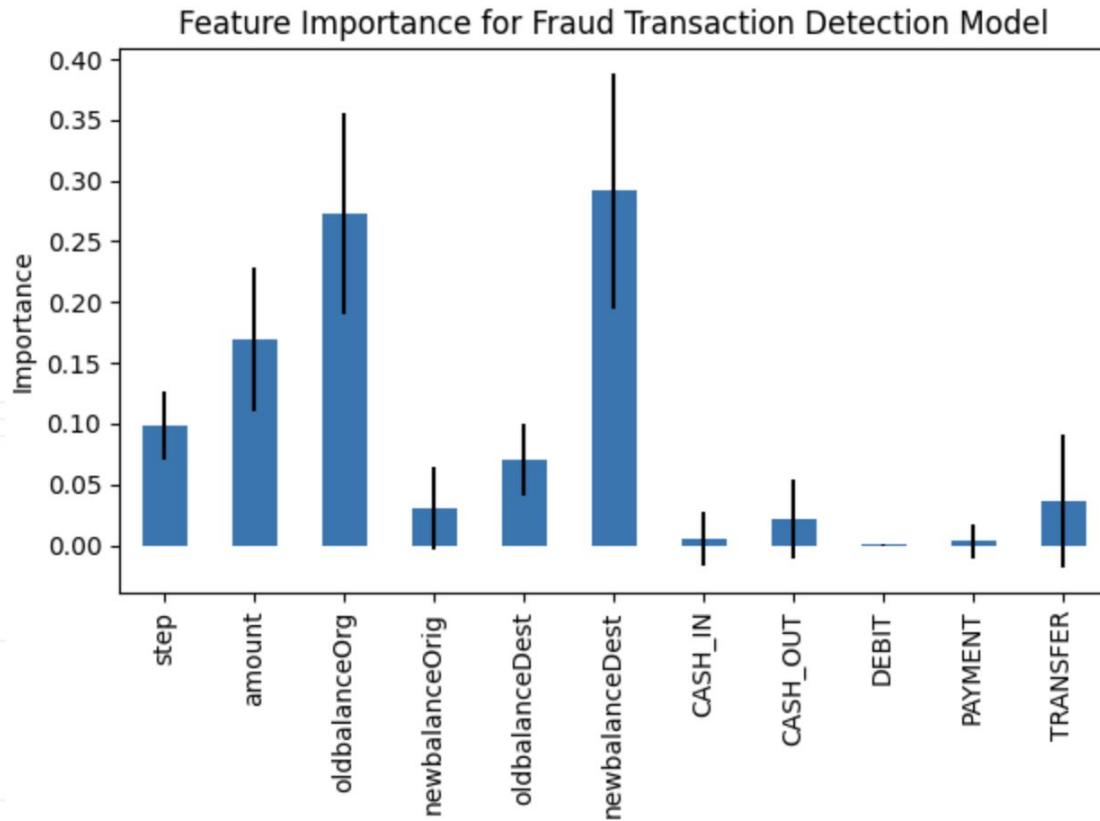
MLU-EXPLAIN

PRECISION & RECALL

Accuracy Is Not Enough

[Jared Wilber](#), March 2022

Analyze results



Save the model to Cloud Storage

Save your model to a pickle file and then, upload your model to Cloud Storage bucket. The uploaded model path is later used for creating a model in the Vertex AI Model Registry.

Note: You can also upload the model to Vertex AI Model Registry from your local environment using the latest Vertex AI SDK for Python.

Check your bucket

Cloud Storage ← Bucket details REFRESH LEA

Buckets your-bucket-name-caltech-class-unique

Monitoring Location: us-central1 (Iowa) Storage class: Standard Public access: Subject to object ACLs Protection: None

Settings

OBJECTS CONFIGURATION PERMISSIONS PROTECTION LIFECYCLE OBSERVABILITY INVENTORY REPORTS

Buckets > your-bucket-name-caltech-class-unique

UPLOAD FILES UPLOAD FOLDER CREATE FOLDER TRANSFER DATA MANAGE HOLDS DOWNLOAD DELETE

Filter by name prefix only Filter objects and folders Show deleted data

Name	Size	Type	Created	Storage class	Last modified	Public access	Version history	Encryption
fraud-detect-model-path-unique/	—	Folder	—	—	—	—	—	—

Name	Size	Type	Created	S
model.pkl	21.8 MB	application/octet-stream	Oct 31, 2023, 5:34:48 AM	S

Create a model in Vertex AI

Set the parameters required for model creation in Vertex AI Model Registry.

```
# set model display name
MODEL_DISPLAY_NAME = "fraud-detection-model-unique" # @param {type:"string"}
# set the GCS path to the model artifact
ARTIFACT_GCS_PATH = f"{BUCKET_URI}/{BLOB_PATH}"
# set the prediction container uri
SERVING_CONTAINER_IMAGE_URI = (
    "us-docker.pkg.dev/vertex-ai/prediction/sklearn-cpu.1-2:latest"
)
```

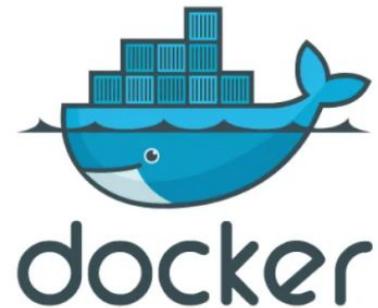
What the heck is a container?

What is a container?

→ Stores your application's code with files/libraries needed to run it on any machine

But why?

- More portable
- Solves the “works on my machine” situation
- Scales very easily when your app gets popular
- Fault tolerant → isolated from other infrastructure



“What Docker really does is separate the application code from infrastructure requirements and needs. It does this by running each application in an isolated environment called a ‘container.’”

Source: [Chinmay Shah](#)

"It works on my computer"

The notebook won't run on his colleagues computer!

```
In [2]: import sys
!{sys.executable} -m pip install --user numpy

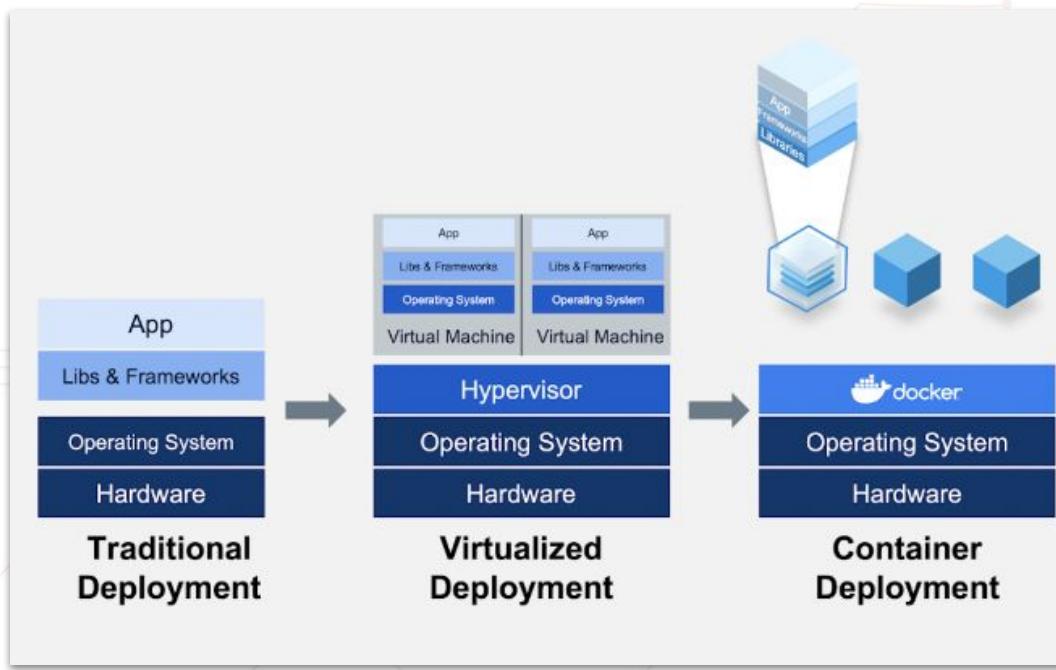
de0e77425/numpy/distutils/command/build_src.py", line 378, in generate_sources
    source = func(extension, build_dir)
File "/private/var/folders/5z/6mpx11111jx9spjrrht8k9y80000gn/T/pip-install-dduq83ml(numpy_8396f5e1bbf1412cae468cb
de0e77425/numpy/core/setup.py", line 456, in generate_config_h
    moredefs, ignored = cocache.check_types(config_cmd, ext, build_dir)
File "/private/var/folders/5z/6mpx11111jx9spjrrht8k9y80000gn/T/pip-install-dduq83ml(numpy_8396f5e1bbf1412cae468cb
de0e77425/numpy/core/setup.py", line 50, in check_types
    out = check_types(*a, **kw)
File "/private/var/folders/5z/6mpx11111jx9spjrrht8k9y80000gn/T/pip-install-dduq83ml(numpy_8396f5e1bbf1412cae468cb
de0e77425/numpy/core/setup.py", line 311, in check_types
    raise SystemError(
SystemError: Cannot compile 'Python.h'. Perhaps you need to install python-dev|python-devel.

ERROR: Failed building wheel for numpy
Failed to build numpy
ERROR: Could not build wheels for numpy, which is required to install pyproject.toml-based projects
WARNING: You are using pip version 21.3.1; however, version 22.0.3 is available.
You should consider upgrading via the '/Library/Frameworks/Python.framework/Versions/3.10/bin/python3 -m pip install
--upgrade pip' command.
```

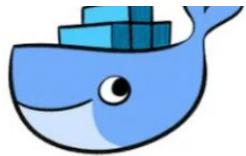


Basic libraries are failing to install

Docker architecture



Lots of container registry options



DockerHub



Amazon ECR



DigitalOcean
Container Registry



Azure
Container Registry



Harbor
Container Registry



GitLab
Container Registry



Google
Container Registry



IBM
Container Registry



Alibaba Cloud
Container Registry

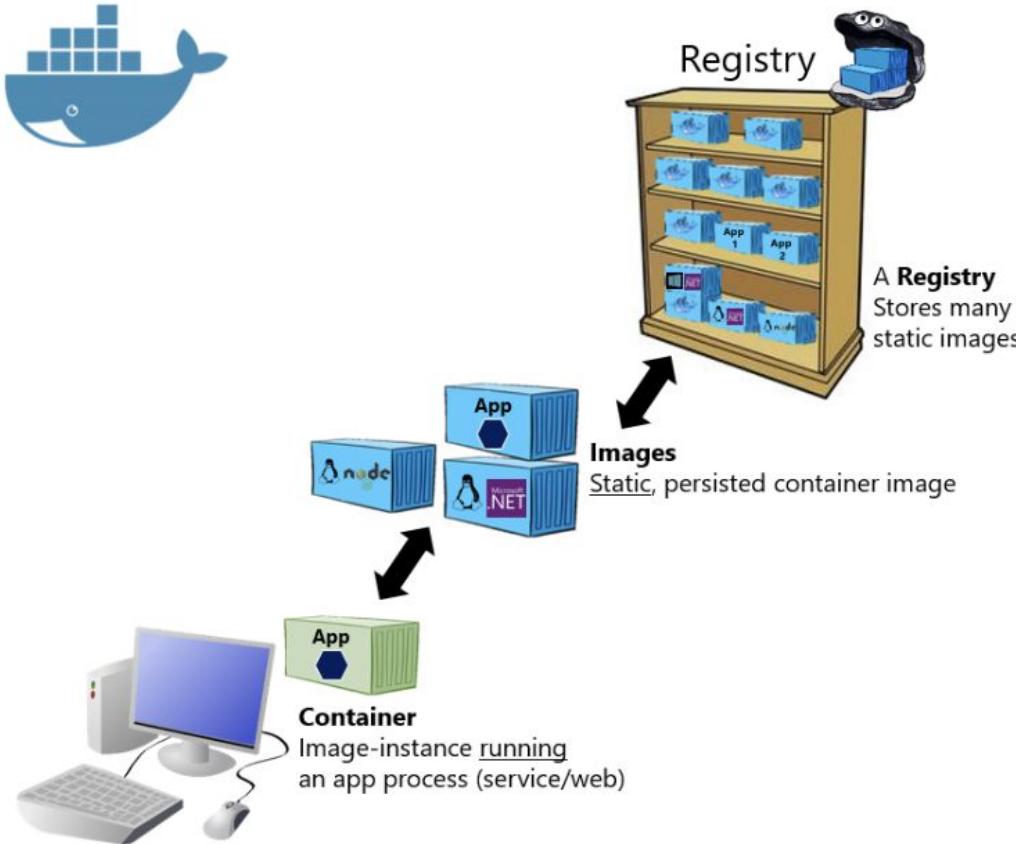


JFrog
Container Registry



Quay
Container Registry

Container Registry



Example container registry

AWS Amazon ECR Public Gallery Find artifact repositories About Share & manage images

Amazon ECR Public Gallery

Share and deploy container images, publicly and privately

Sort by: Popularity 1 2 3 4 5 >

Filters Clear all

Popular Registries

- Amazon
- Docker

Verification Info

- Verified

Operating Systems

- Linux
- Windows

Architectures

- ARM
- ARM 64
- x86
- x86-64

Repositories Showing 1 - 20 results (of 264104)

 cloudwatch-agent/cloudwatch-agent (8.5B+ downloads) by Amazon Cloudwatch Agent  Amazon Cloudwatch Agent OS/Arch: Linux, x86-64, ARM 64	 eks-distro/kubernetes-csi/livenessprobe (5.8B+ downloads) by Amazon EKS Distribution  Amazon EKS Distro Kubernetes CSI Livenessprobe OS/Arch: Linux, x86-64, ARM 64
 eks-distro/kubernetes-csi/node-driver-registrar (5.7B+ downloads) by Amazon EKS Distribution  Amazon EKS Distro Kubernetes CSI Node Driver Registrar OS/Arch: Linux, x86-64, ARM 64	 eks-distro/kubernetes-csi/external-provisioner (4.8B+ downloads) by Amazon EKS Distribution  Amazon EKS Distro Kubernetes CSI External Provisioner OS/Arch: Linux, x86-64, ARM 64
 xray/aws-xray-daemon (4.3B+ downloads) by AWS X-Ray  The AWS X-Ray daemon gathers raw segment data and relays it to the AWS X-Ray API. OS/Arch: Linux	 datadog/agent (3.9B+ downloads) by datadog  Docker container for the new Datadog Agent OS/Arch: Linux, Windows, x86-64, ARM 64
 aws-observability/aws-for-fluent-bit (3.7B+ downloads) by AWS Observability Toolkits  Official AWS distribution of Fluent Bit OS/Arch: Linux, Windows, ARM 64, x86-64	 eks-distro/kubernetes-csi/external-attacher (1.8B+ downloads) by Amazon EKS Distribution  Amazon EKS Distro Kubernetes CSI External Attacher OS/Arch: Linux, x86-64, ARM 64
 eks-distro/kubernetes-csi/external-resizer (1.8B+ downloads) by Amazon EKS Distribution  Amazon EKS Distro Kubernetes CSI External Resizer	 eks-distro/kubernetes/pause (1.6B+ downloads) by Amazon EKS Distribution  Amazon EKS Distro Kubernetes pause image

What is Model Registry?

The Vertex AI Model Registry is a central repository where you can manage the lifecycle of your ML models. From the Model Registry, you have an overview of your models so you can better organize, track, and train new versions. When you have a model version you would like to deploy, you can assign it to an endpoint directly from the registry, or using aliases, deploy models to an endpoint.

Common workflow

There are many valid workflows for working in the Model Registry. To get started, you might want to follow the guidelines below to understand what you can do in the Model Registry and at what stage in your model-training journey.

- Import models to the Model Registry.
- Create new models, assign a model version the default alias, ready for production.
- Add other aliases, or labels to help you manage and organize your models and model versions.
- Deploy your models to an endpoint.
- Run batch prediction, and start your model evaluation pipeline.
- View your model details and view performance metrics from the model details page.

Model Registry

Create a model resource in Vertex AI using the `Model.upload` method.



Learn more about [Vertex AI Model Registry](#).

```
# create a Vertex AI model resource
model = aiplatform.Model.upload(
    display_name=MODEL_DISPLAY_NAME,
    artifact_uri=ARTIFACT_GCS_PATH,
    serving_container_image_uri=SERVING_CONTAINER_IMAGE_URI,
)
# print the model's display name
print("Display name:\n", model.display_name)
# print the model's resource name
print("Resource name:\n", model.resource_name)
```

What is Model Registry?

MODEL DEVELOPMENT

-  Training

-  Experiments

-  Metadata

DEPLOY AND USE

-  Model Registry

-  Online prediction

-  Batch predictions

-  Vector Search

Source: XXX

What is Model Registry?

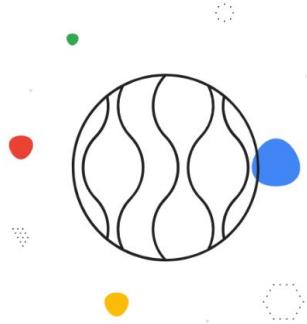
The screenshot shows the Vertex AI Model Registry interface. On the left, there is a sidebar with sections for MODEL DEVELOPMENT (Datasets, Labeling tasks), DEPLOY AND USE (Model Registry, Online prediction, Batch predictions, Vector Search), and MANAGE (Ray on Vertex AI, Marketplace). The main area is titled "Model Registry" and includes "CREATE" and "IMPORT" buttons. A descriptive text block states: "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](#)". Below this is a "Region" dropdown set to "us-central1 (Iowa)". A "Filter" input field is present. A table lists models, with one entry visible: "fraud-detection-model..." (Name), Deployment status: -, Description: -, Default version: 1, Type: Imported (Custom training), Source: Imported, Updated: Oct 31, 2023, 5:38:38 AM.

Name	Deployment status	Description	Default version	Type	Source	Updated	Labels
fraud-detection-model...	-	-	1	Imported	Custom training	Oct 31, 2023, 5:38:38 AM	-

What is Model Registry?

← fraud-detection-model-unique → Version 1 ▾ [EXPORT](#)

EVALUATE DEPLOY & TEST BATCH PREDICT VERSION DETAILS



No model evaluations created yet

Evaluations help you understand your model's performance with metrics, like precision and recall.

[CREATE EVALUATION](#) [LEARN MORE](#)

Create an endpoint

← fraud-detection-model-unique > Version 1 ▾ [EXPORT](#)

EVALUATE

DEPLOY & TEST

BATCH PREDICT

VERSION DETAILS



This model was imported on Oct 31, 2023, 5:38:25 AM.

Model ID 5982387791147827200

Version description —

Created Oct 31, 2023, 5:38:25 AM

Region us-central1

Encryption type Google-managed

Dataset No managed dataset

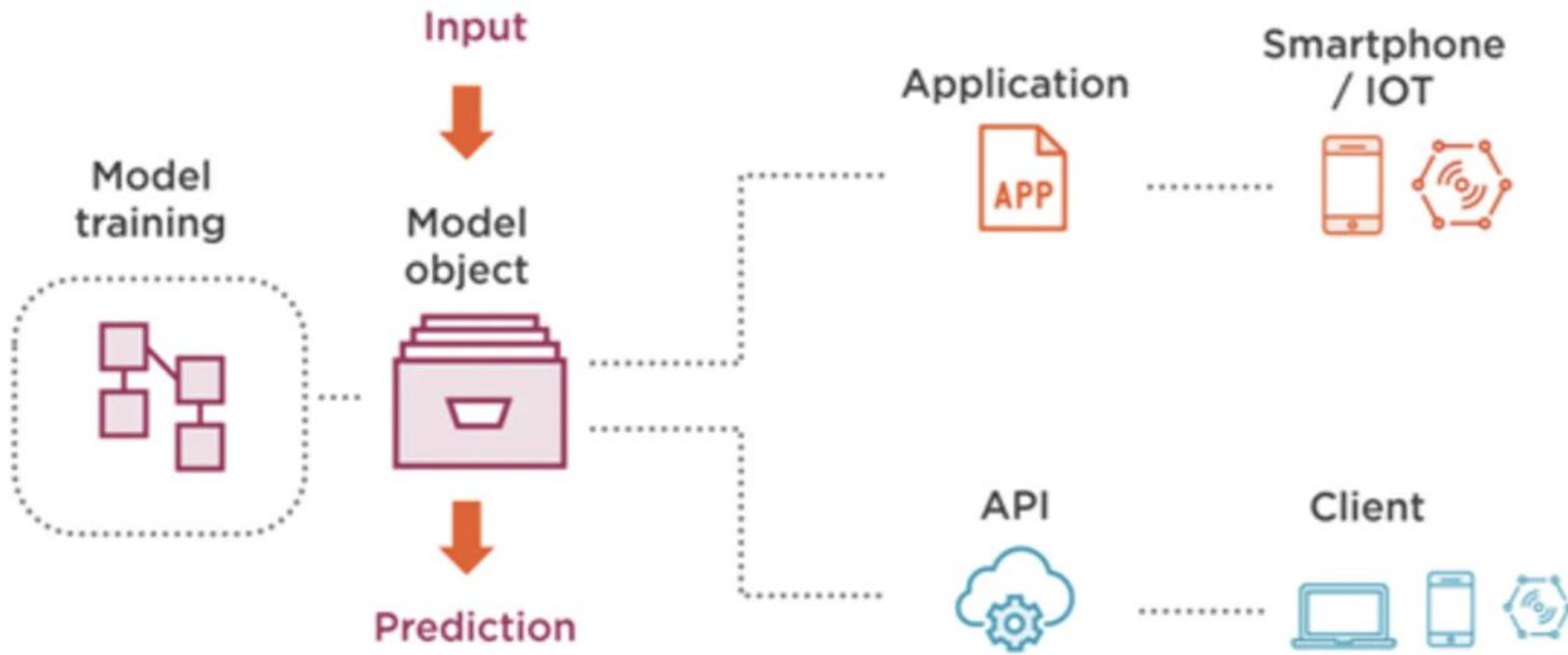
Objective Custom

Source Custom training

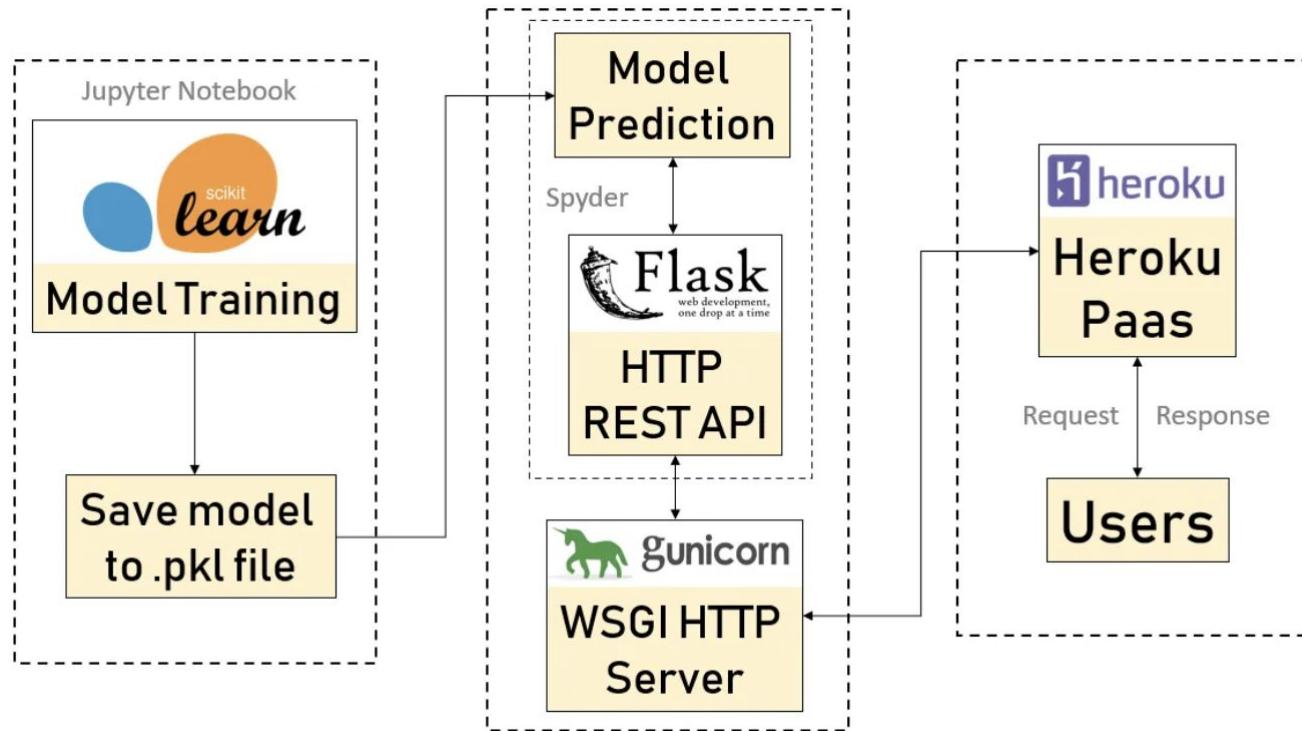
Container image [us-docker.pkg.dev/vertex-ai/prediction/sklearn-cpu.1-2:latest](#)

Model artifact location [gs://your-bucket-name-caltech-class-unique/fraud-detect-model-path-unique](#)

How do users interact with my model?



What is an endpoint?



Endpoints

```
# set the endpoint display name
ENDPOINT_DISPLAY_NAME = "fraud-detect-endpoint-unique" # @param {type:"string"}
# create the Endpoint
endpoint = aiplatform.Endpoint.create(display_name=ENDPOINT_DISPLAY_NAME)
# print the endpoint display name
print("Display name:\n", endpoint.display_name)
# print the endpoint resource name
print("Resource name:\n", endpoint.resource_name)
```

Creating Endpoint

Create Endpoint backing LRO: projects/961778279858/locations/us-central1/endpoints/8169584699631468544/operations/3881366
687039094784

Endpoint created. Resource name: projects/961778279858/locations/us-central1/endpoints/8169584699631468544

To use this Endpoint in another session:

```
endpoint = aiplatform.Endpoint('projects/961778279858/locations/us-central1/endpoints/8169584699631468544')
```

Display name:

fraud-detect-endpoint-unique

Resource name:

projects/961778279858/locations/us-central1/endpoints/8169584699631468544

Endpoint



Deploy model to the endpoint

Set the following parameters for endpoint deployment:

- `endpoint` : The Vertex AI Endpoint resource created in the last step.
- `deployed_model_display_name` : Display name for the model. If not provided, model's display name is used.
- `machine_type` : Machine type required for serving the model on the endpoint.

```
# set the display name for the deployed model  
DEPLOYED_MODEL_NAME = "fraud-detection-deployed-model"  
# set the machine type for the endpoint  
MACHINE_TYPE = "n1-standard-2"
```

Testing the endpoint



Delete resources

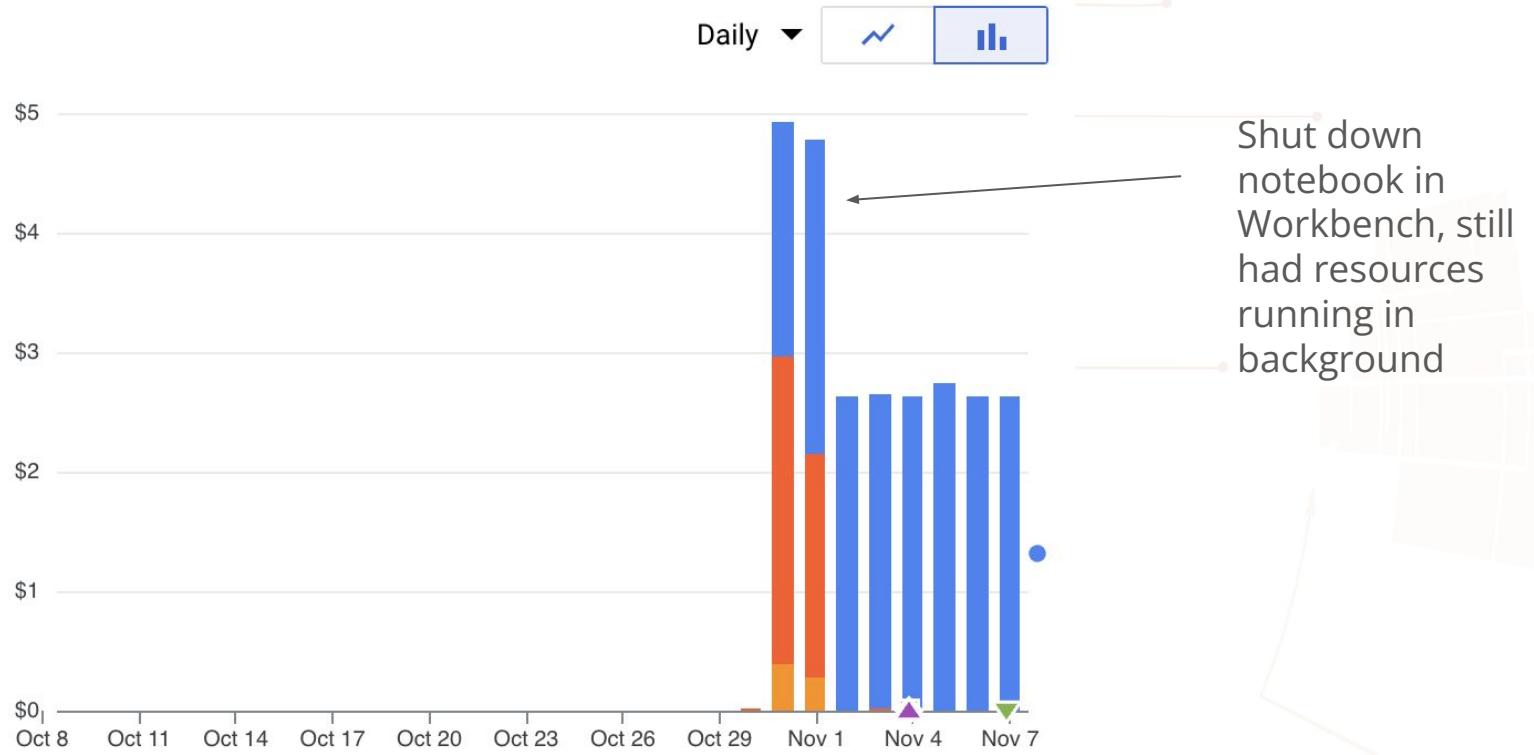
- Undeploy the model
- Delete endpoint
- Delete model from Model Registry
- Delete from Online Predictions
- Delete bucket in Cloud Storage
- Delete Compute Engine instance
- Delete Workbench instance (notebook)

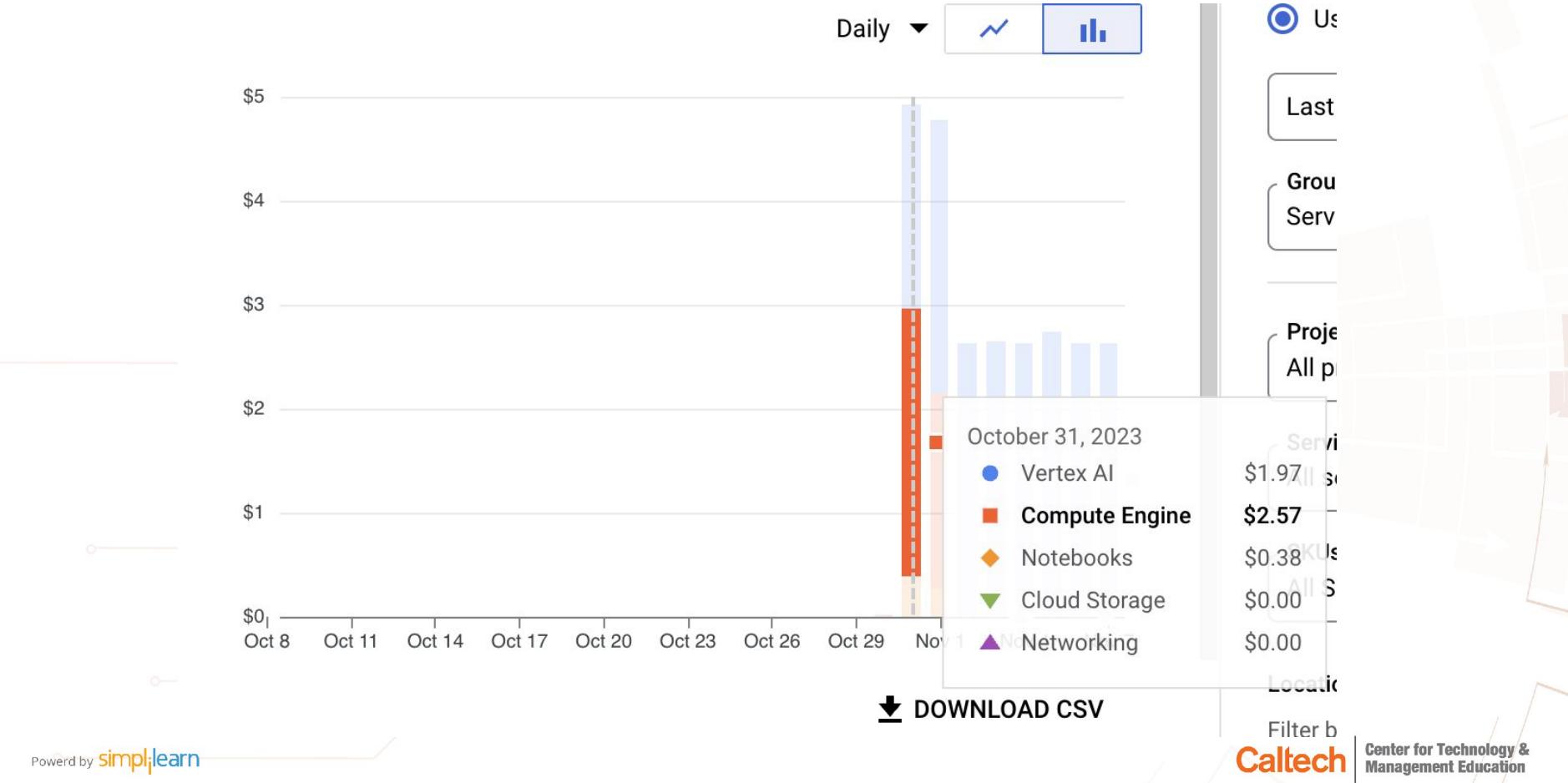
Learning Objectives for Today Review

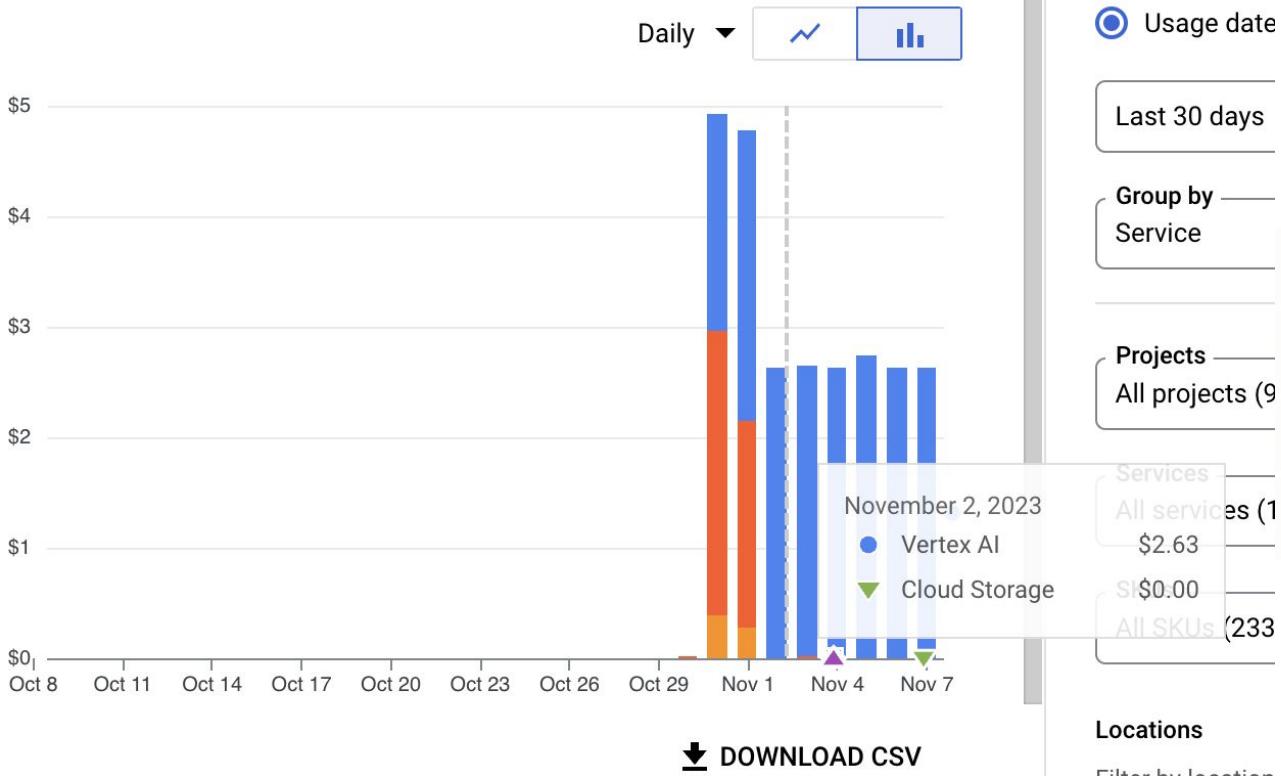
1. Learn about Cloud History and its practical use cases for Machine Learning and AI applications
2. Understand the Machine Learning Lifecycle as it pertains to google cloud platform
3. Apply understanding in a live demo detecting fraud for a large financial institution using Machine Learning

Further Questions?









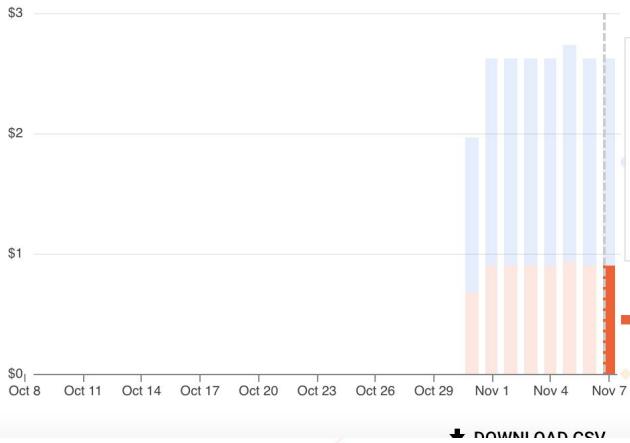
Reports

[PRINT](#)[SHARE](#)[SAVE VIEW](#)[LEARN](#)

Saved views

Free trial credits [?](#)
\$0.00 remaining
out of \$300.00

October 9 – November 7, 2023 (total cost) [?](#)
\$20.47
includes \$0.00 in credits

Daily [▼](#)

Filters

 Usage date Invoice monthLast 30 days [▼](#)

Group by

SKU [▼](#)

Projects

All projects (9) [▼](#)

Services

All services (18) [▼](#)

SKUs

1 out of 222 SKUs [▼](#)

Select the key and values of the labels you want to filter.

Credits

 Discounts [?](#) Sustained use discounts [?](#) Spending based discounts [\(contractual\)](#)

SKU

- Vertex AI:
Online/Batch
Prediction N1
Predefined
Instance
Core running
in Americas
for AI
Platform

- Vertex AI:
Online/Batch
Prediction N1
Predefined
Instance Ram
running in
Americas for
AI Platform

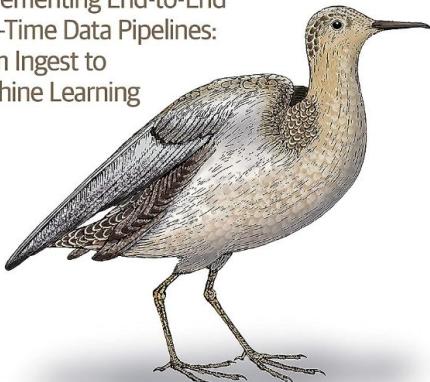
Resources to continue your journey

O'REILLY®

2nd Edition

Data Science on the Google Cloud Platform

Implementing End-to-End
Real-Time Data Pipelines:
From Ingest to
Machine Learning



Valliappa Lakshmanan

[Amazon Link](#)

Powered by  simplilearn



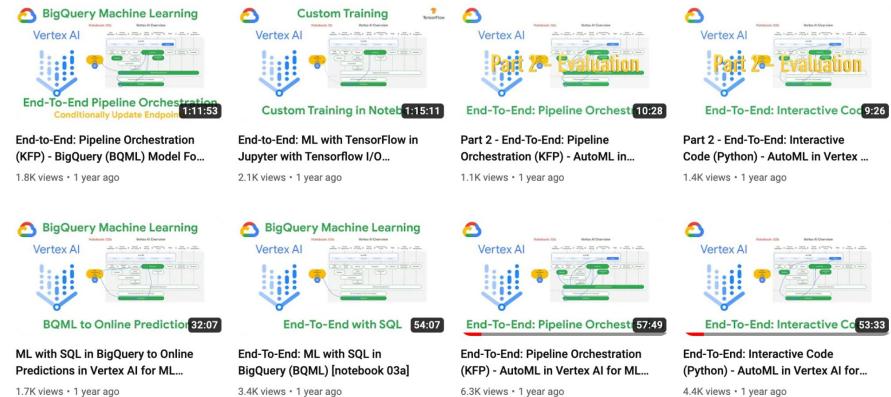
StatMike

@statmike-channel · 1.7K subscribers · 12 videos

Hi, I'm Mike 🙌 >

[github.com/statmike](#) and 1 more link

 Subscribed ▾



[YT Link](#)

Caltech

Center for Technology &
Management Education

Level up!



Next Steps: Pipelines

Intro

<https://medium.com/google-cloud/google-vertex-ai-the-easiest-way-to-run-ml-pipelines-3a41c5ed153>

More advanced

<https://blog.ml6.eu/a-general-framework-for-machine-learning-pipelines-on-gcp-b57e234f7d12>

Google Vertex AI supports two different types of pipelines:

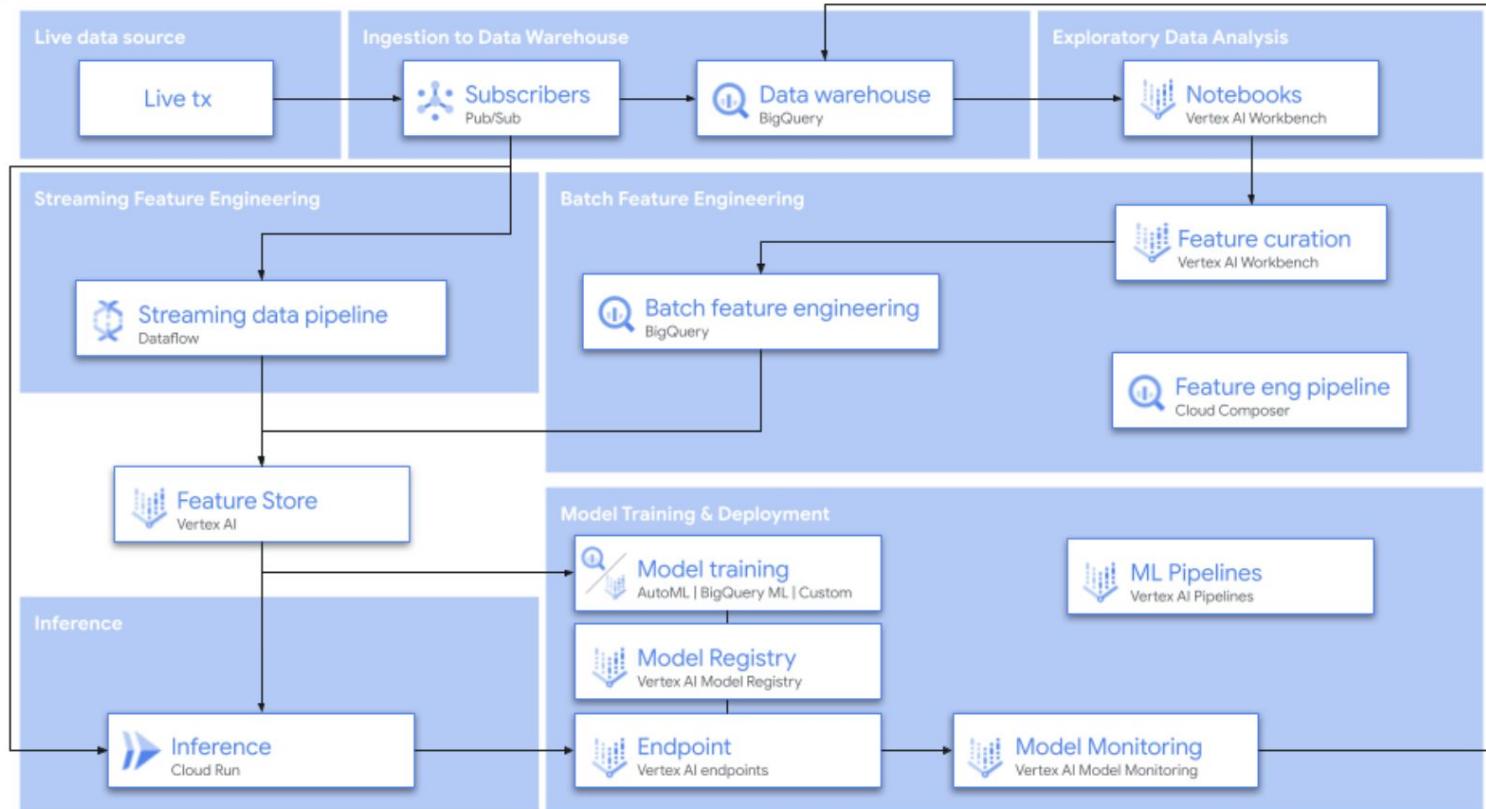
- Kubeflow Pipelines using the Kubeflow SDK ≥ 1.6
- TensorFlow Extended Pipelines using the TFX SDK $\geq 0.30.0$

Pipelining Basics

- Each component has its own compute resources
- A data preprocessing component doesn't need to have a GPU, but a deep learning model training component does
- Can customize to create base images. Think if you need to run a specific algorithm and utilize all cores. You wouldn't need this for the preprocessing step, but for XGBoost you would want to set `jobs=-1`
- Can monitor resource consumption on the component level

Next Steps: MLOps

Data to AI: real-time fraud detection architecture on Google Cloud



Source:

<https://www.googlecloudcommunity.com/gc/Architecture-Framework-Community/Build-an-end-to-end-data-to-AI-solution-on-Google-Cloud-with/ba-p/595682>

Example MLOps Workflow

Setting up the Pipeline

