

DEPLOYMENT TRAINING AND CODE

STEP 1: Enable APIs

- Enable **Compute Engine API** (Need this to create notebook instance)
- Enable **Vertex AI API**
- Enable **Container Registry API** (Need this to create container for custom training job)

STEP 2: Vertex AI Workbench Configuration

Create Vertex AI Workbench Instance

- Go to **Vertex AI** and select workbench instance
- Go to **USER-MANAGED NOTEBOOKS**
- Click **NEW NOTEBOOK** button and select **Pytorch 1.11** with **1 NVIDIA Tesla T4 GPU**

Vertex AI

Dashboard

Datasets

Features

Labelling tasks

Workbench

Pipelines

Training

Experiments

Models

Endpoints

Batch predictions

Metadata

Matching engine

Manage resources

Marketplace

Workbench

MANAGED NOTEBOOKS

Customize...

Python 3
Includes scikit-learn, pandas and more

Python 3 (CUDA Toolkit 11.0)
Optimized for NVIDIA GPUs

TensorFlow Enterprise
Includes Keras, scikit-learn, pandas, NLTK and more

PyTorch 1.11
Includes scikit-learn, pandas, NLTK and more

R 4.1
Includes basic R packages, scikit-learn, pandas, NLTK and more

RAPIDS 0.18 [EXPERIMENTAL]
Optimized for NVIDIA GPUs

Kaggle Python [BETA]
Python image for Kaggle Notebooks, supporting hundreds of machine learning libraries popular on Kaggle

Theia IDE [EXPERIMENTAL]
IDE with notebook support including scikit-learn, pandas, and more

Smart Analytics Frameworks
BigQuery, Apache Beam, Apache Spark, Apache Hive and more

Filter

Enter property

SHOW INFO PANEL

LEARN

Auto-upgrade

Environment

Machine type

GPUs

Owner

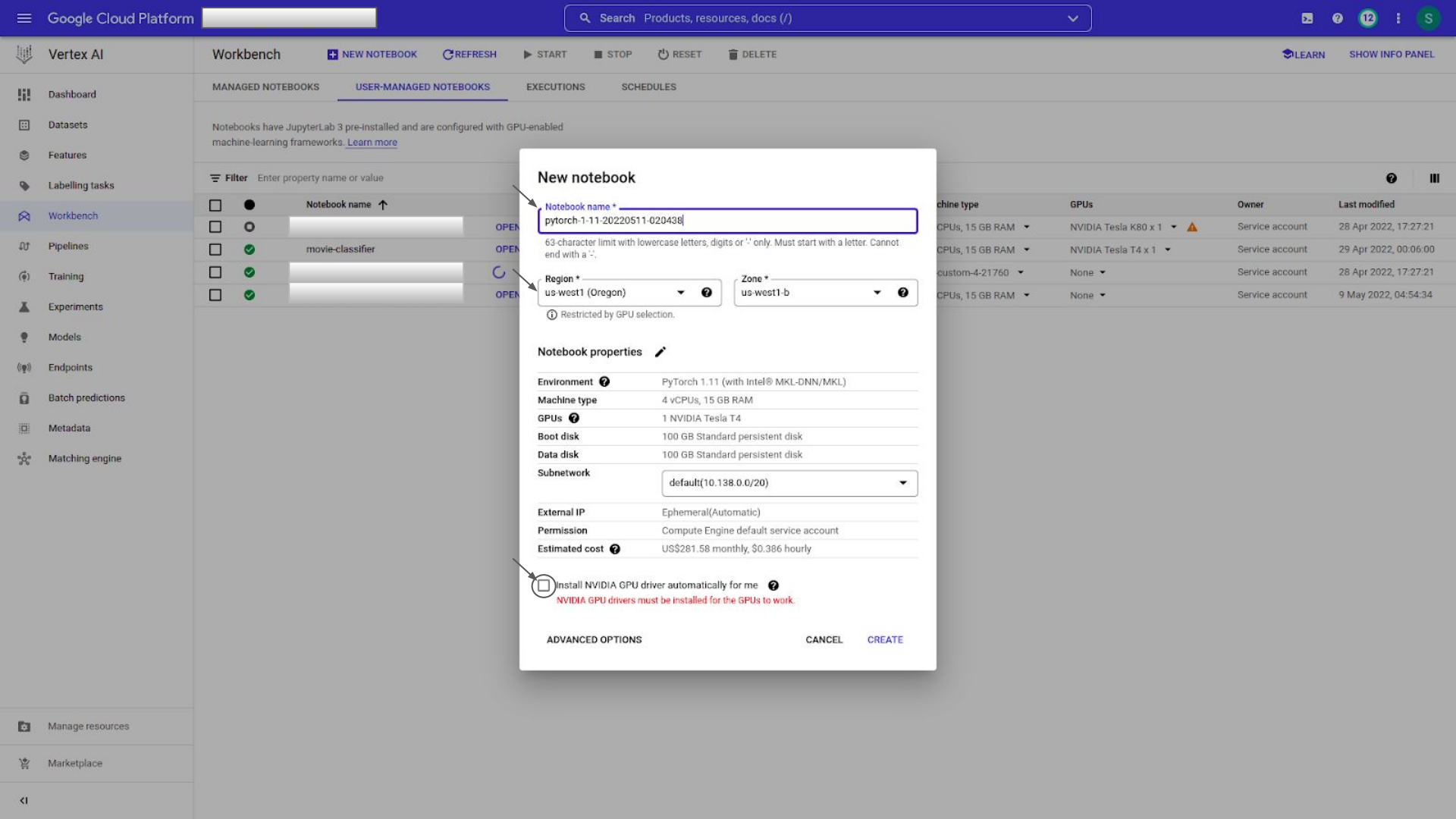
Last modified

Without GPUs

With 1 NVIDIA Tesla T4

ph1.8	4 vCPUs, 15 GB RAM	NVIDIA Tesla K80 x 1	Service account	28 Apr 2022, 17:27:21
ph1.11	4 vCPUs, 15 GB RAM	NVIDIA Tesla T4 x 1	Service account	29 Apr 2022, 00:06:00
ph1.4	e2-custom-4-21760	None	Service account	28 Apr 2022, 17:27:21
TensorFlow/2.3	4 vCPUs, 15 GB RAM	None	Service account	9 May 2022, 04:54:34

- A popup will appear where you have to provide the **notebook name** and select the **region**.
You can give the notebook multiple properties according to your requirements
Enable the checkbox (install **NVIDIA GPU driver automatically for me**) before creating the notebook.



New notebook

Notebook name *

63 character limit with lowercase letters, digits or '-' only. Must start with a letter. Cannot end with a '-'.
?

Region *

ⓘ Restricted by GPU selection.

Zone *

?

Notebook properties

Environment ⓘ PyTorch 1.11 (with Intel® MKL-DNN/MKL)

Machine type 4 vCPUs, 15 GB RAM

GPUs ⓘ 1 NVIDIA Tesla T4

Boot disk 100 GB Standard persistent disk

Data disk 100 GB Standard persistent disk

Subnetwork

External IP Ephemeral(Automatic)

Permission Compute Engine default service account

Estimated cost ⓘ US\$281.58 monthly, \$0.386 hourly

☐ Install NVIDIA GPU driver automatically for me ⓘ
NVIDIA GPU drivers must be installed for the GPUs to work.

ADVANCED OPTIONS

CANCEL

CREATE

Machine type	GPUs	Owner	Last modified
CPUs, 15 GB RAM	NVIDIA Tesla K80 x 1	Service account	28 Apr 2022, 17:27:21
CPUs, 15 GB RAM	NVIDIA Tesla T4 x 1	Service account	29 Apr 2022, 00:06:00
custom-4-21760	None	Service account	28 Apr 2022, 17:27:21
CPUs, 15 GB RAM	None	Service account	9 May 2022, 04:54:34

On the notebook is created, A button (**OPEN JUPYTERLAB**) will appear at the right side of notebook name.

- Click on it and the notebook will open in a new tab.

- Vertex AI
- Dashboard
- Datasets
- Features
- Labelling tasks
- Workbench
- Pipelines
- Training
- Experiments
- Models
- Endpoints
- Batch predictions
- Metadata
- Matching engine
- Manage resources
- Marketplace

Workbench

MANAGED NOTEBOOKS USER-MANAGED NOTEBOOKS EXECUTIONS SCHEDULES

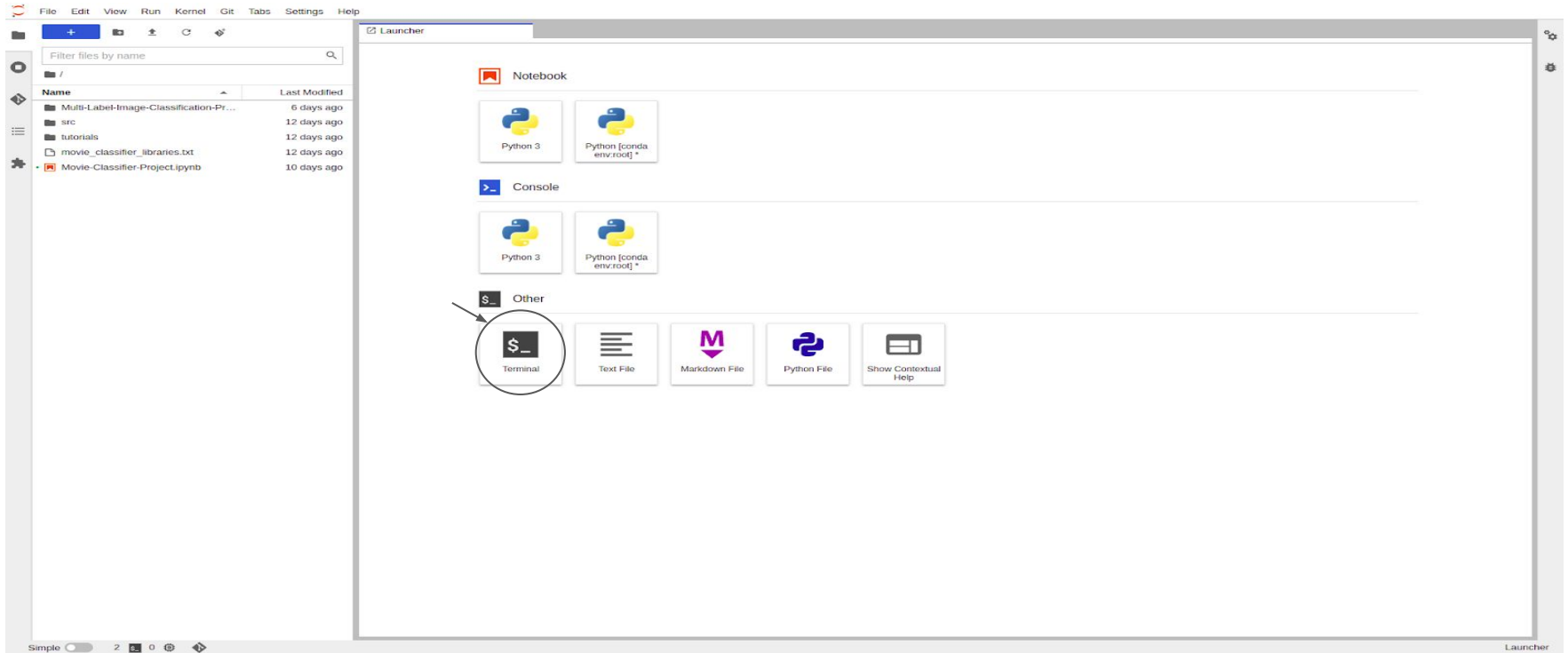
Notebooks have JupyterLab 3 pre-installed and are configured with GPU-enabled machine-learning frameworks. [Learn more](#)

Filter Enter property name or value

		Notebook name ↑		Zone	Auto-upgrade	Environment	Machine type	GPUs	Owner	Last modified
<input type="checkbox"/>	●		OPEN JUPYTERLAB	us-central1-a	—	PyTorch:1.8	4 vCPUs, 15 GB RAM	NVIDIA Tesla K80 x 1	Service account	28 Apr 2022, 17:27:21
<input type="checkbox"/>	●		OPEN JUPYTERLAB	us-central1-a	—	PyTorch:1.11	4 vCPUs, 15 GB RAM	NVIDIA Tesla T4 x 1	Service account	29 Apr 2022, 00:06:00
<input type="checkbox"/>	✓	movie-classifier		us-central1-a	—	PyTorch:1.4	e2-custom-4-21760	None	Service account	28 Apr 2022, 17:27:21
<input type="checkbox"/>	✓		OPEN JUPYTERLAB	us-west1-b	—	TensorFlow:2.3	4 vCPUs, 15 GB RAM	None	Service account	9 May 2022, 04:54:34

STEP 3: Move your code to Vertex AI

- Select **Terminal** under Other's tab.



- Inside the terminal, go to the `movie-classifier` directory and cd into it.
Then go to the terminal and create two new variables with the URI of your container image.
`IMAGE_URI_FOR_TRAINING = "gcr.io/$PROJECT_ID/movie-classifier-training_image:v1"`
`IMAGE_URI_FOR_DEPLOYMENT = "gcr.io/$PROJECT_ID/movie-classifier-deployment_image:v1"`
- Replace `$PROJECT_ID` with your Project ID

STEP 4: Create Docker files

- In order to create Training Docker file, tap on **New Launcher button (+)** and create a text file (under Other's tab).
- Rename it from untitled.txt to **DockerfileForTraining** and put the following code inside the file:

```
FROM gcr.io/cloud-aiplatform/training/pytorch-gpu.1-7
```

```
WORKDIR /
```

```
COPY ./movie_classifier_libraries.txt ./
```

```
RUN pip install -r movie_classifier_libraries.txt
```

```
RUN mkdir model
```

```
RUN mkdir images
```

```
COPY . .
```

```
# Sets up the entry point to invoke the trainer.
```

```
ENTRYPOINT ["python", "-m", "train"]
```

Filter files by name

/ Multi-Label-Image-Classification-Project / trainer /

Name	Last Modified
train.py	16 days ago

Multi-Label-Image-Classification-Project/trainer



Notebook



Python 3

Python [conda
env:root] *

Console



Python 3

Python [conda
env:root] *

Other



Terminal



Text File



Markdown File



Python File

Show Contextual
Help

- Similarly create the Docker file for deployment, rename it to **DockerfileForDeployment** and put following code inside it:

```
FROM python:3.7
```

```
WORKDIR /home/model-server/
```

```
COPY ./movie_classifier_libraries_deployment.txt /home/model-server/
```

```
RUN pip install -r movie_classifier_libraries_deployment.txt
```

```
RUN apt-get update
```

```
RUN apt-get install ffmpeg libsm6 libxext6 -y
```

```
ADD . /home/model-server/
```

```
RUN mkdir model
```

```
RUN python download_model_and_test_df.py
```

```
CMD exec gunicorn -b :5000 --max-requests 1 --graceful-timeout 300 -t 600 main:app
```

STEP 5: Dockerize your Code

- Now that we have docker files, we can create docker images and push them to container registries (e.g. Docker Hub or Amazon ECR). We will be ready to kick off a customer model training job as pushed to container registry.

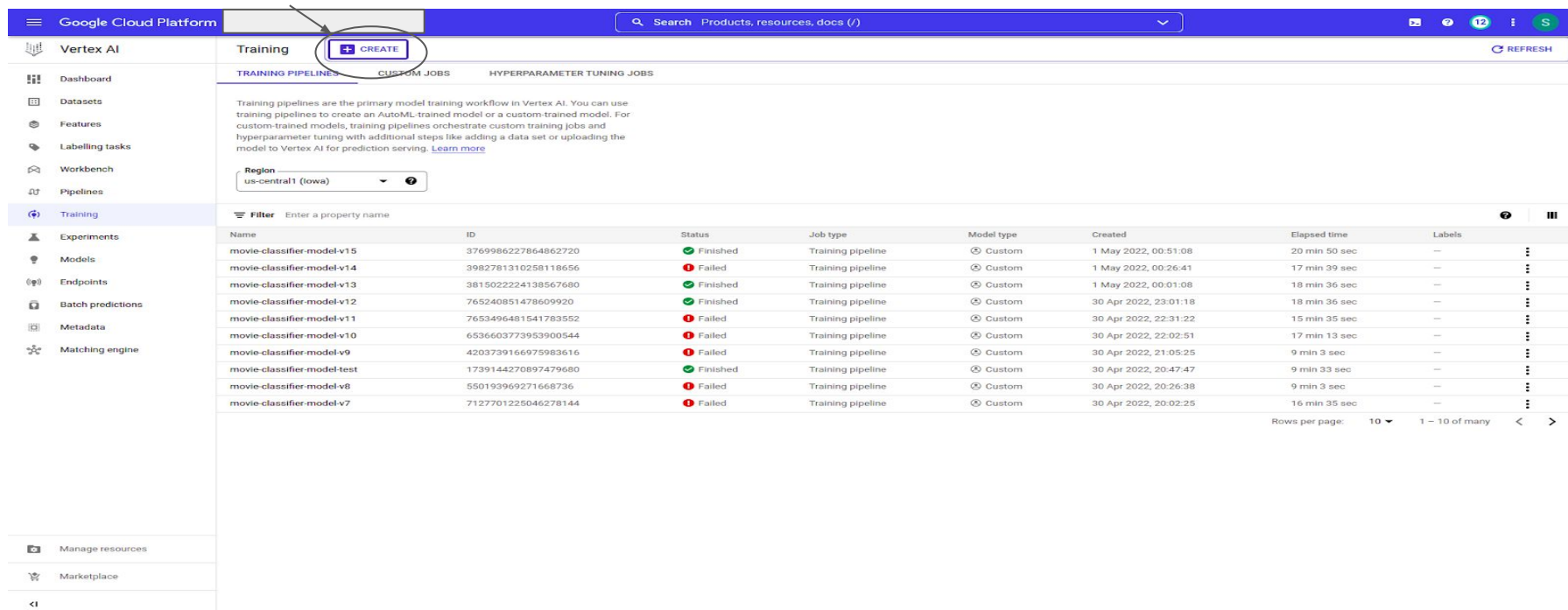
STEP 6: Cloud Storage Bucket Creation

Before starting the training job, you need to create a **Cloud Storage Bucket** and move your training data into it. You also need to create a **Cloud Storage Bucket** to store your trained model. Vertex will later use this bucket to create an exported Storage Bucket to deploy your model. Please make sure that the region selected while creating the model bucket should be same where you deploy your model.

- Run the following command in your **Workbench Terminal** to create a new bucket in your project:
- **PROJECT_ID="your-cloud-project-id"**
- **BUCKET_NAME="gs://\${PROJECT_ID}-bucket"**
- **gsutil mb -l us-central1 \$BUCKET_NAME**

STEP 7: Training Process and Configurations

Navigate to the **TRAINING** section in the Vertex AI of Cloud console and click **CREATE** button to enter the parameters for the training job.



The screenshot displays the Google Cloud Platform Vertex AI Training interface. The left sidebar shows the navigation menu with 'Training' selected. The main content area is titled 'Training' and includes a 'CREATE' button (circled in red) and tabs for 'TRAINING PIPELINES', 'CUSTOM JOBS', and 'HYPERPARAMETER TUNING JOBS'. Below the tabs, there is a description of training pipelines and a 'Region' dropdown menu set to 'us-central1 (Iowa)'. A table lists various training pipelines with their details.

Name	ID	Status	Job type	Model type	Created	Elapsed time	Labels
movie-classifier-model-v15	3769986227864862720	Finished	Training pipeline	Custom	1 May 2022, 00:51:08	20 min 50 sec	—
movie-classifier-model-v14	3982781310258118656	Failed	Training pipeline	Custom	1 May 2022, 00:26:41	17 min 39 sec	—
movie-classifier-model-v13	381502224138567680	Finished	Training pipeline	Custom	1 May 2022, 00:01:08	18 min 36 sec	—
movie-classifier-model-v12	765240851478609920	Finished	Training pipeline	Custom	30 Apr 2022, 23:01:18	18 min 36 sec	—
movie-classifier-model-v11	7653496481541783552	Failed	Training pipeline	Custom	30 Apr 2022, 22:31:22	15 min 35 sec	—
movie-classifier-model-v10	6536603773953900544	Failed	Training pipeline	Custom	30 Apr 2022, 22:02:51	17 min 13 sec	—
movie-classifier-model-v9	4203739166975983616	Failed	Training pipeline	Custom	30 Apr 2022, 21:05:25	9 min 3 sec	—
movie-classifier-model-test	1739144270897479680	Finished	Training pipeline	Custom	30 Apr 2022, 20:47:47	9 min 33 sec	—
movie-classifier-model-v8	550193969271668736	Failed	Training pipeline	Custom	30 Apr 2022, 20:26:38	9 min 3 sec	—
movie-classifier-model-v7	7127701225046278144	Failed	Training pipeline	Custom	30 Apr 2022, 20:02:25	16 min 35 sec	—

Rows per page: 10 1 – 10 of many

- Under **Data Custom training (managed dataset)** select your training method and click **Continue**.

Google Cloud Platform

Vertex AI

Training

TRAINING PIPELINES

Training pipeline training pipelines custom-trained model hyperparameter model to Vertex

Region us-central1 (los)

Filter Enter

Name

movie-classifier-r

movie-classifier-r

movie-classifier-r

movie-classifier-r

movie-classifier-r

movie-classifier-r

movie-classifier-r

movie-classifier-r

movie-classifier-r

movie-classifier-r

Manage resources

Marketplace

Train new model

- 1 Training method
- 2 Model details
- 3 Training container
- 4 Hyperparameters (optional)
- 5 Compute and pricing
- 6 Prediction container (optional)

START TRAINING CANCEL

Dataset * No managed dataset

Annotation set -

Objective Custom

Please refer to the pricing guide for more details (and available deployment options) for each method.

AutoML options are only available when you train with a managed data set.

Model training method

☐ AutoML
Train high quality models with minimal effort and machine learning expertise. Just specify how long you want to train. [Learn more](#)

☐ AutoML Edge
Train a model that can be exported for on-prem/on-device use. Typically has lower accuracy. [Learn more](#)

☒ Custom training (advanced)
Run your TensorFlow, scikit-learn and XGBoost training applications in the cloud. Train with one of Google Cloud's pre-built containers or use your own. [Learn more](#)

CONTINUE

- **Continue** In the next step, enter **Model name** (mandatory) and provide description (optional) and click

Google Cloud Platform

Vertex AI

Training

TRAINING PIPELINE

Training pipeline training pipelines custom-trained model hyperparameter model to Vertex.

Region us-central1 (los)

Filter Enter

Name

movie-classifier-m

movie-classifier-m

movie-classifier-m

movie-classifier-m

movie-classifier-m

movie-classifier-m

movie-classifier-m

movie-classifier-m

movie-classifier-m

movie-classifier-m

movie-classifier-m

Manage resources

Marketplace

Train new model

- ☒ Training method
- ☒ Model details
- 3 Training container
- 4 Hyperparameters (optional)
- 5 Compute and pricing
- 6 Prediction container (optional)

START TRAINING CANCEL

☒ Train new model
Creates a new model group and assigns the trained model as version 1

☐ Train new version
Trains model as a version of an existing model

Name *
Test Model Name

Description

ADVANCED OPTIONS

CONTINUE

- Now select **Custom container** option and browse your container image (that you pushed earlier) from container registry
- Also provide your model bucket directory where your model will be saved and click **Continue**.
- **Tuning checks** are not doing hyperparameter tuning so leave the **Enable hyperparameter tuning** checkbox unchecked and click **Continue**

- training Step 1: Select the **Machine type** of Accelerator Cloud and click **Continue**

Google Cloud Platform

Vertex AI

Train new model

- Training method
- Model details
- Training container
- Hyperparameters (optional)
- 5 Compute and pricing**
- 6 Prediction container (optional)

START TRAINING **CANCEL**

Model training pricing is based on the length of time spent training, machine types and any accelerators used. [Learn more](#)

Region: us-central1 (Iowa)

Compute settings

Select the type of virtual machine to use for your worker pool. You can add up to 4 worker pools. To learn about compute costs and how to map your ML framework's roles to specific worker pools, consult the [documentation](#)

Worker pool 0

Machine type *

Filter Type to filter

Standard	n1-standard-4 4 vCPUs, 15 GiB memory
High-memory	
High-CPU	n1-standard-8 8 vCPUs, 30 GiB memory
Custom	
High-GPU	n1-standard-16 16 vCPUs, 60 GiB memory
Mega-GPU	n1-standard-32 32 vCPUs, 120 GiB memory

CLEAR SELECTION

100

ADD MORE WORKER POOLS (OPTIONAL)

CONTINUE

Train new model

- ✓ Training method
- ✓ Model details
- ✓ Training container
- ✓ Hyperparameters (optional)
- 5 Compute and pricing
- 6 Prediction container (optional)

START TRAINING

CANCEL

Model training pricing is based on the length of time spent training, machine types and any accelerators used. [Learn more](#)

Region
us-central1 (Iowa)

Compute settings

Select the type of virtual machine to use for your worker pool. You can add up to 4 worker pools. To learn about compute costs and how to map your ML framework's roles to specific worker pools, consult the [documentation](#)

Worker pool 0

Machine type *
n1-standard-8, 8 vCPUs, 30 GiB memory

Accelerator type

- NVIDIA_TESLA_K80
- NVIDIA_TESLA_P100
- NVIDIA_TESLA_V100
- NVIDIA_TESLA_P4
- NVIDIA_TESLA_T4

Disk type
SSD

Disk size (GB)
100

ADD MORE WORKER POOLS (OPTIONAL)

CONTINUE

Train new model

- ✓ Training method
- ✓ Model details
- ✓ Training container
- ✓ Hyperparameters (optional)
- 5 Compute and pricing
- 6 Prediction container (optional)

START TRAINING

CANCEL

Model training pricing is based on the length of time spent training, machine types and any accelerators used. [Learn more](#)

Region
us-central1 (Iowa)

Compute settings

Select the type of virtual machine to use for your worker pool. You can add up to 4 worker pools. To learn about compute costs and how to map your ML framework's roles to specific worker pools, consult the [documentation](#)

Worker pool 0

Machine type *
n1-standard-8, 8 vCPUs, 30 GiB memory

Accelerator type
NVIDIA_TESLA_K80

Accelerators can speed up model training that involves intensive compute tasks. [Learn more](#)

Accelerator count
1

Disk size (GB)
100

ADD MORE WORKER POOLS (OPTIONAL)

CONTINUE

- Leave the **No prediction container** option selected and hit **START TRAINING** button to kick off the training job.

The screenshot shows the Google Cloud Platform interface for training a new model. The left sidebar contains the 'Vertex AI' menu with options like Dashboard, Datasets, Features, Labelling tasks, Workbench, Pipelines, Training (selected), Experiments, Models, Endpoints, Batch predictions, Metadata, and Matching engine. The main panel is titled 'Train new model' and shows a progress bar with steps: Training method, Model details, Training container, Hyperparameters (optional), Compute and pricing, and Prediction container (optional). The 'Prediction container (optional)' step is currently active and highlighted. Below the progress bar, there is a 'START TRAINING' button (circled in red) and a 'CANCEL' button. To the right, the 'No prediction container' option is selected with a radio button. A text box explains that you can associate a custom-trained model with a container for prediction requests, and provides links for more information. Other options include 'Pre-built container' (with a link to supported runtimes) and 'Custom container' (with a note about storing in Container Registry or Artifact Registry).

Google Cloud Platform

Vertex AI

Train new model

- Training method
- Model details
- Training container
- Hyperparameters (optional)
- Compute and pricing
- Prediction container (optional)**

START TRAINING CANCEL

You can associate your custom-trained model with a container in order to serve prediction requests using Vertex AI. [Learn more about getting predictions.](#)

☒ No prediction container
You can always import your model artifact later to serve prediction requests

☐ Pre-built container
View the list of [supported runtimes](#) including TensorFlow, scikit-learn and PyTorch versions

☐ Custom container
Build a custom Docker container. Must be stored in [Container Registry](#) or [Artifact Registry](#)

= The training will take around 15-20 minutes to train job. After changes from **Training** to **Finished** and then you can deploy your trained model to Vertex AI.

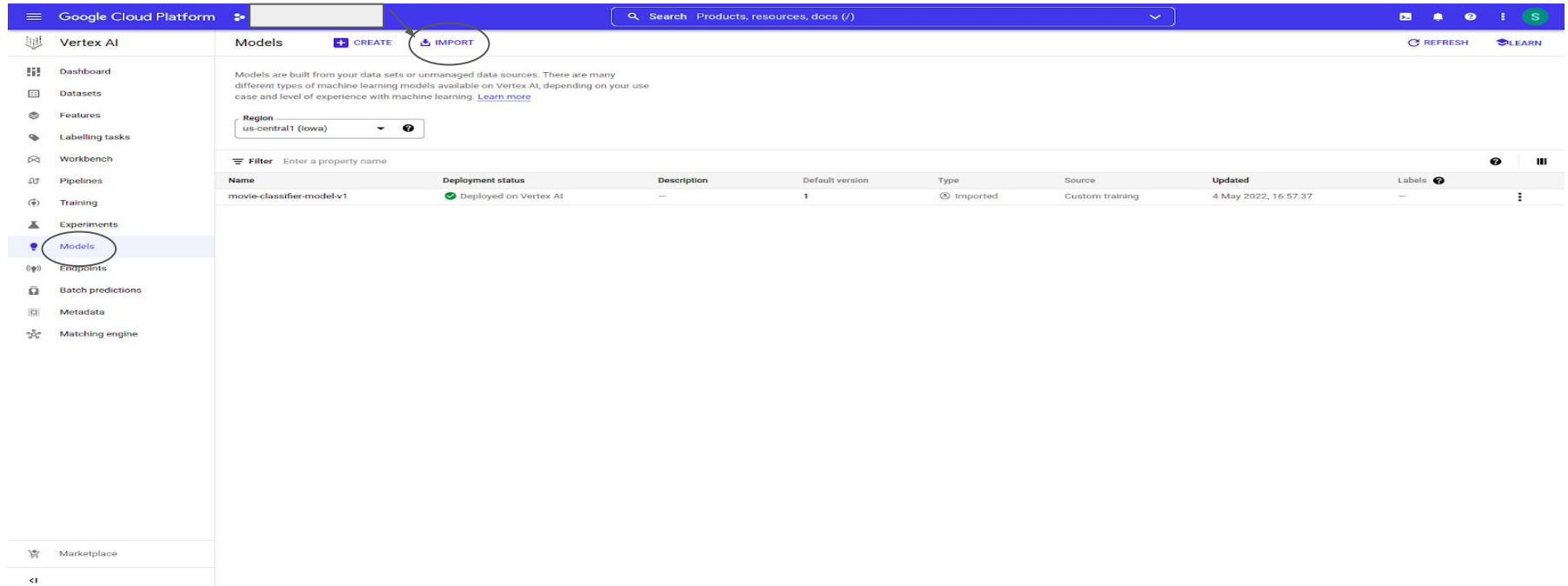
The screenshot shows the Google Cloud Platform interface for Vertex AI. The left sidebar contains navigation links: Vertex AI, Dashboard, Datasets, Features, Labelling tasks, Workbench, Pipelines, Training (highlighted), Experiments, Models, Endpoints, Batch predictions, Metadata, Matching engine, and Marketplace. The main content area is titled 'Training' and includes a 'CREATE' button. Below this, there are tabs for 'TRAINING PIPELINES', 'CUSTOM JOBS', and 'HYPERPARAMETER TUNING JOBS'. A text block explains that training pipelines are the primary model training workflow in Vertex AI. A 'Region' dropdown menu is set to 'us-central1 (Iowa)'. A filter bar shows 'Name: movie-classifier-model-new'. Below the filter is a table with the following data:

Name	ID	Status	Job type	Model type	Created	Elapsed time	Labels
movie-classifier-model-new	5855569570244329472	Finished	Training pipeline	Custom	11 May 2022, 03:37:13	19 min 35 sec	—

Annotations in the image include an arrow pointing to the 'Training' link in the sidebar and another arrow pointing to the 'Finished' status in the table.

STEP 8: Model Deployment

- To deploy the model on Vertex AI, go to the **MODEL** section of Vertex AI and click **IMPORT** button.



The screenshot shows the Google Cloud Platform Vertex AI console. The left sidebar contains a navigation menu with the following items: Vertex AI, Dashboard, Datasets, Features, Labelling tasks, Workbench, Pipelines, Training, Experiments, Models (highlighted with a red circle), Endpoints, Batch predictions, Metadata, and Matching engine. The main content area is titled 'Models' and includes a 'CREATE' button and an 'IMPORT' button (also highlighted with a red circle). Below the buttons, there is a 'Region' dropdown menu set to 'us-central1 (Iowa)'. A table lists the models, with one model 'movie-classifier-model-v1' shown. The table has columns for Name, Deployment status, Description, Default version, Type, Source, Updated, and Labels.

Name	Deployment status	Description	Default version	Type	Source	Updated	Labels
movie-classifier-model-v1	Deployed on Vertex AI	—	1	Imported	Custom training	4 May 2022, 16:57:37	—

- Enter **Model name** (mandatory) and provide **Description** (optional) and click **Continue**.

Google Cloud Platform

Vertex AI

Models

+ CREATE

IMPORT

Models are built from your data sets 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

Name	Deployment status	Description	Default
movie-classifier-model-v1	Deployed on Vertex AI	—	1

Import model

- Name and region
- Model settings
- Explainability (optional)

IMPORT **CANCEL**

You can import model artifacts that have been trained outside of Google Cloud. Once your model has been imported, you can serve it for online or batch predictions and compare it against your other Cloud AI models. [More info](#)

☒ Import as new model
Creates a new model group and assigns the imported model as version 1

☐ Import as new version
Imports the model as a version of an existing model

Name *
movie-classifier-model-v2

Description

Region: us-central1 (Iowa)

ADVANCED OPTIONS

CONTINUE

- ≡ **Next we select Import or existing custom container option and click on Docker Container**
Here we will add the Docker container following routes in Prediction route and Health route text
Prediction route=/get_movie_genres
Health route=/health
(We have specified these routes in our inference.py file and we will use them when we
- **Also we have set the port as 8080 to 5000 in Docker Container**
Once the container is up and running, leave the Docker container selected and hit **IMPORT button.**

Google Cloud Platform

Search

Products, resources, docs

Vertex AI

Models

CREATE

IMPORT

Dashboard

Datasets

Features

Labelling tasks

Workbench

Pipelines

Training

Experiments

Models

Endpoints

Batch predictions

Metadata

Matching engine

Marketplace

Models are built from your data sets 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

Name	Deployment status	Description	Default
movie-classifier-model-v1	Deployed on Vertex AI	—	1

Import model

- 1 Name and region
- 2 Model settings
- 3 Explainability (optional)

IMPORT

CANCEL

☐ Import model artifacts into a new pre-built container
View the list of [supported runtimes](#) including TensorFlow, scikit-learn and XGBoost versions

☒ Import an existing custom container
Build a custom Docker container. Must be stored in [Container Registry](#) or [Artifact Registry](#)

Custom container settings

Container image
gcr.io/pronto-project/movie-classifier-image@sha256:ace98e9b96604e **BROWSE**

The URI for the Docker image in Container Registry or Artifact Registry. You must have permission to access the image. [Learn more](#)

Command

Model artifact location (Cloud storage path)
☒ gs:// movie-classifier-model-bucket **BROWSE**

Path to the Cloud Storage directory where the exported model file is stored (not the path to the model file itself).

Arguments

Optional. Add arguments for the command that runs when the container starts. Overrides the container's CMD instruction. Enter one parameter and its argument per line.

```
--flag_a=xxxx  
-flag2  
flag3
```

Environment variables

Optional. Specifies the command that runs when the container starts. Overrides the container's ENTRYPOINT instructions. Enter one variable per line, separating the key and value with an equal (=) sign.

```
ENV_VAR_1=MY_ENV_VALUE  
ENV_VAR_2=MY_SECOND_ENV_VALUE
```

Prediction route

An HTTP path to send prediction requests to. Vertex AI will forward requests to the entered path on the container's IP address and port. If not set, defaults to a value set when the model is deployed to an endpoint. [Learn more about predict routes](#).

Google Cloud Platform

Search Products, resources, docs

Vertex AI

Dashboards

Datasets

Features

Labelling tasks

Workbench

Pipelines

Training

Experiments

Models

Endpoints

Batch predictions

Metadata

Matching engine

Marketplace

Models

CREATE

IMPORT

Models are built from your data sets 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

Name	Deployment status	Description	Default
movie-classifier-model-v1	<div>Deployed on Vertex AI</div>	—	1

Import model

1

Name and region

2

Model settings

3

Explainability (optional)

IMPORT

CANCEL

Environment variables

Optional. Specifies the command that runs when the container starts. Overrides the container's ENTRYPOINT instructions. Enter one variable per line, separating the key and value with an equal (=) sign.

```
ENV_VAR_1=MY_ENV_VALUE
ENV_VAR_2=MY_SECOND_ENV_VALUE
```

Prediction route

/get_movie_genres

An HTTP path to send prediction requests to. Vertex AI will forward requests to the entered path on the container's IP address and port. If not set, defaults to a value set when the model is deployed to an endpoint. [Learn more about predict routes](#).

Health route

/health

An HTTP path to send health checks to. Vertex AI occasionally sends GET requests to this path on the container's IP address and port to check that the container is healthy. If not set, defaults to a value set when the model is deployed to an endpoint. [Learn more about health routes](#).

Port

5000

Port to expose from the container. Prediction requests and health checks will be sent to the port. If left blank, the default port is 8080.

Predict schemata

Optional. [Learn more about the predict schemata](#)

gs:// Instances

BROWSE

Cloud Storage location to a YAML file that defines the format of a single instance used in prediction and explanation requests.

gs:// Parameters

BROWSE

Cloud Storage location to a YAML file that defines the prediction and explanation parameters.

gs:// Predictions

BROWSE

Cloud Storage location to a YAML file that defines the format of a single prediction or explanation.

CONTINUE

Google Cloud Platform

Vertex AI

Dashboard

Datasets

Features

Labelling tasks

Workbench

Pipelines

Training

Experiments

Models

Endpoints

Batch predictions

Metadata

Matching engine

Models

CREATEIMPORT

Models are built from your data sets 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

Name	Deployment status	Description	Default
movie-classifier-model-v1	Deployed on Vertex AI	—	1

Import model

- ✓ Name and region
- ✓ Model settings
- 3 Explainability (optional)

IMPORT

CANCEL

Explainability options

LEARN

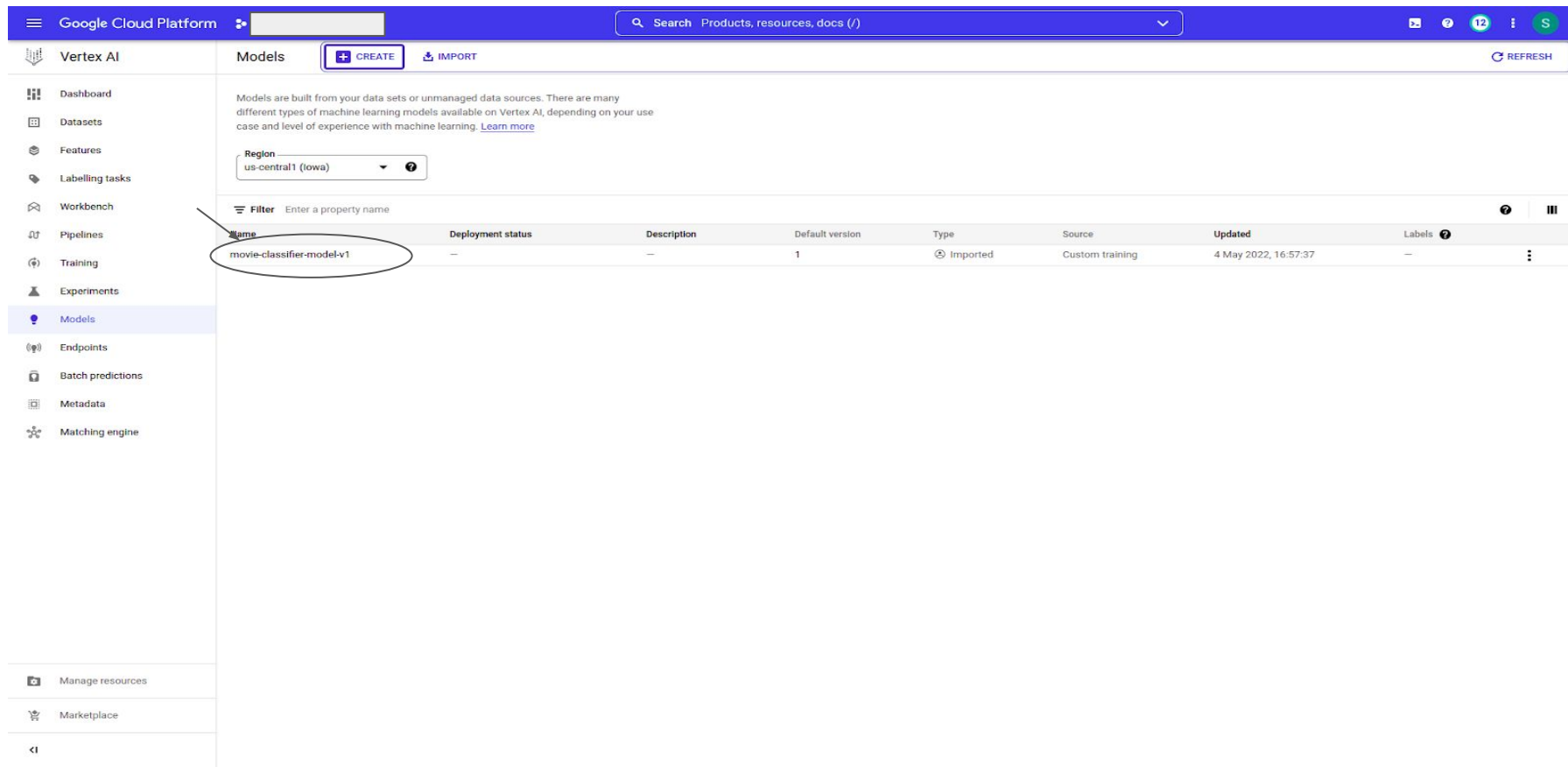
In Vertex AI, models are made explainable through feature attribution, which tells you how much each feature contributed to the predicted result. You can use this information to verify that the model is behaving as expected, recognise bias in your models and get ideas for ways to improve your model and your training data. Explainability will incur a minor additional cost. [Learn more](#)

Select a feature attribution method

Your model's data type determines which attribution methods are available to use. [Learn more about attribution methods](#)

- ☒ None
- ☐ Sampled Shapley (for tabular models)
- ☐ Integrated gradients (for tabular models)
- ☐ Integrated gradients (for image classification models)
- ☐ XRAI (for image classification models)

- After few seconds, the model will appear in **MODEL** section of Vertex AI and now you can deploy this model to an endpoint.



The screenshot shows the Google Cloud Platform interface for Vertex AI Models. The left sidebar contains navigation links: Vertex AI, Dashboard, Datasets, Features, Labelling tasks, Workbench, Pipelines, Training, Experiments, Models (highlighted), Endpoints, Batch predictions, Metadata, and Matching engine. The main content area is titled 'Models' and includes a 'CREATE' button and an 'IMPORT' button. Below this, there is a 'Region' dropdown set to 'us-central1 (Iowa)'. A table lists the models, with the first row 'movie-classifier-model-v1' circled. The table has columns for Name, Deployment status, Description, Default version, Type, Source, Updated, and Labels. A 'Filter' input is located above the table.

Name	Deployment status	Description	Default version	Type	Source	Updated	Labels
movie-classifier-model-v1	—	—	1	Imported	Custom training	4 May 2022, 16:57:37	—

STEP 9: Endpoint Deployment

- When your training job ends, Vertex AI creates a model endpoint for you. In order to use this model, you need to deploy a model endpoint. You can create a model endpoint for your model.
- Go to the **ENDPOINT** section of Vertex AI and then click on **CREATE AN ENDPOINT** button.
 - Give the **Endpoint name** like “movie-classifier-endpoint” and hit the **Continue** button.

Google Cloud Platform

Search

Products, resources, docs

Vertex AI

Dashboard

Datasets

Features

Labelling tasks

Workbench

Pipelines

Training

Experiments

Models

Endpoints

Batch predictions

Metadata

Matching engine

Manage resources

Marketplace

Endpoints

CREATE ENDPOINT

Endpoints are machine learning models made available for online prediction requests. Endpoints are useful for timely predictions from many users (for example, in response to an application request). You can also request batch predictions if you don't need immediate results.

To create an endpoint, you need at least one machine-learning model. [Learn more](#)

Region

us-central1 (Iowa)

Filter

Enter a property name

	Name	ID	Status	Models
<input type="checkbox"/>				0
<input type="checkbox"/>				0
<input type="checkbox"/>				0
<input type="checkbox"/>				0
<input type="checkbox"/>				0
<input type="checkbox"/>				0
<input type="checkbox"/>				0
<input type="checkbox"/>				0
<input type="checkbox"/>				0
<input type="checkbox"/>				0
<input type="checkbox"/>				0
<input type="checkbox"/>				0
<input type="checkbox"/>				0
<input type="checkbox"/>				0

New endpoint

1

Define your endpoint

2

Model settings

CREATE

CANCEL

Endpoint name *

movie-classifier-endpoint

Location

Region

us-central1 (Iowa)

Access

Determines how your endpoint can be accessed. By default, endpoints are available for prediction serving through a REST API. Endpoint access can't be changed after the endpoint is created.

☒ Standard

Makes the endpoint available for prediction serving through a REST API. AutoML and custom-trained models can be added to standard endpoints.

☐ Private

Create a private connection to this endpoint using a VPC network and [private services access](#). Only custom-trained and tabular models can be added to private endpoints. [Learn more](#)

ADVANCED OPTIONS

CONTINUE

Under **Model settings**, there is a model name drop down where all your created model list appears. Select the model that you created in the previous step and select its version. Once you do that, some more options will appear.

- Leave **Traffic split** at 100 and enter 1 for **Minimum number of compute nodes**.
- Under **Machine type**, select **n1-standard-2** (or any machine type you'd like).

From **Service account** drop-down, select a service account that has all the permissions that your code requires.

- Leave the rest of the defaults selected and then click **Continue**.

Google Cloud Platform

Search Products, resources, docs

Vertex AI

Dashboard

Datasets

Features

Labelling tasks

Workbench

Pipelines

Training

Experiments

Models

Endpoints

Batch predictions

Metadata

Matching engine

Manage resources

Marketplace

Endpoints

CREATE ENDPOINT

Endpoints are machine learning models made available for online prediction requests. Endpoints are useful for timely predictions from many users (for example, in response to an application request). You can also request batch predictions if you don't need immediate results.

To create an endpoint, you need at least one machine-learning model. [Learn more](#)

Region
us-central1 (Iowa)

Filter Enter a property name

	Name	ID	Status	Models
<input type="checkbox"/>				0
<input type="checkbox"/>				0
<input type="checkbox"/>				0
<input type="checkbox"/>				0
<input type="checkbox"/>				0
<input type="checkbox"/>				0
<input type="checkbox"/>				0
<input type="checkbox"/>				0
<input type="checkbox"/>				0
<input type="checkbox"/>				0
<input type="checkbox"/>				0
<input type="checkbox"/>				0
<input type="checkbox"/>				0
<input type="checkbox"/>				0

New endpoint

Define your endpoint

Model settings

Model monitoring

CREATE CANCEL

Add model

Model name *
movie-classifier-model-v1

Version
Version 1

Traffic split *
100 % ?

Compute resources

Choose how compute resources will serve prediction traffic to your model

Autoscaling: If you set a minimum and maximum, compute nodes will scale to meet traffic demand within those boundaries

No scaling: If you only set a minimum, then that number of compute nodes will always run regardless of traffic demand (the maximum will be set to minimum)

Once scaling settings are set, they can't be changed unless you redeploy the model. [Pricing guide](#)

Minimum number of compute nodes *
1

Default is 1. If set to 1 or more, then compute resources will continuously run even without traffic demand. This can increase cost but avoid dropped requests due to node initialisation.

Maximum number of compute nodes (optional)

Enter a number equal to or greater than the minimum nodes. Can reduce costs but may cause reliability issues for high traffic.

ADVANCED SCALING OPTIONS

Machine type *
n1-standard-2, 2 vCPUs, 7.5 GiB memory

Accelerator type
NVIDIA_TESLA_K80

Accelerator count
1

Service account
App Engine default service account

A service account determines what Google Cloud resources your service code can access. By default, a Google-managed service account is used with permissions appropriate for most models. You can also use a user-managed service account to customise permissions. [Learn more](#).

= We won't enable model monitoring for this model, so next click **Create** to kick off the endpoint deployment and you can then use it as the prediction, or we appear in Vertex AI Endpoints

Google Cloud Platform

Search Products, resources, docs (/)

REFRESH

LEARN

Vertex AI

Endpoints [+ CREATE ENDPOINT](#)

Dashboard

Datasets

Features

Labelling tasks

Workbench

Pipelines

Training

Experiments

Models

Endpoints

Batch predictions

Metadata

Matching engine

Endpoints are machine learning models made available for online prediction requests. Endpoints are useful for timely predictions from many users (for example, in response to an application request). You can also request batch predictions if you don't need immediate results.

To create an endpoint, you need at least one machine-learning model [Learn more](#)

Endpoint "movie-classifier-endpoint" encountered an error and can't be created. It will automatically be deleted in 30 days. [DISMISS](#)

Region
us-central1 (Iowa)

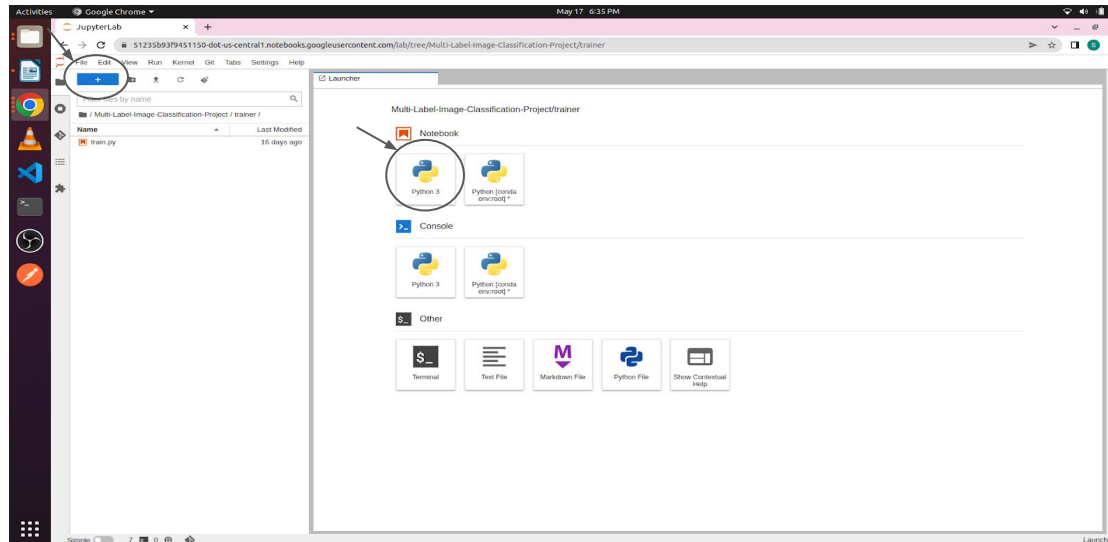
Filter Enter a property name

	Name	ID	Status	Models	Region	Monitoring	Most recent alerts	Last updated	API	Notification	Labels
<input type="checkbox"/>	movie-classifier-endpoint	7602102559280463872	Active	0	us-central1	Disabled	—	12 May 2022, 03:20:54	Sample request		
<input type="checkbox"/>				0	us-central1	Disabled	—	3 May 2022, 17:49:20	—		
<input type="checkbox"/>				0	us-central1	Disabled	—	2 May 2022, 21:21:41	—		
<input type="checkbox"/>				0	us-central1	Disabled	—	2 May 2022, 20:59:52	—		
<input type="checkbox"/>				0	us-central1	Disabled	—	2 May 2022, 19:59:24	—		
<input type="checkbox"/>				0	us-central1	Disabled	—	2 May 2022, 19:27:31	—		
<input type="checkbox"/>				0	us-central1	Disabled	—	2 May 2022, 08:00:25	—		
<input type="checkbox"/>				0	us-central1	Disabled	—	2 May 2022, 07:34:16	—		
<input type="checkbox"/>				0	us-central1	Disabled	—	2 May 2022, 07:08:50	—		
<input type="checkbox"/>				0	us-central1	Disabled	—	1 May 2022, 07:04:00	—		
<input type="checkbox"/>				0	us-central1	Disabled	—	1 May 2022, 06:01:50	—		
<input type="checkbox"/>				0	us-central1	Disabled	—	1 May 2022, 05:24:10	—		
<input type="checkbox"/>				0	us-central1	Disabled	—	1 May 2022, 05:24:08	—		
<input type="checkbox"/>				0	us-central1	Disabled	—	1 May 2022, 02:32:41	—		

Rows per page: 20 1 – 14 of 14

STEP 10: Get predictions on the deployed model

- We'll get predictions on our trained model from a Python notebook, using the **Vertex Python API**.
- Go back to your notebook instance and create a Python 3 notebook from the Launcher. **Also upload the image that you want to use for prediction.**



- In your notebook, run the following in a cell to install the Vertex AI SDK:

```
!pip3 install google-cloud-aiplatform --upgrade --user
```

- Then add a cell in your notebook to import the SDK and create a reference to the endpoint you just deployed:

```
from google.cloud import aiplatform
```

```
PROJECT_ID = "XXXXXXXXX"
```

```
ENDPOINT_ID = "4325213563224324234"
```

```
endpoint = aiplatform.Endpoint(
```

```
    endpoint_name=f"projects/{PROJECT_ID}/locations/us-central1/endpoints/{ENDPOINT_ID}")
```

- **Replace PROJECT_ID with your project id and ENDPOINT_ID with your endpoint id.**

- You can find your endpoint ID in the Endpoints section of the Vertex AI.

Google Cloud Platform

Search Products, resources, docs (/)

Vertex AI

Endpoints [+ CREATE ENDPOINT](#)

Endpoints are machine learning models made available for online prediction requests. Endpoints are useful for timely predictions from many users (for example, in response to an application request). You can also request batch predictions if you don't need immediate results.

To create an endpoint, you need at least one machine-learning model. [Learn more](#)

Region: us-central1 (Iowa)

Filter: Enter a property name

<input type="checkbox"/>	Name	ID	Status	Models	Region	Monitoring	Most recent alerts	Last updated ↓	API	Notification	Labels
<input type="checkbox"/>	movie-classifier-endpoint	7602102559280463872	Active	1	us-central1	Disabled	—	12 May 2022, 03:35:02	Sample request		

Marketplace

- Finally, make a prediction to your endpoint by copying and running the code below in a new cell:

```
import base64

encoded_string = ""

with open("movie-poster-image.jpg", "rb") as image_file:

    encoded_string = base64.b64encode(image_file.read())

instance = [{"b64_string": encoded_string.decode('utf-8')}]

prediction = endpoint.predict(instances=instance)

print(prediction.predictions[0])
```

- Run this cell, and you should see the following prediction output:

```
{'predicted_genres': ['Comedy', 'Drama', 'Romance']}
```



You've learned how to use Vertex AI to:

- Train a model by providing the training code in a custom container. You used a Pytorch model in this example, but you can train a model built with any framework using custom containers.
- Deploy a Pytorch model using a custom container as part of the same workflow you used for training.
- Create a model endpoint and generate a prediction.