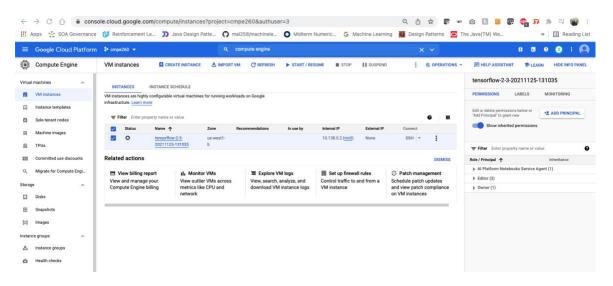# Vertex AI: Training and serving a custom model
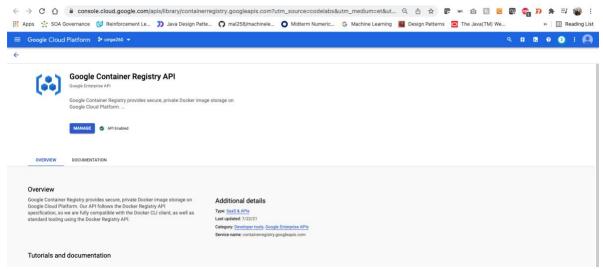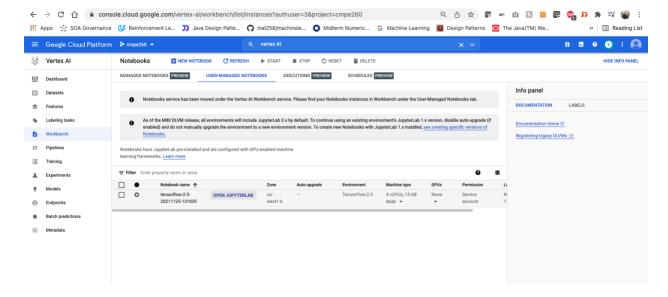
Setting up the environment

1) Enable Compute Engine:
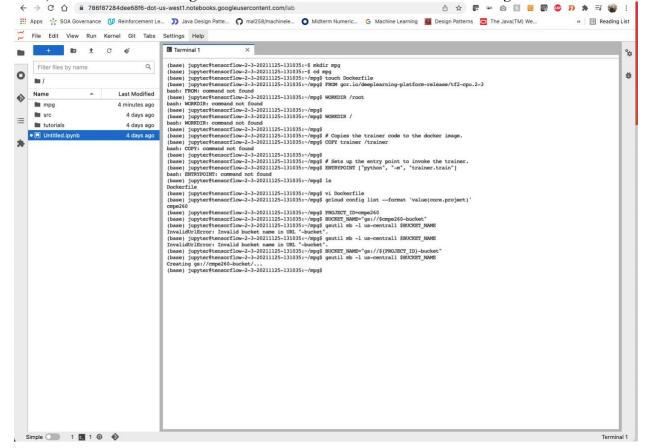


2) Enable Google Container Registry API



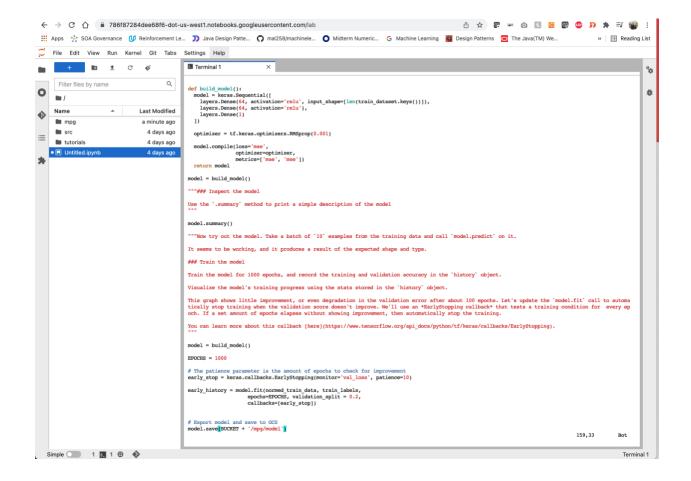3) Enable workbench

4) Containerize training code and create a Docker file and create a storage bucket



5) Create training code:

786f87284dee68f6-dot-us-west1.notebooks.googleusercontent.com/lab

Apps   SOA Governance   Reinforcement Le...   Java Design Patte...   mal258/machinele...   Midterm Numeric...   Machine Learning   Design Patterns   The Java(TM) We...   »   Reading List

File   Edit   View   Run   Kernel   Git   Tabs   Settings   Help

Terminal 1                                                ✕

```python
def build_model():
    model = keras.Sequential([
        layers.Dense(64, activation='relu', input_shape=[len(train_dataset.keys())]),
        layers.Dense(64, activation='relu'),
        layers.Dense(1)
    ])

    optimizer = tf.keras.optimizers.RMSprop(0.001)

    model.compile(loss='mse',
                  optimizer=optimizer,
                  metrics=['mae', 'mse'])
    return model

model = build_model()

"""### Inspect the model

Use the `.summary` method to print a simple description of the model
"""

model.summary()

"""Now try out the model. Take a batch of `10` examples from the training data and call `model.predict` on it.

It seems to be working, and it produces a result of the expected shape and type.

### Train the model

Train the model for 1000 epochs, and record the training and validation accuracy in the `history` object.

Visualize the model's training progress using the stats stored in the `history` object.

This graph shows little improvement, or even degradation in the validation error after about 100 epochs. Let's update the `model.fit` call to automatically stop training when the validation score doesn't improve. We'll use an *EarlyStopping callback* that tests a training condition for  every epoch. If a set amount of epochs elapses without showing improvement, then automatically stop the training.

You can learn more about this callback [here](https://www.tensorflow.org/api_docs/python/tf/keras/callbacks/EarlyStopping).
"""

model = build_model()

EPOCHS = 1000

# The patience parameter is the amount of epochs to check for improvement
early_stop = keras.callbacks.EarlyStopping(monitor='val_loss', patience=10)

early_history = model.fit(normed_train_data, train_labels,
                          epochs=EPOCHS, validation_split = 0.2,
                          callbacks=[early_stop])


# Export model and save to GCS
model.save(BUCKET + '/mpg/model')
```

                                                                    159,33        Bot

## 6) **Build and test the container locally**

7) Run a training job in vertex AI

Apps   SOA Governance   Reinforcement Le...   Java Design Patte...   mal258/machinele...   Midterm Numeric...   Machine Learning   Design Patterns   The Java(TM) We...   »   Reading List

**Google Cloud Platform**

**Vertex AI**

Dashboard
Datasets
Features
Labeling tasks
Workbench
Pipelines
Training
Experiments
Models
Endpoints
Batch predictions
Metadata

**Train new model**

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

START TRAINING   CANCEL

Select a pre-built container or build a custom container using ML frameworks (as well as non-ML dependencies, libraries and binaries) that are not otherwise supported. Learn more

○ Pre-built container
View the list of supported runtimes including TensorFlow and scikit-learn versions

◉ Custom container
Build a custom Docker container. Must be stored in Container Registry

**Custom container settings**

Container image *                                    BROWSE
❶ Container image URL is required.

gs:// Model output directory                         BROWSE

Your model artifacts and other data needed for training will be stored on Cloud Storage. You should specify a path here if you do not set an output directory in your application code or arguments.

**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

**Select container image**                           ✕

CONTAINER REGISTRY   ARTIFACT REGISTRY

Project: cmpe260   CHANGE

▼ gcr.io/cmpe260/mpg

d41cc7cdde   v1                                       3 minutes ago

SELECT   CANCEL

---

Apps   SOA Governance   Reinforcement Le...   Java Design Patte...   mal258/machinele...   Midterm Numeric...   Machine Learning   Design Patterns   The Java(TM) We...   »   Reading List

**Google Cloud Platform**

**Vertex AI**

Dashboard
Datasets
Features
Labeling tasks
Workbench
Pipelines
Training
Experiments
Models
Endpoints
Batch predictions
Metadata

**Train new model**

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

**Pre-built container settings**

Vertex AI provides Docker container images for serving predictions. To use a pre-built container, your trained model code must be in Python 3.7. Learn more about pre-built containers

In order to run in a pre-built container, your code needs to be in Python 3.7

Model framework *
TensorFlow                                           ▼

Model framework version *
2.1                                                  ▼

Accelerator type *
None                                                 ▼

Model directory *
gs:// cmpe260-bucket/mpg                             BROWSE
Cloud Storage location containing the model artifact and any supporting files

**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.

Deploy a model endpoint and get prediction on deployed model