November 5, 2018

# Managed ML Model Training and Serving

Lisa Quera, ML Specialist

Google Cloud

# Machine Learning at Scale

**Cloud ML Engine**

- Serverless, no-ops, ML training and serving platform

- Distributed training infrastructure that supports CPUs, GPUs, and TPUs

- Automatic hyperparameter tuning

- Train, tune, and serve TensorFlow models (batch and online prediction)

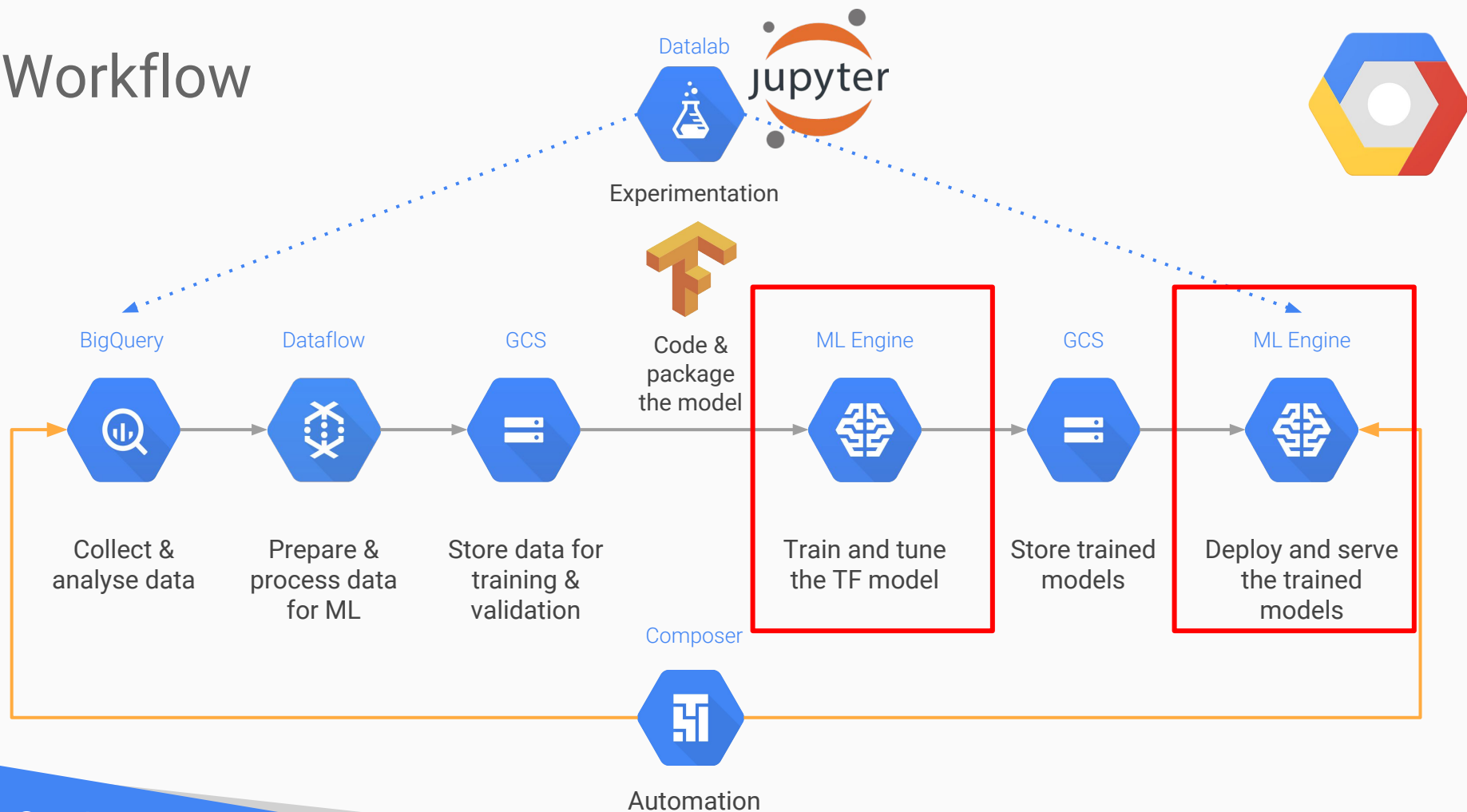- Train and serve Scikit-learn and XGBoost models for online predictions

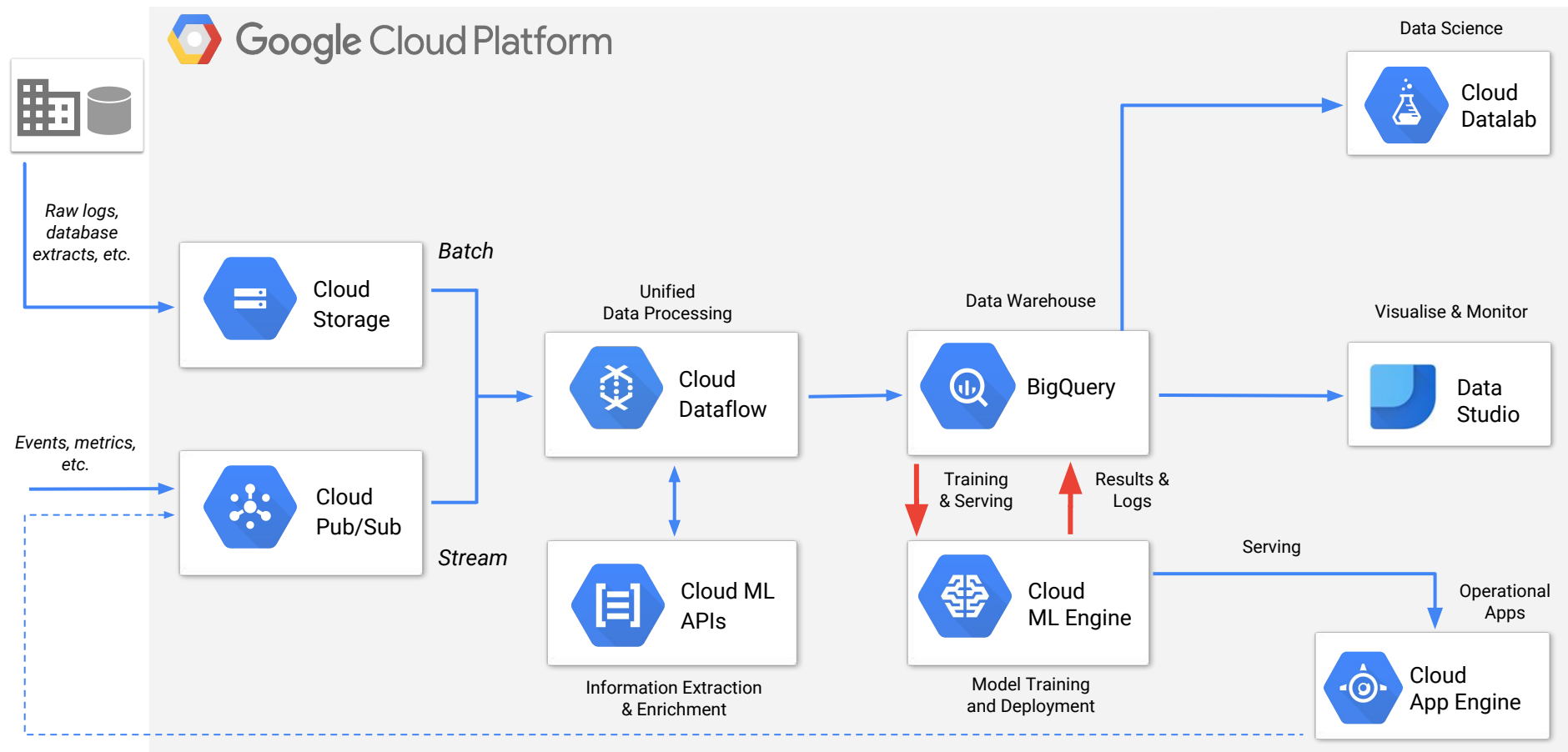- TPUs for training TensorFlow models **(Beta)**

Google Cloud

Workflow

# Example Architecture: BigQuery to deployed model API

# Training locally

*train locally*

*Local path*

*output directory*

```
gcloud ml-engine local train \
    --module-name trainer.task --package-path trainer/ \
    -- \
    --train-files $TRAIN_DATA --eval-files $EVAL_DATA --job-dir $MODEL_DIR
```

*training data*

*evaluation data*

Google Cloud

# Training in the cloud

## with single node

*train in the cloud*

*GCS location*

*Can be a package in GCS*

*region*

```
gcloud ml-engine jobs submit training $JOB_NAME --job-dir $OUTPUT_PATH \
    --runtime-version 1.10 --module-name trainer.task --package-path trainer --region $REGION \
    --scale-tier BASIC
    --train-files $TRAIN_DATA --eval-files $EVAL_DATA --num-epoch 1000 --learning-rate 0.01
```

*single worker*

*GCS locations*

*model-specific params*

Google Cloud

# Training in the cloud at scale

with multiple workers

```
gcloud ml-engine jobs submit training $JOB_NAME --job-dir $OUTPUT_PATH \
    --runtime-version 1.10 --module-name trainer.task --package-path trainer --region $REGION \
    --scale-tier STANDARD_1
    -- \
    --train-files $TRAIN_DATA --eval-files $EVAL_DATA
```

*distributed*

Google Cloud

# Manually Distributing the Training

```python
with tf.device("/job:ps/task:0"):
  weights_1 = tf.Variable(...)
  biases_1 = tf.Variable(...)

with tf.device("/job:ps/task:1"):
  weights_2 = tf.Variable(...)
  biases_2 = tf.Variable(...)

with tf.device("/job:worker/task:7"):
  input, labels = ...
  layer_1 = tf.nn.relu(tf.matmul(input, weights_1) + biases_1)
  logits = tf.nn.relu(tf.matmul(layer_1, weights_2) + biases_2)
  # ...
  train_op = ...

with tf.Session("grpc://worker7.example.com:2222") as sess:
  for _ in range(10000):
    sess.run(train_op)
```

**tf.train.ClusterSpec** construction

```python
tf.train.ClusterSpec({"local": ["localhost:2222", "localhost:2223"]

tf.train.ClusterSpec({
    "worker": [
        "worker0.example.com:2222",
        "worker1.example.com:2222",
        "worker2.example.com:2222"
    ],
    "ps": [
        "ps0.example.com:2222",
        "ps1.example.com:2222"
    ]})
```

# Training in the cloud at scale

## with GPUs (K80/P100/V100 - *availability by region*)

```
gcloud ml-engine jobs submit training $JOB_NAME --job-dir $OUTPUT_PATH \
    --runtime-version 1.10 --module-name trainer.task --package-path trainer --region $REGION \
    --scale-tier BASIC_GPU
    -- \
    --train-files $TRAIN_DATA --eval-files $EVAL_DATA
```

*single GPU*

Google Cloud

# CMLE Machine Types

| Cloud ML Engine scale tier | |
|---|---|
| **BASIC** | A single worker instance. This tier is suitable for learning how to use Cloud ML Engine and for experimenting with new models using small datasets.<br><br>**Compute Engine machine name**: n1-standard-4 |
| **STANDARD_1** | One master instance, plus four workers and three parameter servers.<br><br>**Compute Engine machine name, master**: n1-highcpu-8, **workers**: n1-highcpu-8, **parameter servers**: n1-standard-4 |
| **PREMIUM_1** | One master instance, plus 19 workers and 11 parameter servers.<br><br>**Compute Engine machine name, master**: n1-highcpu-16, **workers**: n1-highcpu-16, **parameter servers**: n1-highmem-8 |
| **BASIC_GPU** | A single worker instance with a single NVIDIA Tesla K80 GPU. To learn more about graphics processing units (GPUs), see the section on training with GPUs.<br><br>**Compute Engine machine name**: n1-standard-8 with one k80 GPU |
| **BASIC_TPU** *(Beta)* | A master VM and a Cloud TPU. See how to use TPUs for your training job.<br><br>**Compute Engine machine name, master**: n1-standard-4, **workers**: Cloud TPU |
| **CUSTOM** | The CUSTOM tier is not a set tier, but rather enables you to use your own cluster specification. When you use this tier, set values to configure your processing cluster according to these guidelines: |

Google

# CMLE Custom Options

| Cloud ML Engine machine name | |
|---|---|
| `standard` | A basic machine configuration suitable for training simple models with small to moderate datasets.<br><br>**Compute Engine machine name:** n1-standard-4 |
| `large_model` | A machine with a lot of memory, specially suited for parameter servers when your model is large (having many hidden layers or layers with very large numbers of nodes).<br><br>**Compute Engine machine name:** n1-highmem-8 |
| `complex_model_s` | A machine suitable for the master and workers of the cluster when your model requires more computation than the standard machine can handle satisfactorily.<br><br>**Compute Engine machine name:** n1-highcpu-8 |
| `complex_model_m` | A machine with roughly twice the number of cores and roughly double the memory of complex_model_s.<br><br>**Compute Engine machine name:** n1-highcpu-16 |
| `complex_model_l` | A machine with roughly twice the number of cores and roughly double the memory of complex_model_m.<br><br>**Compute Engine machine name:** n1-highcpu-32 |
| `standard_gpu` | A machine equivalent to standard that also includes a single NVIDIA Tesla K80 GPU.<br><br>**Compute Engine machine name:** n1-standard-8 with one k80 GPU |
| `complex_model_m_gpu` | A machine equivalent to complex_model_m that also includes four NVIDIA Tesla K80 GPUs.<br><br>**Compute Engine machine name:** n1-standard-16-k80x4 |

Google Cloud

# Cloud TPUs

**Cloud TPUs**

$ Excellent performance / $

Train in days instead of weeks

No more fighting with drivers

Flexibility and scale

Fully-managed in the cloud

Google Cloud

https://github.com/tensorflow/tpu-demos

# Supported models for TPUs

## Image recognition & object detection

**Image recognition:**
AmoebaNet-D
ResNet-50/101/152/200
Inception v2/v3/v4
DenseNet

**Object detection:**
RetinaNet

**Low-resource models:**
MobileNet
SqueezeNet

## Machine translation and language modeling

**Models:**
Machine translation
Language modeling
Sentiment analysis
Question-answering
(all transformer-based)

## Speech recognition

**Model:**
ASR Transformer
(LibriSpeech)

## Image generation

**Models:**
Image Transformer
DCGAN

Google Cloud

# Training Scikit-learn & XGBoost models

*Key parameters to scikit-learn and XGBoost training on CMLE*

```
gcloud ml-engine jobs submit training $JOB_NAME --job-dir $OUTPUT_PATH \
    --runtime-version 1.9 --python-version 2.7 --scale-tier BASIC \
    --module-name sklearn_trainer.task --package-path sklearn_trainer --region $REGION \
    -- \
    --train-files $TRAIN_DATA --eval-files $EVAL_DATA --num-epoch 1000 --learning-rate 0.01
```

Google Cloud

XGBoost
scikit learn

# Keras: what is happening?

- Compatibility module introduced in TensorFlow: tf.keras

- Write your custom estimator model_fn using tf.keras.layers and/or tf.layers (mix-and-match)

- Convert your compiled keras model to tf.estimator using tf.keras.estimator.model_to_estimator

- With TensorFlow, Keras users gain access to new features:
  - Distributed training
  - Multiple GPUs
  - Cloud ML
  - Hyperparameter tuning
  - TF-Serving

K Keras

# Keras example

```python
y_train = ...
y_test = ...

model = Sequential()
model.add(Conv2D(32, kernel_size=(3, 3), activation='relu', input_shape=input_shape))
model.add(Conv2D(64, (3, 3), activation='relu'))
model.add(MaxPooling2D(pool_size=(2, 2)))
model.add(Dropout(0.25))
model.add(Flatten())
model.add(Dense(128, activation='relu'))
model.add(Dropout(0.5))
model.add(Dense(num_classes, activation='softmax'))

model.compile(loss=keras.losses.categorical_crossentropy,
            optimizer=keras.optimizers.Adadelta(), metrics=['accuracy'])

model.fit(x_train, y_train, batch_size=batch_size, epochs=epochs,
        validation_data=(x_test, y_test))
score = model.evaluate(x_test, y_test, verbose=0)
```

Google Cloud

K Keras

# Premade Estimators

**TensorFlow**

Google Cloud

LinearClassifier

LinearRegressor

DNNClassifier

DNNRegressor

DNNLinearCombinedClassifier

(tf.contrib.estimator)RNNClassifier

DNNLinearCombinedRegressor

(tf.contrib.estimator)RNNEstimator

(tf.contrib.kernel_methods.) KernelLinearClassifier

(tf.contrib.factorization.) KMeansClustering

(tf.contrib.timeseries.) ARRegressor

BoostedTreesClassifier

BoostedTreesRegressor

# Hyperparameter tuning

```
gcloud ml-engine jobs submit training $JOB_NAME --job-dir $OUTPUT_PATH \

    --runtime-version 1.7 --module-name trainer.task --package-path trainer/ --region $REGION \

    --scale-tier PREMIUM_1 --config hyperparams.yaml

    -- \

    --train-files $TRAIN_DATA --eval-files $EVAL_DATA
```

*hypertuning*

# Hyperparameter tuning

## hyperparams.yaml

```yaml
trainingInput:
 hyperparameters:
   goal: MAXIMIZE
   hyperparameterMetricTag: accuracy
   maxTrials: 40
   enableTrialEarlyStopping: True
   maxParallelTrials: 2
   algorithm: UNSPECIFIED
   params:
    - parameterName: learning-rate
      type: FLOAT
      minValue: 0.001
      maxValue: 0.1
      scaleType: UNIT_LOG_SCALE
...
```

## task.py

```python
...
        # Initialise the optimizer for the DNN
        optimizer = tf.train.AdagradOptimizer(
            learning_rate=hparams.learning_rate)
...

parser.add_argument(
    '--learning-rate',
    help='Learning rate used by the DNN optimizer',
    default=0.01,
    type=float
 )
...
```

# Distributed Training

Consuming the deployed ML model API for predictions

Google Cloud

# Deploy the trained TF model

gcloud command line tool:

```
# Creating model
NAME=demo_classifier
gcloud ml-engine models create $NAME --regions $REGION

# Creating versions\
VERSION=v2.3
MODEL_DIR=gs://ksalama-gcs/trained_models/demo_classifier_output
gcloud ml-engine versions create $VERSION --model $NAME --origin $MODEL_DIR \
  --runtime-version 1.7 --config config.yaml

# List deployed models
gcloud ml-engine models list
```

```
description: A free-form description of the version.
deploymentUri: gs://path/to/source
runtimeVersion: '1.7'
manualScaling:
  nodes: 10
autoScaling:
  minNodes: 0
```

Google Cloud

# Predicting with TF Model

## Online versus Batch Prediction

Cloud ML Engine provides two ways to get predictions from trained models: *online prediction* (sometimes called HTTP prediction), and *batch prediction*. In both cases, you pass input data to a cloud-hosted machine-learning model and get inferences for each data instance. The differences are shown in the following table:

| Online prediction | Batch prediction |
| --- | --- |
| Optimized to minimize the latency of serving predictions. | Optimized to handle a high volume of instances in a job and to run more complex models. |
| Can process one or more instances per request. | Can process one or more instances per request. |
| Predictions returned in the response message. | Predictions written to output files in a Cloud Storage location that you specify. |
| Input data passed directly as a JSON string. | Input data passed indirectly as one or more URIs of files in Cloud Storage locations. |
| Returns as soon as possible. | Asynchronous request. |
| Accounts with the following IAM roles can request online predictions:<br><br>• Legacy Editor or Viewer<br>• Cloud ML Engine Admin or Developer | Accounts with the following IAM roles can request batch predictions:<br><br>• Legacy Editor<br>• Cloud ML Engine Admin or Developer |
| Runs on the runtime version and in the region selected when you deploy the model. | Can run in any available region, using any available runtime version. Though you should run with the defaults for deployed model versions. |

Google Cloud

# Predicting with TF Model

gcloud command line tool:

```
gcloud ml-engine predict --model $NAME --version $VERSION --json-instances test.json
```

gcloud **batch prediction**:

```
gcloud ml-engine job submit prediction \

$JOB_NAME --model $NAME --version $VERSION \

data-format TEXT \

input-paths $GCS_DATA_DIR \

output-path $GCS_OUT_DIR \
```

Google Cloud

python code - REST API call - **online prediction**:

```python
def estimate(project, model_name, version, instances):

    credentials = GoogleCredentials.get_application_default()
    api = discovery.build('ml', 'v1',
            credentials=credentials,
            discoveryServiceUrl =
            'https://storage.googleapis.com/cloud-ml/discovery/ml_v1_discovery.json')

    request_data = {'instances': instances}

    model_url = 'projects/{}/models/{}/versions/{}'.format(
                project,model_name, version)

    response = api.projects().predict(body=request_data, name=model_url).execute()

    estimates = list(map(lambda item: item["scores"]
            ,response["predictions"]
    ))

    return estimates
```

# Deploy XGBoost & Scikit-learn models

```python
from sklearn.externals import joblib
joblib.dump(estimator, model.joblib)
```

```
gsutil cp ./model.joblib ${MODEL_PATH}/model.joblib
```

*Save and dump model to GCS*

*deploy the model to CMLE*

```
gcloud ml-engine models create ${MODEL_NAME} --regions=${REGION}

gcloud ml-engine versions create ${VERSION} --model=${MODEL_NAME} \
    --origin=${MODEL_PATH} \
    --runtime-version="1.4" \
    --framework="SCIKIT_LEARN"
    --pythonVersion="2.7"
```
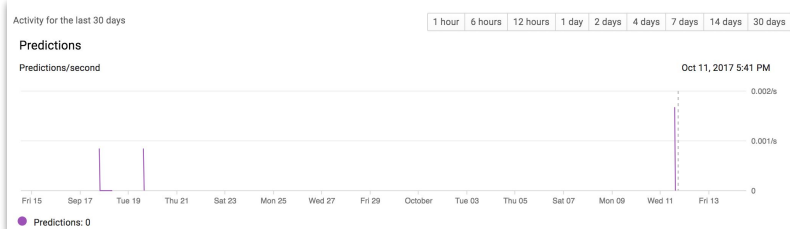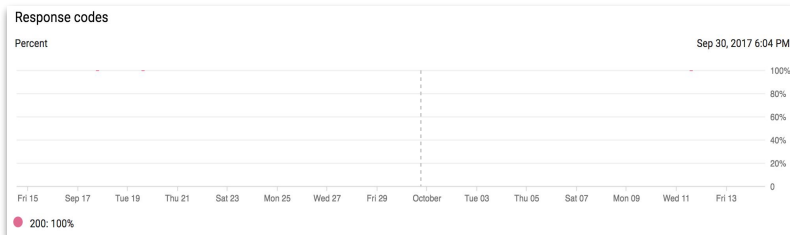
https://cloud.google.com/blog/big-data/2018/04/serving-real-time-scikit-learn-and-xgboost-predictions

Google Cloud

# Monitoring model API serving health



**Response codes**

Percent      Sep 30, 2017 6:04 PM

● 200: 100%

**Activity for the last 30 days** | 1 hour | 6 hours | 12 hours | 1 day | 2 days | 4 days | 7 days | 14 days | 30 days

**Predictions**

Predictions/second      Oct 11, 2017 5:41 PM

● Predictions: 0

**Prediction errors**

Percent      Oct 1, 2017 6:06 AM

● Prediction errors: 0

● v1

| Description | |
|---|---|
| Model | taxifare_estimator |
| Model location | gs://ksalama-gcs-cloudml/ml-models/taxifare/dnn-combined-regression-small/ |
| Creation time | Aug 22, 2017, 11:52:04 PM |
| Last use time | Oct 11, 2017, 1:20:01 PM |

**Requests**

Requests/second      Sep 24, 2017 1:07 AM

● Requests: 0

**Model latency**

Duration      Sep 29, 2017 5:47 AM

● 50th percentile: 163.840ms    ● 95th percentile: 252.314ms    ● 99th percentile: 260.178ms

Google Cloud

# Example Production Workflow

Message Queueing
System

Pre/Post Prediction
Data Processing

Send
Data Points
to Topic

Pubsub

Consume
Data Points
from Topic

Dataflow

Store Data
Points +
Predictions

BigQuery

Stream Data
Source

Store & Serve
Information

Send Data Point
to ML Model API
and Receive
Prediction

ML Engine

Deployed ML
Model APIs

Google Cloud

# Monitor ML Engine job on Cloud Console

# Monitor ML Engine job on Cloud Console

**Stackdriver**
**Logging**

- **Logs**
- Logs-based metrics
- Exports
- Logs ingestion

**CREATE METRIC**  **CREATE EXPORT**

text:average_loss ⊗

| Cloud ML Job, lab3a_181009_204525 ▾ | All logs ▾ | Any log level ▾ | 🕐 No limit ▾ | Jump to now ▾ |

Showing logs from **all time** (PST)    Download logs  View Options ▾

↓        No older entries found matching current filter.        ↓

▸ ℹ 2018-10-09 13:51:40.778 PDT  `master-replica-0` Saving dict for global step 18: average_loss = 212.268, globa…

▸ ℹ 2018-10-09 14:01:39.560 PDT  `master-replica-0` Saving dict for global step 47742: average_loss = 89.53, glob…

▸ ℹ 2018-10-09 14:11:42.629 PDT  `master-replica-0` Saving dict for global step 96544: average_loss = 89.4802, gl…

▸ ℹ 2018-10-09 14:13:40.509 PDT  `master-replica-0` Saving dict for global step 100007: average_loss = 89.4804, g…

↑        Load newer logs        ↑

Thank you!