

Sign up for a Google Cloud Platform free trial (<https://cloud.google.com>) to follow along or try at home.

<https://goo.gl/nqsCo5> (this deck)

End-to-end machine learning on structured data

Lak Lakshmanan, Tech Lead, Data & ML Professional Services, Google Cloud

www.vlakshman.com lak_gcp@

2017-10-12
World Summit AI,
Amsterdam



Agenda

Machine Learning on Google Cloud

Explore the data

Create the dataset

Build the model

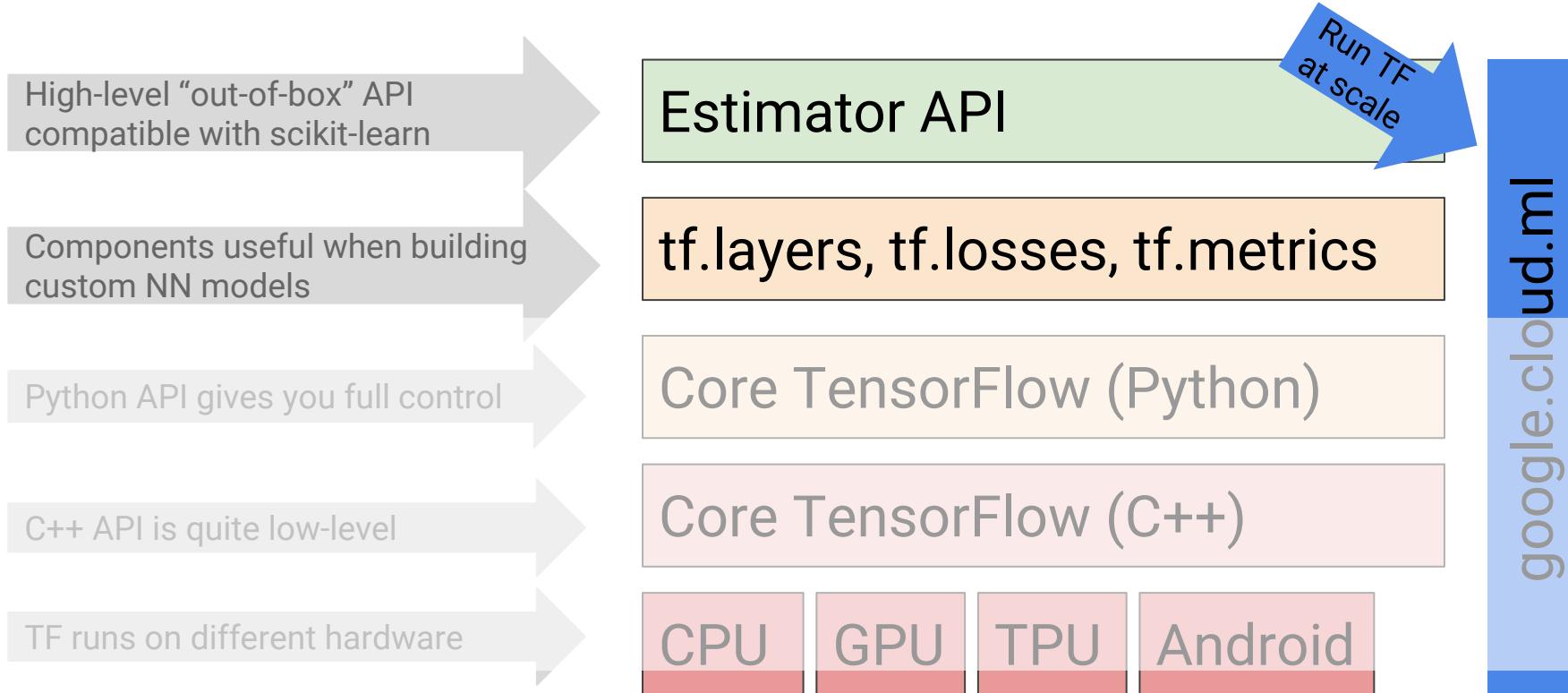
Operationalize model

The most common ML models at Google are models that operate on structured data

Type of network	# of network layers	# of weights	% of deployed models
MLP0	5	20M	61%
MLP1	4	5M	
LSTM0	58	52M	29%
LSTM1	56	34M	
CNN0	16	8M	5%
CNN1	89	100M	

<https://cloud.google.com/blog/big-data/2017/05/an-in-depth-look-at-googles-first-tensor-processing-unit-tpu>

We will use distributed TensorFlow on Cloud ML Engine



Many Machine Learning frameworks can handle toy problems



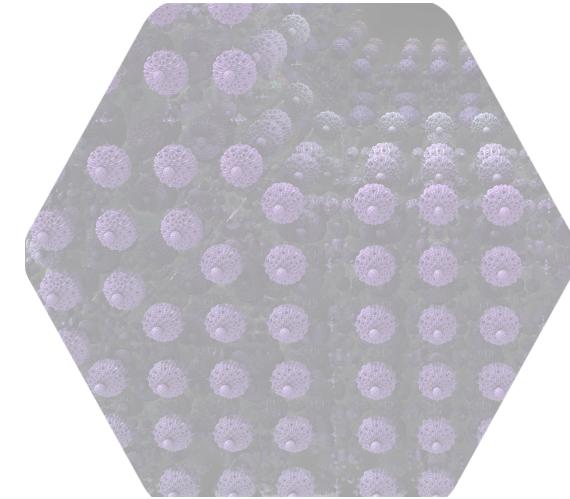
To build effective ML, you need:



Big Data



Feature
Engineering



Model
Architectures

As your data size increases, batching and distribution become important



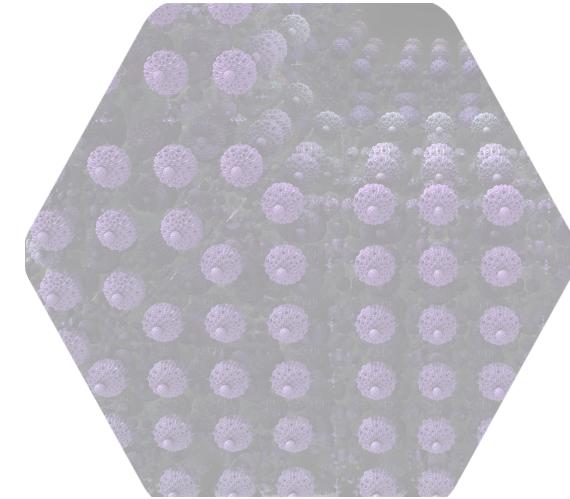
To build effective ML, you need:



Big Data

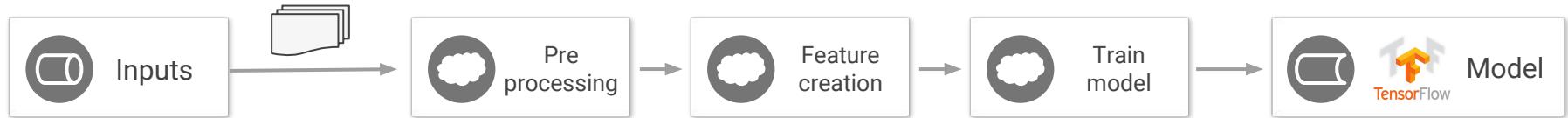


Feature
Engineering



Model
Architectures

Input necessary transformations



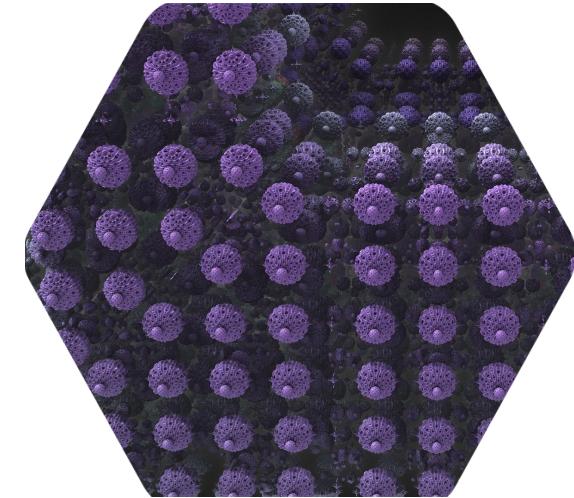
To build effective ML, you need:



Big Data



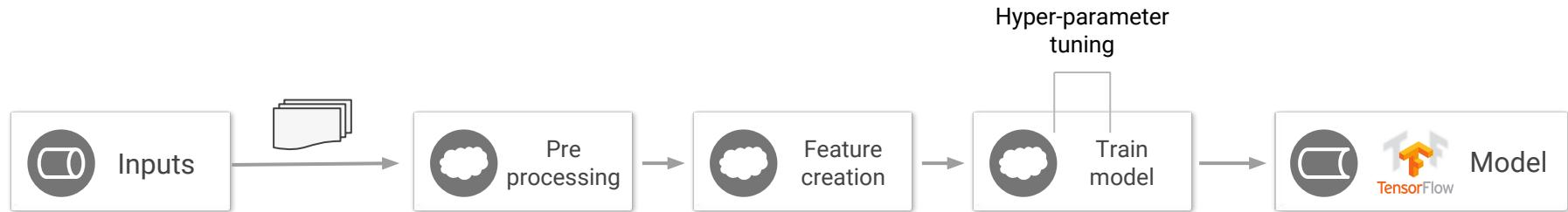
Feature
Engineering



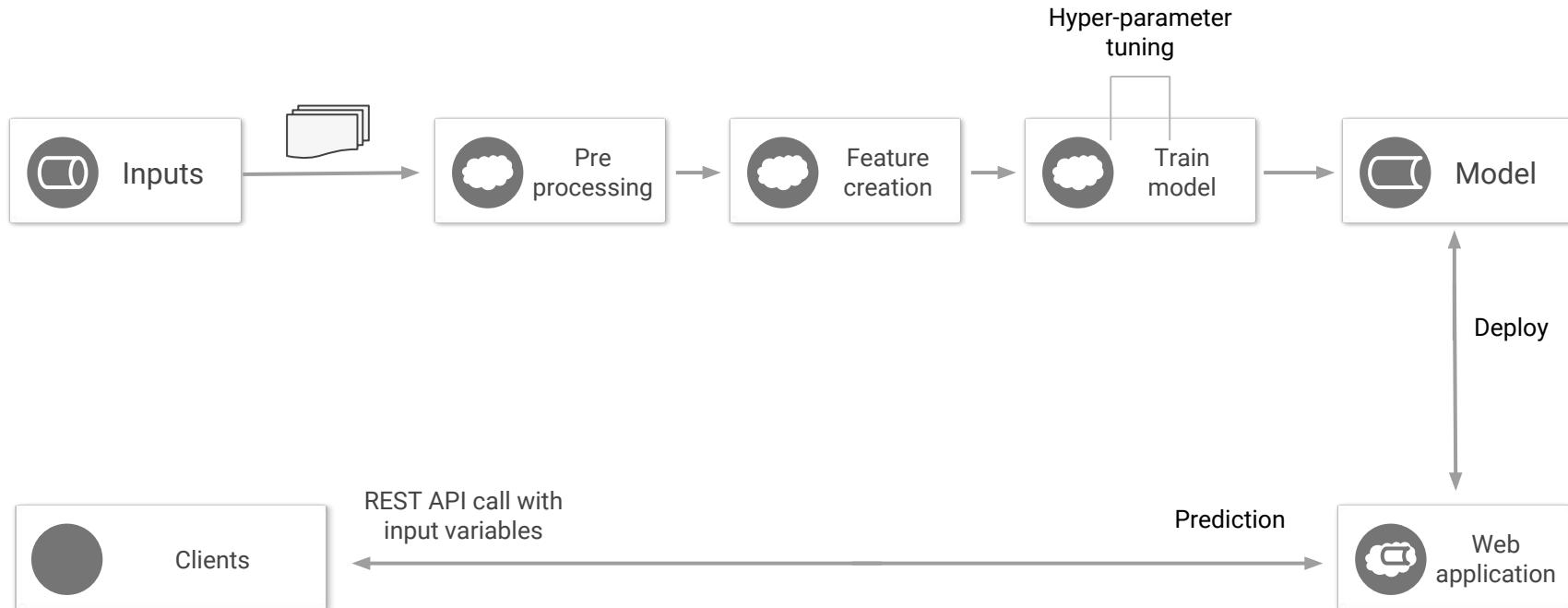
Model
Architectures

What else does a ML framework need to provide?

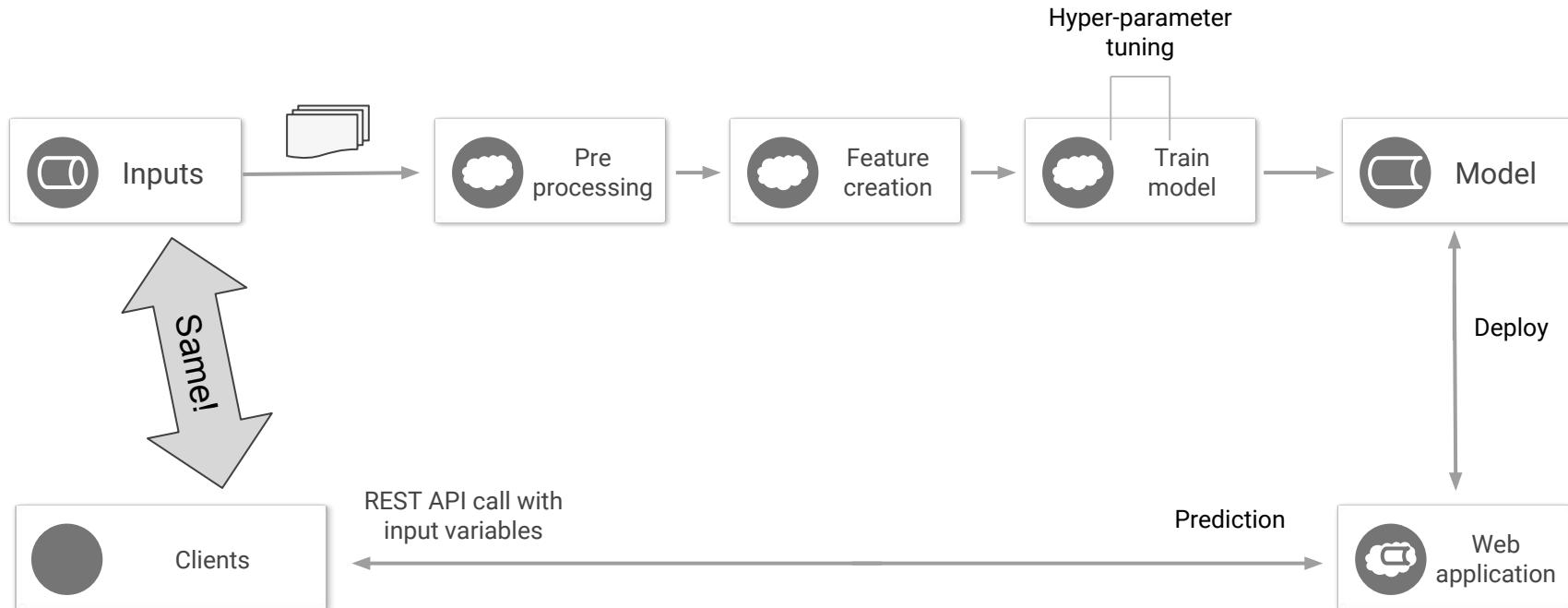
Hyperparameter tuning might be nice



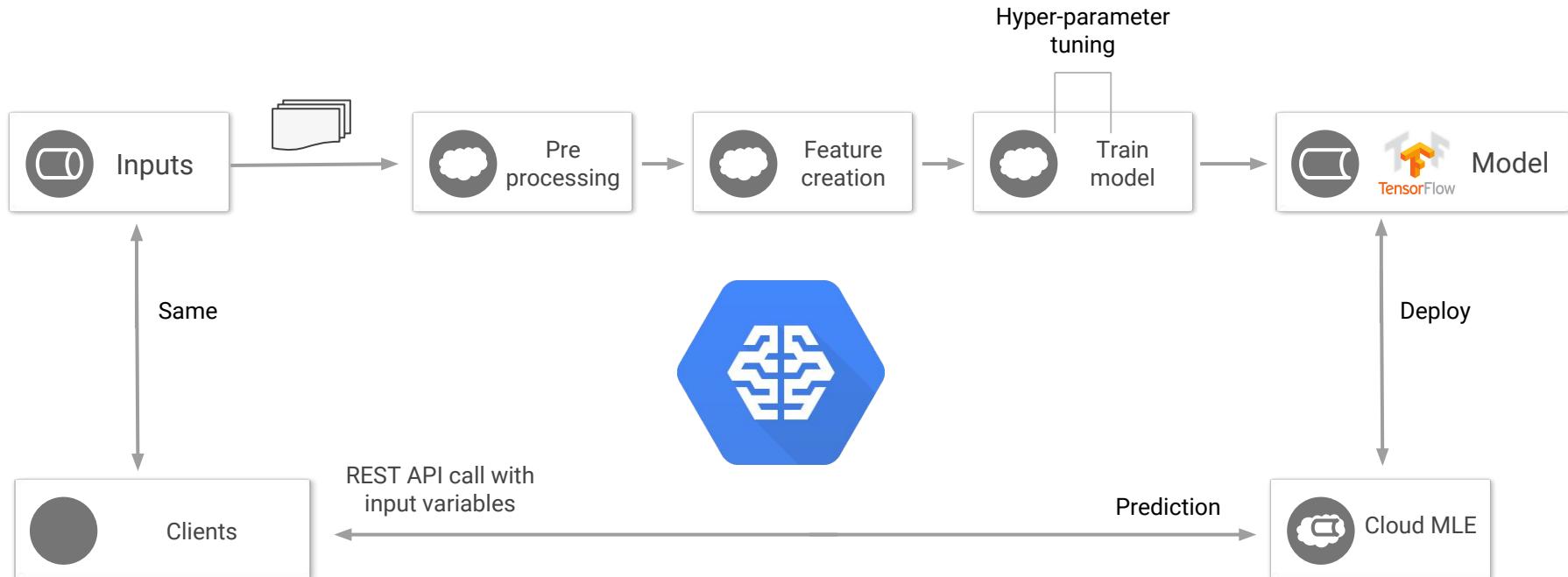
Need to autoscale prediction code



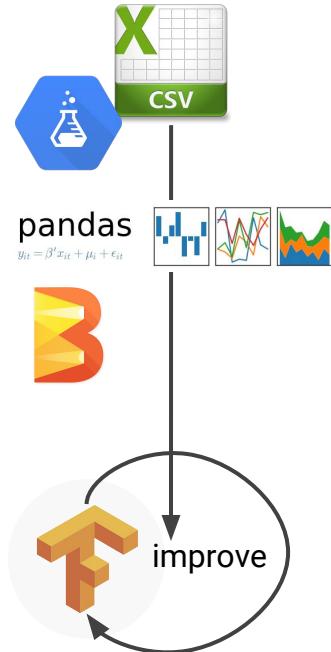
Who does the preprocessing?



Cloud Machine Learning—repeatable, scalable, tuned



In Datalab, start locally on sampled dataset



This notebook is Lab2b of CPB 102, Google's course on Machine Learning using Cloud ML.

In this notebook, we will work with relatively low-level TensorFlow functions to implement a linear regression model. We will use this notebook to demonstrate early stopping -- a technique whereby training is stopped once the error on the validation dataset starts to increase.

```
import databab,bigquery as bq
import tensorflow as tf
import pandas as pd
import numpy as np
import shutil

Code to read data and compute error is the same as Lab2a.

def read_dataset(filename):
    return pd.read_csv(filename, header=None, names=['pickuplon','pickuplat','dropofflon','dropofflat','passengers','fare_amount'])

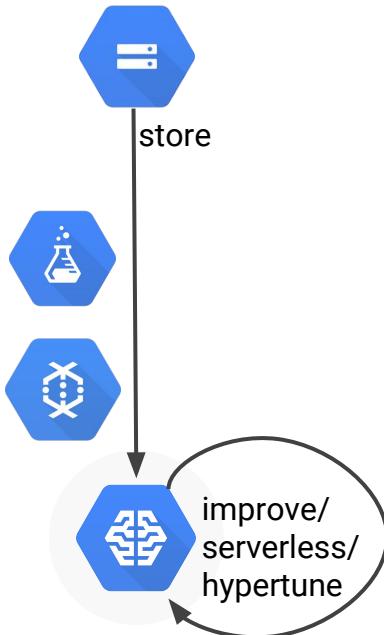
df_train = read_dataset("../lab1a/taxi-train.csv")
df_valid = read_dataset("../lab1a/taxi-valid.csv")
df_test = read_dataset("../lab1a/taxi-test.csv")
df_train[:5]

FEATURE_COLS = np.arange(0,5)
TARGET_COL   = 'fare_amount'

def compute_rmse(actual, predicted):
    return np.sqrt(np.mean((actual-predicted)**2))

def print_rmse(model):
    print "Train RMSE = {}".format(compute_rmse(df_train[TARGET_COL], model.predict(df_train.iloc[:,FEATURE_COLS].values)))
    print "Valid RMSE = {}".format(compute_rmse(df_valid[TARGET_COL], model.predict(df.valid.iloc[:,FEATURE_COLS].values)))
```

Then, scale it out to GCP using serverless technology



Training on cloud

In order to train on the cloud, we have to copy the model and data to our bucket on Google Cloud Storage (GCS).

```
%bash
rm -rf taxifare.tar.gz taxi_trained
tar cvfz taxifare.tar.gz taxifare
gsutil cp taxifare.tar.gz gs://${BUCKET}/taxifare/source/taxifare.tar.gz
gsutil cp ./lab1a/*.csv gs://${BUCKET}/taxifare/input/
gsutil -m rm -r -f gs://${BUCKET}/taxifare/taxi_preproc
gsutil -m rm -r -f gs://${BUCKET}/taxifare/taxi_trained
```

Running.

When you run your preprocessor, you have to change the input and output to be on GCS.

Using DirectPipelineRunner runs Dataflow locally, but the inputs & outputs are on the cloud. Using BlockingDataflowPipelineRunner will use Cloud Dataflow (and take much longer because of the overhead involved for such a small dataset). To see the status of your BlockingDataflowPipelineRunner job, visit <https://console.cloud.google.com/dataflow>

```
# imports
import apache_beam as beam
import google.cloud.ml as ml
import google.cloud.ml.dataflow.io.tfrecordio as tfrecordio
import google.cloud.ml.io as io
import os

# Change as needed
#RUNNER = 'DirectPipelineRunner' #
RUNNER = 'BlockingDataflowPipelineRunner'

# defines
feature_set = TaxifareFeatures()
OUTPUT_DIR = 'gs://{}/taxifare/taxi_preproc'.format(BUCKET)
```

Agenda

Machine Learning on Google Cloud

Explore the data

Create the dataset

Build the model

Operationalize model

Goal: predict weight of a newborn



An open dataset of births is available in BigQuery

Births recorded in the 50 states of the USA from
1969 to 2008

Table ID	bigquery-public-data:samples.natality
Table Size	21.9 GB
Long Term Storage Size	21.9 GB
Number of Rows	137,826,763

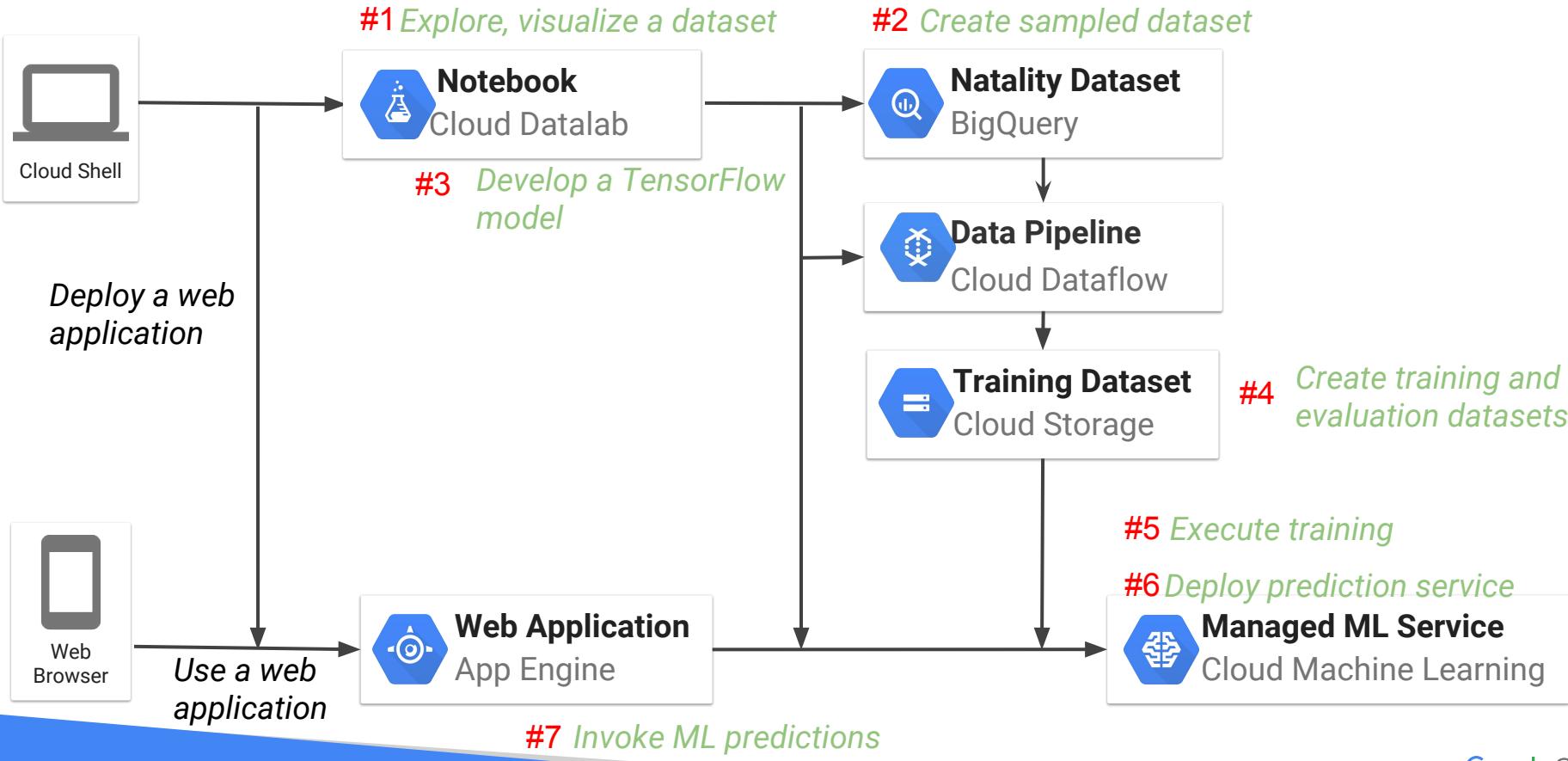


<https://bigquery.cloud.google.com/table/bigquery-public-data:samples.natality>

The data set includes details about the pregnancy

Date of birth	year	INTEGER	NULLABLE	Four-digit year of the birth. Example: 1975.
	month	INTEGER	NULLABLE	Month index of the date of birth, where 1=January.
	day	INTEGER	NULLABLE	Day of birth, starting from 1.
	wday	INTEGER	NULLABLE	Day of the week, where 1 is Sunday and 7 is Saturday.
Location of birth (US state)	state	STRING	NULLABLE	The two character postal code for the state. Entries after 2004 do not include this value.
Baby's birth weight (lbs)	weight_pounds	FLOAT	NULLABLE	Weight of the child, in pounds.
Mother's age at birth	mother_age	INTEGER	NULLABLE	Reported age of the mother when giving birth.
Duration of pregnancy	gestation_weeks	INTEGER	NULLABLE	The number of weeks of the pregnancy.
Mother's use of cigarette or alcohol	cigarette_use	BOOLEAN	NULLABLE	True if the mother smoked cigarettes. Available starting 2003.
Mother's weight gain (lbs)	alcohol_use	BOOLEAN	NULLABLE	True if the mother used alcohol. Available starting 1989.
	weight_gain_pounds	INTEGER	NULLABLE	Number of pounds gained by the mother during pregnancy.

The end-to-end machine learning set of labs



<https://goo.gl/nqsCo5> (this deck)

<codelabs.developers.google.com/cloud-quest-scientific-data>

<github.com/GoogleCloudPlatform/training-data-analyst/blob/master/blogs/babyweight/>

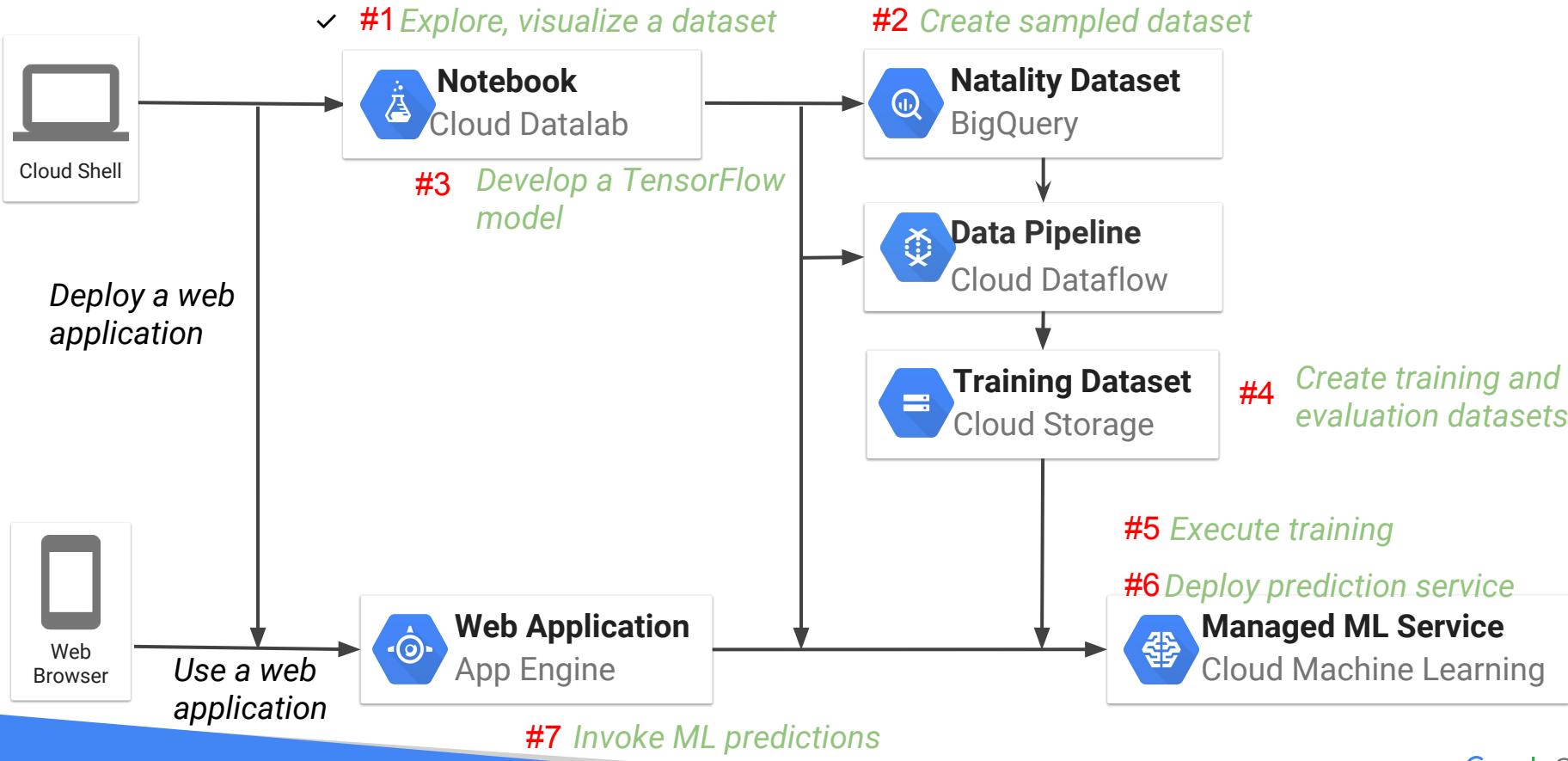
Lab #1: Exploring a BigQuery dataset to find features to use in model

Lab #1: Exploring a BigQuery dataset to find features to use in model

BigQuery: publicdata.samples.natality

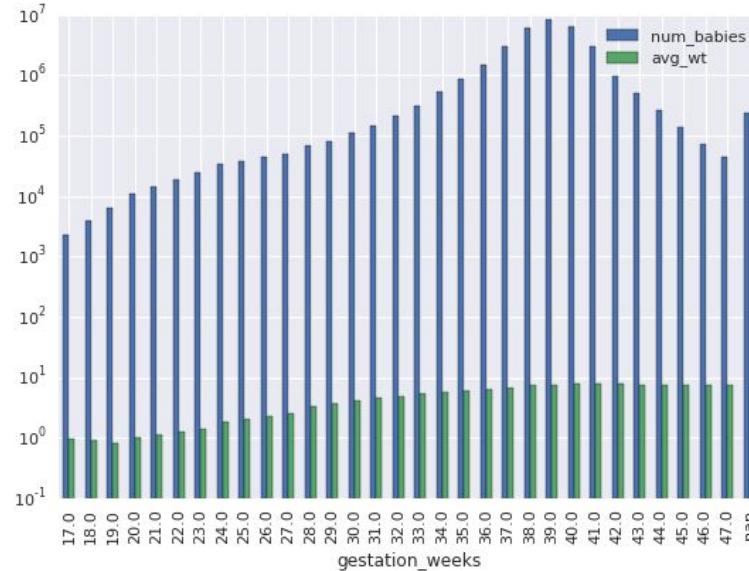


The End-to-End process



Debrief

What features did you use? Why?



Is it fair practice to use mother's race as a feature?

**How would we check whether
the use/non-use of mother's
race amplifies injustice?**



Agenda

Machine Learning on Google Cloud

Explore the data

Create the dataset

Build the model

Operationalize model

Building a ML model involves:



Create the dataset



Build the model



Operationalize the
model

<https://goo.gl/nqsCo5> (this deck)

<codelabs.developers.google.com/cloud-quest-scientific-data>

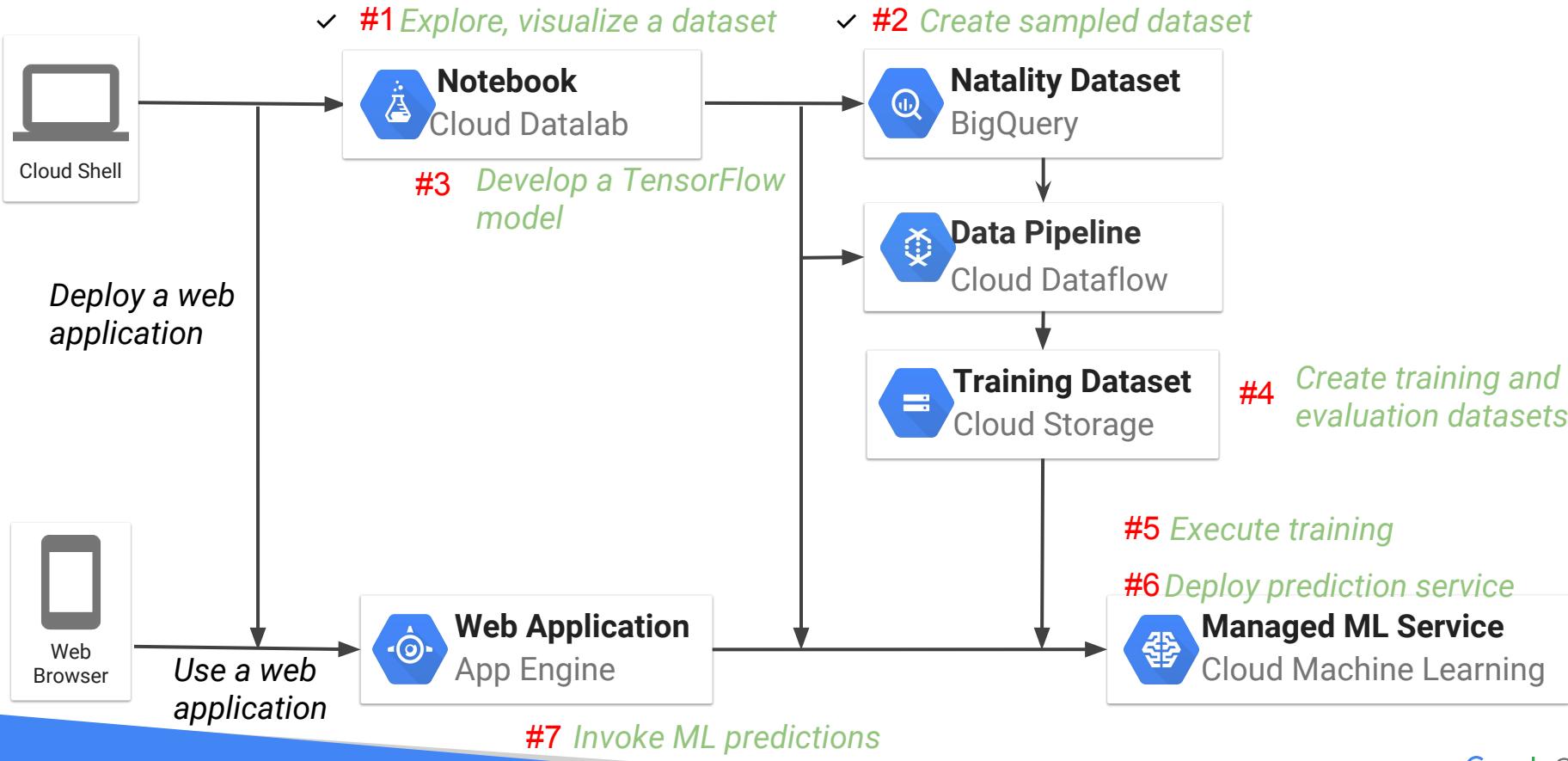
<github.com/GoogleCloudPlatform/training-data-analyst/blob/master/blogs/babyweight/>

Lab #2: Creating a sampled dataset

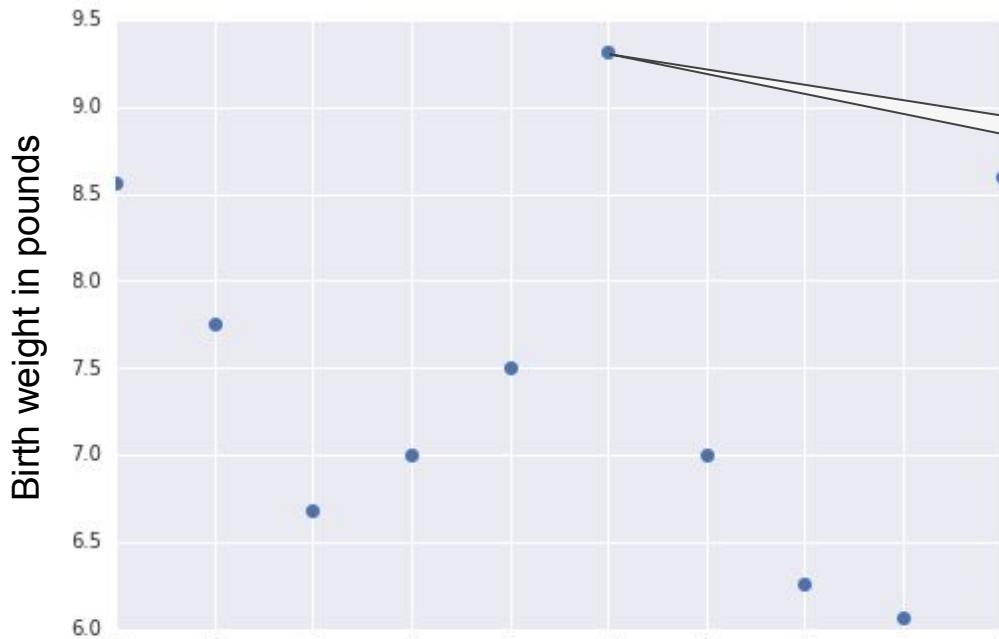
Lab #2: Creating a sampled dataset

1. Sample the full BigQuery, creating a smaller dataset so that you can use it for model development and local training
2. Make sure that this smaller dataset is representative of the full dataset.

The End-to-End process



Debrief: The simplest option is to sample rows randomly



- Each data point is a birth record from the natality dataset

weight	year	mother_age	gestation_weeks	cigarette_use	alcohol_use
6.03	2004	29	39	false	false

- Random sampling eliminates potential biases due to order of the training examples but ...

But what about triplets?

3 rows with essentially the same data!

How can we make this data unique?

How can we solve this?

Use FARM_FINGERPRINT so that you get a repeatable, nonoverlapping sample.

Export train.csv and eval.csv (2 files)



Agenda

Machine Learning on Google Cloud

Explore the data

Create the dataset

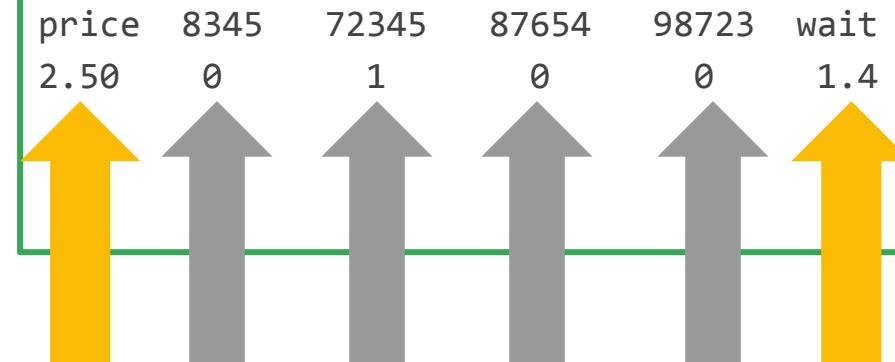
Build the model

Operationalize model

Two types of features: Dense & sparse

```
{  
    "transactionId": 42,  
    "name": "Ice Cream",  
    "price": 2.50,  
    "tags": ["cold", "dessert"],  
    "servedBy": {  
        "employeeId": 72365,  
        "waitTime": 1.4,  
        "customerRating": 4  
    },  
    "storeLocation": {  
        "latitude": 35.3,  
        "longitude": -98.7  
    }  
},
```

price	8345	72345	87654	98723	wait
2.50	0	1	0	0	1.4

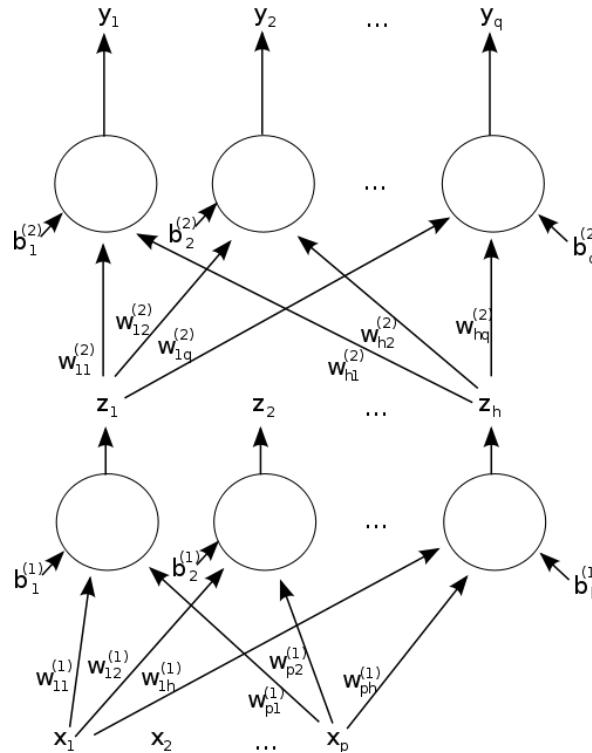


DNNs good for dense, highly correlated

pixel_values (



)



1024^2 input
nodes

10 hidden
nodes \rightarrow 10
image
features

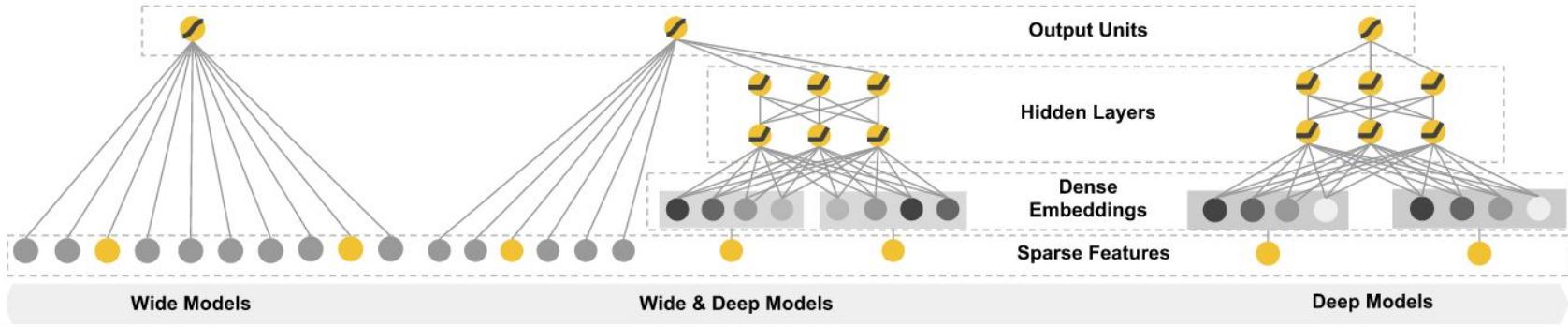
https://commons.wikimedia.org/wiki/File:Two_layer_ann.svg

Linear for sparse, independent features

One-hot encoding

[0.	1.	0.	0.	0.	0.	0.	0.	0.	0.]
[0.	0.	0.	0.	0.	0.	0.	1.	0.	0.]
[0.	0.	0.	0.	0.	0.	0.	0.	0.	1.]
[0.	1.	0.	0.	0.	0.	0.	0.	0.	0.]
[0.	0.	0.	0.	0.	0.	0.	1.	0.	0.]
[0.	0.	0.	1.	0.	0.	0.	0.	0.	0.]
[0.	0.	0.	0.	0.	0.	1.	0.	0.	0.]
[0.	1.	0.	0.	0.	0.	0.	0.	0.	0.]
[0.	0.	0.	0.	1.	0.	0.	0.	0.	0.]
[0.	1.	0.	0.	0.	0.	0.	0.	0.	0.]

Wide-and-deep models let you handle both





Wide & Deep



memorization
relevance

generalization
diversity

Wide-and-deep network in Estimator API

```
model = tf.contrib.learn.DNNLinearCombinedClassifier(  
    model_dir=...,  
    linear_feature_columns=wide_columns,  
    dnn_feature_columns=deep_columns,  
    dnn_hidden_units=[100, 50])
```



<https://goo.gl/nqsCo5> (this deck)

<codelabs.developers.google.com/cloud-quest-scientific-data>

<github.com/GoogleCloudPlatform/training-data-analyst/blob/master/blogs/babyweight/>

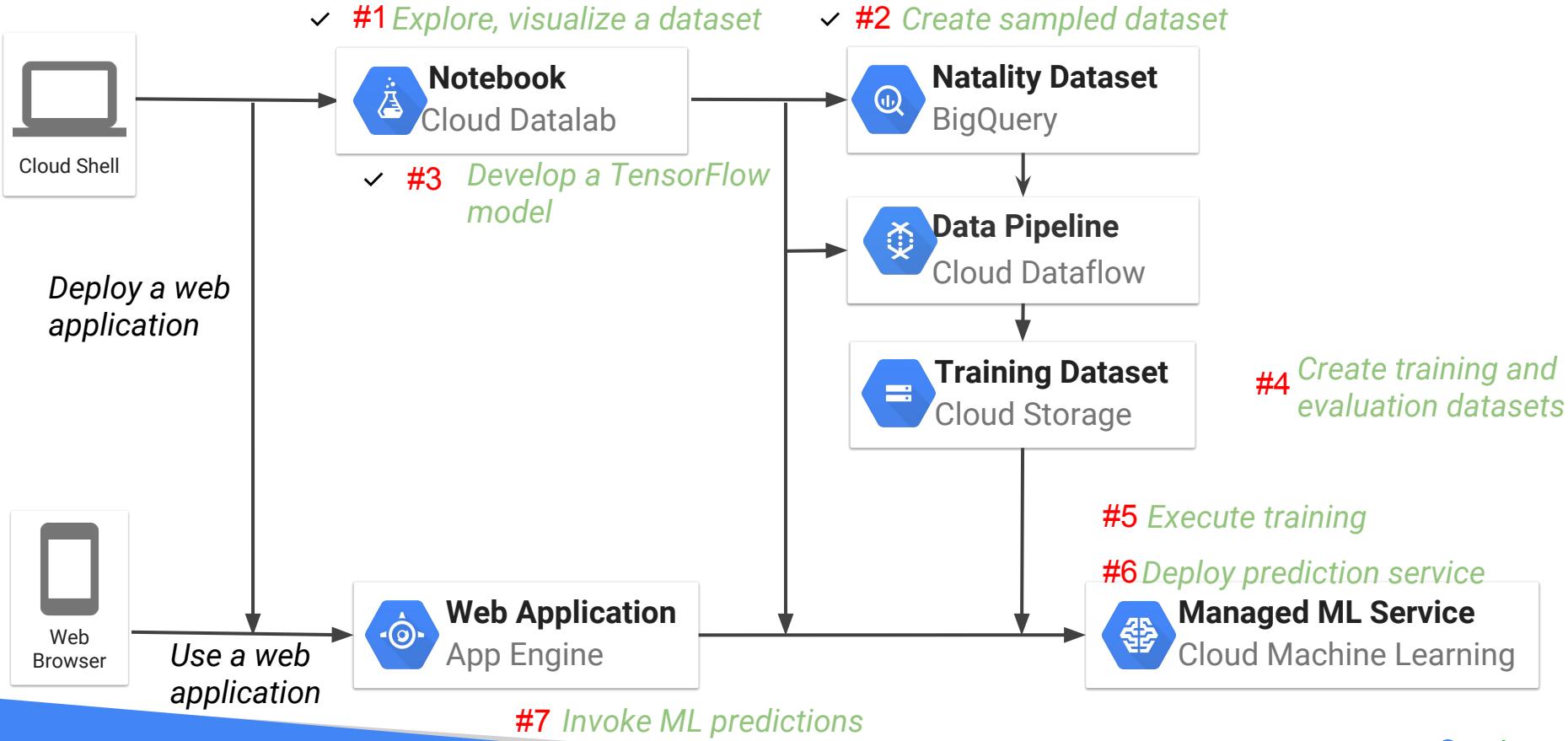
Lab #3: Create Tensorflow model

Lab #3: Create Tensorflow model

In this lab, you use the high-level Estimator API for wide-and-deep model

1. Determine your wide columns and deep columns
2. Determine if you need any embeddings
3. Feel free to iterate and add further preprocessing of features

The End-to-End process



Agenda

Machine Learning on Google Cloud

Explore the data

Create the dataset

Build the model

Operationalize model

<https://goo.gl/nqsCo5> (this deck)

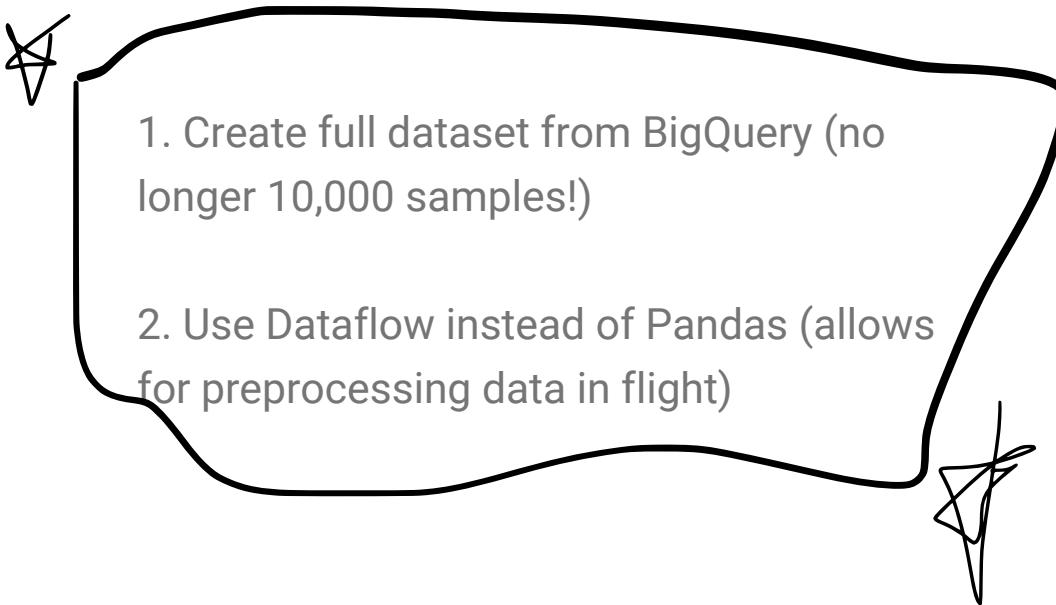
<codelabs.developers.google.com/cloud-quest-scientific-data>

<github.com/GoogleCloudPlatform/training-data-analyst/blob/master/blogs/babyweight/>

Lab #4: Preprocess using Dataflow

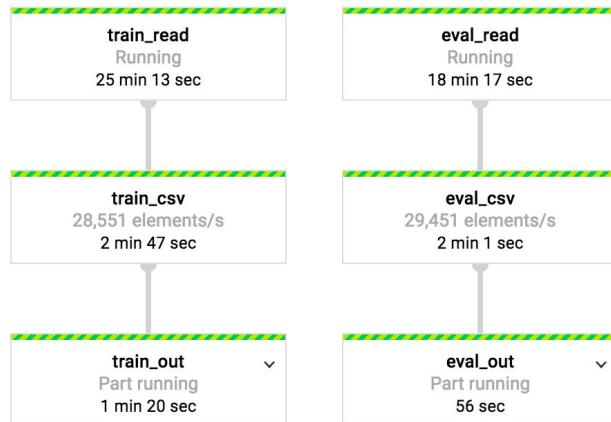
Lab #4: Preprocess using Dataflow

While Pandas is fine for experimenting, for operationalization of your workflow, it is better to do preprocessing in Apache Beam (Dataflow)

- 
1. Create full dataset from BigQuery (no longer 10,000 samples!)
 2. Use Dataflow instead of Pandas (allows for preprocessing data in flight)

Split full dataset into train/eval and do preprocessing

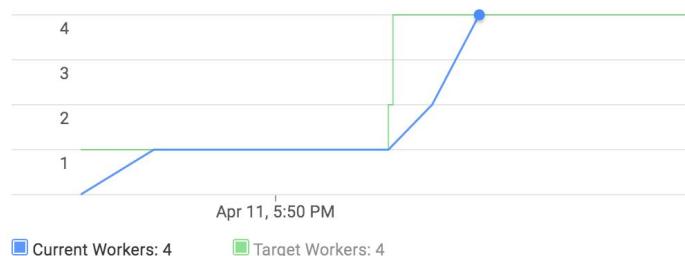
BigQuery -> Dataflow -> CSV



Autoscaling

Workers	15
Current State	Autoscaling: Raised the number of workers to 4 based on the currently running step(s).

Worker History

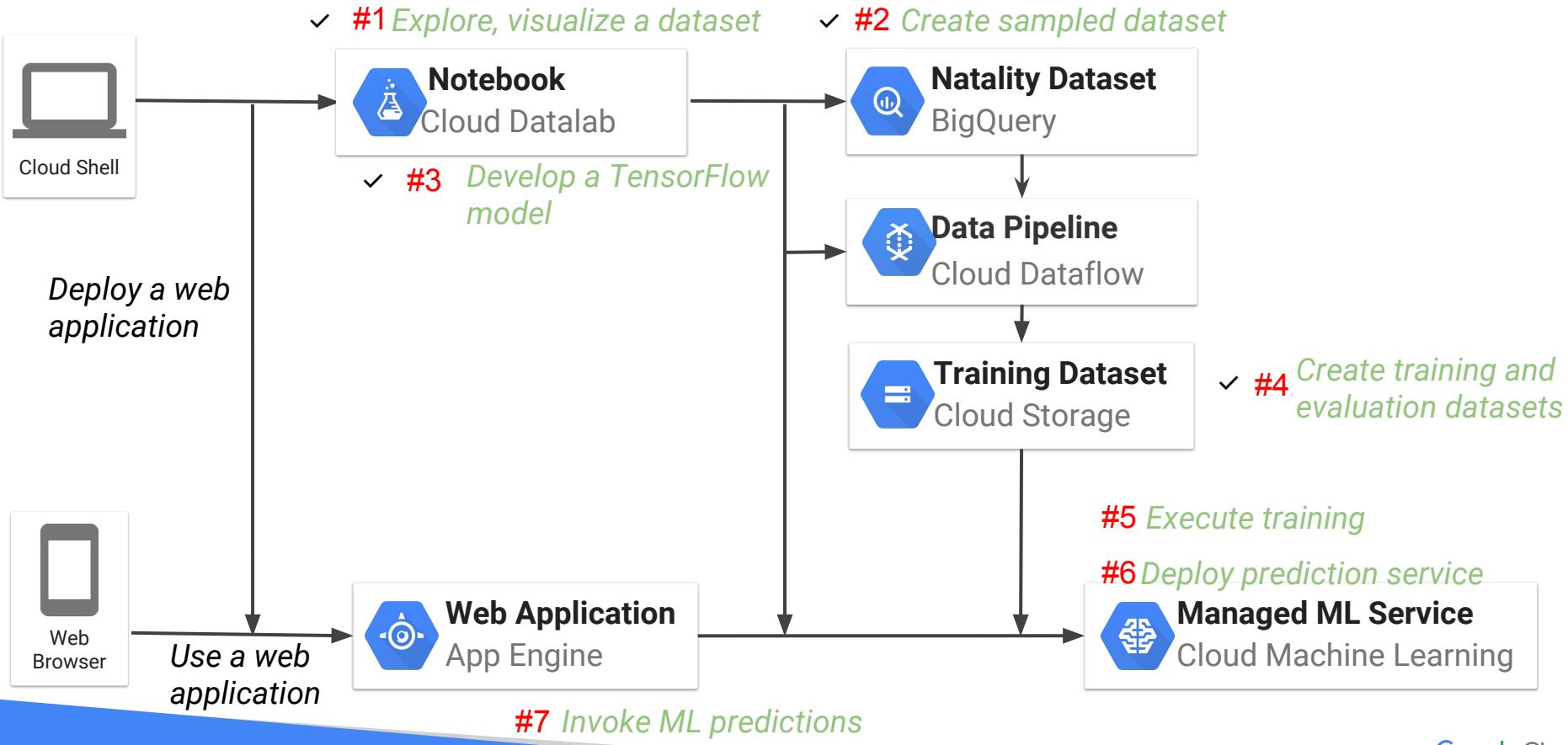


See More History

Resource Metrics

Current vCPUs	15
Total vCPU Time	1.561 vCPU hr

The End-to-End process



<https://goo.gl/nqsCo5> (this deck)

codelabs.developers.google.com/cloud-quest-scientific-data

github.com/GoogleCloudPlatform/training-data-analyst/blob/master/blogs/babyweight/

Lab #5: Train on Cloud ML Engine

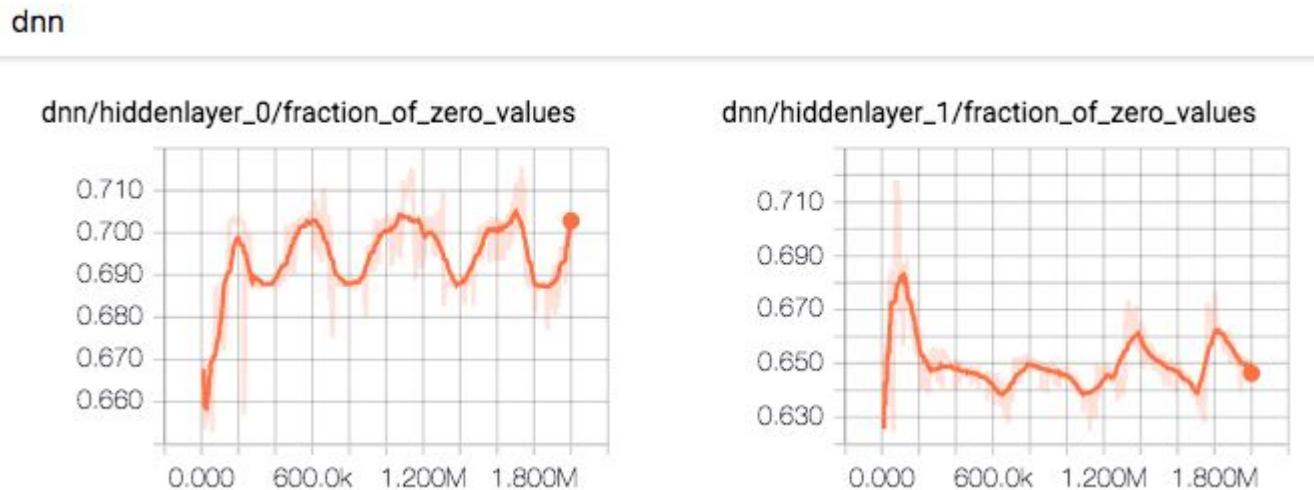
Lab #5: Train on Cloud ML Engine

This lab illustrates distributed training and hyperparameter tuning on Cloud ML Engine

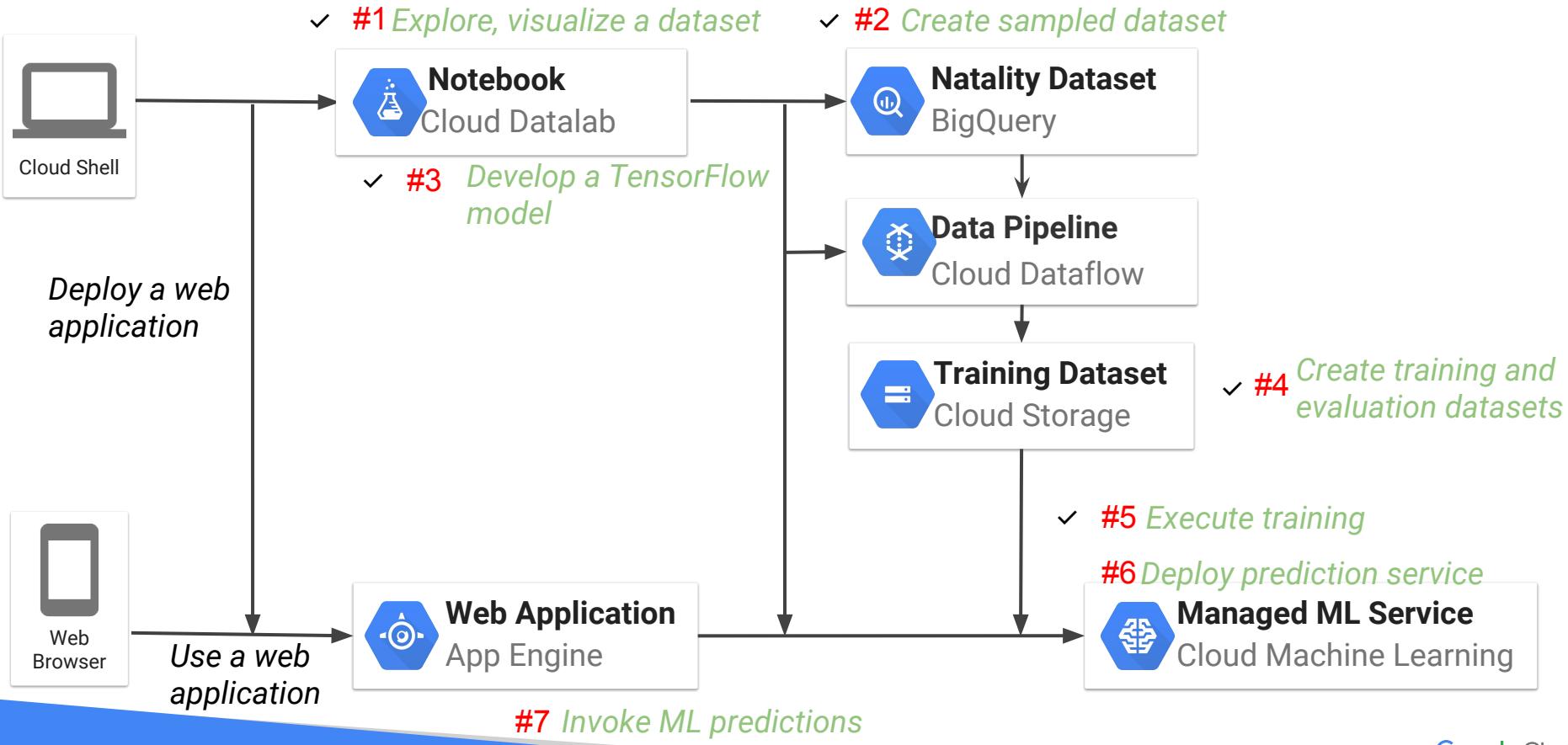
1. change batchsize if necessary
2. use train_steps instead of num_epochs

Train on Cloud ML Engine

Submit training job on full dataset, monitor using TensorBoard



The End-to-End process



<https://goo.gl/nqsCo5> (this deck)

<codelabs.developers.google.com/cloud-quest-scientific-data>

<github.com/GoogleCloudPlatform/training-data-analyst/blob/master/blogs/babyweight/>

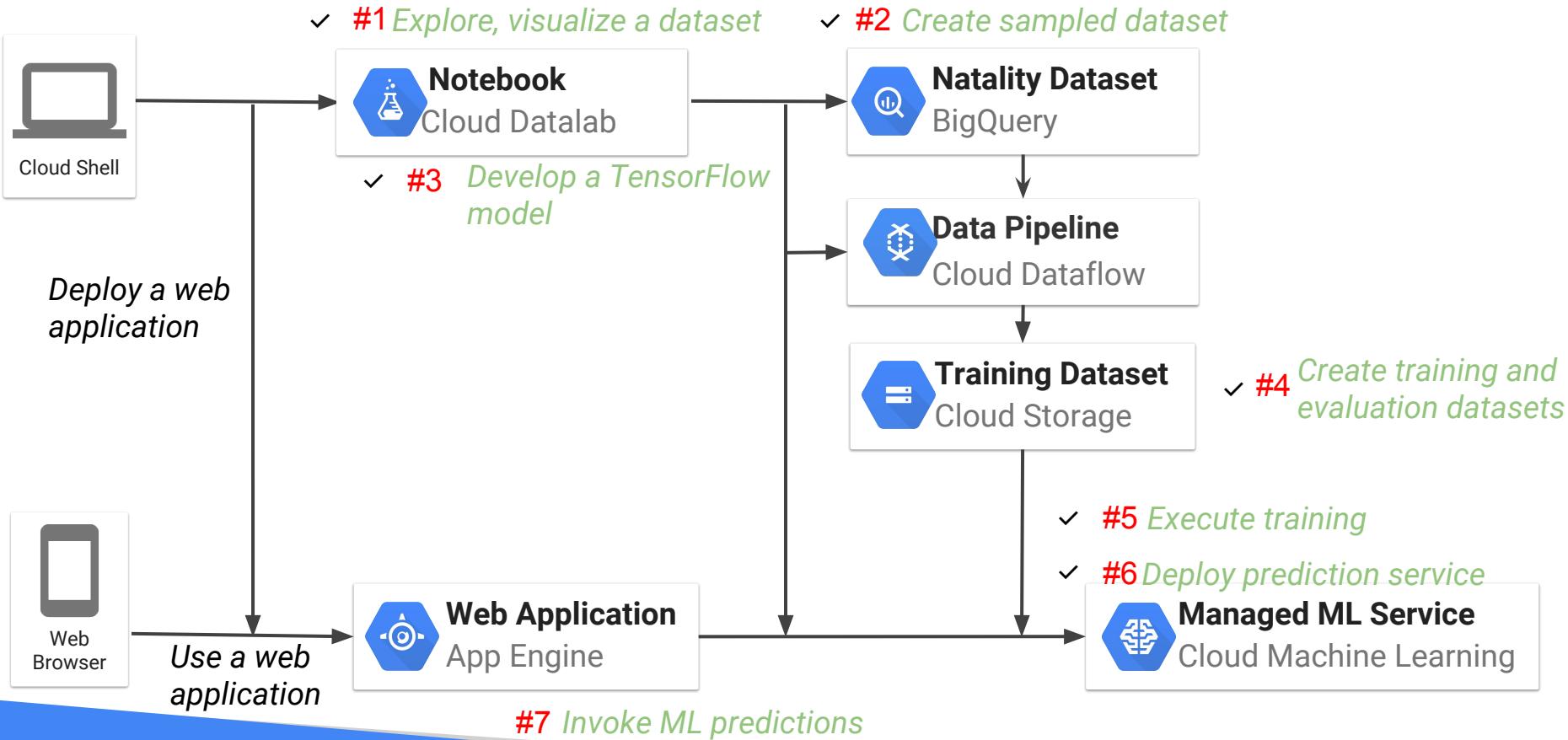
Lab #6: Deploy and predict

Lab #6: Deploy and predict

In this lab, you will:

1. Deploy trained model to Cloud ML Engine
2. Send JSON request to model to get predictions

The End-to-End process



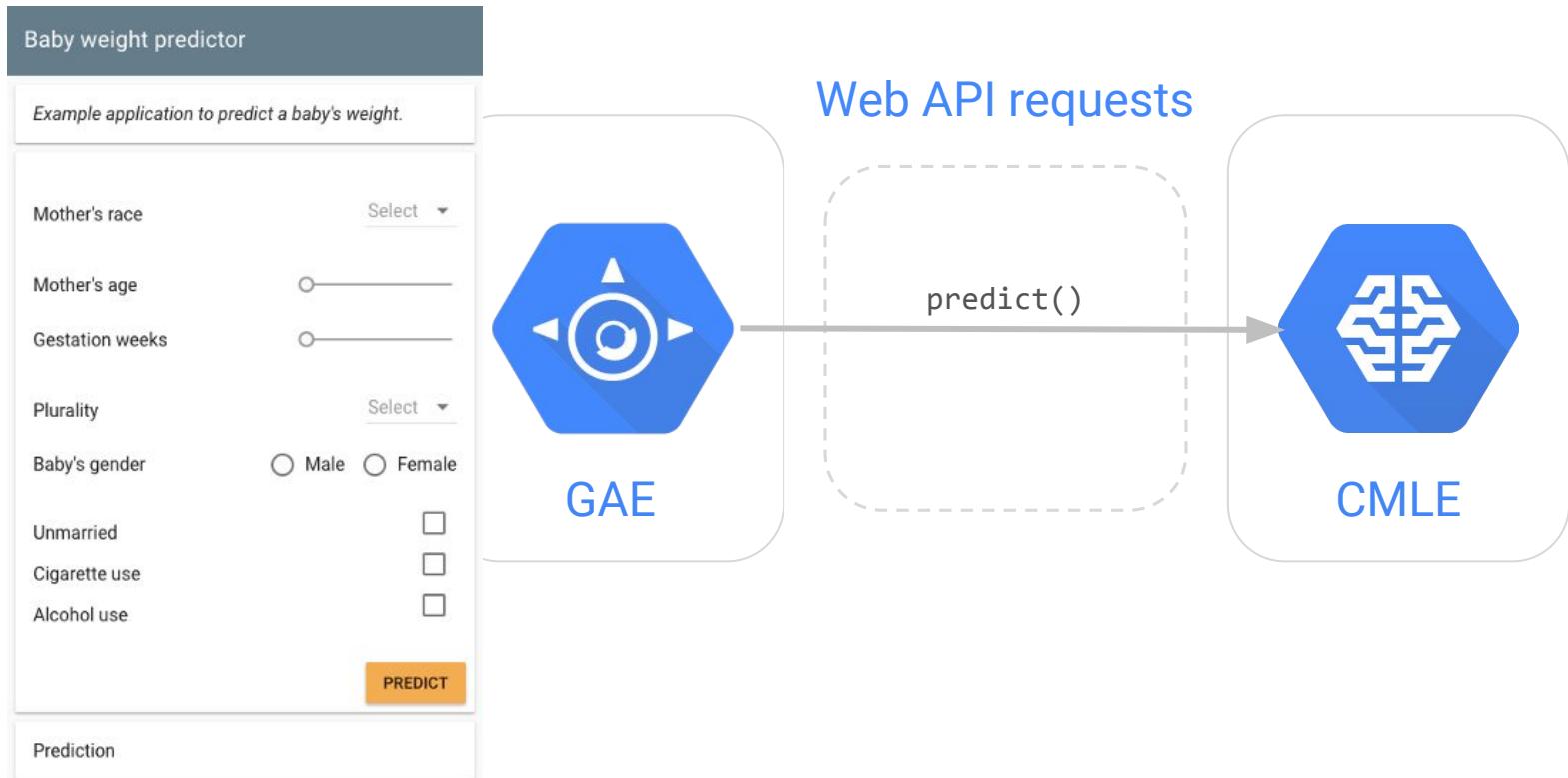
<https://goo.gl/nqsCo5> (this deck)

<codelabs.developers.google.com/cloud-quest-scientific-data>

<github.com/GoogleCloudPlatform/training-data-analyst/blob/master/blogs/babyweight/>

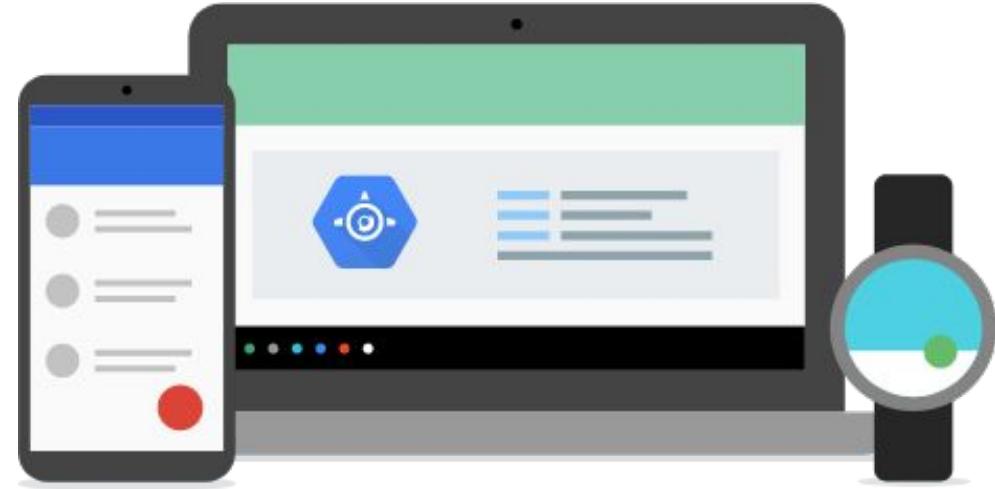
Lab #7: Build an AppEngine app to serve ML predictions

In this lab, we will use AppEngine to invoke ML predictions

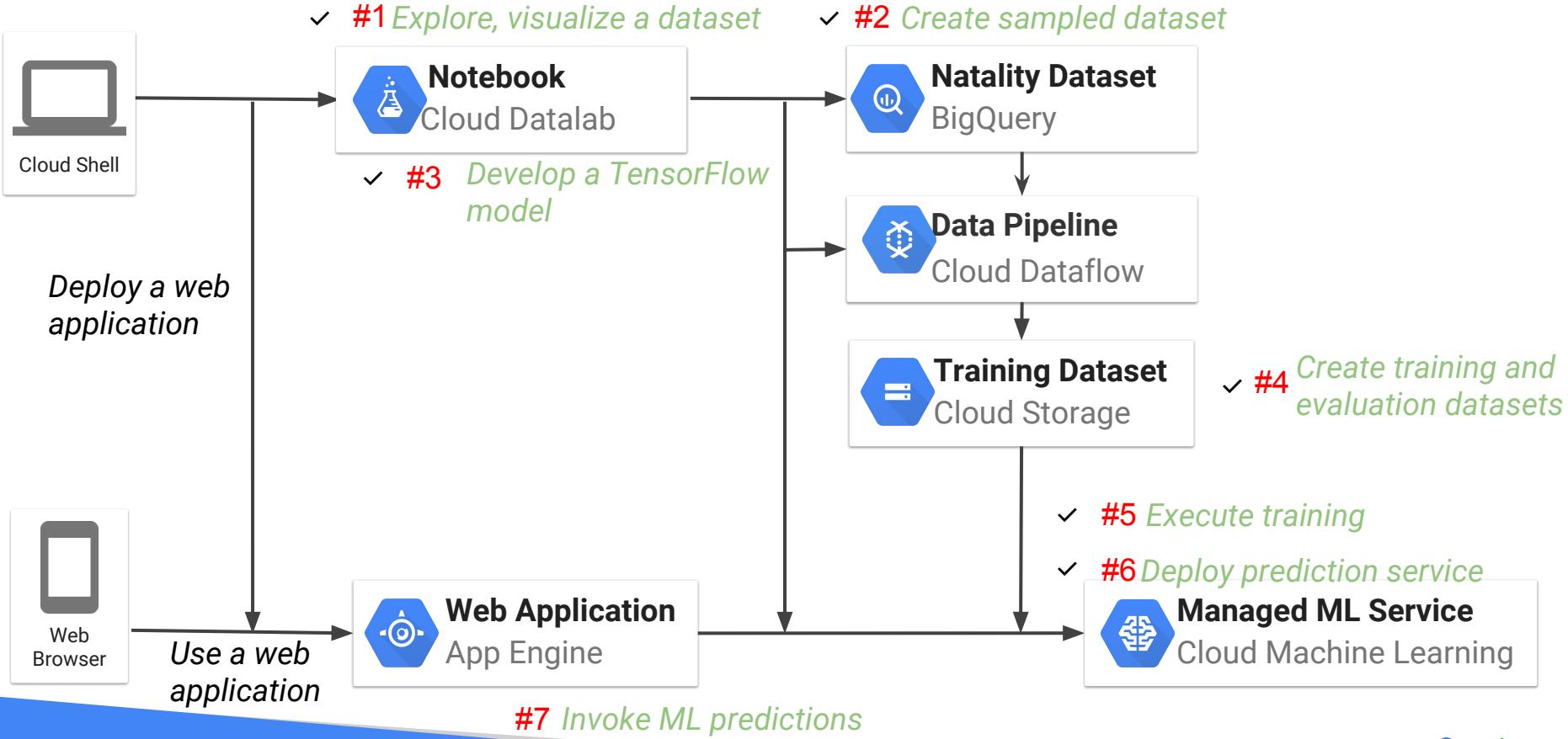


Google App Engine is a fully-managed service for building web backends

Supports Java, Node.js,
Ruby, C#, Go, Python, and
PHP



Summary: An end-to-end process to operationalize ML models



To learn more about TensorFlow on Google Cloud

Hands-on Tensorflow (next session)

goo.gl/nqsCo5

Serverless Machine Learning with TensorFlow bootcamp in Amsterdam

Thursday, October 19 (9:00 – 17:00) at GoDataDriven, Wibautstraat 202

www.godatatest.com/google

Serverless Machine Learning with TensorFlow on Coursera

<https://www.coursera.org/learn/serverless-machine-learning-gcp/>

cloud.google.com

