# Basics of
# Machine Learning

Dmitry Ryabokon, github.com/dryabokon

# Lesson 21
# Deployments in GCP

# Copying files cloudconsole ↔ VM

Copy one file from cloudshell to VM

gcloud compute scp a.txt dmytro_ryabokon@instance-test-dr:~ --project=ml-ops-poc-695

Copy one file from VM to cloudshell

gcloud compute scp dmytro_ryabokon@instance-test-dr:~/b.txt ./b.txt --project=ml-ops-poc-695

Copy folder file from VM to cloudshell

gcloud compute scp --recurse dmytro_ryabokon@instance-test-dr:~/sources/prj_console/* ~/sources/prj_console/ --project=ml-ops-poc-695

# Push docker image to GCR

1) https://console.cloud.google.com/iam-admin/serviceaccount

2) https://cloud.google.com/container-registry/docs/advanced-authentication
gcloud auth configure-docker

3) sudo docker tag hello_dima gcr.io/ml-ops-poc-695/hello_dima
4) sudo docker push gcr.io/ml-ops-poc-695/hello_dima:latest

# Push docker image to GCR

## Create Service Account

# Push docker image to GCR

## Create Service Account: add key, store it

# Push docker image to GCR

## Configure authentication

# Push docker image to GCR

## Configure authentication

```
root@instance-test-dr:/home/rsa-key-20200330# ls
ml-ops-poc-695-0e4f48ea77f8.json  sources
root@instance-test-dr:/home/rsa-key-20200330# gcloud auth activate-service-account gcr-acc@ml-ops-poc-695.iam.gserviceaccount.com --key-file=./ml-ops-poc-695-0e4f48ea77f8.json
Activated service account credentials for: [gcr-acc@ml-ops-poc-695.iam.gserviceaccount.com]
root@instance-test-dr:/home/rsa-key-20200330#
root@instance-test-dr:/home/rsa-key-20200330#
root@instance-test-dr:/home/rsa-key-20200330#
root@instance-test-dr:/home/rsa-key-20200330#
root@instance-test-dr:/home/rsa-key-20200330# gcloud auth configure-docker
WARNING: Your config file at [/root/.docker/config.json] contains these credential helper entries:

{
  "credHelpers": {
    "gcr.io": "gcloud",
    "us.gcr.io": "gcloud",
    "eu.gcr.io": "gcloud",
    "asia.gcr.io": "gcloud",
    "staging-k8s.gcr.io": "gcloud",
    "marketplace.gcr.io": "gcloud"
  }
}
Adding credentials for all GCR repositories.
WARNING: A long list of credential helpers may cause delays running 'docker build'. We recommend passing the registry name to configure only the registry you are using.
gcloud credential helpers already registered correctly.
root@instance-test-dr:/home/rsa-key-20200330#
```

# Push docker image to GCR

## Push docker image

```
root@instance-test-dr:/home/rsa-key-20200330# sudo docker images
REPOSITORY                       TAG          IMAGE ID        CREATED         SIZE
prj_flask_nginx_nginx            latest       aac06a4109ae    10 hours ago    109MB
prj_flask_nginx_flask_app        latest       dfc41c0075e3    10 hours ago    928MB
hello_dima                       latest       032962ec3886    31 hours ago    5.6MB
gcr.io/ml-ops-poc-695/hello_dima latest       032962ec3886    31 hours ago    5.6MB
ubuntu                           latest       597ce1600cf4    8 days ago      72.8MB
busybox                          latest       16ea53ea7c65    3 weeks ago     1.24MB
alpine                           latest       14119a10abf4    6 weeks ago     5.6MB
python                           3.6-jessie   890456b21ed5    2 years ago     703MB
nginx                            1.15.8       f09fe80eb0e7    2 years ago     109MB
python                           3.6.7        1ec4d11819ad    2 years ago     918MB
root@instance-test-dr:/home/rsa-key-20200330#
root@instance-test-dr:/home/rsa-key-20200330#
root@instance-test-dr:/home/rsa-key-20200330#
root@instance-test-dr:/home/rsa-key-20200330#
root@instance-test-dr:/home/rsa-key-20200330#
root@instance-test-dr:/home/rsa-key-20200330# sudo docker push gcr.io/ml-ops-poc-695/hello_dima:latest
The push refers to repository [gcr.io/ml-ops-poc-695/hello_dima]
e2eb06d8af82: Layer already exists
latest: digest: sha256:50f64478c42a993af03592591f1e7ba1435267ac8a1a25814ff71113545e31fd size: 528
root@instance-test-dr:/home/rsa-key-20200330#
```

# Push docker image to GCR

## Docker image appears at GCR

# Flask – nginx docker

https://towardsdatascience.com/how-to-deploy-ml-models-using-flask-gunicorn-nginx-docker-9b32055b3d0
https://github.com/ivanpanshin/flask_gunicorn_nginx_docker

# Flask – nginx docker

```python
#https://www.w3schools.com/tags/ref_httpmethods.asp
#https://towardsdatascience.com/how-to-deploy-ml-models-using-flask-gunicorn-nginx-docker-9b32055b3d0
#https://github.com/ivanpanshin/flask_gunicorn_nginx_docker
# ------------------------------------------------------------------------
from flask import Flask
from flask import request, jsonify
# ------------------------------------------------------------------------
server = Flask(__name__)
# ------------------------------------------------------------------------
#request data from a specified resource
def run_request_GET():
    return 'Get response OK\n'
# ------------------------------------------------------------------------
#send data to a server to create/update a resource
#POST method is called when you have to add a child resource
def run_request_POST():

    data_dct = request.json
    response = jsonify(data_dct)
    return response
# ------------------------------------------------------------------------
#send data to a server to create/update a resource
#calling the same PUT request multiple times will always produce the same result (PUT requests are idempotent)
#PUT method is called when you have to modify a single resource
def run_request_PUT():
    data_dct = request.json
    response = jsonify(data_dct)
    return response
# ------------------------------------------------------------------------
@server.route('/', methods=['GET', 'POST', 'PUT'])
def hello_world():
    if   request.method == 'GET': return run_request_GET()
    elif request.method == 'POST':return run_request_POST()
    elif request.method == 'PUT' :return run_request_PUT()
# ------------------------------------------------------------------------
if __name__ == "__main__":
    server.run(debug=True,port=8000)
    #curl -X GET  http://127.0.0.1:8000/
    #curl -X POST http://127.0.0.1:8000/ -d {\"key1\":\"value1\",\"key2\":\"value2\"} -H "Content-Type: application/json"

    #curl -w "\n" -s https://api.ipify.org
    #curl -X GET  http://34.122.156.46
    #curl -X POST http://34.122.156.46 -d {\"key1\":\"value1\",\"key2\":\"value2\"} -H "Content-Type: application/json"
```

```dockerfile
FROM python:3.6.7
pip install flask gunicorn
COPY . .
```

```yaml
version: '3'

services:
  flask_app:
    container_name: flask_app
    restart: always
    build: ./flask_app
    ports:
      - "8000:8000"
    command: gunicorn -w 1 -b 0.0.0.0:8000 main:server

  nginx:
    container_name: nginx
    restart: always
    build: ./nginx
    ports:
      - "80:80"
    depends_on:
      - flask_app
```
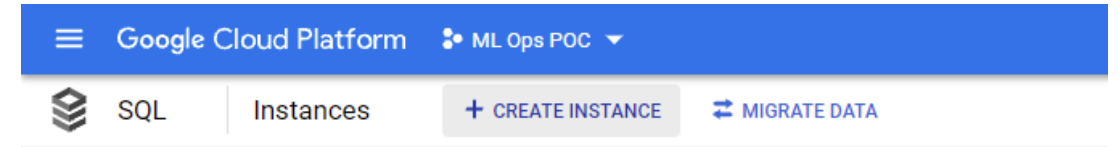
# Dataproc, BigQuery, and Apache Spark for ML

https://cloud.google.com/dataproc/docs/tutorials/bigquery-sparkml#spark-ml-tutorial_regression-console

# Dataproc, BigQuery, and Apache Spark for ML
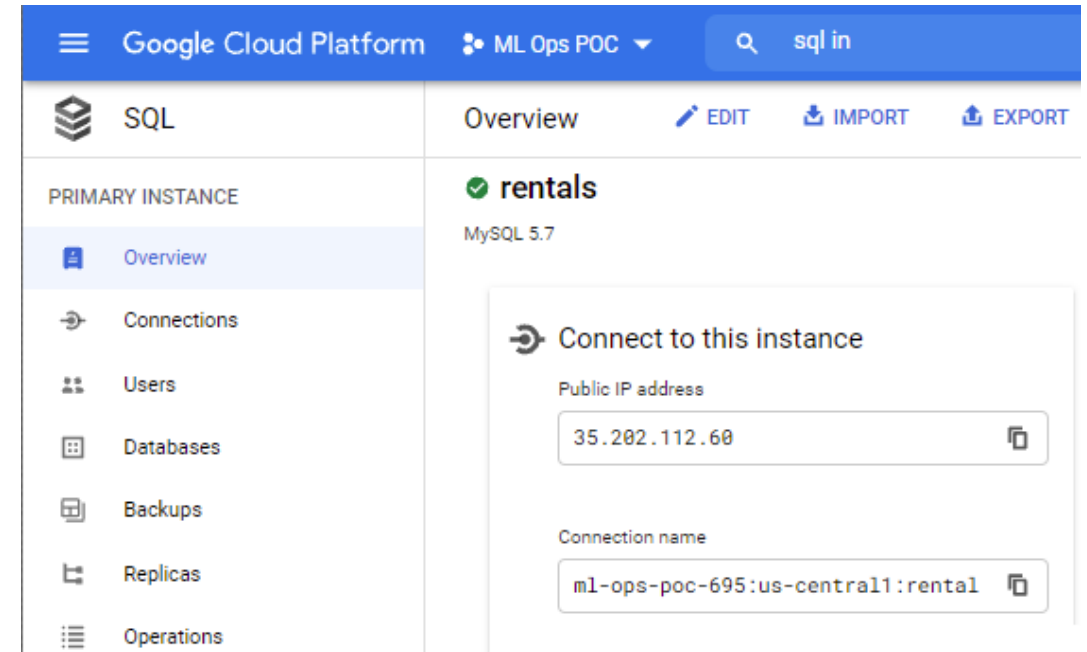
## 1. Create Cloud MySQL instance

```
#sudo apt-get install mysql\*
#sudo gcloud sql connect rentals
```

# Dataproc, BigQuery, and Apache Spark for ML

## 2. Create Tables in Cloud MySQL

```
MySQL [(none)]> CREATE DATABASE IF NOT EXISTS
recommendation_spark;
DROP TABLE IF EXISTS Accommodation;
CREATE TABLE IF NOT EXISTS Accommodation
(
  id varchar(255),
  title varchar(255),
  location varchar(255),
  price int,
  rooms int,
  rating float,
  type varchar(255),
  PRIMARY KEY (ID)
);
CREATE TABLE  IF NOT EXISTS Rating
(
  userId varchar(255),
  accoId varchar(255),
  rating int,
  PRIMARY KEY(accoId, userId),
  FOREIGN KEY (accoId)
    REFERENCES Accommodation(id)
);
CREATE TABLE  IF NOT EXISTS Recommendation
(
  userId varchar(255),
  accoId varchar(255),
  prediction float,
  PRIMARY KEY(userId, accoId),
  FOREIGN KEY (accoId)
    REFERENCES Accommodation(id)
);
```



15

# Dataproc, BigQuery, and Apache Spark for ML

## 3. Export data to Cloud Storage Bucket

# Dataproc, BigQuery, and Apache Spark for ML

## 4. Import data from Bucket to Cloud SQL

## 5. Explore Cloud SQL data

```
sudo gcloud sql connect rentals
SHOW DATABASES;
USE recommendation_spark;
SELECT * FROM Accommodation limit 100;
```

# Dataproc, BigQuery, and Apache Spark for ML

## 6. Dataproc cluster setup



```
echo "Authorizing Cloud Dataproc to connect with Cloud SQL"
CLUSTER=rentals
CLOUDSQL=rentals
ZONE=us-central1-f
NWORKERS=2
machines="$CLUSTER-m"
for w in `seq 0 $(($NWORKERS - 1))`; do
    machines="$machines $CLUSTER-w-$w"
done
echo "Machines to authorize: $machines in $ZONE ... finding
their IP addresses"
ips=""
for machine in $machines; do
    IP_ADDRESS=$(gcloud compute instances describe $machine -
-zone=$ZONE --
format='value(networkInterfaces.accessConfigs[].natIP)' | sed
"s/\['//g" | sed "s/'\]//g" )/32
    echo "IP address of $machine is $IP_ADDRESS"
    if [ -z  $ips ]; then
        ips=$IP_ADDRESS
    else
        ips="$ips,$IP_ADDRESS"
    fi
done
echo "Authorizing [$ips] to access cloudsql=$CLOUDSQL"
gcloud sql instances patch $CLOUDSQL --authorized-networks
$ips
```

# Dataproc, BigQuery, and Apache Spark for ML

## 7. Prepare PySpark script at GS

```
gsutil cp gs://cloud-training/bdml/v2.0/model/train_and_apply.py train_and_apply.py

#patch file with credentials
#cloudshell edit train_and_apply.py

gsutil cp train_and_apply.py gs://$DEVSHELL_PROJECT_ID
```

# Dataproc, BigQuery, and Apache Spark for ML

## 8. Run PySpak script

# Dataproc, BigQuery, and Apache Spark for ML

## 9. Explore results

```
sudo gcloud sql connect rentals
SHOW DATABASES;
USE recommendation_spark;
SELECT * FROM Recommendation limit 10;
```

```
rsa-key-20200330@instance-test-dr:~$ sudo gcloud sql connect rentals
Allowlisting your IP for incoming connection for 5 minutes...done.
Connecting to database with SQL user [root].Enter password:
Welcome to the MariaDB monitor.  Commands end with ; or \g.
Your MySQL connection id is 10269
Server version: 5.7.34-google-log (Google)

Copyright (c) 2000, 2018, Oracle, MariaDB Corporation Ab and others.

Type 'help;' or '\h' for help. Type '\c' to clear the current input statement.

MySQL [(none)]> USE recommendation_spark;
Reading table information for completion of table and column names
You can turn off this feature to get a quicker startup with -A

Database changed
MySQL [recommendation_spark]> show tables;
+--------------------------------+
| Tables_in_recommendation_spark |
+--------------------------------+
| Accommodation                  |
| Rating                         |
| Recommendation                 |
+--------------------------------+
3 rows in set (0.003 sec)

MySQL [recommendation_spark]> SELECT * FROM Recommendation limit 10;
+--------+--------+------------+
| userId | accoId | prediction |
+--------+--------+------------+
| 6      | 30     |   4.289363 |
| 6      | 12     |  4.2010007 |
| 6      | 38     |  4.0971465 |
| 18     | 61     |  2.1708128 |
| 18     | 33     |  2.1591156 |
| 6      | 75     |  3.8684156 |
| 7      | 34     |   2.206352 |
| 19     | 59     |  2.7065306 |
| 7      | 54     |   2.006525 |
| 19     | 66     |  2.5661232 |
+--------+--------+------------+
10 rows in set (0.003 sec)

MySQL [recommendation_spark]> █
```
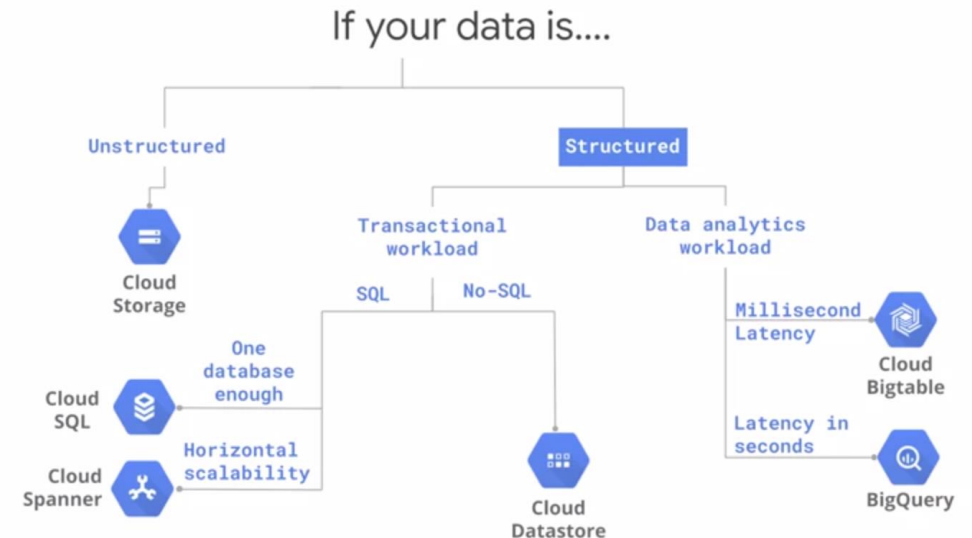
# References

https://cloud.google.com/dataproc/docs/tutorials/bigquery-sparkml#spark-ml-tutorial_regression-console

# Dataproc, BigQuery, and Apache Spark for ML



Choose your solutions based on access pattern

| | Cloud Storage | Cloud SQL | Datastore | Bigtable | BigQuery |
|---|---|---|---|---|---|
| Capacity | Petabytes + | Gigabytes | Terabytes | Petabytes | Petabytes |
| Access metaphor | Like files in a file system | Relational database | Persistent Hashmap | Key-value(s), HBase API | Data warehouse |
| Read | Have to copy to local disk | SELECT rows | filter objects on property | scan rows | SELECT rows |
| Write | One file | INSERT row | put object | put row | Batch/stream |
| Update granularity | An object (a "file") | Field | Attribute | Row | Field |
| Usage | Store blobs | No-ops SQL database on the cloud | Structured data from AppEngine apps | No-ops, high throughput, scalable, flattened data | Interactive SQL* querying fully managed warehouse |

# Dataproc, BigQuery, and Apache Spark for ML



Rich open-source ecosystem for big data

1 — Hadoop is the canonical open-source MapReduce framework.

2 — Pig provides a convenient scripting language that can be compiled into Hadoop MapReduce jobs.

3 — Hive is a data warehousing system and query language.

4 — Spark is a fast, interactive, general-purpose framework for SQL, streaming, machine learning, etc.