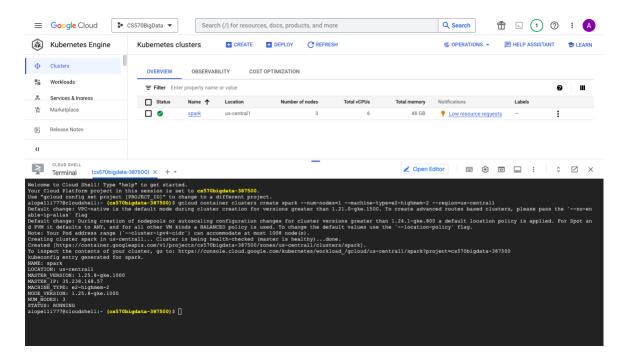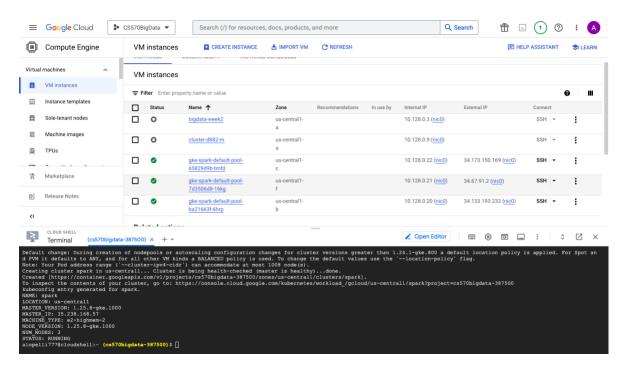1. Create a cluster on Google Kubernetes Engine(GKE) by running the below command on the cloud shell on GCP

```
gcloud container clusters create spark --num-nodes=1 --machine-type=e2-highmem-2 --region=us-central1
```



We can see the 3 nodes that are created

2. Create image and deploy spark to Kubernetes

- Install the NFS Server Provisioner

```
helm repo add stable https://charts.helm.sh/stable
```



```
helm repo update
```



```
helm install nfs stable/nfs-server-provisioner \
--set persistence.enabled=true,persistence.size=5Gi
```

3. To create a persistent disk volume and a pod to use NFS - create a yaml file with name spar-pvc.yaml and insert the code

```
vi spark-pvc.yaml
```

```
kind: PersistentVolumeClaim
apiVersion: v1
metadata:
  name: spark-data-pvc
spec:
  accessModes:
    - ReadWriteMany
  resources:
    requests:
      storage: 2Gi
  storageClassName: nfs
---
apiVersion: v1
kind: Pod
metadata:
  name: spark-data-pod
spec:
  volumes:
    - name: spark-data-pv
      persistentVolumeClaim:
        claimName: spark-data-pvc
  containers:
    - name: inspector
      image: bitnami/minideb
      command:
        - sleep
        - infinity
      volumeMounts:
        - mountPath: "/data"
          name: spark-data-pv
```

We can see this code on the cloud shell with the command:

```
cat spark-pvc.yaml
```

```
alopelli777@cloudshell:~ (cs570bigdata-387500)$ vi spark-pvc.yaml
alopelli777@cloudshell:~ (cs570bigdata-387500)$ cat spark-pvc.yaml
kind: PersistentVolumeClaim
apiVersion: v1
metadata:
  name: spark-data-pvc
spec:
  accessModes:
    - ReadWriteMany
  resources:
    requests:
      storage: 2Gi
  storageClassName: nfs
---
apiVersion: v1
kind: Pod
metadata:
  name: spark-data-pod
spec:
  volumes:
    - name: spark-data-pv
      persistentVolumeClaim:
        claimName: spark-data-pvc
  containers:
    - name: inspector
      image: bitnami/minideb
      command:
        - sleep
        - infinity
      volumeMounts:
        - mountPath: "/data"
          name: spark-data-pv
alopelli777@cloudshell:~ (cs570bigdata-387500)$
```

4. Apply the above yaml descriptor

```
kubectl apply -f spark-pvc.yaml
```

```
alopelli777@cloudshell:~ (cs570bigdata-387500)$ kubectl apply -f spark-pvc.yaml
persistentvolumeclaim/spark-data-pvc created
pod/spark-data-pod created
alopelli777@cloudshell:~ (cs570bigdata-387500)$
```

5. Create and prepare your application JAR file

```
docker run -v /tmp:/tmp -it bitnami/spark -- find /opt/bitnami/spark/examples/jars/ -name
spark-examples* -exec cp {} /tmp/my.jar \;
```

```
alopelli777@cloudshell:~ (cs570bigdata-387500)$ docker run -v /tmp:/tmp -it bitnami/spark -- find /opt/bitnami/spark/examples/jars/ -name spark-examples* -exec cp {} /tmp/my.jar \;
spark 19:35:22.13
spark 19:35:22.13 Welcome to the Bitnami spark container
spark 19:35:22.13 Subscribe to project updates by watching https://github.com/bitnami/containers
spark 19:35:22.13 Submit issues and feature requests at https://github.com/bitnami/containers/issues
spark 19:35:22.13
```

6. Add a test file with a line of words that we will be using later for the word count test

```
echo "the quick brown fox the fox ate the mouse how now brown cow" > /tmp/test.txt
```

```
alopelli777@cloudshell:~ (cs570bigdata-387500)$ echo "the quick brown fox the fox ate the mouse how now brown cow" > /tmp/test.txt
alopelli777@cloudshell:~ (cs570bigdata-387500)$
```

7. Copy the JAR file containing the application, and any other required files, to the PVC using the mount point.

```
kubectl cp /tmp/my.jar spark-data-pod:/data/my.jar
kubectl cp /tmp/test.txt spark-data-pod:/data/test.txt
```

```
alopelli777@cloudshell:~ (cs570bigdata-387500)$ kubectl cp /tmp/my.jar spark-data-pod:/data/my.jar
alopelli777@cloudshell:~ (cs570bigdata-387500)$ kubectl cp /tmp/test.txt spark-data-pod:/data/test.txt
alopelli777@cloudshell:~ (cs570bigdata-387500)$
```

8. Make sure the files a inside the persistent volume

```
kubectl exec -it spark-data-pod -- ls -al /data
```

```
alopelli777@cloudshell:~ (cs570bigdata-387500)$ kubectl exec -it spark-data-pod -- ls -al /data
total 1540
drwxrwsrwx 2 root root     4096 Jun 27 19:35 .
drwxr-xr-x 1 root root     4096 Jun 27 19:14 ..
-rw-r--r-- 1 1001 root 1564259 Jun 27 19:35 my.jar
-rw-r--r-- 1 1000 1000      60 Jun 27 19:35 test.txt
alopelli777@cloudshell:~ (cs570bigdata-387500)$
```

9. Deploy Apache Spark on Kubernetes using the shared volume spark-chart.yaml:

```
nano spark-chart.yaml
cat spark-chart.yaml
```

```
alopelli777@cloudshell:~ (cs570bigdata-387500)$ nano spark-chart.yaml
alopelli777@cloudshell:~ (cs570bigdata-387500)$ cat spark-chart.yaml
service:
  type: LoadBalancer
worker:
  replicaCount: 3
  extraVolumes:
    - name: spark-data
      persistentVolumeClaim:
        claimName: spark-data-pvc
  extraVolumeMounts:
    - name: spark-data
      mountPath: /data
alopelli777@cloudshell:~ (cs570bigdata-387500)$
```

10. Check the pods is running:

```
kubectl get pods
```

```
alopelli777@cloudshell:~ (cs570bigdata-387500)$ kubectl get pods
NAME                             READY   STATUS    RESTARTS   AGE
nfs-nfs-server-provisioner-0     1/1     Running   0          80m
spark-data-pod                   1/1     Running   0          46m
alopelli777@cloudshell:~ (cs570bigdata-387500)$
```

11. Deploy Apache Spark on the Kubernetes cluster using the Bitnami Apache Spark Helm chart and supply it with the configuration file above

```
helm repo add bitnami https://charts.bitnami.com/bitnami
helm install spark bitnami/spark -f spark-chart.yaml
```



12. Get the external IP of the running pod

```
kubectl get svc -l "app.kubernetes.io/instance=spark,app.kubernetes.io/name=spark"
```



13. Open the external ip on your browser( I did by pasting the 34.171.254.91 in a separate browser)

Word Count on Spark

1. Submit a word count task and you see the below content after running the command

```
kubectl run --namespace default spark-client --rm --tty -i --restart='Never' \
    --image docker.io/bitnami/spark:3.4.1-debian-11-r3 \
    -- spark-submit --master spark://34.27.61.122:7077 \
    --deploy-mode cluster \
    --class org.apache.spark.examples.JavaWordCount \
    /data/my.jar /data/test.txt
```



2. And on your browser, you should see this task finished

View the output of the completed jobs

1. On the browser, you should see the worker node ip address of the finished task



2. Get the name of the worker node( my worker node address is 10.52.0.4)

```
kubectl get pods -o wide | grep WORKER-NODE-ADDRESS
```



3. Execute this pod and see the result of the finished tasks

```
kubectl exec -it <Worker node name> -- bash
```



```
cd /opt/bitnami/spark/work
cat <task-name>/stdout
```

The task name here is the Submission ID in the completed Drivers section of the URL



Running python PageRank onPySpark on the pods

1. Execute the spark master pods  and Go to the directory where pagerank.py located

```
kubectl exec -it spark-master-0 – bash
cd /opt/bitnami/spark/examples/src/main/python
```

```
alopelli777@cloudshell:~ (cs570bigdata-387500)$ kubectl exec -it spark-master-0 -- bash
I have no name!@spark-master-0:/opt/bitnami/spark$ cd /opt/bitnami/spark/examples/src/main/python
I have no name!@spark-master-0:/opt/bitnami/spark/examples/src/main/python$ ls
__init__.py  avro_inputformat.py  logistic_regression.py  mllib        parquet_inputformat.py  sort.py  status_api_demo.py  transitive_closure.py
als.py       kmeans.py            ml                      pagerank.py  pi.py                   sql      streaming           wordcount.py
I have no name!@spark-master-0:/opt/bitnami/spark/examples/src/main/python$
```

2. Run the pagerank using pyspark

spark-submit pagerank.py /opt 2

Notice, /opt is an example directory, you can enter any directory you like, and 2 is the number of iterations you want the pagerank to run, you can also change to any numbers you like

```
        file:/opt/bitnami/spark/venv/lib/python3.9/site-packages/pandas/compat/numpy
        file:/opt/bitnami/spark/examples/src/main/java/org/apache/spark/examples/sql
        file:/opt/bitnami/spark/venv/lib/python3.9/site-packages/botocore/data/sagemaker-edge/2020-09-23
        file:/opt/bitnami/python/lib/python3.9/site-packages/virtualenv/discovery/windows
        file:/opt/bitnami/spark/venv/lib/python3.9/site-packages/botocore/data/cognito-idp/2016-04-18
        file:/opt/bitnami/spark/venv/lib/python3.9/site-packages/awscli/examples/robomaker
        file:/opt/bitnami/spark/venv/lib/python3.9/site-packages/botocore/data/elb/2012-06-01
        file:/opt/bitnami/spark/venv/lib/python3.9/site-packages/botocore/data/verifiedpermissions/2021-12-01
        file:/opt/bitnami/spark/venv/lib/python3.9/site-packages/numpy/f2py/tests/src/return_real
        file:/opt/bitnami/spark/venv/lib/python3.9/site-packages/botocore/data/cloudfront/2016-08-20
        file:/opt/bitnami/spark/licenses
        file:/opt/bitnami/spark/venv/lib/python3.9/site-packages/pandas/io/sas
        file:/opt/bitnami/python/lib/python3.9/config-3.9-x86_64-linux-gnu
        file:/opt/bitnami/spark/venv/lib/python3.9/site-packages/awscli/customizations/ec2
        file:/opt/bitnami/spark/data/mllib
        file:/opt/bitnami/spark/venv/lib/python3.9/site-packages/botocore/data/directconnect/2012-10-25
        file:/opt/bitnami/python/lib/python3.9/multiprocessing
        file:/opt/bitnami/spark/venv/lib/python3.9/site-packages/rsa
        file:/opt/bitnami/python/lib/python3.9/test/test_import/data/unwritable
        file:/opt/bitnami/spark/venv/lib/python3.9/site-packages/setuptools/config
        file:/opt/bitnami/spark/examples/src/main/scala/org/apache/spark/examples/pythonconverters
        file:/opt/bitnami/spark/venv/lib/python3.9/site-packages/botocore/data/ecr/2015-09-21
        file:/opt/bitnami/spark/venv/lib/python3.9/site-packages/awscli/examples/organizations
        file:/opt/bitnami/spark/venv/lib/python3.9/site-packages/awscli/examples/grafana
        file:/opt/bitnami/spark/python/pyspark/sql/connect
        file:/opt/bitnami/python/lib/python3.9/http
        file:/opt/bitnami/java/legal/jdk.accessibility
        file:/opt/bitnami/java/legal/jdk.internal.opt
        file:/opt/bitnami/spark/venv/lib/python3.9/site-packages/botocore/data/sqs/2012-11-05
        file:/opt/bitnami/spark/venv/lib/python3.9/site-packages/botocore/data/autoscaling/2011-01-01
        file:/opt/bitnami/spark/examples/src/main/scala/org/apache/spark/examples/mllib
        file:/opt/bitnami/spark/venv/lib/python3.9/site-packages/awscli/examples/xray
        file:/opt/bitnami/spark/venv/lib/python3.9/site-packages/awscli/examples/acm-pca
        file:/opt/bitnami/spark/venv/lib/python3.9/site-packages/awscli/examples/ecs/wait
        file:/opt/bitnami/spark/python/pyspark/testing
        file:/opt/bitnami/spark/venv/lib/python3.9/site-packages/awscli/examples/globalaccelerator
        file:/opt/bitnami/spark/venv/lib/python3.9/site-packages/dateutil
        file:/opt/bitnami/spark/venv/lib/python3.9/site-packages/pytz/zoneinfo/arctic
        file:/opt/bitnami/spark/venv/lib/python3.9/site-packages/botocore/data/oam/2022-06-10
        file:/opt/bitnami/java/legal/jdk.jcmd


23/07/13 03:27:22 INFO SparkContext: Invoking stop() from shutdown hook
23/07/13 03:27:22 INFO SparkContext: SparkContext is stopping with exitCode 0.
23/07/13 03:27:22 INFO SparkUI: Stopped Spark web UI at http://spark-master-0.spark-headless.default.svc.cluster.local:4040
23/07/13 03:27:22 INFO MapOutputTrackerMasterEndpoint: MapOutputTrackerMasterEndpoint stopped!
23/07/13 03:27:22 INFO MemoryStore: MemoryStore cleared
23/07/13 03:27:22 INFO BlockManager: BlockManager stopped
23/07/13 03:27:22 INFO BlockManagerMaster: BlockManagerMaster stopped
23/07/13 03:27:22 INFO OutputCommitCoordinator$OutputCommitCoordinatorEndpoint: OutputCommitCoordinator stopped!
23/07/13 03:27:22 INFO SparkContext: Successfully stopped SparkContext
23/07/13 03:27:22 INFO ShutdownHookManager: Shutdown hook called
23/07/13 03:27:22 INFO ShutdownHookManager: Deleting directory /tmp/spark-01f7afe7-6ee4-4cdf-8c5f-5faaae65574e
23/07/13 03:27:22 INFO ShutdownHookManager: Deleting directory /tmp/spark-3fc97a6b-a86c-4361-8f81-7654985d201c
23/07/13 03:27:22 INFO ShutdownHookManager: Deleting directory /tmp/spark-3fc97a6b-a86c-4361-8f81-7654985d201c/pyspark-67b7452b-d114-440d-8f78-dbf7148c3fc0
I have no name!@spark-master-0:/opt/bitnami/spark/examples/src/main/python$
```