## 1. Moving files to bucket



## 2. Creating Code



```python
from pyspark.sql import SparkSession

def map_time(s):
    val=('x',1)
    if s!='Time':
        final = int(s.replace(":",""))
        if final >=0 and final<=600:
            val=("00-06",1)
        elif final > 600 and final<=1200:
            val=("06-12",1)
        elif final > 1200 and final <=1800:
            val=("12-18",1)
        elif final>1800 and final <=2400:
            val=("18-24",1)
    return val

spark= SparkSession.builder.appName("myFile").getOrCreate()

df = spark.read.text("gs://my_bucket_20/hash_file.txt")

rdd = df.rdd

sep = rdd.map(lambda x :x[0].split("\t"))

time = sep.map(lambda x:x[1])

time_sep = time.map(lambda x: map_time(x))

sorting = time_sep.filter(lambda x:x[0] != 'x').sortBy(lambda x: x[0])

groupdata= sorting.reduceByKey(lambda a,b: a+b)

print(groupdata.collect())
```
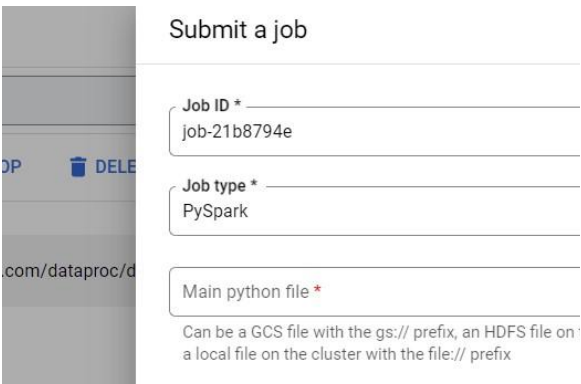
3. The TEXT File



4. Creating Clusters



5. Creating Job

## 6. The final Output

← Job details    ⬚ CLONE    🗑 DELETE    ■ STOP    ↻ REFRESH

| | |
|---|---|
| Job ID | GA3-job-059ce8cd |
| Job UUID | c1d7ccb1-e89f-46f6-b88b-cf6f0e372a43 |
| Type | Dataproc Job |
| Status | ✔ Succeeded |

**MONITORING**    CONFIGURATION

ℹ The charts below represent the metrics from the cluster this job ran on, scoped to the time that this job was running. It is possible for more than one job to run on a cluster a

### Output    LINE WRAP: OFF

ℹ Spark jobs take ~60 seconds to initialize resources.

```
23/04/13 03:49:26 INFO org.apache.hadoop.yarn.client.api.impl.YarnClientImpl: Submitted application application_1681336220566_0008
23/04/13 03:49:27 INFO org.apache.hadoop.yarn.client.RMProxy: Connecting to ResourceManager at my-cluster-6ebf-m/10.128.0.2:8030
23/04/13 03:49:29 INFO com.google.cloud.hadoop.repackaged.gcs.com.google.cloud.hadoop.gcsio.GoogleCloudStorageImpl: Ignoring exception of type GoogleJsonRespons
[('00-06', 4), ('06-12', 8), ('12-18', 12), ('18-24', 6)]
23/04/13 03:49:50 INFO org.sparkproject.jetty.server.AbstractConnector: Stopped Spark@54e29800{HTTP/1.1, (http/1.1)}{0.0.0.0:0}
```