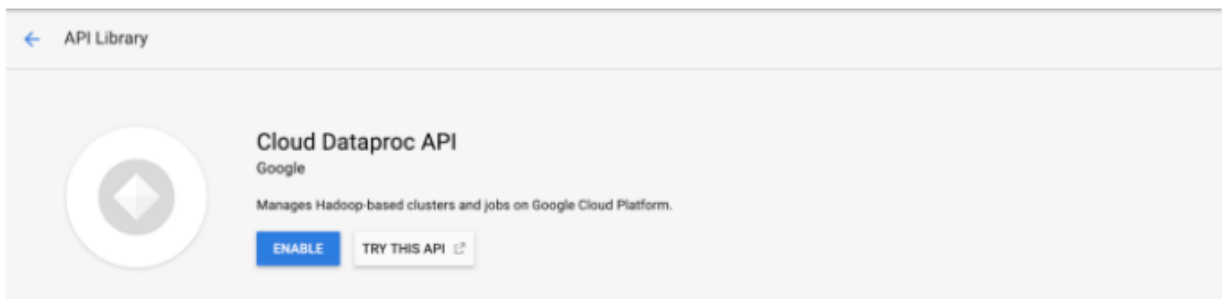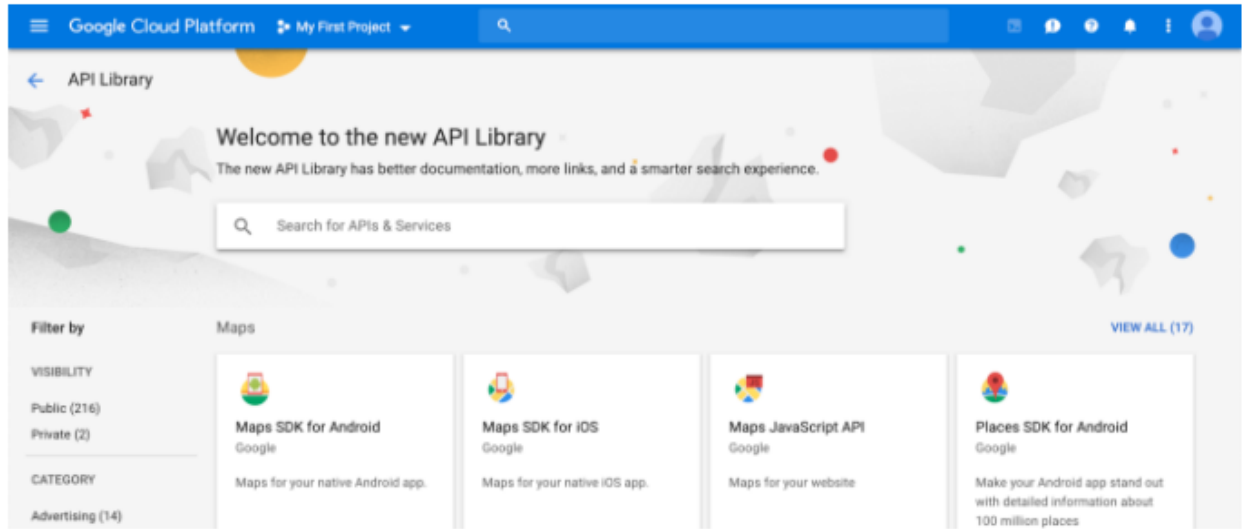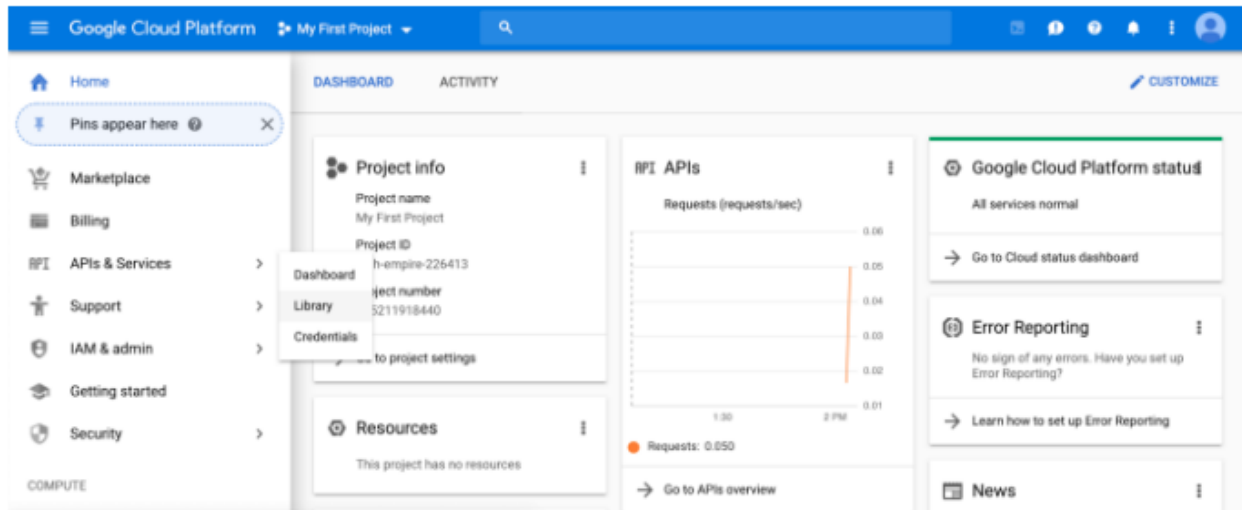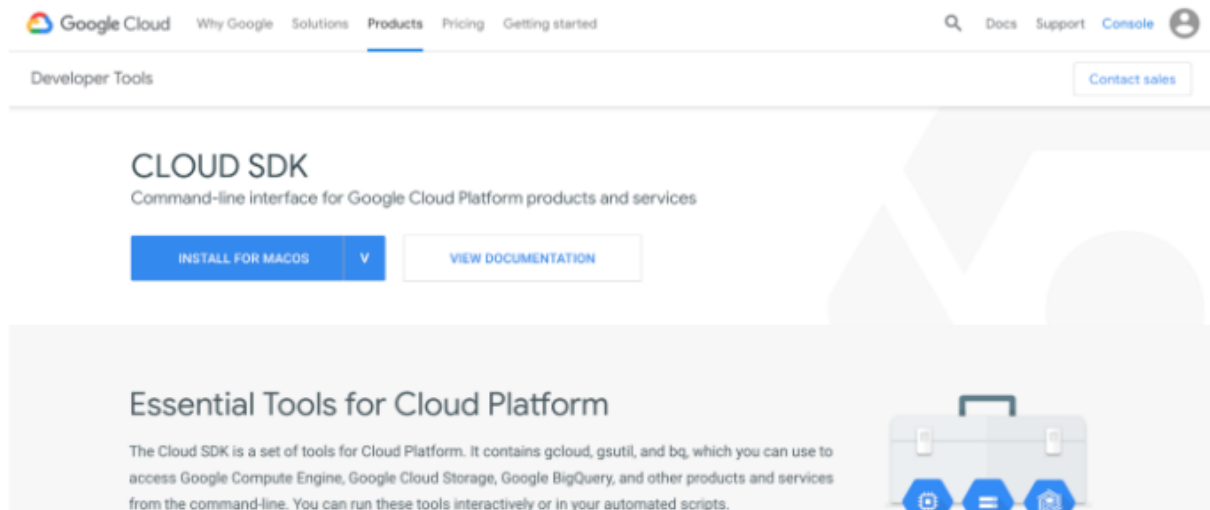# SENTIMENT ANALYSIS ON GOOGLE CLOUD DATAPROC USING PYSPARK

## SERIES OF STEPS PERFORMED IN GOOGLE CLOUD PLATFORM USING CLOUD DATAPRO

# SENTIMENT ANALYSIS ON GOOGLE CLOUD DATAPROC USING PYSPARK
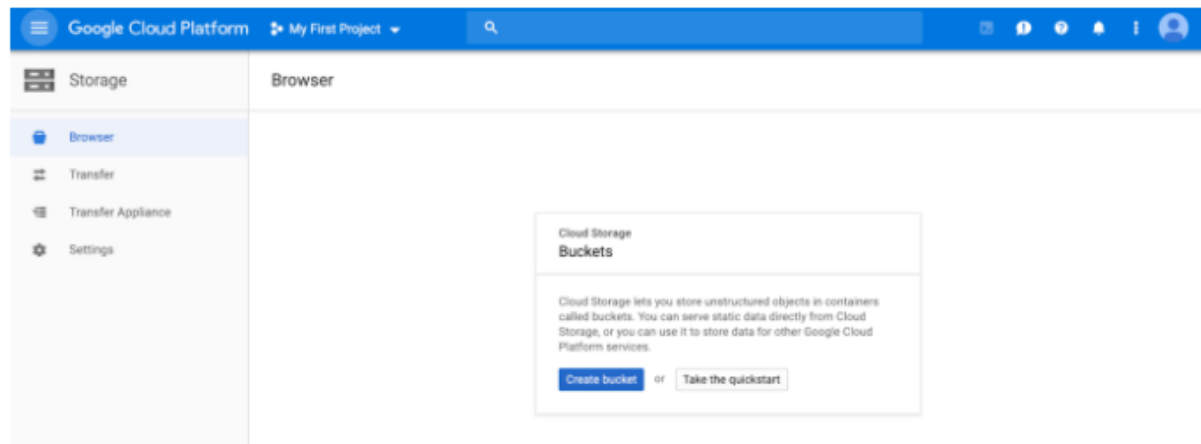


## CLOUD SDK
Command-line interface for Google Cloud Platform products and services

**INSTALL FOR MACOS** | V | **VIEW DOCUMENTATION**

## Essential Tools for Cloud Platform

The Cloud SDK is a set of tools for Cloud Platform. It contains gcloud, gsutil, and bq, which you can use to access Google Compute Engine, Google Cloud Storage, Google BigQuery, and other products and services from the command-line. You can run these tools interactively or in your automated scripts.

Install Google Cloud SDK by following instructions on https://cloud.google.com/sdk/

# SENTIMENT ANALYSIS ON GOOGLE CLOUD DATAPROC USING PYSPARK



```
Warning: wget 1.20 is already installed and up-to-date
To reinstall 1.20, run `brew reinstall wget`
URL transformed to HTTPS due to an HSTS policy
--2018-12-23 20:35:57--  https://cs.stanford.edu/people/alecmgo/trainingandtestdata.zip
Resolving cs.stanford.edu (cs.stanford.edu)... 171.64.64.64
Connecting to cs.stanford.edu (cs.stanford.edu)|171.64.64.64|:443... connected.
HTTP request sent, awaiting response... 200 OK
Length: 81363704 (78M) [application/zip]
Saving to: 'trainingandtestdata.zip'

trainingandtestdata.zip  100%[===============================>]  77.59M  6.54MB/s    in 13s

2018-12-23 20:36:11 (5.76 MB/s) - 'trainingandtestdata.zip' saved [81363704/81363704]

Archive:  trainingandtestdata.zip
  inflating: testdata.manual.2009.06.14.csv
  inflating: training.1600000.processed.noemoticon.csv
Copying file://pyspark_sa_train_data.csv [Content-Type=text/csv]...
==> NOTE: You are uploading one or more large file(s), which would run
significantly faster if you enable parallel composite uploads. This
feature can be enabled by editing the
"parallel_composite_upload_threshold" value in your .boto
configuration file. However, note that if you do this large files will
be uploaded as `composite objects
<https://cloud.google.com/storage/docs/composite-objects>`_,which
means that any user who downloads such objects will need to have a
compiled crcmod installed (see "gsutil help crcmod"). This is because
without a compiled crcmod, computing checksums on composite objects is
so slow that gsutil disables downloads of composite objects.

| [1 files][219.2 MiB/219.2 MiB]  304.0 KiB/s
Operation completed over 1 objects/219.2 MiB.
Copying file://pyspark_sa_test_data.csv [Content-Type=text/csv]...
/ [1 files][  2.2 MiB/  2.2 MiB]
Operation completed over 1 objects/2.2 MiB.
```

# SENTIMENT ANALYSIS ON GOOGLE CLOUD DATAPROC USING PYSPARK

# SENTIMENT ANALYSIS ON GOOGLE CLOUD DATAPROC USING PYSPARK

# SENTIMENT ANALYSIS ON GOOGLE CLOUD DATAPROC USING PYSPARK