# CPR E 419 - Lab 6
# Spark

**1.**

Below are screenshots of the commands run and successful completion of program execution.



Output Screenshot: (from 2 files)

```
Java 462182 elasticsearch/elasticsearch 15698
Ruby 363801 rails/rails 29825
Python 331883 jakubroztocil/httpie 21283
PHP 273999 zurb/foundation 22636
C++ 159831 rogerwang/node-webkit 27350
C 145354 neovim/neovim 17395
C# 116155 dotnet/corefx 9176
```

**2.** Verified and attached the results of the data on a smaller dataset. The dataset and results are provided in the submissions folder. Will include results of running the experiment on patents dataset in the final submission.

```
cpre419@cpre419-VirtualBox:~/Downloads$ spark-submit --class NumOfTriangles triangle2.jar /data/patents.txt /op7
2021-04-11 08:45:46,915 WARN util.Utils: Your hostname, cpre419-VirtualBox resolves to a loopback address: 127.0.1.1; using 10.0.2.15 instead (on interface enp0s3)
2021-04-11 08:45:46,916 WARN util.Utils: Set SPARK_LOCAL_IP if you need to bind to another address
2021-04-11 08:45:49,030 INFO client.RMProxy: Connecting to ResourceManager at /0.0.0.0:8032
2021-04-11 08:45:49,550 INFO yarn.Client: Requesting a new application from cluster with 1 NodeManagers
2021-04-11 08:45:49,745 INFO yarn.Client: Verifying our application has not requested more than the maximum memory capability of the cluster (8192 MB per container)
2021-04-11 08:45:49,745 INFO yarn.Client: Will allocate AM container, with 5632 MB memory including 512 MB overhead
2021-04-11 08:45:49,747 INFO yarn.Client: Setting up container launch context for our AM
2021-04-11 08:45:49,749 INFO yarn.Client: Setting up the launch environment for our AM container
2021-04-11 08:45:49,795 INFO yarn.Client: Preparing resources for our AM container
2021-04-11 08:45:50,068 WARN yarn.Client: Neither spark.yarn.jars nor spark.yarn.archive is set, falling back to uploading libraries under SPARK_HOME.
2021-04-11 08:45:54,102 INFO yarn.Client: Uploading resource file:/tmp/spark-2da60134-7e69-4d2a-a7e4-fb9cd67228b5/__spark_libs__8767342431313587164.zip -> hdfs://localhost:9000/user/cpre419/.sparkStaging/
application_1618148656645_0001/__spark_libs__8767342431313587164.zip
2021-04-11 08:45:58,873 INFO yarn.Client: Uploading resource file:/home/cpre419/Downloads/triangle2.jar -> hdfs://localhost:9000/user/cpre419/.sparkStaging/application_1618148656645_0001/triangle2.jar
2021-04-11 08:45:58,910 WARN hdfs.DFSClient: Caught exception
```