

Project 3 Presentation

Team 2

Overview

Intro to Big Data
Project 3

Analyze Twitter's data using Apache Spark.

- ▢ Spark 2.1.0 has APIs to support different languages. For this project PySpark, the python API, is used.
- ▢ Three different analytical tasks are implemented.
Common theme: **How do trends behave?**
- ▢ Built with a Graphical User Interface to provide a better user experience.

The Big Question

On twitter, trending is an important way for ideas to stand out from a sea of tweets. So we ask the question: How do trends behave?

- **Task 1:** Are users with lots of friends and followers important vectors of hashtag spread?
- **Task 2:** Are retweets important vectors of hashtag spread? How do they compare to regular tweets?
- **Task 3:** How do trends change over time? What about the top 3 trends?

Note: Task 1 & Task 2 conclusions are made from data collected over a 15 minute period, containing ~46,000 tweets. Task 3 is based on data provided via class. They are used to demonstrate the analysis, but should not be generalized for all tweets.

Task 1

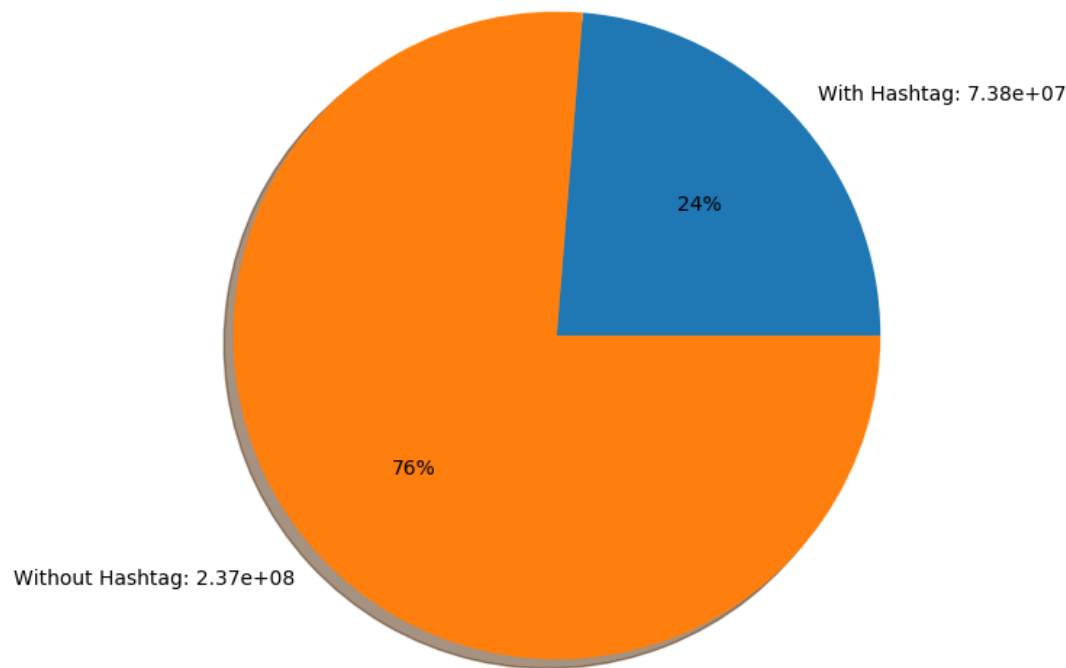
space

Tweets without hashtags have three times as many friends and followers.

Despite this, the correlation between hashtag and follower + friends is weak.

Users with many followers and friends are not important vectors of hashtag spread in this set of data.

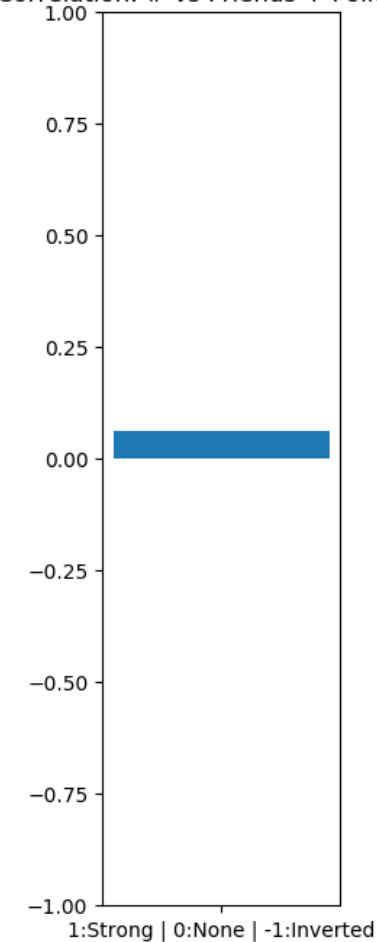
Hashtag Presence vs Followers & Friends (Spark SQL & DataFrames)



Does friend/follower correlate with hashtag usage? Divide tweets into two buckets based on presence or absence of hashtag. Each bucket will hold the cumulative sum of friends and followers. This is shown on the pie chart. A raw count is not enough to establish correlation however. For this, a correlation analysis is done and is shown on the bar chart

Began query at: 2017-05-01 18:24:57

Correlation: # vs Friends + Followers



Task 1

Methodology

From a tweet, four attributes are retrieved:

- user id, hashtag, followers, friends
 - hashtag converted to zero or one, based on absence/presence

For analysis:

- Users that appear multiple times and sometimes do and sometimes don't use hashtags in their tweets are removed.
- Users that appear multiple times but consistently do/don't use hashtags are only counted once.
- Friends + Followers are summed cumulatively by absence/presence of hashtag and displayed in pie chart
- Correlation is done between two lists: the first signifying if tweet has hashtag, and the second holding the sum of friends and followers for that tweet

Task 2

Task 2: Methodology

From a tweet, three attributes are retrieved:

- ▢ hashtag, retweet hashtag, retweet's retweet count
 - ▢ hashtag and retweet hashtag converted to zero or one
 - ▢ RT will stand for retweet, for readability

For analysis:

- ▢ Although collected tweets have a retweet count of zero, if they were a RT, then there is a RT Count on the RT, or RT RT Count
- ▢ If a tweet is a RT that contains a hashtag, see how many times the original RT has been RTed. This tells us how many times the hashtag has been spread via retweeting.
- ▢ Aggregate that into the first category, RTs with #. Compare that with retweet volume that does not contain a hashtag.
- ▢ Finally compare it to new tweet volume (not a RT) with a hashtag.

Task 3

space

Cumulative trend volume, across all trends, varies erratically.

The top three trends, although brief, account for significant portions of surges in trend volume.

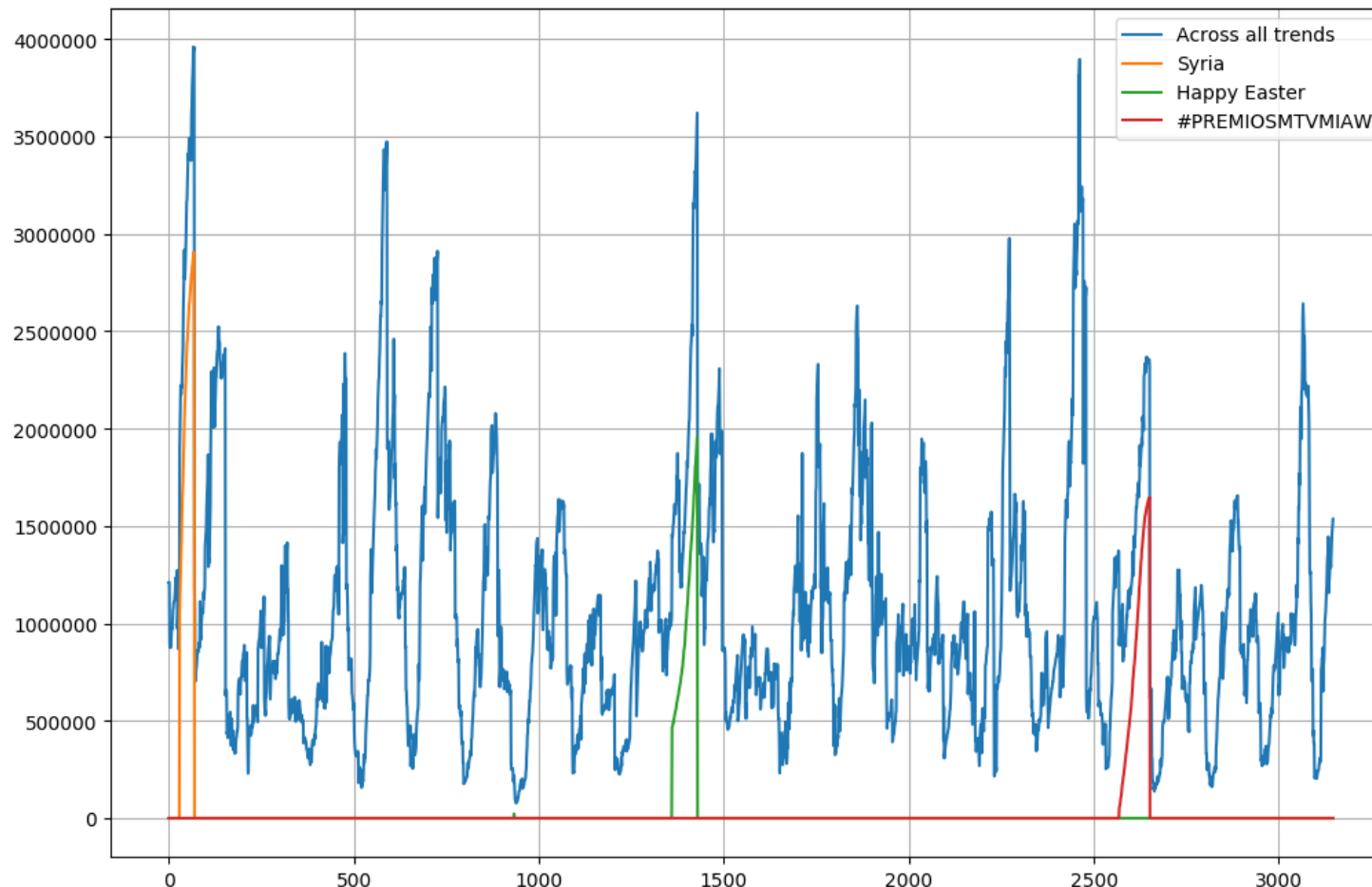
Trends are brief, but can flood the Twitter quickly.

Total Trend Volume over Time - (Spark SQL & DataFrames)

Start: 2017-04-06 20:48:18

End: 2017-04-28 17:47:58

Began query at: 2017-05-01 18:26:11



How do trends behave over time? The main plot shows cumulative volume across all trends over time. In addition, the top three trends are also shown over time. The top three are identified by finding cumulative volume across all time by trend. The top 3 are then plotted, showing behavior.

Task 3: Methodology

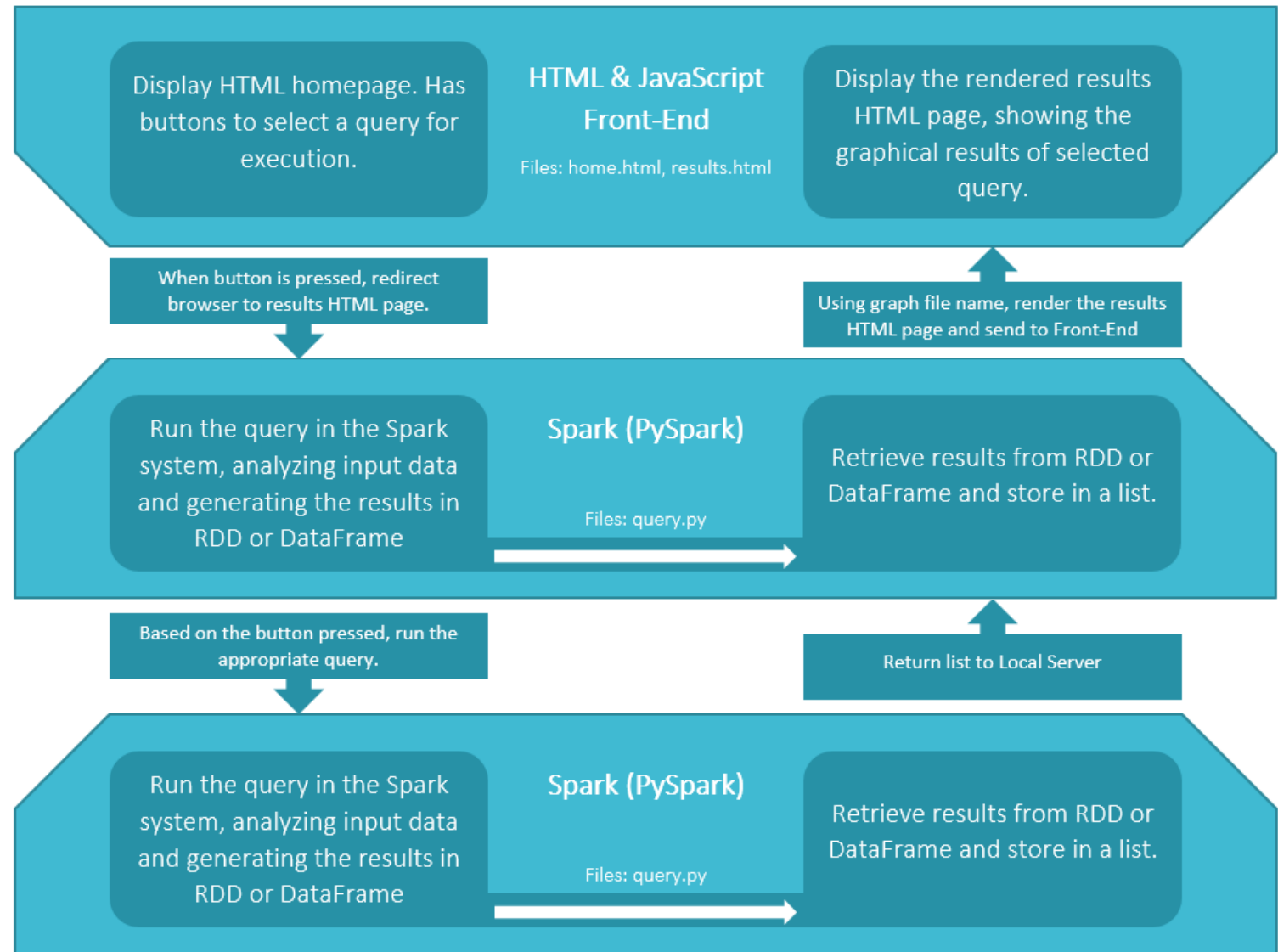
From a query, three attributes are retrieved:

- created_at, name, tweet_volume

For analysis:

- For the main plot, cumulative tweet volume, aggregate the tweet volume across all trends, by time.
- To identify the top 3 trends, aggregate the tweet volume across all time, by trend. The trends with the most tweet volume are identified as the top 3 trends.
- Retrieve tweet volume from the main plot that corresponds to the top 3 trends.
- Plot the top 3 trends over time. Where the top trend is not present, plot a zero.

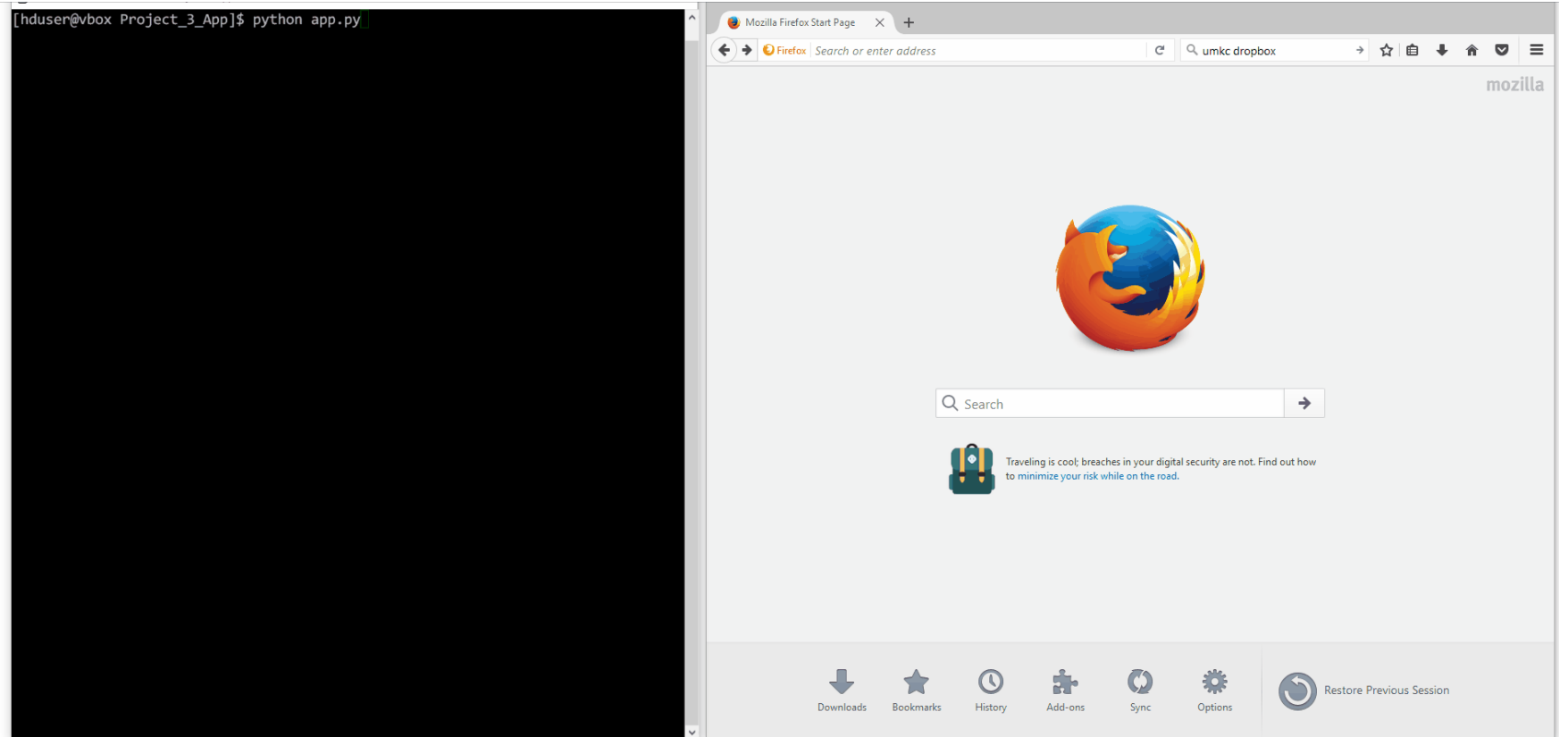
Extra Req. Software Stack



Flask Server and Task 1

Demo

Server console on left
GUI on right



Demo

Server console on left
GUI on right

Task 2

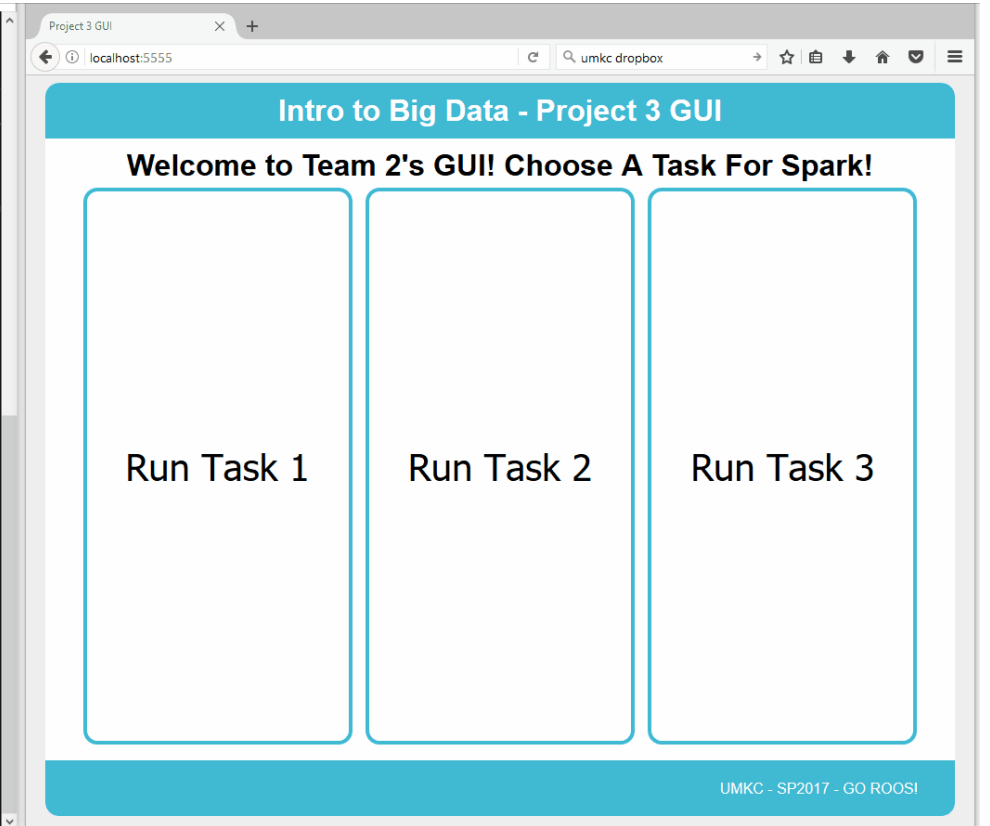
Task 3

Demo

Server console on left
GUI on right

```
[hduser@vbox Project_3_App]$ python app.py
Using Spark's default log4j profile: org/apache/spark/log4j-defaults.properties
Setting default log level to "WARN".
To adjust logging level use sc.setLogLevel(newLevel). For SparkR, use setLogLevel(newLevel).
17/05/01 22:51:45 WARN NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable
17/05/01 22:51:45 WARN Utils: Your hostname, vbox resolves to a loopback address: 127.0.0.1; using 10.0.2.15 instead (on interface enp0s3)
17/05/01 22:51:45 WARN Utils: Set SPARK_LOCAL_IP if you need to bind to another address
* Running on http://0.0.0.0:5000/ (Press CTRL+C to quit)

Rendering home HTML page now...
10.0.2.2 - - [01/May/2017 22:51:50] "GET / HTTP/1.1" 200 -
```



Questions?