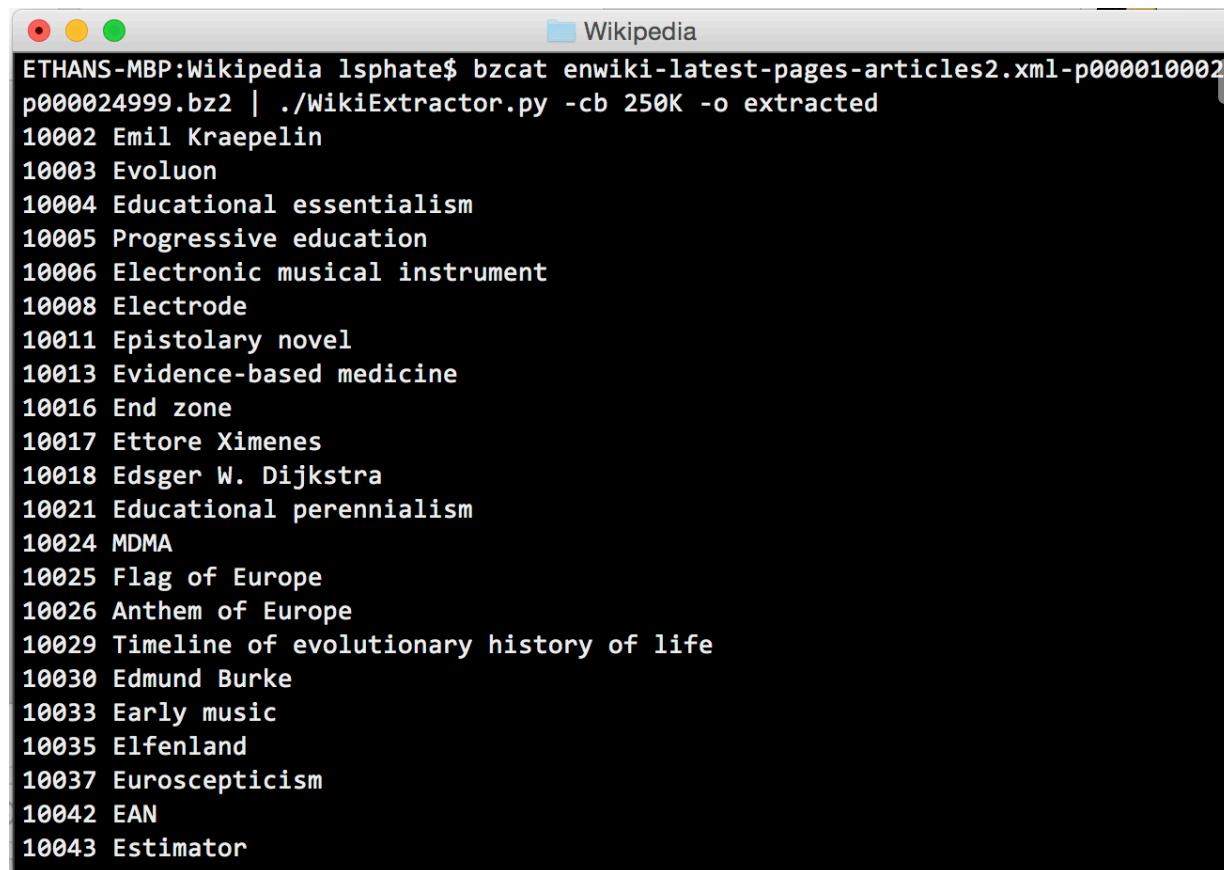


# EECS6895 Adv. Bigdata Analytics HW 1 - Sun-Yi Lin(sl3833)

## Prob. 2 Wikipedia Articles

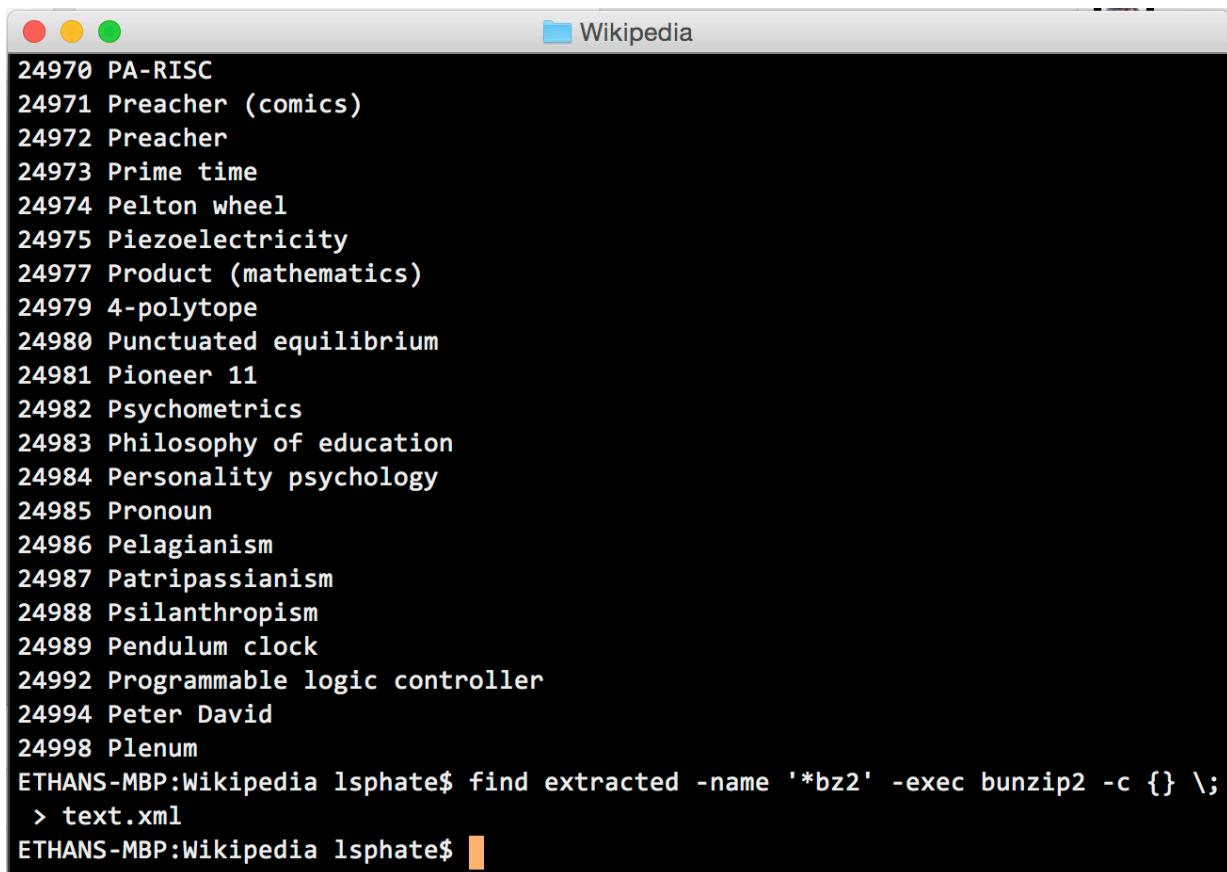
### Step 1. Download & extract Wikipedia articles

I downloaded the latest articles from <http://dumps.wikimedia.org/enwiki/latest/>. I choose the smallest file since we only need 1,000 pages. Then I use [Wikipedia Extractor](#) to extract the text of each articles.



```
ETHANS-MBP:Wikipedia lspbate$ bzcat enwiki-latest-pages-articles2.xml-p000010002
p000024999.bz2 | ./WikiExtractor.py -cb 250K -o extracted
10002 Emil Kraepelin
10003 Evoluon
10004 Educational essentialism
10005 Progressive education
10006 Electronic musical instrument
10008 Electrode
10011 Epistolary novel
10013 Evidence-based medicine
10016 End zone
10017 Ettore Ximenes
10018 Edsger W. Dijkstra
10021 Educational perennialism
10024 MDMA
10025 Flag of Europe
10026 Anthem of Europe
10029 Timeline of evolutionary history of life
10030 Edmund Burke
10033 Early music
10035 Elfenland
10037 Euroscepticism
10042 EAN
10043 Estimator
```

Then we need to combine all the files into one text files by keeping the moderate numbers of outputs and use the command below:



```
24970 PA-RISC
24971 Preacher (comics)
24972 Preacher
24973 Prime time
24974 Pelton wheel
24975 Piezoelectricity
24977 Product (mathematics)
24979 4-polytope
24980 Punctuated equilibrium
24981 Pioneer 11
24982 Psychometrics
24983 Philosophy of education
24984 Personality psychology
24985 Pronoun
24986 Pelagianism
24987 Patrilinearism
24988 Psilanthropism
24989 Pendulum clock
24992 Programmable logic controller
24994 Peter David
24998 Plenum
ETHANS-MBP:Wikipedia lsphate$ find extracted -name '*bz2' -exec bunzip2 -c {} \;
> text.xml
ETHANS-MBP:Wikipedia lsphate$
```

The output format of **Wikipedia Extractor** is like below:

```
<doc id="10002" url="http://en.wikipedia.org/wiki?curid=10002" title="Emil Kraepelin">
Emil Kraepelin

Emil Kraepelin (15 February 1856 – 7 October 1926) was a German psychiatrist. H.J. Eysenck's "Encyclopedia of Psychology" identifies him as the founder of modern scientific psychiatry, as well as of psychopharmacology and psychiatric genetics. Kraepelin believed the chief origin of psychiatric disease to be biological and genetic malfunction. His theories dominated psychiatry at the start of the twentieth century and, despite the later psychodynamic influence of Sigmund Freud and his disciples, enjoyed a revival at century's end.

Family and early life.
Kraepelin, the son of a civil servant, was born in 1856 in Neustrelitz, in the Duchy of Mecklenburg-Strelitz in Germany. He was first introduced to biology by his brother Karl, 10 years older and, later, the director of the Zoological Museum of Hamburg.

...
External links.
For biographies of Kraepelin see:
For English translations of Kraepelin's work see:

</doc>
```

Each element contains the ID, title and text of one page, all the external links and other tags will be eliminated.

## Step 2. Do TF-IDF

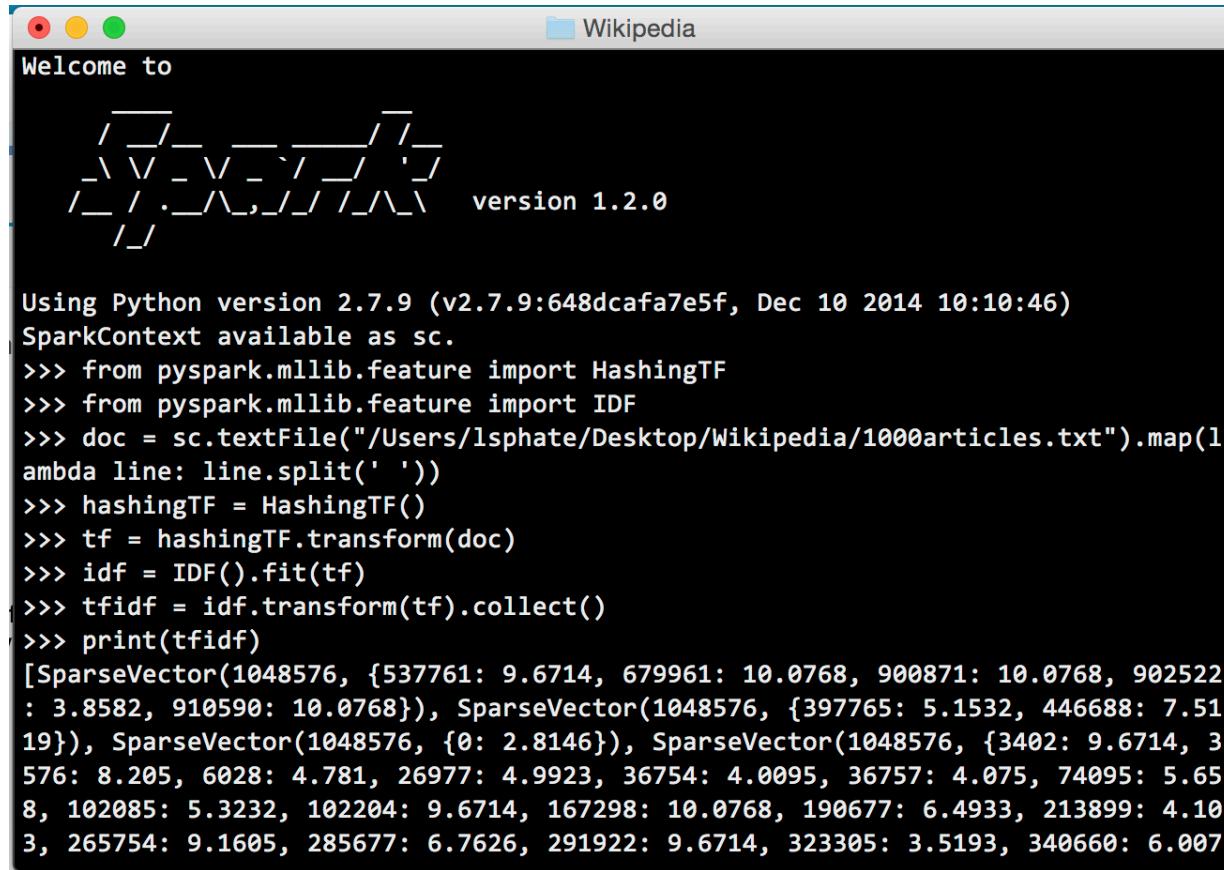
After we got we use **pyspark** to create TD-IDF of these articles.

```
from pyspark.mllib.feature import HashingTF
from pyspark.mllib.feature import IDF

doc = sc.textFile("PATH_OF_INPUT").map(lambda line: line.split(' '))
hashingTF = HashingTF()
tf = hashingTF.transform(doc)
idf = IDF().fit(tf)
tfidf = idf.transform(tf).collect()

print(tfidf)
```

After several minutes of computing, we'll get the result:



The screenshot shows a terminal window titled "Wikipedia". The window title bar has a blue square icon followed by the word "Wikipedia". The main area of the terminal displays the output of a Python script. It starts with a welcome message "Welcome to" followed by a logo consisting of various symbols like slashes and dots. Below the logo, it says "version 1.2.0". Then it prints the Python version and SparkContext availability: "Using Python version 2.7.9 (v2.7.9:648dcafa7e5f, Dec 10 2014 10:10:46)" and "SparkContext available as sc.". The script then imports HashingTF and IDF from pyspark.mllib.feature, reads a file containing 1000 Wikipedia articles, creates a HashingTF object, performs a transform, fits an IDF model, transforms the data, and finally prints the resulting TF-IDF vectors. The printed output is a list of SparseVector objects, each representing a document's features. The first few vectors are shown in full, while the rest are truncated with an ellipsis.

```
Welcome to
   _\ /_ \
  / \ \_ \_ `/_ /'_ \
 /_ / ._/\_,/_/ /_/\_ \
 /_/
Using Python version 2.7.9 (v2.7.9:648dcafa7e5f, Dec 10 2014 10:10:46)
SparkContext available as sc.
>>> from pyspark.mllib.feature import HashingTF
>>> from pyspark.mllib.feature import IDF
>>> doc = sc.textFile("/Users/lspbate/Desktop/Wikipedia/1000articles.txt").map(lambda line: line.split(' '))
>>> hashingTF = HashingTF()
>>> tf = hashingTF.transform(doc)
>>> idf = IDF().fit(tf)
>>> tfidf = idf.transform(tf).collect()
>>> print(tfidf)
[SparseVector(1048576, {537761: 9.6714, 679961: 10.0768, 900871: 10.0768, 902522: 3.8582, 910590: 10.0768}), SparseVector(1048576, {397765: 5.1532, 446688: 7.5119}), SparseVector(1048576, {0: 2.8146}), SparseVector(1048576, {3402: 9.6714, 3576: 8.205, 6028: 4.781, 26977: 4.9923, 36754: 4.0095, 36757: 4.075, 74095: 5.658, 102085: 5.3232, 102204: 9.6714, 167298: 10.0768, 190677: 6.4933, 213899: 4.103, 265754: 9.1605, 285677: 6.7626, 291922: 9.6714, 323305: 3.5193, 340660: 6.007}
```

## Prob. 3 Tweets Regarding 5 Companies

### Step 1. Pull the real-time tweets data

As the request of this problem, we need to choose 5 companies to collect the real-time tweets mentions them. In order to be

allowed to collect the data from Twitter, we need to create a twitter application and obtain the access token to use in the **Twitter Streaming API**.

 Application Management



## Twitter Apps

[Create New App](#)



CU2015S\_ABD\_HW1

Homework Assignment for ADB

And also we need the third-party library called [Tweepy](#) to pull down the tweets. I followed the instruction form [here](#), and wrote the Python application to get the tweets.

```
from tweepy.streaming import StreamListener
from tweepy import OAuthHandler
from tweepy import Stream

access_token = "YOUR_ACCESS_TOKEN"
access_token_secret = "YOUR_ACCESS_SECRET"
consumer_key = "YOUR_KEY"
consumer_secret = "YOUR_SECRET"

class StdOutListener(StreamListener):
    def on_data(self, data):
        print data
        return True
    def on_error(self, status):
        print status

if __name__ == '__main__':
    l = StdOutListener()
    auth = OAuthHandler(consumer_key, consumer_secret)
    auth.set_access_token(access_token, access_token_secret)
    stream = Stream(auth, l)
    stream.filter(track=['KEYWORD'])
```

Where I replace the **KEYWORD** by the names of the 5 companies' name I choose: **BMW, Burberry, Heineken, Nikon** and **Tumblr**, then start these applications to grab the data.

```

dyn-209-2-220-191:Adv Big Data Analytics lsphate$ python twitterlistener_BMW.py > BMW_data.txt
dyn-209-2-220-191:Adv Big Data Analytics lsphate$ python twitterlistener_Nikon.py > Nikon_data.txt
dyn-209-2-220-191:Adv Big Data Analytics lsphate$ python twitterlistener_Tumblr.py > Tumblr_data.txt
dyn-209-2-220-191:Adv Big Data Analytics lsphate$ python twitterlistener_Heineken.py > Heineken_data.txt
dyn-209-2-220-191:Adv Big Data Analytics lsphate$ python twitterlistener_Uniqlo.py > Uniqlo_data.txt

```

The tweets pulled down are in these format:

```

{
  "created_at": "Thu Feb 12 21:23:02 +0000
  2015", "id": 565984447057895424, "id_str": "565984447057895424", "text": "RT @PdeTannhauser: Hoy hace 8
  a\u00f1os, en St Moritz, @NickHeidfeld se lo pasaba pipa con su BMW Sauber F1.07 CC @VirutasF1
  http://\u2026t.co/CDLQu1e\u2026", ...
  "lang": "es", "timestamp_ms": "1423776182318"
}

```

## Step 2. Do TF-IDF

Similar to the former problem, I use **pyspark** to create TF-IDF for each file.

```

from pyspark.mllib.feature import HashingTF
from pyspark.mllib.feature import IDF

doc = sc.textFile("INPUT TWITTER DATA").map(lambda line: line.split(' '))
hashingTF = HashingTF()
hashingTF.indexOf("COMPANY NAME")
tf = hashingTF.transform(doc)
idf = IDF().fit(tf)
tfidf = idf.transform(tf).collect()

print(tfidf)

```

And following are the results of each company:

```
[Q: 928201] ETHANS-MBP:5Twitter lophile$ pyspark
Python 2.7.9 (v2.7.9:648dcfa7e5f, Dec 10 2014, 10:18:46)
[GCC 4.2.1 (Apple Inc. build 5666) (dot 3)] on darwin
Type "help", "copyright", "credits" or "license" for more information.
Spark assembly has been built with Hive, including DataNucleus jars on classpath
15/02/14 14:19:43 WARN NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable
Welcome to
      / \   / \ / \ / \ / \
     /_ \ /_ \/_ \ /_ \ \_ \
    /_ \ /_ \ /_ \ /_ \ \_ \
   /_ \ /_ \ /_ \ /_ \ \_ \
  /_ \ /_ \ /_ \ /_ \ \_ \
 /_ \ /_ \ /_ \ /_ \ \_ \
version 1.2.0

Using Python version 2.7.9 (v2.7.9:648dcfa7e5f, Dec 10 2014 10:10:46)
SparkContext available as sc
-> from pyspark.mllib.feature import HashingTF
-> from pyspark.mllib.feature import IDF
-> doc = sc.textFile("/Users/lophile/Desktop/5Twitter/Heineken_data.txt").map(lambda line: line.split(' '))
-> hashingTF = HashingTF()
-> hashingTF.indexOf("Heineken")
928201
-> tf = hashingTF.transform(doc)
->> idf = IDF().fit(tf)
->> print(idf)
>> [SparseVector(1048576, [3932: 1.8326, 26942: 2.5257, 40140: 2.5257, 56805: 2.5257, 67332: 0.6539, 70933: 2.5257, 85911: 2.1203, 98147: 2.5257, 100335: 2.5257, 168353: 2.5257, 189632: 2.5257, 194731: 2.5257, 2.5257, 252462: 2.5257, 273694: 2.5257, 258589: 2.1203, 293156: 2.5257, 315202: 2.5257, 323305: 2.5257, 351085: 2.5257, 353470: 1.3079, 390422: 2.5257, 394544: 2.5257, 405286: 0.6539, 452545: 2.5257, 513365: 2.5257, 525510: 2.5257, 561081: 2.5257, 673757: 2.5257, 645467: 2.5257, 715640: 5.4977, 719143: 2.1203, 725841: 2.1203, 735901: 2.5257, 745081: 2.5257, 778979: 2.5257, 809211: 2.5257, 866705: 2.5257, 870594: 5.1079, 905151: 2.5257, 955895: 2.5257, 977746: 2.5257, 979681: 2.5257, 980539: 2.5257, 981801: 2.5257, 982011: 2.5257, 982339: 2.5257, 982401: 2.5257, 982479: 2.5257, 982501: 2.5257, 982539: 2.5257, 982561: 2.5257, 982579: 2.5257, 982581: 2.5257, 982591: 2.5257, 982601: 2.5257, 982611: 2.5257, 982621: 2.5257, 982631: 2.5257, 982641: 2.5257, 982651: 2.5257, 982661: 2.5257, 982671: 2.5257, 982681: 2.5257, 982691: 2.5257, 982701: 2.5257, 982711: 2.5257, 982721: 2.5257, 982731: 2.5257, 982741: 2.5257, 982751: 2.5257, 982761: 2.5257, 982771: 2.5257, 982781: 2.5257, 982791: 2.5257, 982801: 2.5257, 982811: 2.5257, 982821: 2.5257, 982831: 2.5257, 982841: 2.5257, 982851: 2.5257, 982861: 2.5257, 982871: 2.5257, 982881: 2.5257, 982891: 2.5257, 982901: 2.5257, 982911: 2.5257, 982921: 2.5257, 982931: 2.5257, 982941: 2.5257, 982951: 2.5257, 982961: 2.5257, 982971: 2.5257, 982981: 2.5257, 982991: 2.5257, 983001: 2.5257, 983011: 2.5257, 983021: 2.5257, 983031: 2.5257, 983041: 2.5257, 983051: 2.5257, 983061: 2.5257, 983071: 2.5257, 983081: 2.5257, 983091: 2.5257, 983101: 2.5257, 983111: 2.5257, 983121: 2.5257, 983131: 2.5257, 983141: 2.5257, 983151: 2.5257, 983161: 2.5257, 983171: 2.5257, 983181: 2.5257, 983191: 2.5257, 983201: 2.5257, 983211: 2.5257, 983221: 2.5257, 983231: 2.5257, 983241: 2.5257, 983251: 2.5257, 983261: 2.5257, 983271: 2.5257, 983281: 2.5257, 983291: 2.5257, 983301: 2.5257, 983311: 2.5257, 983321: 2.5257, 983331: 2.5257, 983341: 2.5257, 983351: 2.5257, 983361: 2.5257, 983371: 2.5257, 983381: 2.5257, 983391: 2.5257, 983401: 2.5257, 983411: 2.5257, 983421: 2.5257, 983431: 2.5257, 983441: 2.5257, 983451: 2.5257, 983461: 2.5257, 983471: 2.5257, 983481: 2.5257, 983491: 2.5257, 983501: 2.5257, 983511: 2.5257, 983521: 2.5257, 983531: 2.5257, 983541: 2.5257, 983551: 2.5257, 983561: 2.5257, 983571: 2.5257, 983581: 2.5257, 983591: 2.5257, 983601: 2.5257, 983611: 2.5257, 983621: 2.5257, 983631: 2.5257, 983641: 2.5257, 983651: 2.5257, 983661: 2.5257, 983671: 2.5257, 983681: 2.5257, 983691: 2.5257, 983701: 2.5257, 983711: 2.5257, 983721: 2.5257, 983731: 2.5257, 983741: 2.5257, 983751: 2.5257, 983761: 2.5257, 983771: 2.5257, 983781: 2.5257, 983791: 2.5257, 983801: 2.5257, 983811: 2.5257, 983821: 2.5257, 983831: 2.5257, 983841: 2.5257, 983851: 2.5257, 983861: 2.5257, 983871: 2.5257, 983881: 2.5257, 983891: 2.5257, 983901: 2.5257, 983911: 2.5257, 983921: 2.5257, 983931: 2.5257, 983941: 2.5257, 983951: 2.5257, 983961: 2.5257, 983971: 2.5257, 983981: 2.5257, 983991: 2.5257, 984001: 2.5257, 984011: 2.5257, 984021: 2.5257, 984031: 2.5257, 984041: 2.5257, 984051: 2.5257, 984061: 2.5257, 984071: 2.5257, 984081: 2.5257, 984091: 2.5257, 984101: 2.5257, 984111: 2.5257, 984121: 2.5257, 984131: 2.5257, 984141: 2.5257, 984151: 2.5257, 984161: 2.5257, 984171: 2.5257, 984181: 2.5257, 984191: 2.5257, 984201: 2.5257, 984211: 2.5257, 984221: 2.5257, 984231: 2.5257, 984241: 2.5257, 984251: 2.5257, 984261: 2.5257, 984271: 2.5257, 984281: 2.5257, 984291: 2.5257, 984301: 2.5257, 984311: 2.5257, 984321: 2.5257, 984331: 2.5257, 984341: 2.5257, 984351: 2.5257, 984361: 2.5257, 984371: 2.5257, 984381: 2.5257, 984391: 2.5257, 984401: 2.5257, 984411: 2.5257, 984421: 2.5257, 984431: 2.5257, 984441: 2.5257, 984451: 2.5257, 984461: 2.5257, 984471: 2.5257, 984481: 2.5257, 984491: 2.5257, 984501: 2.5257, 984511: 2.5257, 984521: 2.5257, 984531: 2.5257, 984541: 2.5257, 984551: 2.5257, 984561: 2.5257, 984571: 2.5257, 984581: 2.5257, 984591: 2.5257, 984601: 2.5257, 984611: 2.5257, 984621: 2.5257, 984631: 2.5257, 984641: 2.5257, 984651: 2.5257, 984661: 2.5257, 984671: 2.5257, 984681: 2.5257, 984691: 2.5257, 984701: 2.5257, 984711: 2.5257, 984721: 2.5257, 984731: 2.5257, 984741: 2.5257, 984751: 2.5257, 984761: 2.5257, 984771: 2.5257, 984781: 2.5257, 984791: 2.5257, 984801: 2.5257, 984811: 2.5257, 984821: 2.5257, 984831: 2.5257, 984841: 2.5257, 984851: 2.5257, 984861: 2.5257, 984871: 2.5257, 984881: 2.5257, 984891: 2.5257, 984901: 2.5257, 984911: 2.5257, 984921: 2.5257, 984931: 2.5257, 984941: 2.5257, 984951: 2.5257, 984961: 2.5257, 984971: 2.5257, 984981: 2.5257, 984991: 2.5257, 985001: 2.5257, 985011: 2.5257, 985021: 2.5257, 985031: 2.5257, 985041: 2.5257, 985051: 2.5257, 985061: 2.5257, 985071: 2.5257, 985081: 2.5257, 985091: 2.5257, 985101: 2.5257, 985111: 2.5257, 985121: 2.5257, 985131: 2.5257, 985141: 2.5257, 985151: 2.5257, 985161: 2.5257, 985171: 2.5257, 985181: 2.5257, 985191: 2.5257, 985201: 2.5257, 985211: 2.5257, 985221: 2.5257, 985231: 2.5257, 985241: 2.5257, 985251: 2.5257, 985261: 2.5257, 985271: 2.5257, 985281: 2.5257, 985291: 2.5257, 985301: 2.5257, 985311: 2.5257, 985321: 2.5257, 985331: 2.5257, 985341: 2.5257, 985351: 2.5257, 985361: 2.5257, 985371: 2.5257, 985381: 2.5257, 985391: 2.5257, 985401: 2.5257, 985411: 2.5257, 985421: 2.5257, 985431: 2.5257, 985441: 2.5257, 985451: 2.5257, 985461: 2.5257, 985471: 2.5257, 985481: 2.5257, 985491: 2.5257, 985501: 2.5257, 985511: 2.5257, 985521: 2.5257, 985531: 2.5257, 985541: 2.5257, 985551: 2.5257, 985561: 2.5257, 985571: 2.5257, 985581: 2.5257, 985591: 2.5257, 985601: 2.5257, 985611: 2.5257, 985621: 2.5257, 985631: 2.5257, 985641: 2.5257, 985651: 2.5257, 985661: 2.5257, 985671: 2.5257, 985681: 2.5257, 985691: 2.5257, 985701: 2.5257, 985711: 2.5257, 985721: 2.5257, 985731: 2.5257, 985741: 2.5257, 985751: 2.5257, 985761: 2.5257, 985771: 2.5257, 985781: 2.5257, 985791: 2.5257, 985801: 2.5257, 985811: 2.5257, 985821: 2.5257, 985831: 2.5257, 985841: 2.5257, 985851: 2.5257, 985861: 2.5257, 985871: 2.5257, 985881: 2.5257, 985891: 2.5257, 985901: 2.5257, 985911: 2.5257, 985921: 2.5257, 985931: 2.5257, 985941: 2.5257, 985951: 2.5257, 985961: 2.5257, 985971: 2.5257, 985981: 2.5257, 985991: 2.5257, 986001: 2.5257, 986011: 2.5257, 986021: 2.5257, 986031: 2.5257, 986041: 2.5257, 986051: 2.5257, 986061: 2.5257, 986071: 2.5257, 986081: 2.5257, 986091: 2.5257, 986101: 2.5257, 986111: 2.5257, 986121: 2.5257, 986131: 2.5257, 986141: 2.5257, 986151: 2.5257, 986161: 2.5257, 986171: 2.5257, 986181: 2.5257, 986191: 2.5257, 986201: 2.5257, 986211: 2.5257, 986221: 2.5257, 986231: 2.5257, 986241: 2.5257, 986251: 2.5257, 986261: 2.5257, 986271: 2.5257, 986281: 2.5257, 986291: 2.5257, 986301: 2.5257, 986311: 2.5257, 986321: 2.5257, 986331: 2.5257, 986341: 2.5257, 986351: 2.5257, 986361: 2.5257, 986371: 2.5257, 986381: 2.5257, 986391: 2.5257, 986401: 2.5257, 986411: 2.5257, 986421: 2.5257, 986431: 2.5257, 986441: 2.5257, 986451: 2.5257, 986461: 2.5257, 986471: 2.5257, 986481: 2.5257, 986491: 2.5257, 986501: 2.5257, 986511: 2.5257, 986521: 2.5257, 986531: 2.5257, 986541: 2.5257, 986551: 2.5257, 986561: 2.5257, 986571: 2.5257, 986581: 2.5257, 986591: 2.5257, 986601: 2.5257, 986611: 2.5257, 986621: 2.5257, 986631: 2.5257, 986641: 2.5257, 986651: 2.5257, 986661: 2.5257, 986671: 2.5257, 986681: 2.5257, 986691: 2.5257, 986701: 2.5257, 986711: 2.5257, 986721: 2.5257, 986731: 2.5257, 986741: 2.5257, 986751: 2.5257, 986761: 2.5257, 986771: 2.5257, 986781: 2.5257, 986791: 2.5257, 986801: 2.5257, 986811: 2.5257, 986821: 2.5257, 986831: 2.5257, 986841: 2.5257, 986851: 2.5257, 986861: 2.5257, 986871: 2.5257, 986881: 2.5257, 986891: 2.5257, 986901: 2.5257, 986911: 2.5257, 986921: 2.5257, 986931: 2.5257, 986941: 2.5257, 986951: 2.5257, 986961: 2.5257, 986971: 2.5257, 986981: 2.5257, 986991: 2.5257, 987001: 2.5257, 987011: 2.5257, 987021: 2.5257, 987031: 2.5257, 987041: 2.5257, 987051: 2.5257, 987061: 2.5257, 987071: 2.5257, 987081: 2.5257, 987091: 2.5257, 987101: 2.5257, 987111: 2.5257, 987121: 2.5257, 987131: 2.5257, 987141: 2.5257, 987151: 2.5257, 987161: 2.5257, 987171: 2.5257, 987181: 2.5257, 987191: 2.5257, 987201: 2.5257, 987211: 2.5257, 987221: 2.5257, 987231: 2.5257, 987241: 2.5257, 987251: 2.5257, 987261: 2.5257, 987271: 2.5257, 987281: 2.5257, 987291: 2.5257, 987301: 2.5257, 987311: 2.5257, 987321: 2.5257, 987331: 2.5257, 987341: 2.5257, 987351: 2.5257, 987361: 2.5257, 987371: 2.5257, 987381: 2.5257, 987391: 2.5257, 987401: 2.5257, 987411: 2.5257, 987421: 2.5257, 987431: 2.5257, 987441: 2.5257, 987451: 2.5257, 987461: 2.5257, 987471: 2.5257, 987481: 2.5257, 987491: 2.5257, 987501: 2.5257, 987511: 2.5257, 987521: 2.5257, 987531: 2.5257, 987541: 2.5257, 987551: 2.5257, 987561: 2.5257, 987571: 2.5257, 987581: 2.5257, 987591: 2.5257, 987601: 2.5257, 987611: 2.5257, 987621: 2.5257, 987631: 2.5257, 987641: 2.5257, 987651: 2.5257, 987661: 2.5257, 987671: 2.5257, 987681: 2.5257, 987691: 2.5257, 987701: 2.5257, 987711: 2.5257, 987721: 2.5257, 987731: 2.5257, 987741: 2.5257, 987751: 2.5257, 987761: 2.5257, 987771: 2.5257, 987781: 2.5257, 987791: 2.5257, 987801: 2.5257, 987811: 2.5257, 987821: 2.5257, 987831: 2.5257, 987841: 2.5257, 987851: 2.5257, 987861: 2.5257, 987871: 2.5257, 987881: 2.5257, 987891: 2.5257, 987901: 2.5257, 987911: 2.5257, 987921: 2.5257, 987931: 2.5257, 987941: 2.5257, 987951: 2.5257, 987961: 2.5257, 987971: 2.5257, 987981: 2.5257, 987991: 2.5257, 988001: 2.5257, 988011: 2.5257, 988021: 2.5257, 988031: 2.5257, 988041: 2.5257, 988051: 2.5257, 988061: 2.5257, 988071: 2.5257, 988081: 2.5257, 988091: 2.5257, 988101: 2.5257, 988111: 2.5257, 988121: 2.5257, 988131: 2.5257, 988141: 2.5257, 988151: 2.5257, 988161: 2.5257, 988171: 2.5257, 988181: 2.5257, 988191: 2.5257, 988201: 2.5257, 988211: 2.5257, 988221: 2.5257, 988231: 2.5257, 988241: 2.5257, 988251: 2.5257, 988261: 2.5257, 988271: 2.5257, 988281: 2.5257, 988291: 2.5257, 988301: 2.5257, 988311: 2.5257, 988321: 2.5257, 988331: 2.5257, 988341: 2.5257, 988351: 2.5257, 988361: 2.5257, 988371: 2.5257, 988381: 2.5257, 988391: 2.5257, 988401: 2.5257, 988411: 2.5257, 988421: 2.5257, 988431: 2.5257, 988441: 2.5257, 988451: 2.5257, 988461: 2.5257, 988471: 2.5257, 988481: 2.5257, 988491: 2.5257, 988501: 2.5257, 988511: 2.5257, 988521: 2.5257, 988531: 2.5257, 988541: 2.5257, 988551: 2.5257, 988561: 2.5257, 988571: 2.5257, 988581: 2.5257, 988591: 2.5257, 988601: 2.5257, 988611: 2.5257, 988621: 2.5257, 988631: 2.5257, 988641: 2.5257, 988651: 2.5257, 988661: 2.5257, 988671: 2.5257, 988681: 2.5257, 988691: 2.5257, 988701: 2.5257, 988711: 2.5257, 988721: 2.5257, 988731: 2.5257, 988741: 2.5257, 988751: 2.5257, 988761: 2.5257, 988771: 2.5257, 988781: 2.5257, 988791: 2.5257, 988801: 2.5257, 988811: 2.5257, 988821: 2.5257, 988831: 2.5257, 988841: 2.5257, 988851: 2.5257, 988861: 2.5257, 988871: 2
```

## Prob. 4 The Stock Prices of 5 Companies

## **Step 1. Pull the real-time stock prices**

We need the [yahoo-finance](#) library to get the shared data of stocks from Yahoo Finance. I downloaded, installed the library, and wrote a Python application to get the data.

```
import time
from yahoo_finance import Share

price = Share('COMPANY CODE')
print price.get_open() + " at open"
print price.get_price() + " at " + price.get_trade_datetime()
for x in range(0, 28):
    time.sleep(60)
    price = Share('COMPANY CODE')
    print price.get_price() + " at " + price.get_trade_datetime()

print "Data pulling completed."
```

My choices of 5 companies are **Amazon**, **HP**, **IBM**, **Microsoft** and **Yahoo**. The datas I pulled down and outputted were in this format:

```

{'Sector': 'Services', 'end': '2015-02-17', 'CompanyName': None, 'symbol': 'AMZN', 'start': '1997-05-16',
'FullTimeEmployees': '154100', 'Industry': 'Catalog & Mail Order Houses'}
378.10 at open
375.2698 at 2015-02-17 18:01:00 UTC+0000
375.12 at 2015-02-17 18:01:00 UTC+0000
375.01 at 2015-02-17 18:03:00 UTC+0000
375.01 at 2015-02-17 18:03:00 UTC+0000
375.29 at 2015-02-17 18:05:00 UTC+0000
375.25 at 2015-02-17 18:06:00 UTC+0000
375.355 at 2015-02-17 18:07:00 UTC+0000
375.19 at 2015-02-17 18:08:00 UTC+0000
375.11 at 2015-02-17 18:09:00 UTC+0000
375.05 at 2015-02-17 18:10:00 UTC+0000
374.90 at 2015-02-17 18:11:00 UTC+0000
374.792 at 2015-02-17 18:12:00 UTC+0000
374.58 at 2015-02-17 18:13:00 UTC+0000
374.65 at 2015-02-17 18:14:00 UTC+0000
375.33 at 2015-02-17 18:14:00 UTC+0000
375.29 at 2015-02-17 18:16:00 UTC+0000
375.14 at 2015-02-17 18:17:00 UTC+0000
375.245 at 2015-02-17 18:18:00 UTC+0000
375.54 at 2015-02-17 18:19:00 UTC+0000
375.18 at 2015-02-17 18:20:00 UTC+0000
375.18 at 2015-02-17 18:20:00 UTC+0000
375.475 at 2015-02-17 18:22:00 UTC+0000
375.68 at 2015-02-17 18:23:00 UTC+0000
375.524 at 2015-02-17 18:24:00 UTC+0000
375.36 at 2015-02-17 18:25:00 UTC+0000
375.14 at 2015-02-17 18:26:00 UTC+0000
375.248 at 2015-02-17 18:27:00 UTC+0000
375.3852 at 2015-02-17 18:28:00 UTC+0000
375.16 at 2015-02-17 18:28:00 UTC+0000
Data pulling completed.

```

## Step 2. Use filter to find outliers

I wrote an application to extract the data points that are 2 standard deviations away against the mean of the prices. The data that imported has been pre-processed to left the price numbers only:

```

import math

def stof(x):
    try:
        return float(x[0])
    except (ValueError, TypeError):
        return 0.0

price = sc.textFile("PRICE_DATA_NAME").map(lambda line: line.split(' '))
pricen = price.map(stof)
stats = pricen.stats()
stddev = stats.stdev()
mean = stats.mean()
outliers = pricen.filter(lambda x: math.fabs(x - mean) > 2 * stddev)
print outliers.collect()

```

The image shows five separate terminal windows side-by-side, each displaying a command-line session. The sessions involve navigating to specific directories and running Python scripts to extract data from files named after companies (Amazon, IBM, Microsoft, HP, and Yahoo). The output of these commands is not visible in the image.

```

Terminal 1 (Top Left):
Last login: Sat Feb 14 00:07:31 on ttys001
ETHANS-MBP:~ lsphate$ cd /Volumes/JetDrive/Columbia/2015\ Spring/Adv\ Big\ Data\ Analytics/HW1/YFinance
ETHANS-MBP:YFinance lsphate$ python price_Amazon.py > Amazon_price.txt

Terminal 2 (Top Right):
Last login: Mon Feb 16 11:15:08 on ttys000
ETHANS-MBP:~ lsphate$ cd /Volumes/JetDrive/Columbia/2015\ Spring/Adv\ Big\ Data\ Analytics/HW1/YFinance
ETHANS-MBP:YFinance lsphate$ python price_IBM.py > IBM_price.txt

Terminal 3 (Second Row, Left):
Last login: Mon Feb 16 11:15:15 on ttys001
ETHANS-MBP:~ lsphate$ cd /Volumes/JetDrive/Columbia/2015\ Spring/Adv\ Big\ Data\ Analytics/HW1/YFinance
ETHANS-MBP:YFinance lsphate$ python price_Microsoft.py > Microsoft_price.txt

Terminal 4 (Second Row, Right):
Last login: Mon Feb 16 11:15:18 on ttys003
ETHANS-MBP:~ lsphate$ cd /Volumes/JetDrive/Columbia/2015\ Spring/Adv\ Big\ Data\ Analytics/HW1/YFinance
ETHANS-MBP:YFinance lsphate$ python price_HP.py > HP_price.txt

Terminal 5 (Bottom):
Last login: Mon Feb 16 11:15:17 on ttys002
ETHANS-MBP:~ lsphate$ cd /Volumes/JetDrive/Columbia/2015\ Spring/Adv\ Big\ Data\ Analytics/HW1/YFinance
ETHANS-MBP:YFinance lsphate$ python price_Yahoo.py > Yahoo_price.txt

```

And we can get the results of the outliers:

The image shows an IPython notebook cell with the title "lsphate". The cell contains Python code for reading a file, mapping it to a float list, calculating statistics, and filtering outliers. The output shows the list of outliers: [374.58, 374.65, 375.68].

```

In [1]: 
    _\ \ \ - \ \ - ` / _/ ' / 
    /_ / . / \_, _/ / / / \ \
    / /           version 1.2.0

Using Python version 2.7.9 (v2.7.9:648dcafa7e5f, Dec 10 2014 10:10:46)
SparkContext available as sc.
>>> import math
>>>
>>> def stof(x):
...     try:
...         return float(x[0])
...     except (ValueError, TypeError):
...         return 0.0
...
>>> price = sc.textFile("/Volumes/JetDrive/Columbia/2015\ Spring/Adv\ Big\ Data\ Analytics/HW1/YFinance/Amazon_price.txt").map(lambda line: line.split(' '))
>>> pricen = price.map(stof)
>>> stats = pricen.stats()
>>> stddev = stats.stdev()
>>> mean = stats.mean()
>>> outliers = pricen.filter(lambda x: math.fabs(x - mean) > 2 * stddev)
>>> print outliers.collect()
[374.58, 374.65, 375.68]
>>>

```

```
lsphate
...
    try:
...
        return float(x[0])
...
    except (ValueError, TypeError):
...
        return 0.0
...
>>> price = sc.textFile("/Volumes/JetDrive/Columbia/2015\ Spring/Adv\ Big\ Data\
Analytics/HW1/YFinance/Amazon_price.txt").map(lambda line: line.split(' '))
>>> pricen = price.map(stof)
>>> stats = pricen.stats()
>>> stddev = stats.stdev()
>>> mean = stats.mean()
>>> outliers = pricen.filter(lambda x: math.fabs(x - mean) > 2 * stddev)
>>> print outliers.collect()
[374.58, 374.65, 375.68]
>>> price = sc.textFile("/Volumes/JetDrive/Columbia/2015\ Spring/Adv\ Big\ Data\
Analytics/HW1/YFinance/HP_price.txt").map(lambda line: line.split(' '))
>>> pricen = price.map(stof)
>>> stats = pricen.stats()
>>> stddev = stats.stdev()
>>> mean = stats.mean()
>>> outliers = pricen.filter(lambda x: math.fabs(x - mean) > 2 * stddev)
>>> print outliers.collect()
[38.475]
>>> 
```

```
lsphate
>>> price = sc.textFile("/Volumes/JetDrive/Columbia/2015\ Spring/Adv\ Big\ Data\
Analytics/HW1/YFinance/IBM_price.txt").map(lambda line: line.split(' '))
>>> pricen = price.map(stof)
>>> stats = pricen.stats()
>>> stddev = stats.stdev()
>>> mean = stats.mean()
>>> outliers = pricen.filter(lambda x: math.fabs(x - mean) > 2 * stddev)
>>> print outliers.collect()
[]
>>> 
```

```
lsphate
>>> price = sc.textFile("/Volumes/JetDrive/Columbia/2015\ Spring/Adv\ Big\ Data\
    Analytics/HW1/YFinance/Microsoft_price.txt").map(lambda line: line.split(' '))
>>> pricen = price.map(stof)
>>> stats = pricen.stats()
>>> stddev = stats.stdev()
>>> mean = stats.mean()
>>> outliers = pricen.filter(lambda x: math.fabs(x - mean) > 2 * stddev)
>>> print outliers.collect()
[43.2801]
>>> 
```

```
lsphate
>>> stats = pricen.stats()
>>> stddev = stats.stdev()
>>> mean = stats.mean()
>>> outliers = pricen.filter(lambda x: math.fabs(x - mean) > 2 * stddev)
>>> print outliers.collect()
[43.2801]
>>> import math
>>>
>>> def stof(x):
...     try:
...         return float(x[0])
...     except (ValueError, TypeError):
...         return 0.0
...
>>> price = sc.textFile("/Volumes/JetDrive/Columbia/2015\ Spring/Adv\ Big\ Data\
    Analytics/HW1/YFinance/Yahoo_price.txt").map(lambda line: line.split(' '))
>>> pricen = price.map(stof)
>>> stats = pricen.stats()
>>> stddev = stats.stdev()
>>> mean = stats.mean()
>>> outliers = pricen.filter(lambda x: math.fabs(x - mean) > 2 * stddev)
>>> print outliers.collect()
[43.77]
>>> 
```