In this tutorial, we will download and analyze the 1996 to 2000 stock market data. The files span over 100,000+ directories and include stocks from around the world. The data does not include the location in the excel file so we will add it after the data is obtained.

1. Open a Git Bash terminal and SSH

```
takum@DESKTOP-CNCNCF8 MINGW64 ~/Desktop
$ ssh mbarrio5@129.150.64.74
-- WARNING -- This system is for the use of authorized users only. Individuals
using this computer system without authority or in excess of their authority
are subject to having all their activities on this system monitored and
recorded by system personnel. Anyone using this system expressly consents to
such monitoring and is advised that if such monitoring reveals possible
evidence of criminal activity system personnel may provide the evidence of such
monitoring to law enforcement officials.

mbarrio5@129.150.64.74's password:
Last login: Mon Nov 16 19:31:05 2020 from 162-204-124-52.lightspeed.irvnca.sbcglobal.net
-bash-4.1$ df -hl
Filesystem            Size  Used Avail Use% Mounted on
/dev/mapper/vg_main-lv_root
                       21G   16G  3.9G  80% /
tmpfs                  30G    0   30G   0% /dev/shm
/dev/xvdb1            477M   72M  376M  17% /boot
/dev/xvdd1           4.8G  130M  4.5G   3% /u01/app/oracle/tools
/dev/mapper/vg_bin-lv_bin
                       50G  7.8G   39G  17% /u01/bdcsce
/dev/mapper/vg_adata-lv_adata
                       34G   48M   32G   1% /adata
/dev/mapper/vg_data-lv_data
                      202G  126G   66G  66% /data
-bash-4.1$ cd ..
-bash-4.1$ cd ..
-bash-4.1$ cd ..
-bash-4.1$ cd dev
-bash-4.1$ cd shm
-bash-4.1$ pwd
/dev/shm
-bash-4.1$ |
```

$ ssh [yourusernamehere]@129.150.64.74

$ cd ..

$ cd ..

$ cd ..

$ cd dev

$ cd shm

$ pwd

## 2. Download the dataset to local Linux system

```
-bash-4.1$ wget --load-cookies /tmp/cookies.txt "https://docs.google.com/uc?export=download&confirm=$(wget --quiet --save-cookies /tmp/co
okies.txt --keep-session-cookies --no-check-certificate 'https://docs.google.com/uc?export=download&id=1fr-FVU0ZddRNuYvmPo79Mpkj403BguOt'
 -O- | sed -rn 's/.*confirm=([0-9A-Za-z_]+).*/\1\n/p')&id=1fr-FVU0ZddRNuYvmPo79Mpkj403BguOt" -O stock_data && rm -rf /tmp/cookies.txt
Connecting to doc-14-1s-docs.googleusercontent.com|172.217.8.193|:443... connected.
HTTP request sent, awaiting response... 200 OK
Length: unspecified [application/x-zip-compressed]
Saving to: "stock_data"

    [                                                                      <=>

2020-11-26 02:08:44 (132 MB/s) - "stock_data" saved [2469842861]
```

$ wget --load-cookies /tmp/cookies.txt "https://docs.google.com/uc?export=download&confirm=$(wget --quiet --save-cookies /tmp/cookies.txt --keep-session-cookies --no-check-certificate 'https://docs.google.com/uc?export=download&id=1fr-FVU0ZddRNuYvmPo79Mpkj403BguOt' -O- | sed -rn 's/.*confirm=([0-9A-Za-z_]+).*/\1\n/p')&id=1fr-FVU0ZddRNuYvmPo79Mpkj403BguOt" -O stock_data && rm -rf /tmp/cookies.txt

**Note: "stock_data" saved will appear if the download is successful**

## 3. Unzipping dataset  (takes approximately 5 minutes)

```
-bash-4.1$ ls
stock_data
-bash-4.1$ unzip stock_data
```

-Unzip in progress

```
 inflating: Data/Data/053270.KQ/053270.KQ.csv
 inflating: Data/Data/053290.KQ/053290.KQ.csv
 inflating: Data/Data/0533-OL.HK/0533-OL.HK.csv
 inflating: Data/Data/0533.HK/0533.HK.csv
 inflating: Data/Data/053300.KQ/053300.KQ.csv
```

-Unzip complete

```
 inflating: Data/Data/ZZZ.TO/ZZZ.TO.csv
 inflating: Data/Data/dir.txt
 inflating: Tickers.xlsx
-bash-4.1$ ls
Data  stock_data  Tickers.xlsx
-bash-4.1$
```

$ ls

$ unzip stock_data

$ hdfs dfs -mkdir dataproj

4. Move files from local linux to HDFS (Moving "Data" to dataproj/ takes 1.1 hours)

```
-bash-4.1$ hdfs dfs -put Tickers.xlsx dataproj/
-bash-4.1$ hdfs dfs -put Data dataproj/
-bash-4.1$ hdfs dfs -ls dataproj/
Found 2 items
drwxr-xr-x   - mbarrio5 hdfs          0 2020-11-26 02:18 dataproj/Data
-rw-r--r--   2 mbarrio5 hdfs    4730902 2020-11-26 02:17 dataproj/Tickers.xlsx
-bash-4.1$ |
```

$ hdfs dfs -put Tickers.xlsx dataproj/

$ hdfs dfs -put Data dataproj/

$ hdfs dfs -ls dataproj/

**Note: remove files to save storage on system**

```
-bash-4.1$ ls
Data  stock_data  Tickers.xlsx
-bash-4.1$ rm -r Data
-bash-4.1$ rm -r stock_data
-bash-4.1$ rm -r Tickers.xlsx
-bash-4.1$ ls
-bash-4.1$ |
```

$ rm -r Data

$ rm -r stock_data

$ rm -r Tickers.xlsx

5. Change file permissions of dataset

```
-bash-4.1$ hdfs dfs -chmod 777 /user/mbarrio5/dataproj/Data/*/*
-bash-4.1$ |
```

$ hdfs dfs -chmod 777 /user/mbarrio5/dataproj/Data/Data/*/*

6. Using Pig to consolidate all 101,000 directories and files (In the step we create a test schema)

```
-bash-4.1$ pig
WARNING: Use "yarn jar" to launch YARN applications.
20/11/26 03:39:31 INFO pig.ExecTypeProvider: Trying ExecType : LOCAL
20/11/26 03:39:31 INFO pig.ExecTypeProvider: Trying ExecType : MAPREDUCE
20/11/26 03:39:31 INFO pig.ExecTypeProvider: Picked MAPREDUCE as the ExecType
2020-11-26 03:39:31,726 [main] INFO  org.apache.pig.Main - Apache Pig version 0.15.0 (r: unknown) compiled Jun 06 2017, 02:55:08
2020-11-26 03:39:31,727 [main] INFO  org.apache.pig.Main - Logging error messages to: /dev/shm/pig_1606361971691.log
2020-11-26 03:39:31,915 [main] INFO  org.apache.pig.impl.util.Utils - Default bootup file /home/mbarrio5/.pigbootup not found
2020-11-26 03:39:32,400 [main] INFO  org.apache.pig.backend.hadoop.executionengine.HExecutionEngine - Connecting to hadoop file system at:
2020-11-26 03:39:33,837 [main] INFO  org.apache.pig.PigServer - Pig Script ID for the session: PIG-default-14112b91-2961-4add-9759-798ec8a6
2020-11-26 03:39:34,212 [main] INFO  org.apache.hadoop.yarn.client.api.impl.TimelineClientImpl - Timeline service address: http://bigdai-no
8/ws/v1/timeline/
2020-11-26 03:39:34,332 [main] INFO  org.apache.pig.backend.hadoop.ATSService - Created ATS Hook
grunt>
```

$ pig

```
grunt> stock_data = LOAD '/user/mbarrio5/dataproj/Data/Data/0002.HK/0002.HK.csv' USING PigStor
age(',','-tagFile')
>> AS (date:chararray, open:chararray, high:chararray,
>> low:chararray, close:chararray, adjclose:chararray,
>> volume:chararray);
grunt> DESCRIBE stock_data;
stock_data: {date: chararray,open: chararray,high: chararray,low: chararray,close: chararray,a
djclose: chararray,volume: chararray}
grunt> stock_data_subset = limit stock_data 50;
grunt> DUMP stock_data_subset;
```

$ stock_data = LOAD '/user/mbarrio5/dataproj/Data/Data/0002.HK/0002.HK.csv' USING PigStorage(',','-tagFile')
       AS (date:chararray, open:chararray, high:chararray,
       low:chararray, close:chararray, adjclose:chararray,
       volume:chararray);

$ Describe stock_data;

$ stock_data_subset = limit stock_data 50;

$ DUMP stock_data_subset;

**Output of "DUMP stock_data_subset;":**

```
(0002.HK.csv,2000-02-24,33.000000,33.000000,32.599998,33.000000,13.371803)
(0002.HK.csv,2000-02-25,32.900002,33.099998,32.099998,33.000000,13.371803)
(0002.HK.csv,2000-02-28,33.000000,33.299999,32.299999,32.799999,13.290762)
(0002.HK.csv,2000-02-29,32.799999,34.500000,32.799999,34.299999,13.898570)
(0002.HK.csv,2000-03-01,34.599998,34.599998,32.700001,32.700001,13.250243)
(0002.HK.csv,2000-03-02,33.099998,34.000000,32.599998,33.299999,13.493366)
(0002.HK.csv,2000-03-03,33.400002,33.799999,32.799999,33.099998,13.412320)
(0002.HK.csv,2000-03-06,33.000000,33.099998,32.299999,32.299999,13.088157)
(0002.HK.csv,2000-03-07,32.299999,33.599998,32.099998,33.400002,13.533885)
(0002.HK.csv,2000-03-08,32.900002,34.500000,32.900002,34.400002,13.939089)
(0002.HK.csv,2000-03-09,34.400002,35.200001,33.299999,34.799999,14.101176)
(0002.HK.csv,2000-03-10,34.700001,34.700001,33.599998,33.700001,13.655445)
grunt>
```

## 7. Complete schema to map reduce ALL of the data

```
2020-11-26 21:43:14,814 [main] INFO  org.apache.pig.backend.hadoop.ATSService - Created ATS Ho
ok
grunt> stock_data = LOAD '/user/mbarrio5/dataproj/Data/Data/*/*' USING PigStorage(',','-tagFil
e')
>> AS (date:chararray, open:chararray, high:chararray,
>> low:chararray, close:chararray, adjclose:chararray,
>> volume:chararray);
grunt> stock_data_subset = stock_data;
grunt> store stock_data_subset into 'output/cmpldata' using PigStorage(',');
```

Store function in progress:

```
2020-11-26 22:56:49,611 [JobControl] INFO  org.apache.pig.backend.hadoop.executionengine.util.
MapRedUtil - Total input paths to process : 104124
2020-11-26 22:57:50,129 [JobControl] INFO  org.apache.pig.backend.hadoop.executionengine.util.
MapRedUtil - Total input paths (combined) to process : 99
2020-11-26 22:57:50,530 [JobControl] INFO  org.apache.hadoop.mapreduce.JobSubmitter - number o
f splits:99
2020-11-26 22:57:50,667 [JobControl] INFO  org.apache.hadoop.mapreduce.JobSubmitter - Submitti
ng tokens for job: job_1605128590665_0285
2020-11-26 22:57:50,821 [JobControl] INFO  org.apache.hadoop.mapred.YARNRunner - Job jar is no
t present. Not adding any jar to the list of resources.
2020-11-26 22:57:51,130 [JobControl] INFO  org.apache.hadoop.yarn.client.api.impl.YarnClientIm
pl - Submitted application application_1605128590665_0285
2020-11-26 22:57:51,182 [JobControl] INFO  org.apache.hadoop.mapreduce.Job - The url to track
the job: http://bigdai-nov-bdcsce-3.compute-608214094.oraclecloud.internal:8088/proxy/applicat
ion_1605128590665_0285/
2020-11-26 22:57:51,183 [main] INFO  org.apache.pig.backend.hadoop.executionengine.mapReduceLa
yer.MapReduceLauncher - HadoopJobId: job_1605128590665_0285
2020-11-26 22:57:51,183 [main] INFO  org.apache.pig.backend.hadoop.executionengine.mapReduceLa
yer.MapReduceLauncher - Processing aliases stock_data
2020-11-26 22:57:51,183 [main] INFO  org.apache.pig.backend.hadoop.executionengine.mapReduceLa
yer.MapReduceLauncher - detailed locations: M: stock_data[1,13],stock_data[-1,-1] C:  R:
2020-11-26 22:57:51,203 [main] INFO  org.apache.pig.backend.hadoop.executionengine.mapReduceLa
yer.MapReduceLauncher - 0% complete
2020-11-26 22:57:51,203 [main] INFO  org.apache.pig.backend.hadoop.executionengine.mapReduceLa
yer.MapReduceLauncher - Running jobs are [job_1605128590665_0285]
2020-11-26 22:58:28,372 [main] INFO  org.apache.pig.backend.hadoop.executionengine.mapReduceLa
yer.MapReduceLauncher - 4% complete
2020-11-26 22:58:28,372 [main] INFO  org.apache.pig.backend.hadoop.executionengine.mapReduceLa
yer.MapReduceLauncher - Running jobs are [job_1605128590665_0285]
```

$ stock_data = LOAD '/user/mbarrio5/dataproj/Data/Data/*/*' USING PigStorage(',','-tagFile')
    AS (date:chararray, open:chararray, high:chararray,

    low:chararray, close:chararray, adjclose:chararray,

    volume:chararray);

$ stock_data_subset = stock_data;

$ store stock_data_subset into 'output/cmpldata' using PigStorage(',');

```
2020-11-26 23:08:58,151 [main] INFO  org.apache.hadoop.yarn.client.ConfiguredRMFailoverProxyPr
ovider - Failing over to rm2
2020-11-26 23:08:58,154 [main] INFO  org.apache.hadoop.mapred.ClientServiceDelegate - Applicat
ion state is completed. FinalApplicationStatus=SUCCEEDED. Redirecting to job history server
2020-11-26 23:08:58,209 [main] WARN  org.apache.pig.backend.hadoop.executionengine.mapReduceLa
yer.MapReduceLauncher - Encountered Warning ACCESSING_NON_EXISTENT_FIELD 192457 time(s).
2020-11-26 23:08:58,209 [main] INFO  org.apache.pig.backend.hadoop.executionengine.mapReduceLa
yer.MapReduceLauncher - Success!
grunt> |
```

**Note: Make sure to "quit" after completing the "store" line**

$ quit

8. Change permissions and merge all 100 Pig outputs into 1 CSV file

**INSERT SCREENSHOT OF CODE EXECUTION HERE

$ hdfs dfs -chmod 777 output/*

$ hdfs dfs -cat output/cmpldata/* | hadoop fs -put - testdata.csv

$ hdfs dfs -chmod 777 testdata.csv

Continue to next page

## 9. Download CSV to temp file system due to size restrictions

```
-bash-4.1$ cd ..
-bash-4.1$ cd ..
-bash-4.1$ pwd
/
-bash-4.1$ cd dev
-bash-4.1$ cd shm
-bash-4.1$ pwd
/dev/shm
-bash-4.1$ hdfs dfs -get testdata.csv
-bash-4.1$ ls
testdata.csv
-bash-4.1$ hdfs dfs -ls
Found 5 items
drwx------    - mbarrio5 hdfs             0 2020-11-26 23:08 .staging
drwxr-xr-x    - mbarrio5 hdfs             0 2020-11-26 02:18 dataproj
drwxr-xr-x    - mbarrio5 hdfs             0 2020-11-26 22:57 output
-rwxrwxrwx    2 mbarrio5 hdfs 14710621726 2020-11-26 23:45 testdata.csv
drwxr-xrwx    - mbarrio5 hdfs             0 2020-11-16 21:18 tmp
-bash-4.1$ ls
testdata.csv
-bash-4.1$ chmod 777 testdata.csv
-bash-4.1$ |
```

$ cd ..

$ cd ..

$ cd dev

$ cd shm

$ ls

$ hdfs dfs -get testdata.csv

$ ls


## 10. SCP to retrieve testdata.csv (note the " . " at the end of csv)

$ scp mbarrio5@129.150.64.74:/home/mbarrio5/testdata.csv .

```
takum@DESKTOP-CNCNCF8 MINGW64 ~/Desktop
$ scp mbarrio5@129.150.64.74:/dev/shm/testdata.csv .
-- WARNING -- This system is for the use of authorized users only. Individuals
using this computer system without authority or in excess of their authority
are subject to having all their activities on this system monitored and
recorded by system personnel. Anyone using this system expressly consents to
such monitoring and is advised that if such monitoring reveals possible
evidence of criminal activity system personnel may provide the evidence of such
monitoring to law enforcement officials.

mbarrio5@129.150.64.74's password:
testdata.csv                                 100%   14GB   3.8MB/s 1:00:51
```

11. Note: a successful file SCP will show the file at your desktop



12. Open testdata and create a new column like so. Then copy ALL of column "A"

13. Paste into the NEW column then rename column A and B like so (Stock Ticker and Country)



14. Now copy column A again, make a new sheet, and paste it into the new sheet. This is to ensure that we did not lose the data when we modify the sheet using the Microsoft Excel Replace tool.

15. Now we begin the find and replace section. This is to add a location to each stock ticker. This is due to the fact that the dataset did not come with a location dimension. All of the tickers have an extension such as: HK, F, DE which are codenames for countries. We will use the replace tool to find these extensions and replace all matching cells with the proper country (ex: .HK = Hong Kong, .PA = France, .DE = Germany). The following replace entries **MUST** be done in the order specified. **Note:** all us tickers do NOT have an extension so the remaining .csv are tickers in the United States.

**Example:**
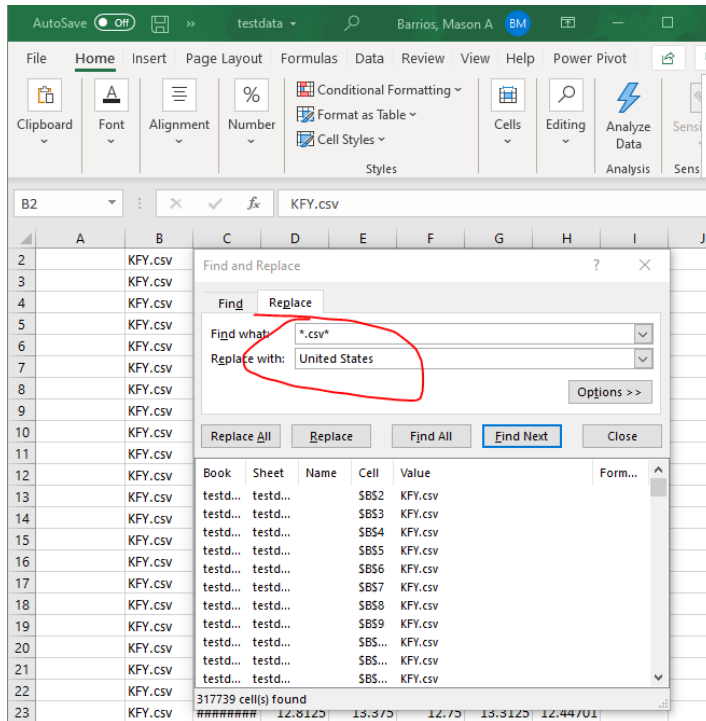CTRL + F
Example:
Find What = *.L*
Replace With = United Kingdom

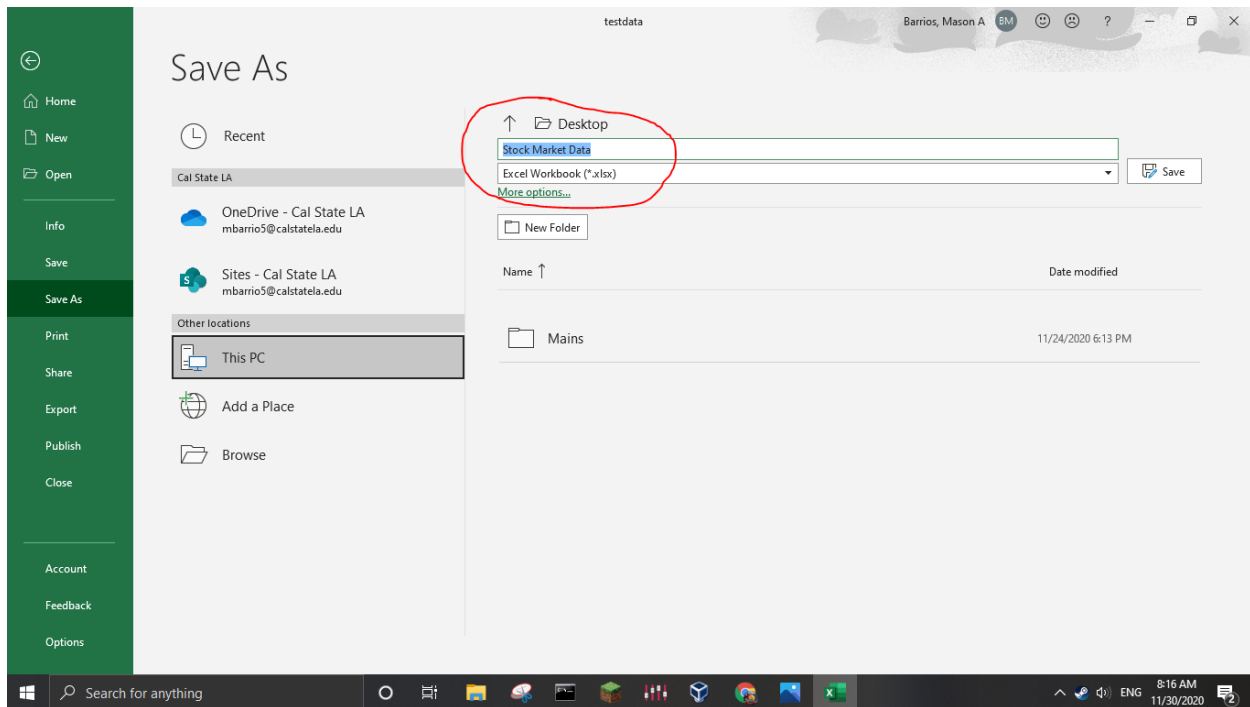**Entries for the Replace tool (complete in this order, do not include "|"):**
*.L* United Kingdom | *.SW* Switzerland | *.ST* Sweden | *.MC* Spain | *.KQ* South Korea
*.OL* Norway | *.AS* Netherlands | *.TA* Israel | *.JK* Indonesia | *.NS* India | *.BO* India |
*.AT* Greece | *.DE* Germany | *.MU* Germany | *.BE* Germany | *.PA* France |
*.F* France | *.CO* Denmark | *.SS* China | *.TO* Canada | *.V* Canada | *.SA* Brazil |
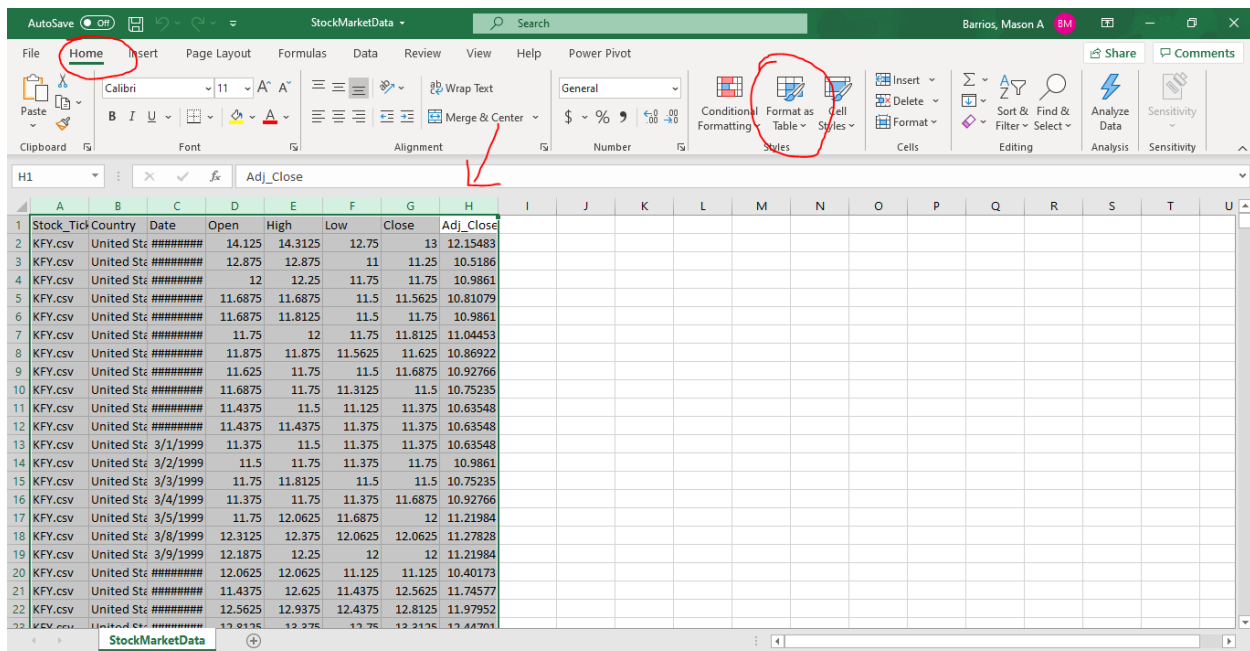*.BR* Belgium | *.AX* Australia | *.CSV* United States

13. Here is the last replace you will do, for the US:
Now go to the other sheet in the excel document to
copy column A in sheet1 to the empty column (column A)
into the sheet test data. Your sheet should like the image on
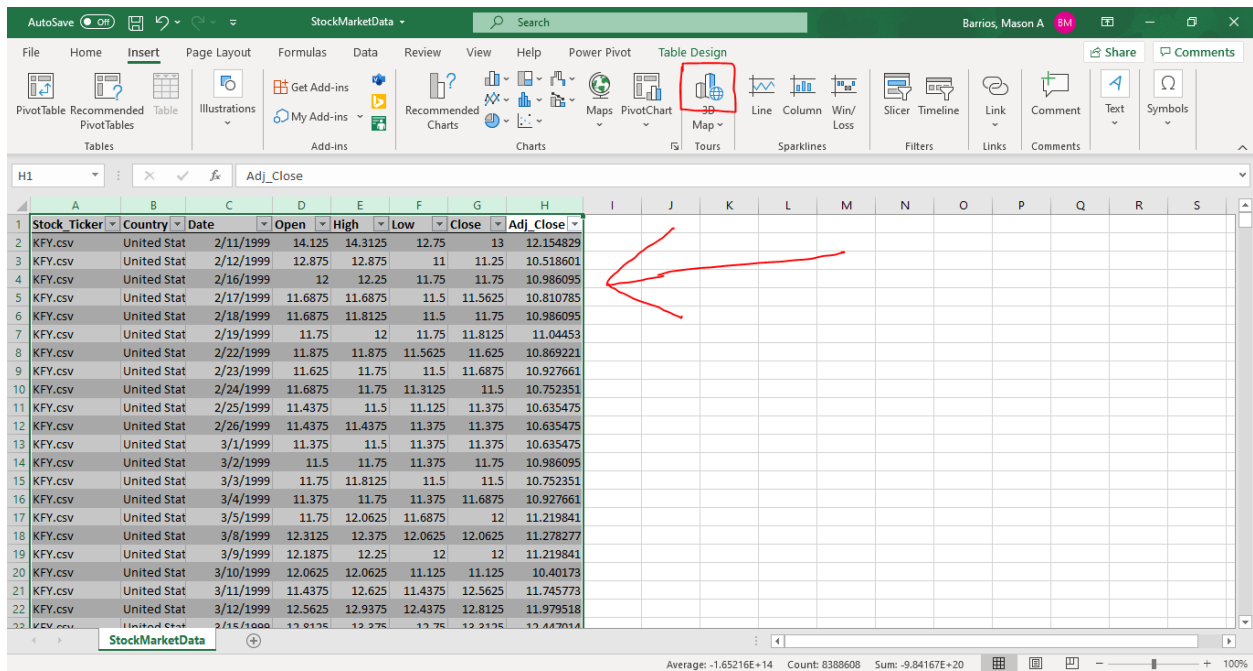the left.

15. Now go to save as and save the document as "Stock Market Data" and the .xlsx extension. It is necessary to be saved in .xlsx to use power maps.
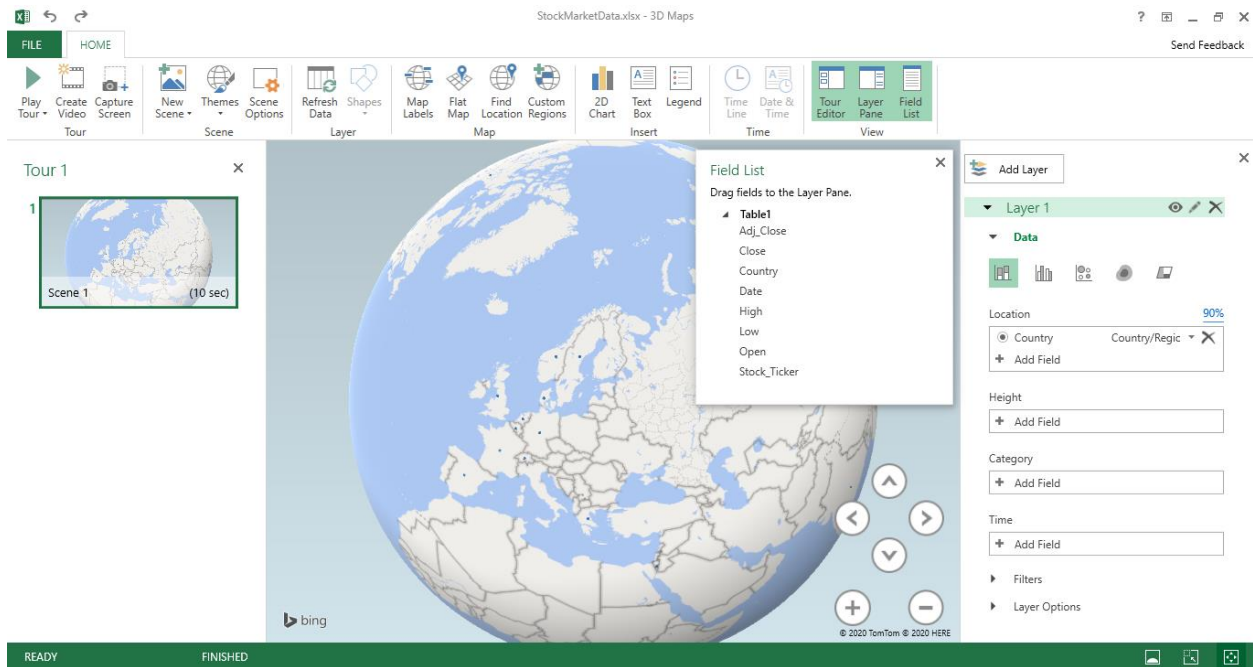


16. Select all then format as table via the home page. Use any design.

17. Now select all again and click on 3D map



Here is how the 3D map looks like upon opening

16. Now we can select "stock ticker" as the category, adj close as the size, and use the pie chart for visualization. We can click on the any section of the chart in a region and find the various SUM of the various stock tickers.