

Lending Club Loan Data Analysis

Group B

Riki Chang, Carolina Rivera, David Sayad, Chris Wallace
California State University Los Angeles

E-mail: cwallac9@calstatela.edu, criver47@calstatela.edu, dsayad2@calstatela.edu,
rchang12@calstatela.edu

Abstract: Our aim is to take this loan data from the Leading Club Company and create a geospatial analysis of each state in the United States of America. This will include state by state information on Home Ownership, Loan Status, and Purpose of Load for each state.

1. Introduction

We took data from kaggle.com, a database of public datasets, tutorials, and machine learning job assets service.

This dataset offers 887379 entries of data totalling at 392 MB. Each line includes a great amount of data that is not used for our aim such as: loan amount, term, interest rate, installment payment, grade, annual income, payment installations, etc. [1]

Narrowing down this information into a form which can be properly analyzed will bring much insight. This will be changed into be in terms of total loans given in a month or state-by-state allocation.

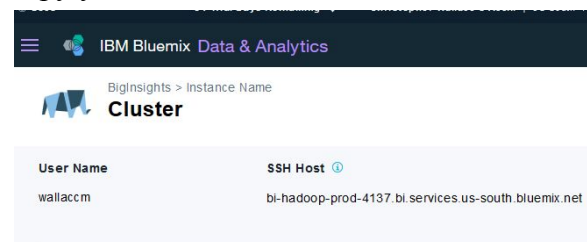
2. General Instructions

2.1 BigInsights

You need to open up your BigInsights account or create one on IBM's website if you do not have one[2]. We have included two links below the first being IBM's Bluemix and the second being the link to our dataset. You can download the data file Lending Club data from Kaggle which is what we did. [1]

2.2 Hadoop

Once you have logged into your BigInsights copy your ssh link.



Start the logging in process in putty or terminal. Once in terminal do as follows:

\$ ssh [username@yourlink.com](#)

You will be prompted to enter your password next.(it will be invisible when you type.. Do not be alarmed it is recording your keystrokes)

2.3 Import Data

After logging in enter the next command to create a directory inside HDFS that will be used to store the data.

\$ hdfs dfs -mkdir tmp/project;

\$ hdfs dfs -mkdir tmp/project/tables;

Next we will manually upload the data through Ambari. Navigate to the same page where you grabbed your login link and click on the link entitled Ambari. Use the same username and password that you used to log into HDFS. After logging into Ambari do as follows:

FileBrowser > User > (username) > tmp > project > tables > Upload > local place you stored Data

2.4 Hive

Start up your hive shell by entering the following command.

\$ hive

In the hive shell we will upload the data into a table stored in the tables directory. Next we will create a table that will extract the columns of the data that we are looking for.

```
hive>CREATE EXTERNAL TABLE IF NOT EXISTS project( json_responce STRING) STORED AS TEXTFILE LOCATION "/tmp/project/tables";
```

```
hive>CREATE TABLE IF NOT EXISTS project( id BIGINT, loan_amnt BIGINT, home_ownership STRING, issue_d STRING, loan_status STRING, purpose STRING, addr_state STRING); INSERT OVERWRITE TABLE project SELECT id, loan_amnt, home_ownership, issue_d, loan_status, purpose, addr_state FROM credit_card_dataset_CSV.csv WHERE id = id
```

To be sure that both command work type in the following command to look at all the tables in your hive shell.

```
hive> show tables;
```

Once it is verified that the tables have been created you can execute a number of commands to query the table further.

Example:

```
hive> Select * From project limit 50;
```

Which translates to select the first fifty rows from the project table.

3. Visualization of Data

3.1 Power View

To visualize the data you must download the csv file from Ambari File Manager. From then you open this file in excel and save it as an excel file(xlsx). From there navigate to the **Insert** tab and select **Power View** to begin a Power View Report. From here you select your **Power View Fields** that you would like to visualize(date, loan amount). Make sure to order and format the date in order to output the correct graph. Next from the **Design** Tab select **Other Charts (Bar)** to display a bar graph of Total Loan Amounts given by the Leading Club (Appendix 1).

3.2 3D Map

Visualizing this data in a state by state basis brings much insight to the needs of the state population. To do this you must have Data Analysis add-in enabled in excel. After that navigate to the Insert Tab and select **3D Map**. From here you must create different layers to show different attributes of data. Select **Add Layer** and select **Bubble** visualization. From here the **Location** of the data will always be *addr_state* and you must select **State/Province**. Then for **Category** select the attribute you would like to to visualize. For our project we selected to have it in terms of Home Ownership (Appendix 2), Loan Status (Appendix 3), and Purpose (Appendix 4).

4. Analysis

4.1 Temporal Analysis

The temporal analysis can show a great amount of knowledge to the loan usage of the general public of the United States. As one can see there are major spikes of total loan amounts in the months of March, July, and October (Appendix 1). They are tapered off greatly after these spikes also until the next month in this cycle is approached. The

end of Winter and the beginning of Spring is typically the time when people buy homes and this could explain the spike in March. Colleges start towards the end of August and many loans for students need to be taken out, this could explain the spike in July. October is the highest total loan amount and one could rationalize this to the holiday season and the expenses involved with them.

4.2 Geospatial Analysis

The geospatial analysis was also showed many trends in the initial states with loan receivers home ownership status, loan status, and purpose of the loan. It is very interesting that in most states about half of the people within this dataset own their house with the notable exception of California and New York (Appendix 2). This can be explained by the well known cost of living in these states and the amount of capital is needed to own a home. In terms of Loan Status per State there does not seem to be much variability between the states (Appendix 3). The only notable exception being that of California having a slightly larger percentage of Fully Paid loans that being 27.8%. Lastly in terms of Geospatial Analysis, we have the Loan Status per State (Appendix 4). There does not seem to be a distinguishable trend from state to state everything is very similar.

5. Key Terms

Ambari: The Apache Ambari project is aimed at making Hadoop management simpler by developing software for provisioning, managing, and monitoring Apache Hadoop clusters. [5]

Analytics: The process of collecting, processing and analyzing data to generate insights that inform fact-based decision-making. In many cases it involves software-based analysis using algorithms. [3]

Apache Pig: a platform for analyzing large data sets that consists of a high-level language for expressing data analysis programs, coupled with infrastructure for evaluating these programs. [4]

BigInsights: It's a software platform designed to help firms discover and analyze business insights hidden in large volumes of a diverse range of data—data that's often ignored or discarded because it's too impractical or difficult to process using traditional means. [7]

Flume: A distributed and reliable way to collect, group, and transfer large amounts of data from many sources to a central data store.[3]

Hadoop: Apache Hadoop is one of the most widely used software frameworks in big data. It is a collection of programs which allow storage, retrieval and analysis of very large data sets using distributed hardware (allowing the data to be spread across many smaller storage devices rather than one very large one). [3]

HCatalog: Makes metadata (metastore) for Hive and merges it with what Pig does.[3]

HDFS: Hadoop Distributed File System; the way that Hadoop structures its files. [3]

Hive: A higher level language that uses HQL, which is similar to SQL (Structured Query Language) in its syntax. [3]

MapReduce: Refers to the software procedure of breaking up an analysis into pieces that can be distributed across different computers in different locations. It first distributes the analysis (map) and then collects the results back into one report (reduce). [3]

Sqoop: Apache Sqoop(TM) is a tool designed for efficiently transferring bulk data between Apache Hadoop and structured datastores such as relational databases. [6]

6. Conclusion

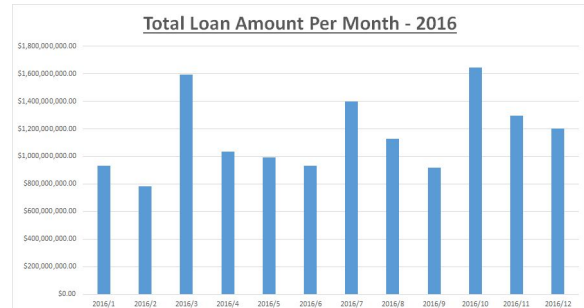
After loading the data into excel we gained insights into how all of this lending data can be visually observed. By looking at the maps provided we can see the total loan amounts(Appendix 1), the loan status(Appendix 3), the purpose of the loan(Appendix 4), and proportionally what kind of loans people are taking out for home ownership(Appendix 2). For example in the Geospatial Analysis we found out that Californians tend to take out more loans for rent where Texans they tend to take out loans for mortgages. We find expect this to be due to the price of owning a house being very high in California as opposed to Texas where it much lower. In conclusion, there are a number of ways this data can be interpreted and we will leave it up to you to continue to explore and find helpful insights.

References

- [1] Kan, W. "Lending Club Loan Data." Retrieved from <https://www.kaggle.com/wendykan/lending-club-loan-data>
- [2]<https://www.ibm.com/cloud-computing/bluemix/>
- [3]"Big Data: The Key Vocabulary Everyone Should Understand"(n.d). Retrieved from <https://www.linkedin.com/pulse/20141203075716-64875646-big-data-the-key-vocabulary-everyone-should-understand>
- [4]"Welcome to Apache Pig!" Retrived from <https://pig.apache.org/>
- [5]"Apache Ambari" Retrieved from <http://ambari.apache.org/>
- [6] "Apache Sqoop" Retrieved from <http://sqoop.apache.org/>
- [7] "Understanding InfoSphere BigInsights" Retrieved from

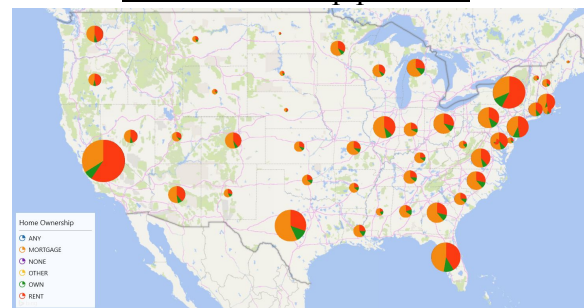
<http://www.ibm.com/developerworks/data/library/techarticle/dm-1110biginsightsintro/index.html>

Appendix 1



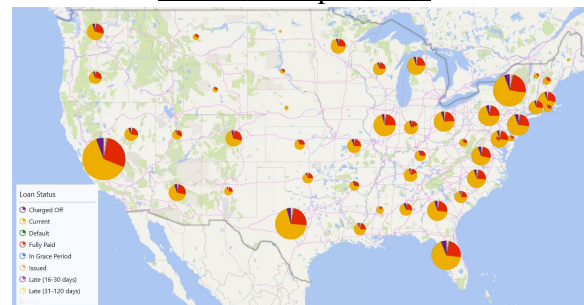
Appendix 2

Home Ownership per State



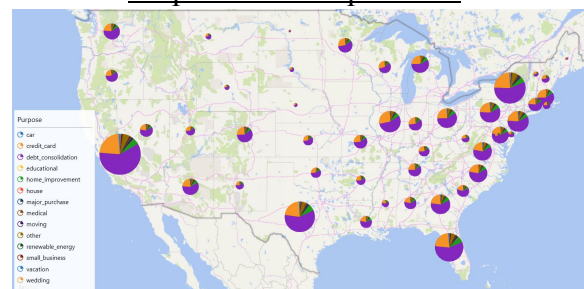
Appendix 3

Loan Status per State



Appendix 4

Purpose of Loan per State



Term Paper Rubric: 100%

It should be almost same as the team presentation. But, mostly, I will take a look at if you revise the content per my comment at the presentation. Thus, any penalty at the presentation can be recovered. You also need to email the instructor the peer evaluation for the term paper. If you don't email me peer evaluation, I assume, all of you contribute the work fairly well. For example, your team score is 95% and your peer evaluation by your team members are 100%, your score is 95 ($= 95 \times 100\%$)

Team Project Tutorial Rubric: 100%

1. Materials Available (30%)

- a. If Data Set can be downloadable per the direction (15%)
- b. If source code is downloadable per the direction (15%)

2. Completeness (70%)

- a. If each step is clear to follow (15%)
- b. If the source code is correct (15%)
- c. If the source codes are executable or possible to copy/paste to execute (15%)
- d. If the visualization using other tools are easy to follow or clearly executable (15%)
- e. If the geo-spatial visualization is clear (10%)

NOTE: You also need to email the instructor the peer evaluation for the term paper (Optional). If you don't email me peer evaluation, I assume, all of you contribute the work fairly well. For example, your team score is 95% and your peer evaluation by your team members are 100%, your score is 95 ($= 95 \times 100\%$)

Project Lab Tutorial Group B

Riki Chang, Carolina Rivera, David Sayad, Chris Wallace
California State University Los Angeles

E-mail: cwallac9@calstatela.edu, criver47@calstatela.edu, dsayad2@calstatela.edu,
rchang12@calstatela.edu

Lending Club Loan Data Analysis

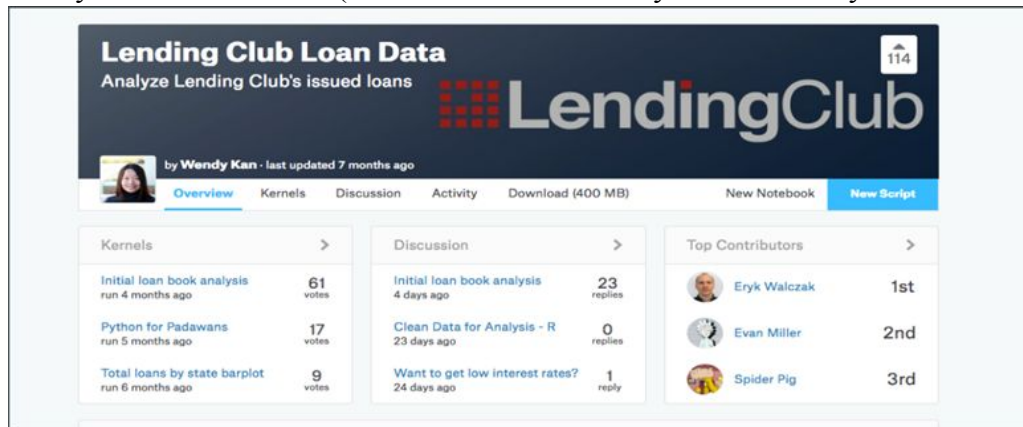
Objectives

In this hands-on lab, you will learn how to:

- Create IBM Bluemix Account
- Create directories in cluster and load data
- Learn how to use Ambari
- Introduced to Hive Shell
- Hive commands to perform the analysis.
- Visualization both Temporal and Geospatial

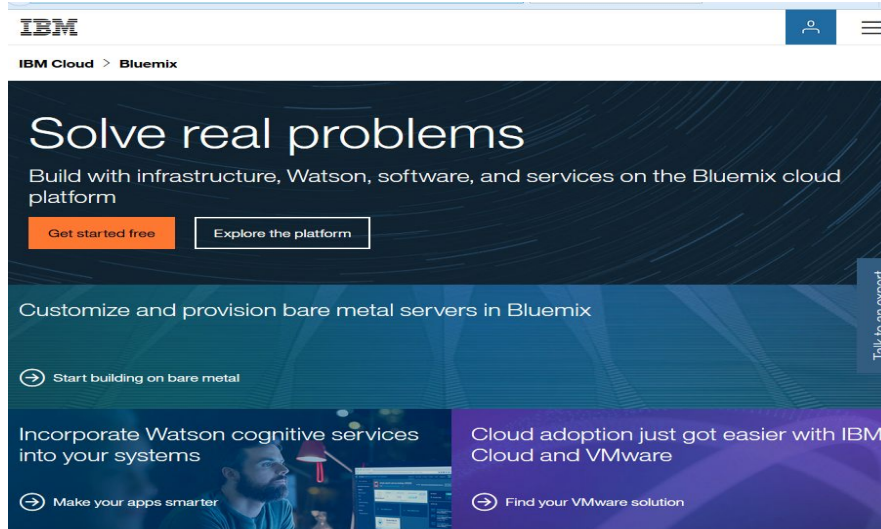
Exercise 1: Get data manually from keggel

1. Go to <https://www.kaggle.com/wendykan/lending-club-loan-data> to download the data and save it to your local machine. (Remember the location you save it too you will need it later)

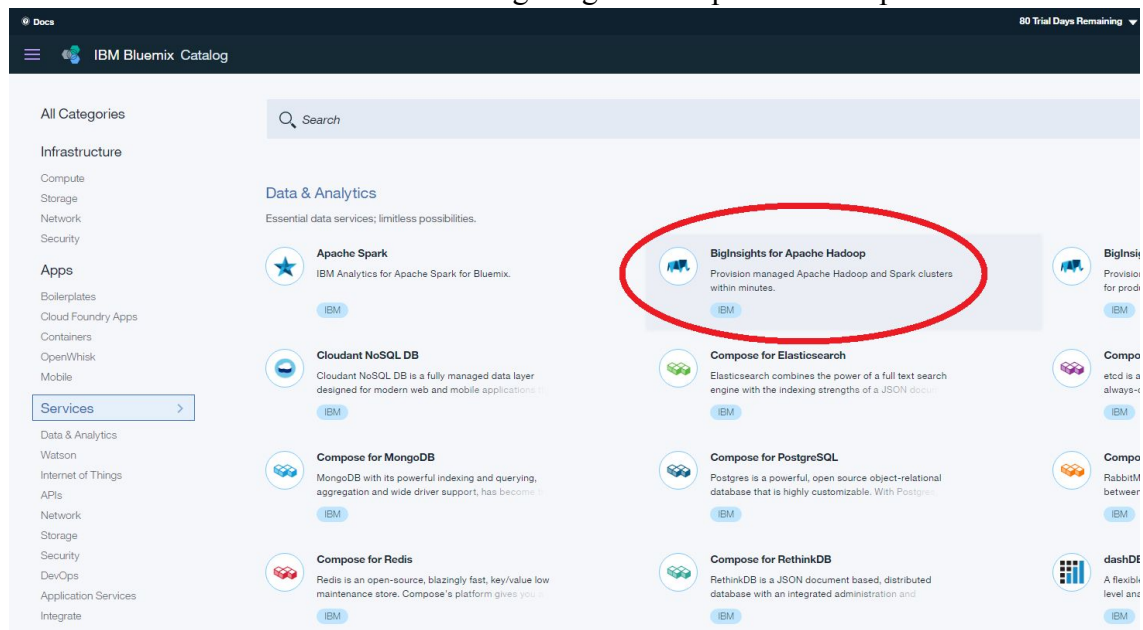


Exercise 2: Create Bluemix Account and a Cluster

1. Next open up a new tab in your browser and go to <https://www.ibm.com/cloud-computing/bluemix/> to create your account.



2. Select create service and select BigInSights for Apache Hadoop



3. Then navigate to **Manage Cluster** then **New Cluster** to create your cluster.

Configuration

Data Storage(GB)

244

Specify the required storage size for the data being analyzed. You can specify a maximum of 732 GB.

Number of Data Nodes *

1

Specify the number of data nodes required for your cluster. You can specify a maximum of 3 data nodes.

IBM Open Platform Version ⓘ

IOP 4.2

Cloud Storage ⓘ

NONE

Click Cloud Storage to configure data stores.

Mandatory Components

✓ HDFS

✓ HBASE

✓ AMBARI METRICS

✓ YARN

✓ ZOOKEEPER

✓ HIVE

✓ MAPREDUCE2

✓ KNOX

Optional Components

☒ SPARK

☒ PIG

☒ SQOOP

☐ OOZIE

☒ FLUME

☐ R

4. Make sure that **SPARK**, **PIG**, **SQOOP**, and **FLUME** are all checked off on the configuration and create your cluster

Exercise 3: Create directories in cluster and load data

1. Open up terminal/Putty on your the machine you are using and we will login to HDFS by typing in this command substituted with your info
\$ssh username@bi-hadoop-prod-4137.bi.services.us-south.ibm.com

BigInsights > Instance Name
Cluster

User Name

wallaccm

SSH Host ⓘ

bi-hadoop-prod-4137.bi.services.us-south.ibm.com

2. After logging in enter the next command to create a directory inside HDFS that will be used to store the data that we downloaded from keggel in the previous step.

```
$ hdfs dfs -mkdir tmp/project;
$ hdfs dfs -mkdir tmp/project/tables;
```

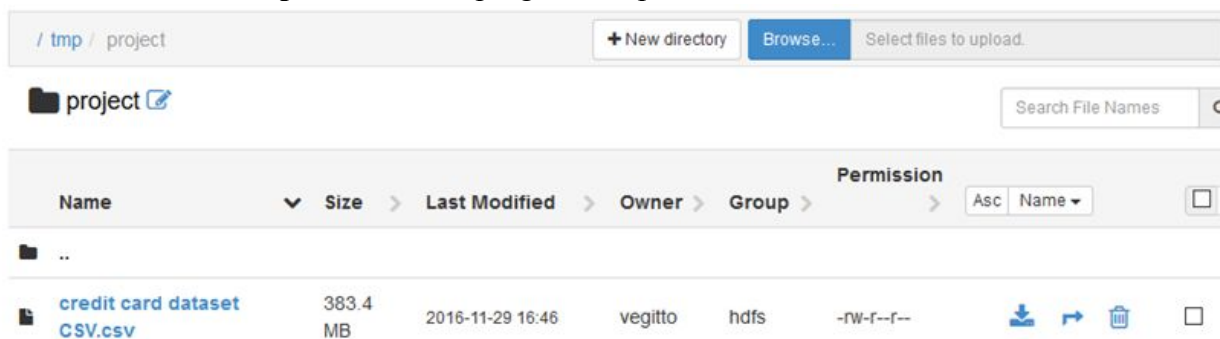


```
bi-hadoop-prod-4152.bi.services.us-south.bluemix.net - PuTTY
bash-4.1$ hdfs dfs -mkdir /tmp/project
mkdir: '/tmp/project': File exists
bash-4.1$
```

3. Next we will manually upload the data through Ambari. Navigate to the same page where you grabbed your login link and click on the link entitled Ambari. Use the same username and password that you used to log into HDFS. After logging into Ambari do as follows:

FileBrowser > User > (username) > tmp > project > tables

4. Then select **Upload** on the top right and upload the data.



Exercise 4: Hive Commands

1. Start up your hive shell by entering the following command.

\$ hive

2. In the hive shell we will upload the data into a table stored in the tables directory. Next we will create a table that will extract the columns of the data that we are looking for.

**hive>CREATE EXTERNAL TABLE IF NOT EXISTS project(json_response STRING)
STORED AS TEXTFILE LOCATION "/tmp/project/tables";**

**hive>CREATE TABLE IF NOT EXISTS project(id BIGINT, loan_amnt BIGINT,
home_ownership STRING, issue_d STRING, loan_status STRING, purpose STRING,
addr_state STRING);
INSERT OVERWRITE TABLE project SELECT id, loan_amnt, home_ownership,
issue_d, loan_status, purpose, addr_state FROM credit_card_dataset_CSV.csv WHERE id
= id**

3. To be sure that both command work type in the following command to look at all the tables in your hive shell.

hive> show tables;

4. Once it is verified that the tables have been created you can execute a number of commands to query the table further.

Example:

hive> Select * From project limit 50;

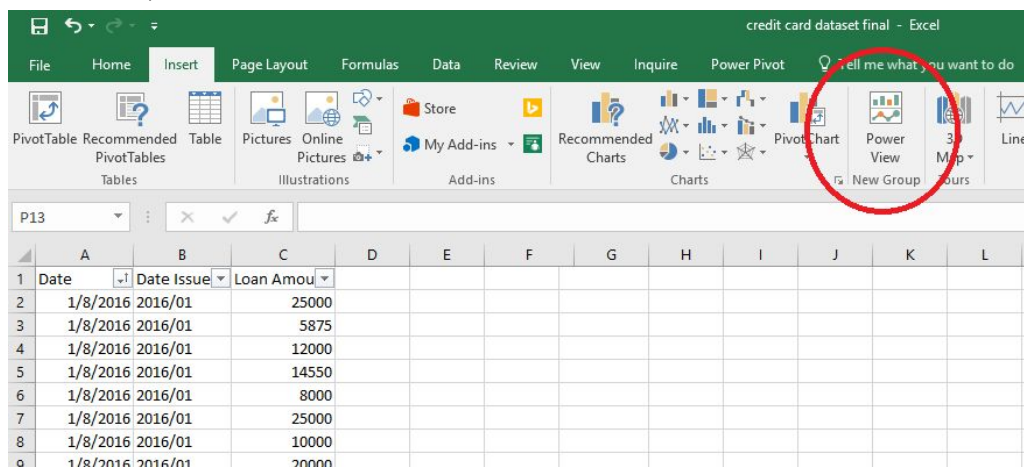
Which translates to select the first fifty rows from the project table.

Exercise 5: Temporal Analysis

1. To visualize the data, you must download the csv file from Ambari File Manager.

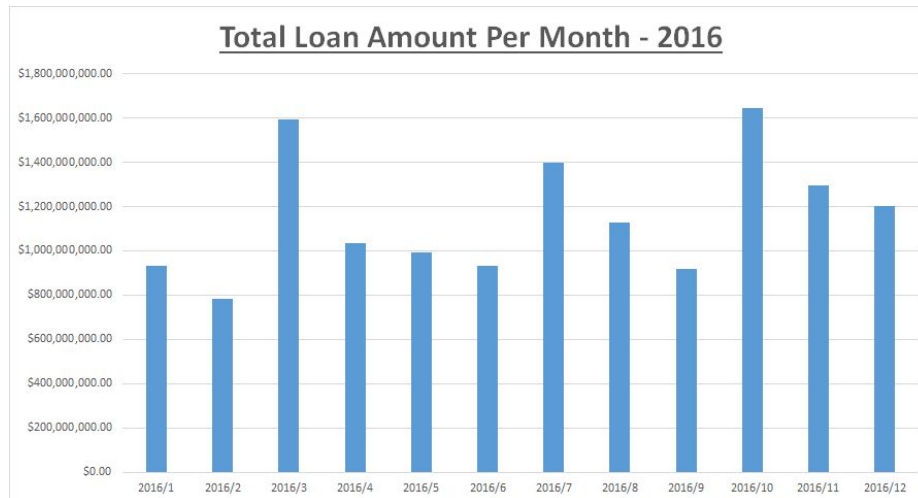
2. From there you open this file in excel and save it as an excel file(xlsx).

3. From there navigate to the Insert tab and select Power View to begin a Power View Report. (Make sure that you enable Power View through the options under all add-ins and add it to the Excel Ribbon)



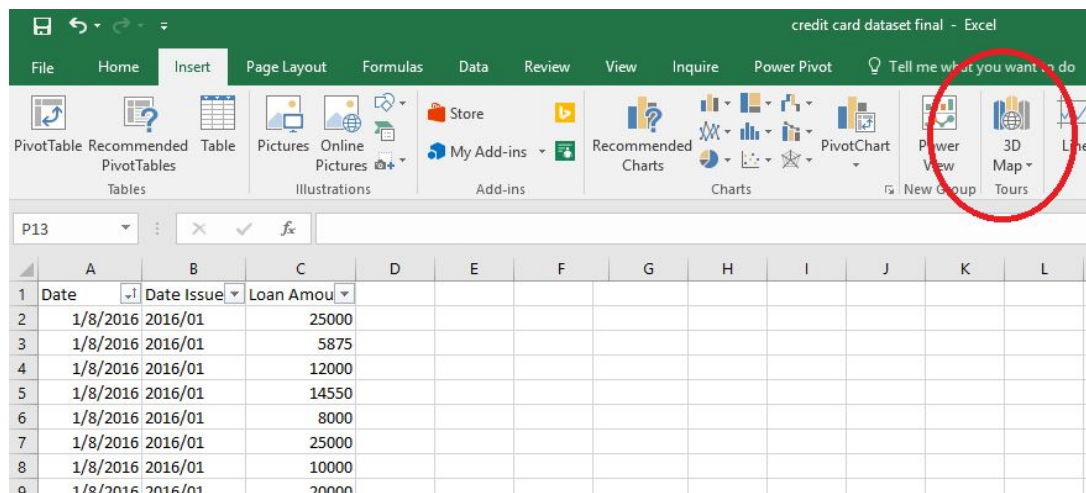
4. From here you select your Power View Fields that you would like to visualize (date, loan amount). Make sure to order and format the date in order to output the correct graph.

5. Next from the Design Tab select Other Charts (Bar) to display a bar graph of Total Loan Amounts given by the Leading Club shown below.

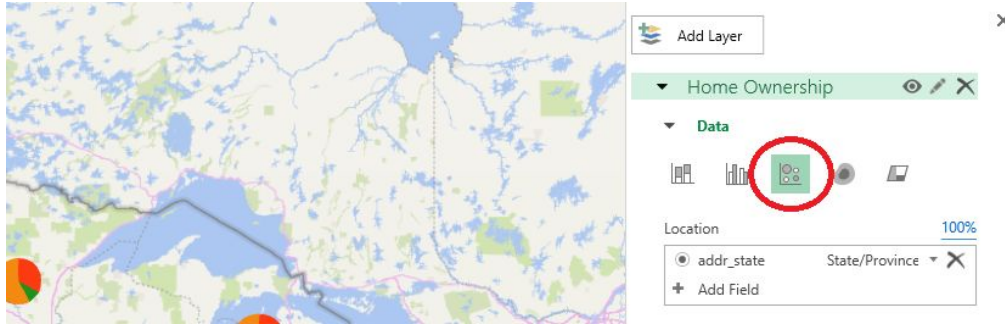


Exercise 6: Geospatial Analysis

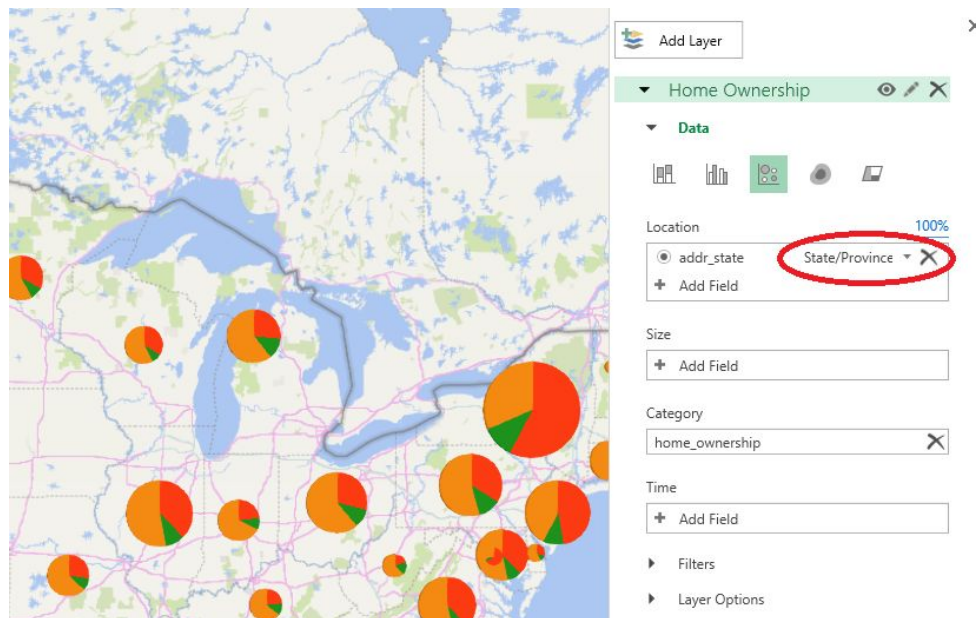
1. First you must have the Data Analysis add-in enabled in excel. After making sure that it is enabled, navigate to the Insert Tab and select **3D Map**.



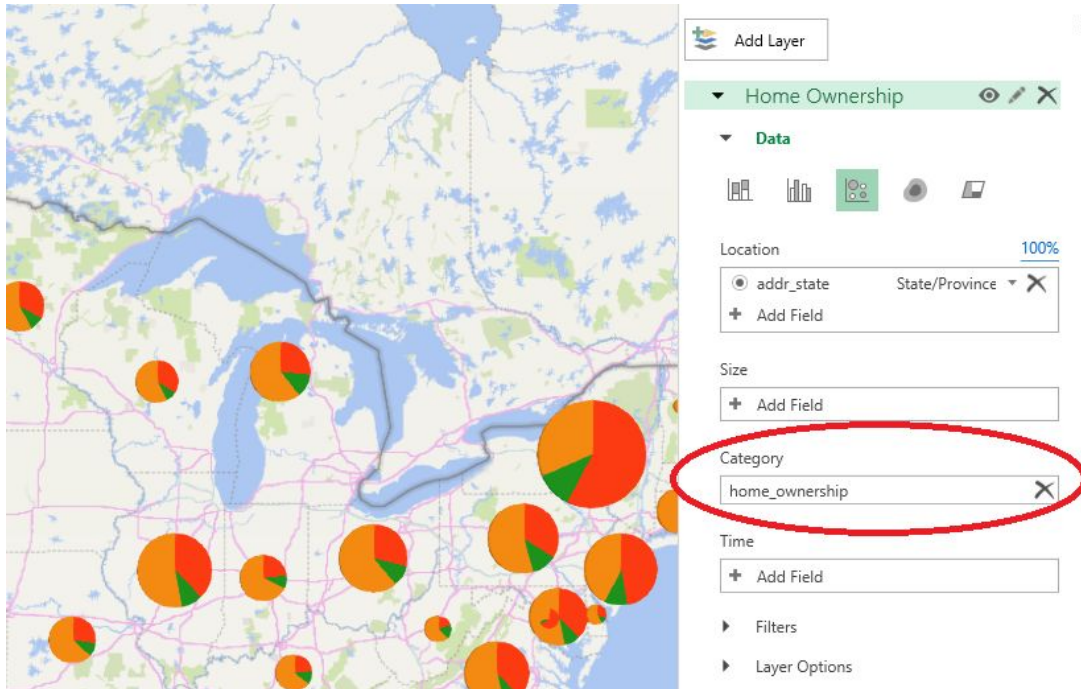
2. From here you must create different layers to show different attributes of data. Select **Add Layer** and select **Bubble** visualization.



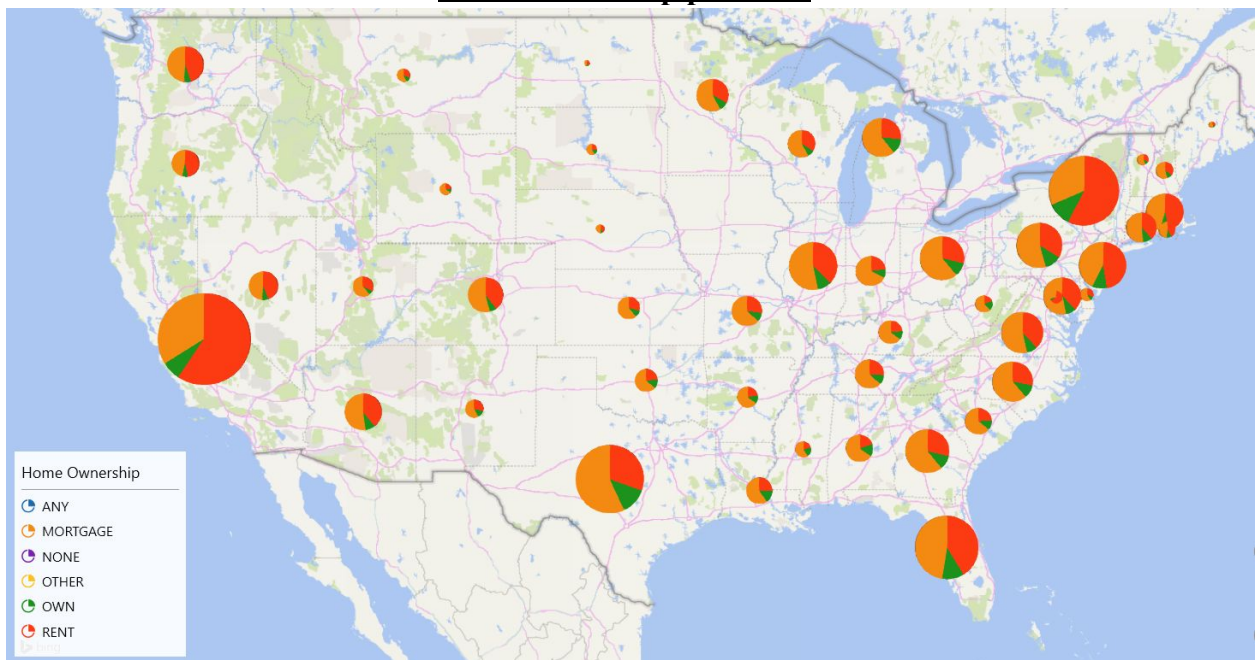
3. From here the **Location** of the data will always be *addr_state* and you must select **State/Province**.



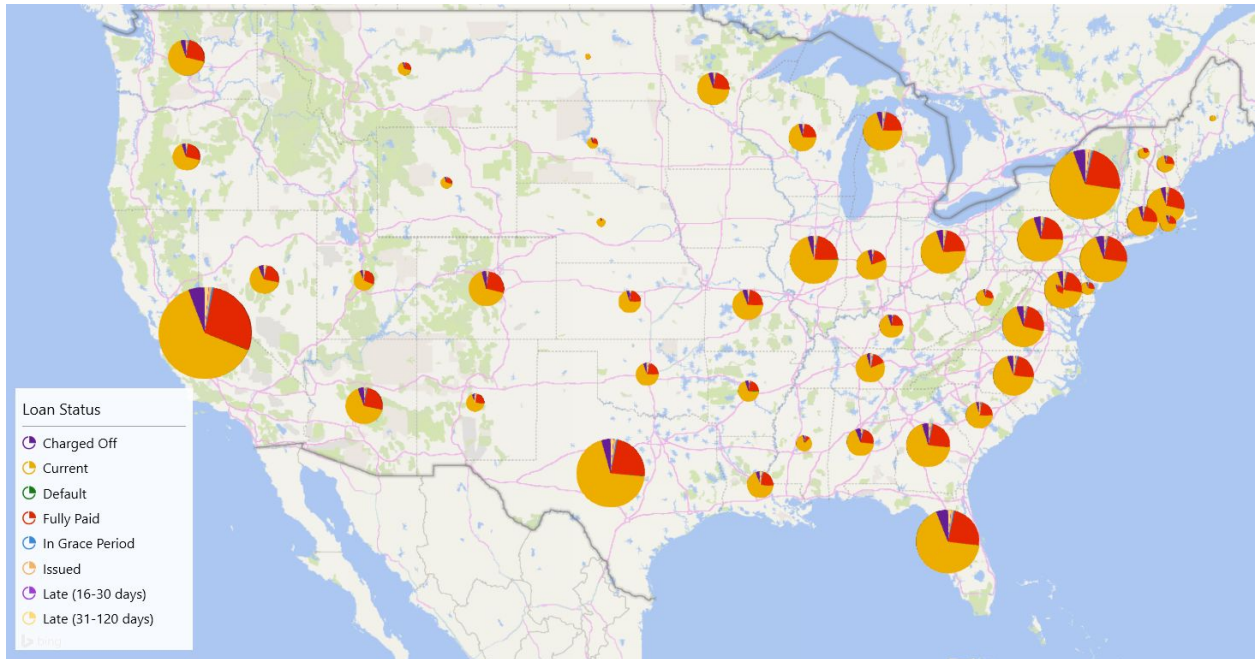
4. Then for **Category** select the attribute you would like to visualize. We selected to have it in terms of Home Ownership, Loan Status, and Purpose. (Shown Below)



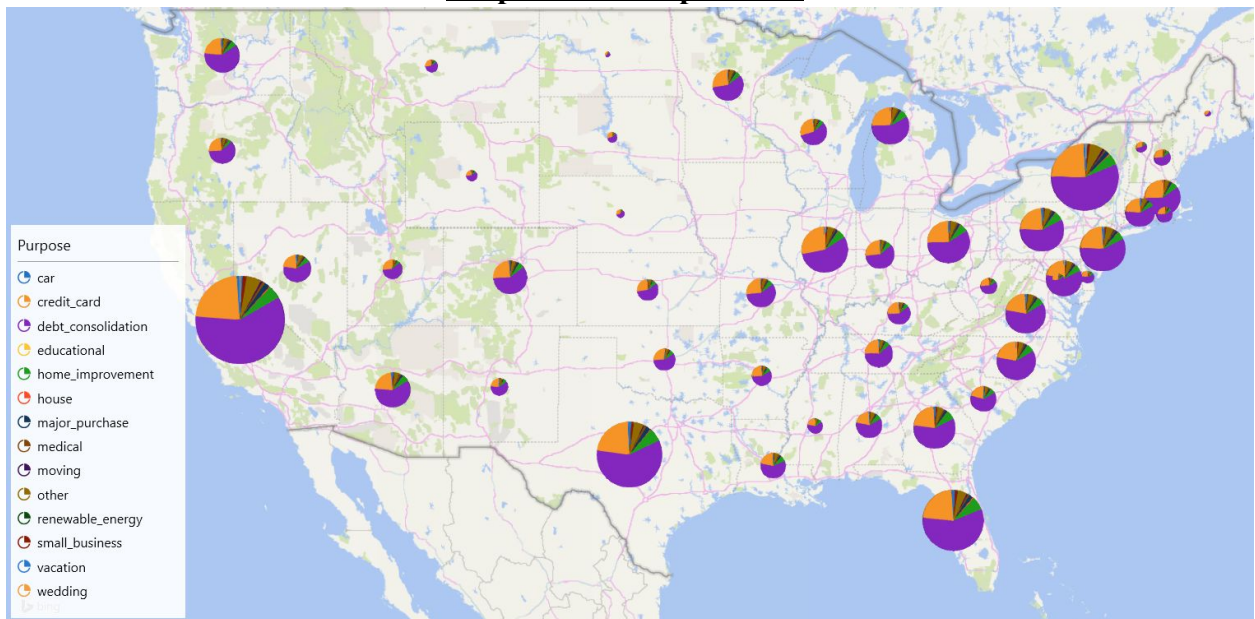
Home Ownership per State



Loan Status per State



Purpose of Loan per State



END OF LAB