

Bank Loan

By: Dawit Abay

INTRO

For my project I combines HDFS data and MapReduce algorithms to evaluate the risk of granting a loan to a person based on their location, loan type, and average risk.



Input Data

| CustomerID | Customer Name | Loan Account Number | Sanctioned Loan Amount | Currency | Disbursed Loan Amount | Loan Status | Risk in % | City | State | Country | Reason For Taking Loan |
|------------|---------------|---------------------|------------------------|-----------|-----------------------|-------------|-----------|-------------------------|---------|--------------|------------------------|
| 1 | Michael | HHLNO1235648 | 300000 | US-Dollar | 250000 | Current | 20 | Atlanta | Georgia | United State | Purchased Home |
| 2 | Christopher | HPLNOD24578 | 150000 | US-Dollar | 130000 | Late | 10 | Augusta-Richmond County | Georgia | United State | Medical Emergency |
| 3 | Jessica | HVLQV458792 | 70000 | US-Dollar | 50000 | Current | 7 | Columbus | Georgia | United State | Purchased Vechel |
| 4 | Matthew | HRLQF854586 | 200000 | US-Dollar | 20000 | Closed | 0 | Macon-Bibb County | Georgia | United State | Purchased Furniture |
| 5 | Ashley | HHLNO1235649 | 300000 | US-Dollar | 250000 | Current | 20 | Savannah | Georgia | United State | Purchased Home |

This is CSV data that show collections of CustomerID ,Customer Name ,Loan Account Number, Sanctioned Loan Amount, Currency ,Disbursed Loan Amount ,Loan Status,Risk in % ,City, State, Country, Reason For Taking Loan.

There are over 500 account in the csv data

Hadoop Code Mapper

Mapper maps input key/value pairs to a set of intermediate key/value pairs.

Reducer reduces a set of intermediate values which share a key to a smaller set of values.

```
package com.AvgRiskCat;
import java.io.IOException;
import org.apache.hadoop.io.DoubleWritable;
import org.apache.hadoop.io.LongWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Mapper;
public class MyMapper extends Mapper<LongWritable,Text,Text,DoubleWritable> {
    public void map(LongWritable key, Text value, Context con) throws IOException, InterruptedException
    {
        String line = value.toString();
        String[] linePart = line.split(",");
        //problem -1 - avg risk
        Double risk = Double.parseDouble(linePart[7]);
        //problem-2 - avg risk per location
        //String loc = linePart[8].toString();
        String cat = linePart[2].toString().substring(1, 3);
        if(cat.equals("HL"))
        {
            cat = "Home Loan";
        }
        else if(cat.equals("PL"))
        {
            cat = "Personal Loan";
        }
        else if(cat.equals("VL"))
        {
            cat = "Vischel Loan";
        }
        else
        {
            cat = "Retailer Loan";
        }
        con.write(new Text(cat), new DoubleWritable(risk));
    }
}
```

Hive

Hive is a data warehouse infrastructure that facilitates querying and managing large data sets which resides in distributed storage system

After starting hdfs,yarn,mapreduce put input file in hdfs.

```
>hadoop fs -put source_path destination_path_with_file_name.
```

then start hive

```
>hive
```

To see database

```
>show database;
```

to create database

```
>create database bank;
```

To use a database

```
>use bank;
```

To create table

```
>create table
```

```
bankloan(cos_id,cos_name,loan_no,amount,currency,amount_paid,status,risk,local,reason)
```

```
>row format delimited
```

```
>fields terminated by ','
```

```
>stored as textfile;
```

Hadoop Code Reducer

This code shows us how to calculate the risk of each loan

```
double glSum = 0;
int glcount = 0;
//double avg = 0;
public void reduce(Text key,Iterable<DoubleWritable> val,Context con) throws IOException, InterruptedException
{
    for (DoubleWritable values : val) {
        glSum = glSum + values.get();
        glcount = glcount +1;
    }
    con.write(new Text(key), new DoubleWritable(glSum/glcount));
}
}
```

Run HDFS

Applications that run on HDFS need streaming access to their data sets. They are not general purpose applications that typically run on general purpose file systems. HDFS is designed more for batch processing rather than interactive use by users. The emphasis is on high throughput of data access rather than low latency of data access. POSIX imposes many hard requirements that are not needed for applications that are targeted for HDFS. POSIX semantics in a few key areas has been traded to increase data throughput rates.

```
king-one@kingone:~/Homework/Bank_Project/mapreduce/categoryRisk$ hadoop jar Avgcat.jar /user/dawit/data/input-data_new.txt /user/dawit/bank-project-folder
2021-07-19 17:40:23.810 INFO impl.MetricsConfig: Loaded properties from hadoop-metrics2.properties
king-one@kingone:~/Homework/Bank_Project/mapreduce/categoryRisk$ hdfs dfs -ls /user/dawit/bank-project-folder

Found 3 items
drwxr-xr-x - king-one supergroup          0 2021-07-04 17:06 /user/dawit/bank-project-folder/output-bankAVG
drwxr-xr-x - king-one supergroup          0 2021-07-04 17:08 /user/dawit/bank-project-folder/output-bankCategory
drwxr-xr-x - king-one supergroup          0 2021-07-04 17:04 /user/dawit/bank-project-folder/output-bankLocal
king-one@kingone:~/Homework/Bank_Project/mapreduce/categoryRisk$
king-one@kingone:~/Homework/Bank_Project/mapreduce/categoryRisk$ hdfs dfs -ls /user/dawit/bank-project-folder/output-bankAVG

Found 2 items
-rw-r--r-- 1 king-one supergroup          0 2021-07-04 17:06 /user/dawit/bank-project-folder/output-bankAVG/_SUCCESS
-rw-r--r-- 1 king-one supergroup          7 2021-07-04 17:06 /user/dawit/bank-project-folder/output-bankAVG/part-r-00000
king-one@kingone:~/Homework/Bank_Project/mapreduce/categoryRisk$
king-one@kingone:~/Homework/Bank_Project/mapreduce/categoryRisk$ hdfs dfs -ls /user/dawit/bank-project-folder/output-bankAVG/part-r-00000
-rw-r--r-- 1 king-one supergroup          7 2021-07-04 17:06 /user/dawit/bank-project-folder/output-bankAVG/part-r-00000
king-one@kingone:~/Homework/Bank_Project/mapreduce/categoryRisk$ hadoop fs -cat /user/dawit/bank-project-folder/output-bankAVG/part-r-00000
10.15
king-one@kingone:~/Homework/Bank_Project/mapreduce/categoryRisk$ hadoop fs -cat /user/dawit/bank-project-folder/output-bankCategory/part-r-00000
Home Loan      18.048
Personal Loan  13.968
Retailer Loan  12.261333333333333
Viechel Loan   10.15
king-one@kingone:~/Homework/Bank_Project/mapreduce/categoryRisk$ hadoop fs -cat /user/dawit/bank-project-folder/output-bankLocal/part-r-00000
Alpharetta      10.485714285714286
Atlanta         10.324137931034484
Augusta-Richmond County 10.513636363636364
Columbus        10.362068965517242
Dunwoody        10.222222222222221
Macon-Bibb County 10.104651162790697
Savannah       10.15
king-one@kingone:~/Homework/Bank_Project/mapreduce/categoryRisk$
```

Hadoop OutPut

Avg risk

The risk of each loan for home,medical
retailer and vehicle

The risk for each location in Georgia

```
10.15
Home Loan      18.048
Personal Loan  13.968
Retailer Loan  12.261333333333333
Vehicle Loan   10.15
Alpharetta     10.485714285714286
Atlanta        10.324137931034484
Augusta-Richmond County 10.513636363636364
Columbus       10.362068965517242
Dunwoody       10.222222222222221
Macon-Bibb County 10.104651162790697
Savannah      10.15
```


R Studio code

```
library(readr)
```

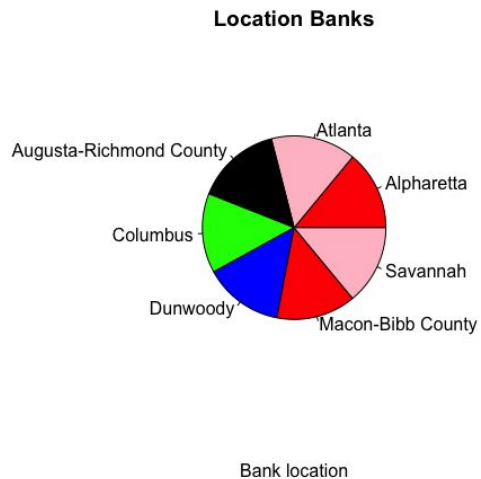
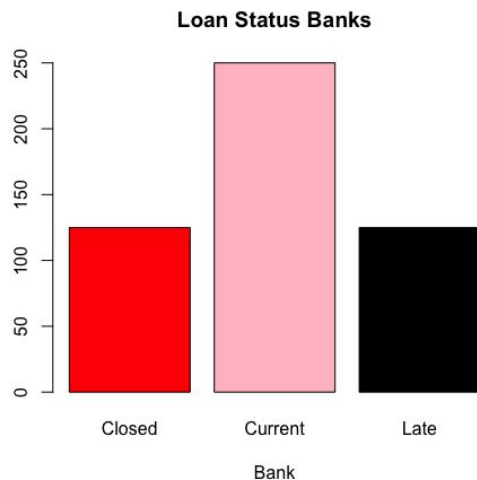
```
library(plotrix)
```

```
input_data <-  
read.csv("~/Desktop/hadoop_bank_project/input-data.  
csv",TRUE, sep = ",")
```

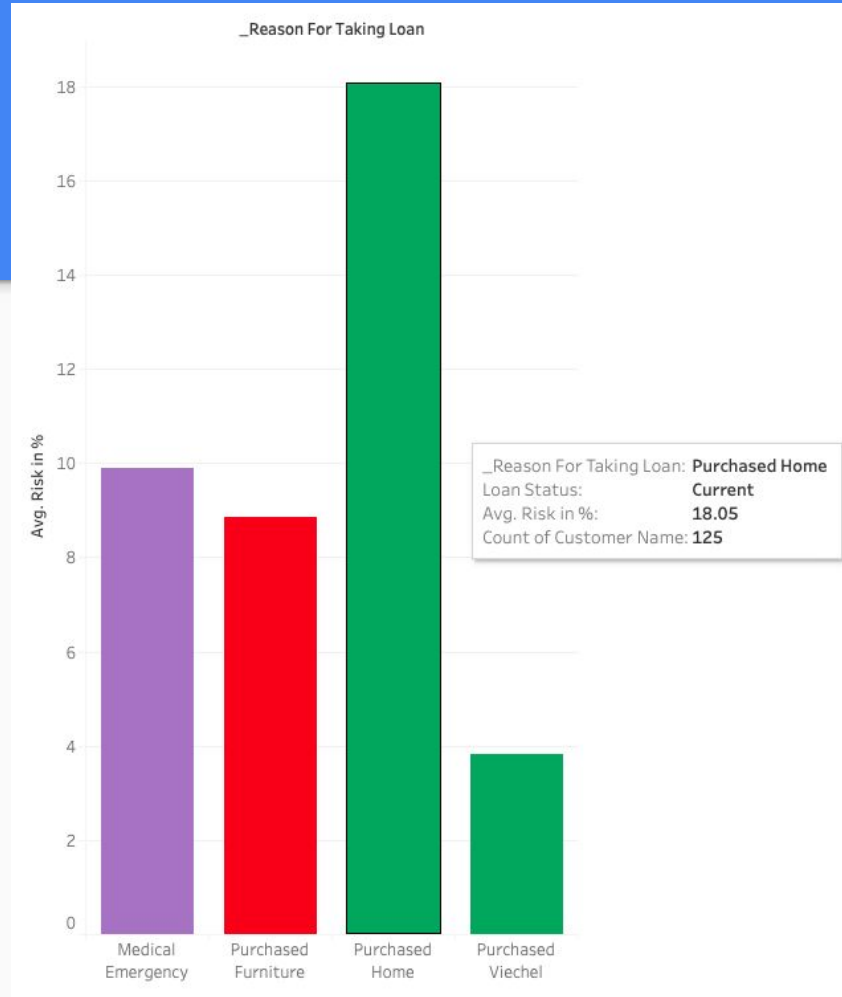
```
require("RColorBrewer")
```

```
pie(xtabs(~input_data$Location),main="Location  
Banks",xlab="Bank location",col =  
c("red","pink","black","green","blue"))
```

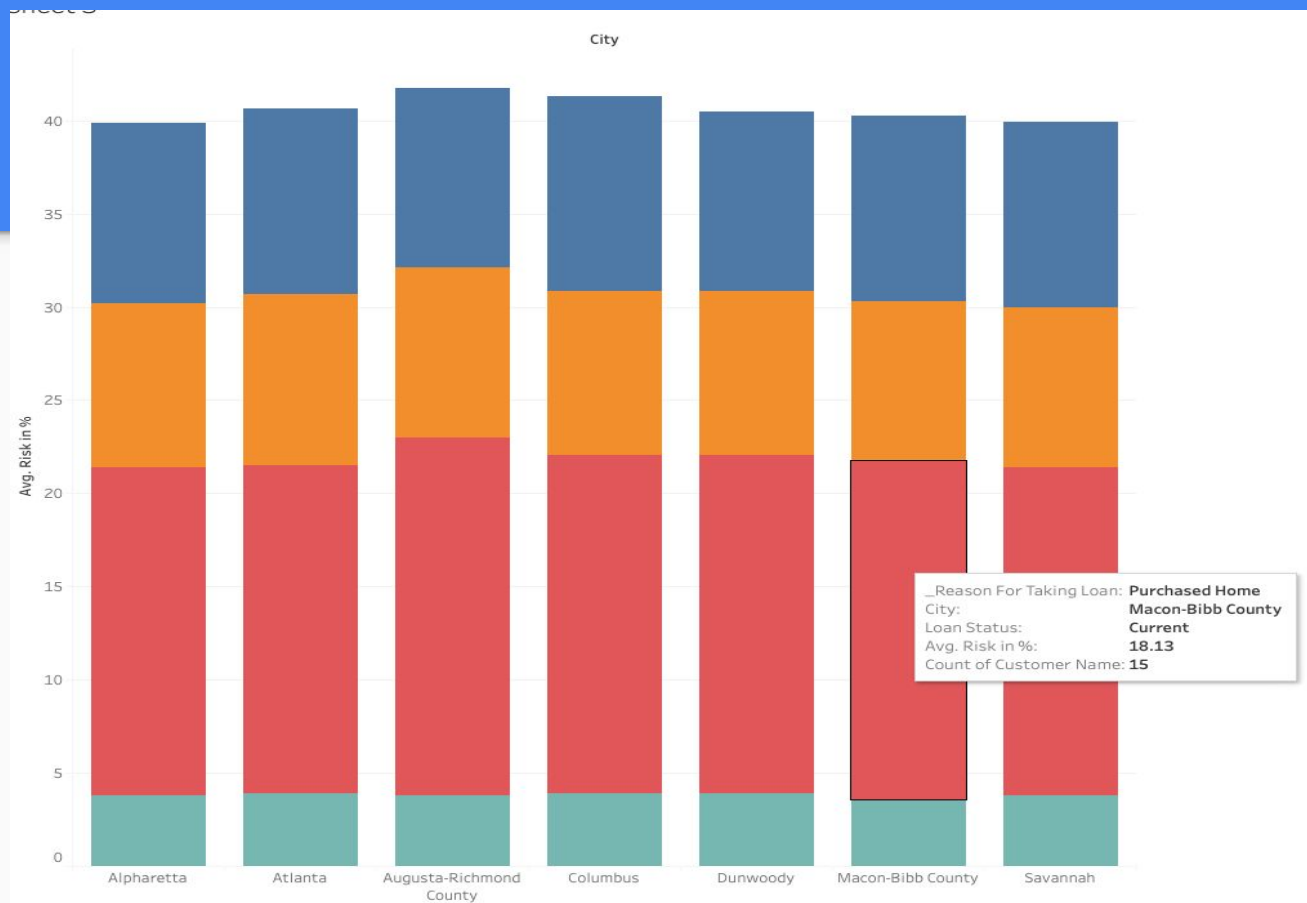
```
#barplot(xtabs(~input_data$`Loan  
Status`),main="Loan Status Banks",xlab="Bank",col =  
c("red","pink","black","green","blue"))
```



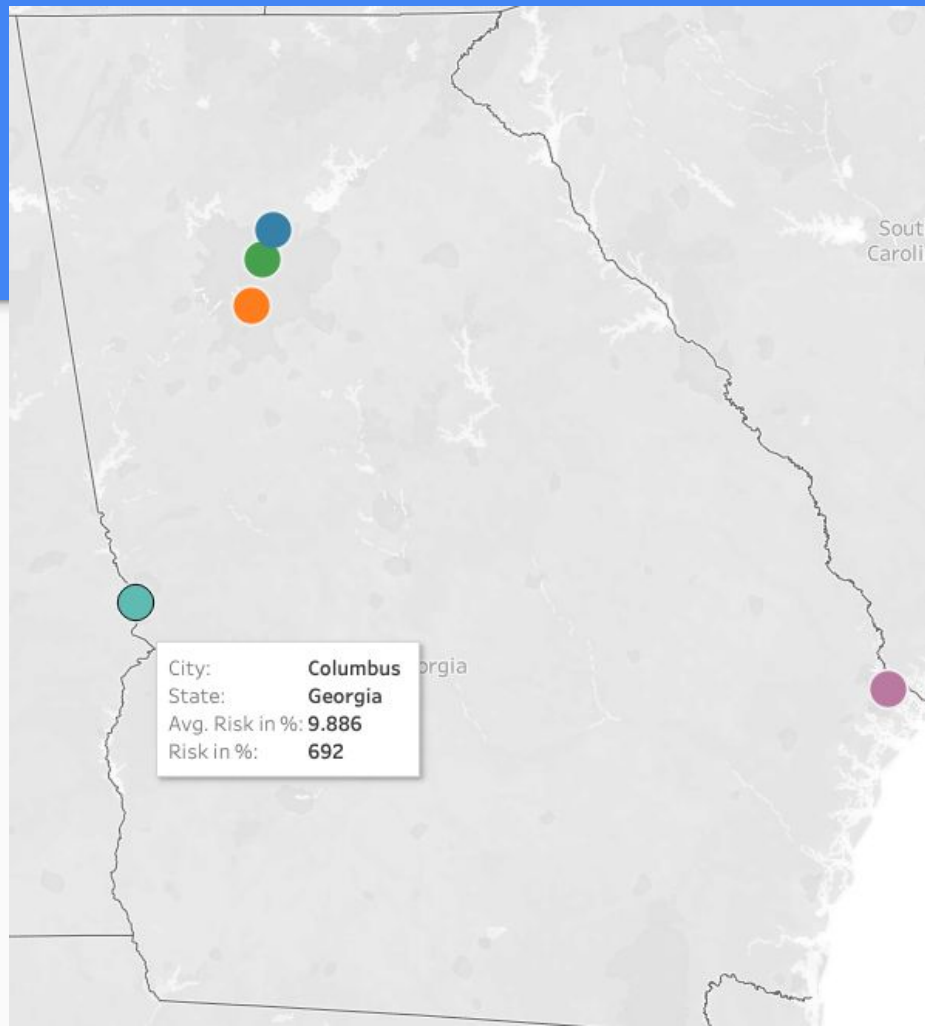
Graph's



Graph's



Map



Run Time

Do to size of the data the run time was fast

Resource

key algorithms/technology:

- Eclipse
- Hadoop
- Linux OS(Ubuntu 20.04)

Language - Java

Which part of Hadoop used in this project?

- HDFS
- MapReduce
- hadoop-common-3.1.2.jar
- hadoop-mapreduce-client-core-3.1.2.jar

data source: <https://data.world/datasets/loan>

<https://www.kaggle.com/zaurbegiev/my-dataset>

<https://www.kaggle.com/panamby/bank-loan-status-dataset>