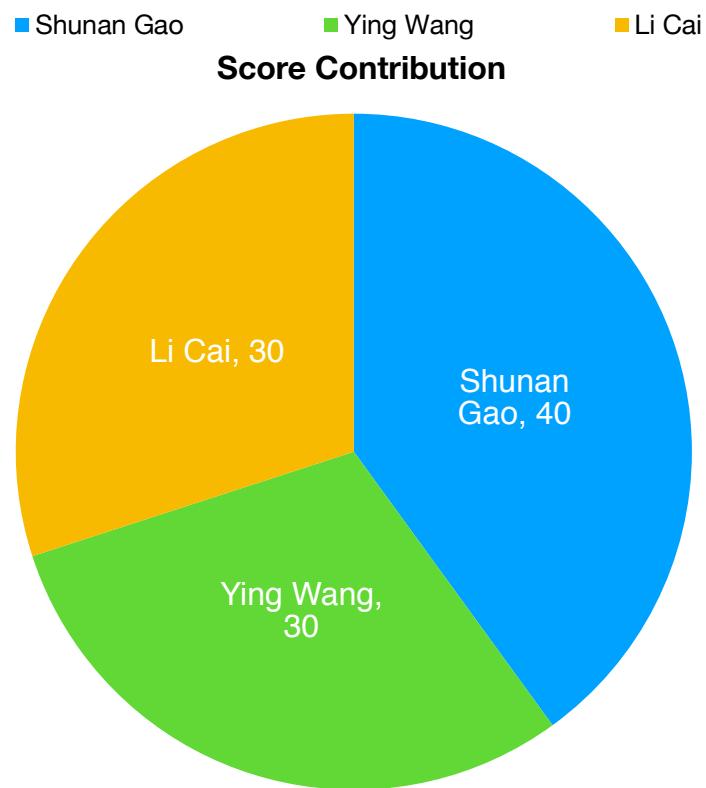


Assignment Report



Part I Data Wrangling:

docker hub link:

<https://hub.docker.com/r/shunangao/hw3-p1/>

Then

docker pull shunangao/hw3-p1

Then

docker run shunangao/hw3-p1

1. Download data from website:

```

USERNAME='gao.shu@husky.neu.edu'
PASSWORD='gj{vx<Io'
print('username='+uname)
print('password='+pwd)

payload = {
    "username": USERNAME,
    "password": PASSWORD
}

session_requests = requests.session()

login_url = "https://freddiemac.embs.com/FLoan/secure/auth.php"

```

To store the user credential, we need to store them in the request session so that user didn't redirect back to the login page whenever he/she required to download a file from the Freddie Mac posted dataset.

To programmatically downloading the file, first the user should create username and password. Once logged in, the user can download all the file required for analysis. We have used the python requests library for this purpose. To store the user credential, we need to store them in the request session so that user didn't redirect back to the login page whenever he/she required to download a file from the Freddie Mac posted dataset.

To store the user credential, we need to store them in the request session so that user didn't redirect back to the login page.

```

result = session_requests.post(
    login_url,
    data = payload,
    headers = dict(referer=login_url)
)
tree = html.fromstring(result.content)
all_links = tree.findall("//a")

```

Once we log in. We use `html.fromstring()` to get the url content and use `findall()` function to get all the link.

```
##### Download zips #####
for link in all_links:
    href=link.get("href")
    if "sample" in href:
        if int(link.text[-8:-4]) >= 2005:
            url= 'https://freddiemac.embs.com/FLoan/Data/' + href
            r = session_requests.get(url, stream=True)
            with open(os.path.join('downloaded_zips',link.text), 'wb') as f:
                for chunk in r.iter_content(chunk_size=1024):
                    if chunk: # filter out keep-alive new chunks
                        f.write(chunk)
```

The data need to be downloaded is in the ‘herf’ part of the link and for this assignment

we need to download all the sample file after 2005.

```
##### Unzip and extract text files #####
try:
    zip_files = os.listdir('downloaded_zips')
    for f in zip_files:
        z = zipfile.ZipFile(os.path.join('downloaded_zips', f), 'r')
        for file in z.namelist():
            if file.endswith('.txt'):
                z.extract(file, r'downloadable_zips_unzipped')
    logging.info('Zip files successfully extracted to folder: downloadable_zips_unzipped.')
except Exception as e:
    logging.error(str(e))
    exit()
```

Unzip all the sample file and save it to a folder

2. Data preprocessing:

After download the data we need to clean the data for each file and save it to a new file.

```
for file in orig_filelists:
    data = pd.read_csv(file, sep='|', names= col_names_orig, low_memory=False)

    data = data.drop(columns = ['super_conforming_flag','pre_harp_loan_sequence_no'])
    # remove rows contain more than 4 NaN
    data.dropna(thresh=4)
    #remove rows doesn't have loan_sequence_no
    data = data[pd.notnull(data['loan_sequence_no'])]
```

For orig file, we first remove some columns which has too many NaN value. Then we will remove the row which has more than 4 missing value and the row don't have loan

```
# replace credit_score with mode
data.credit_score=data.credit_score.replace(r'\s+', np.nan, regex=True).astype('float64')
cs = pd.DataFrame(data['credit_score'])
mode=cs.mode()
data['credit_score'] = data['credit_score'].fillna(mode.iloc[0]['credit_score'])

# replace unknown value of msa with 0000
data.msa=data.msa.replace(r'\s+', np.nan, regex=True)
data['msa'] = data['msa'].fillna('00000').astype('float64')

# replace unknown value of mi_percentage with 0
data.mi_percentage=data.mi_percentage.replace(r'\s+', np.nan, regex=True).astype('float64')
data['mi_percentage'] = data['mi_percentage'].fillna(0)

# replace unknown value of no_of_units with 1
data.no_of_units=data.no_of_units.replace(r'\s+', np.nan, regex=True).astype('float64')
data['no_of_units'] = data['no_of_units'].fillna(1)
```

sequence
number.

```
for file in svcg_filelists:
    data=pd.read_csv(file, sep='|', names= col_names_svcg, low_memory=False)
    #remove rows contain NaN loan_sequence_no
    data = data[pd.notnull(data['loan_sequence_no'])]

    #remove rows contain NaN monthly_report_period
    data = data[pd.notnull(data['monthly_reporting_period'])]

    #remove rows contain NaN current_actual_upb
    data = data[pd.notnull(data['current_actual_upb'])]

    #remove rows contain NaN current_loan_delinquency_status
    data = data[pd.notnull(data['current_loan_delinquency_status'])]

    #remove rows contain NaN loan_age
    data = data[pd.notnull(data['loan_age'])]

    #remove rows contain NaN remaning_months_on_legal_maturity
    data = data[pd.notnull(data['remaning_months_on_legal_maturity'])]
```

After that we will fill some missing value with the appropriate value.

For svcg file. First remove rows which has the NaN value in some columns which can not be none.

```

#remove cols which has too many NaN
data = data.drop(columns = ['repurchase_flag','modification_flag','zero_bal_eff_date','ddlpi',
'mi_recoveries', 'net_sales_proceeds','non_mi_recoveries','expenses','legal_costs',
'maintenance_preservation_cost','taxes_insurance','misc_expenses','actual_loss_calc',
'modification_cost','step_modification_flag','deferred_payment_modification',])

#replace unknown with 0
data.estimated_loan_to_value=data.estimated_loan_to_value.replace(r'\s+', np.nan, regex=True)
data['estimated_loan_to_value'] = data['estimated_loan_to_value'].fillna('0').astype('float64')

```

And the we will remove some columns which has too many NaN data .

All the cleaned orig and svvg files will be save to the new folder for further processing.

Summary File:

To summarize the orig file, we create a new data frame and add some columns. The new columns is to store the summary information of the orig file and for each year is

```

summ_df['year'] = file[-8:-4]
summ_df["aveCreditScore"] = data["credit_score"].mean()
summ_df["loanCount"] = np.count_nonzero(data["loan_sequence_no"])
summ_df["aveMI"] = data["mi_percentage"].mean()
summ_df["aveUnitNumber"] = data["no_of_units"].mean()
summ_df["aveCLTV"] = data["original_cltv"].mean()
summ_df["aveDTI"] = data["original_dti_ratio"].mean()
summ_df["aveLTV"] = data["original_ltv"].mean()
summ_df["totalUPB"] = data["original_upb"].sum()
summ_df["aveUPB"] = data["original_upb"].mean()
summ_df["aveInterestRate"] = data["original_interest_rate"].mean()
summ_df["mostFrequentState"] = pd.DataFrame(data.groupby('property_state').size().rename('cnt')).idxmax()[0]
summ_df["mostFrequentType"] = pd.DataFrame(data.groupby('property_type').size().rename('cnt')).idxmax()[0]
summ_df["mostFrequentPostalCode"] = pd.DataFrame(data.groupby('postal_code').size().rename('cnt')).idxmax()[0]
summ_df["aveNumberOfBorrowers"] = data["no_of_borrowers"].mean()
summ_df["mostFrequentSeller"] = pd.DataFrame(data.groupby('seller_name').size().rename('cnt')).idxmax()[0]
summ_df["mostFrequentServicer"] = pd.DataFrame(data.groupby('servicer_name').size().rename('cnt')).idxmax()[0]

```

one row. We combine all the year and the save the df to a new csv file.

```

summ_df = OrderedDict()
summ_df['year'] = file[-8:-4]
total_loans = data['loan_sequence_no'].nunique()
summ_df['no_distinct_loans_per_year'] = data['loan_sequence_no'].nunique()
summ_df['most_length_loan_months'] = data.groupby('loan_sequence_no').size().max()
total_reos = data[data['current_loan_delinquency_status'].astype(str) == 'R']['Loan_sequence_no'].nunique()
summ_df['total_reos_per_year_percent'] = 100 * total_reos/total_loans
summ_df['mean_upb_year'] = data['current_actual_upb'].mean()
summ_df['avg_loan_age'] = data['loan_age'].mean()
summ_df['avg_current_interest_rate'] = data['current_interest_rate'].mean()

```

And also the svvg file:

These two summary csv file will also be stored in a new folder.

```
concatDf = pd.DataFrame(columns=colNames)
concatDf = pd.concat(dataList ,axis = 0)
concatDf.to_csv(os.path.join('cleanFilesWithSummaries',file_name),index=False)
```

Then
we
com-
bine
all

the cleaned csv file together:

Exploratory Data Analysis:

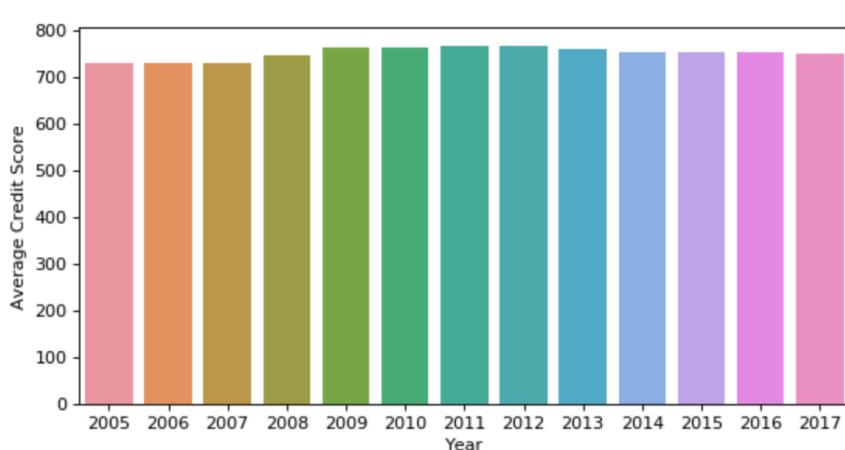
After finishing the data processing we will use a jupyter notebook to visualize the data.

Origination Summary:

```
orig = orig.sort_values(by=[ 'year' ])
orig.head()
```

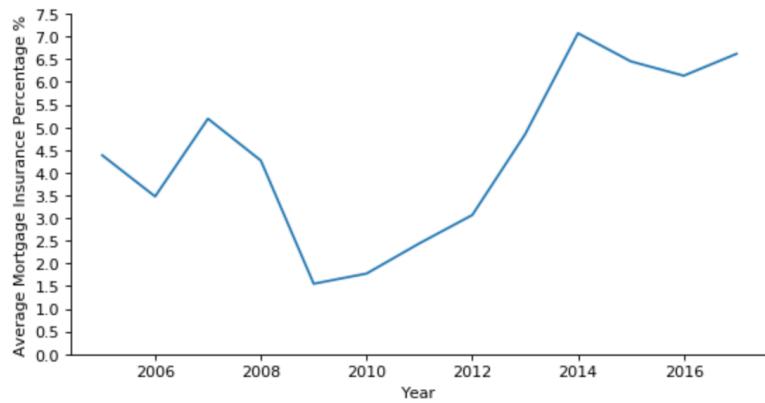
	year	aveCreditScore	loanCount	aveMI	aveUnitNumber	aveCLTV	aveDTI	aveLTV	aveFICO
7	2005	728.49904	50000	4.38436	1.02644	71.21610	65.38164	69.56804	710.0
12	2006	729.95500	50000	3.47626	1.02774	73.13102	57.58518	70.71394	710.0
9	2007	728.46154	50000	5.19054	1.03666	74.53890	61.62562	72.10360	710.0
0	2008	746.36954	50000	4.26810	1.03666	71.50462	57.18436	70.28216	710.0
1	2009	762.33614	50000	1.55242	1.01600	66.85308	34.21672	65.44966	710.0

The origination file
looks like this:



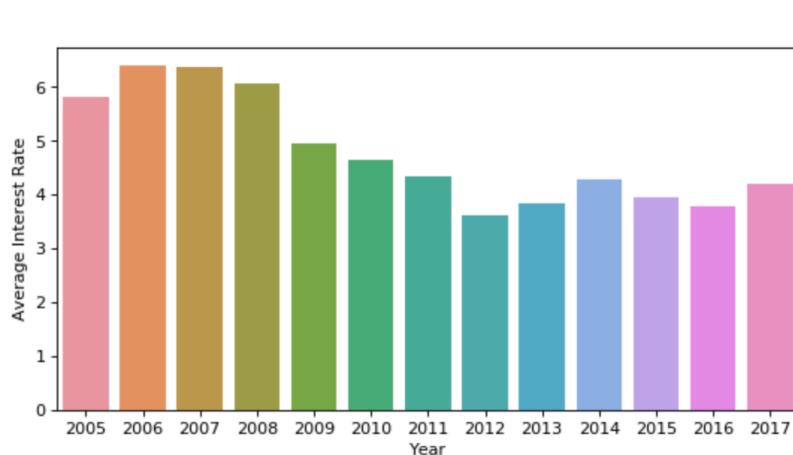
The average credit
score vs year plot:

we can see that average credit score doesn't change that much.



Average mortgage insurance percentage vs year plot:

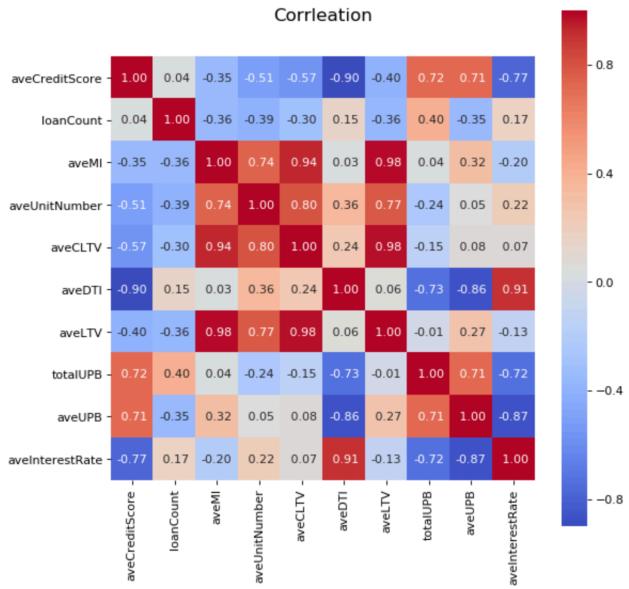
we can see the average MI lowest is 2009 and the peak value is 2014



Average interest rate:

we can see the average interest rate peak value is 2006 and the bottom value is 2012

Correlation:

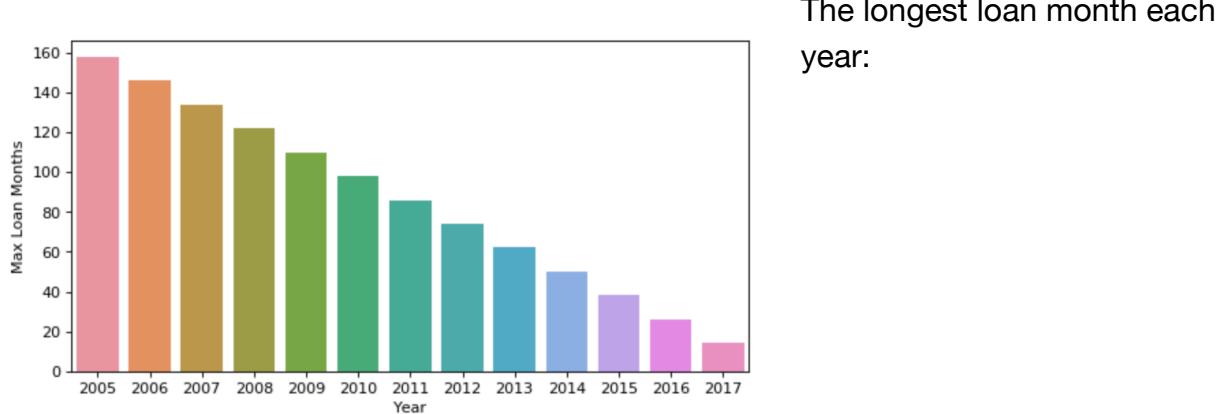


We can see that some of the attribute has the correlation such as credit score and DTI

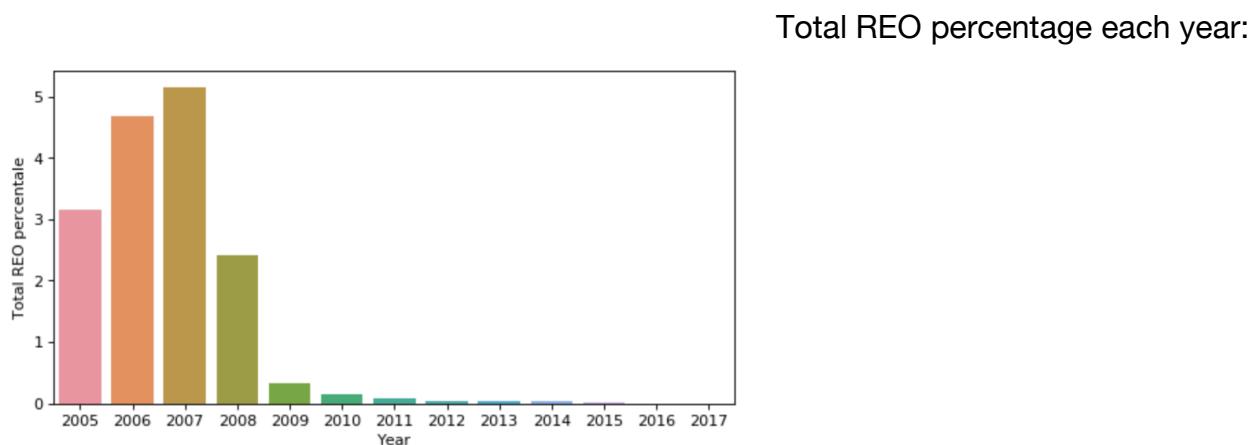
```
svcg = svcg.sort_values(by=[ 'year' ])
svcg.head()
```

	year	no_distinct_loans_per_year	most_lenth_loan_months	total_reos_per_year_percent	mean_upb_year
8	2005	49999	158	3.144063	147470.359979
5	2006	49997	146	4.670280	157320.881169
6	2007	49996	134	5.154412	160294.361364
12	2008	49995	122	2.416242	169830.639024
11	2009	49998	110	0.322013	175750.124093

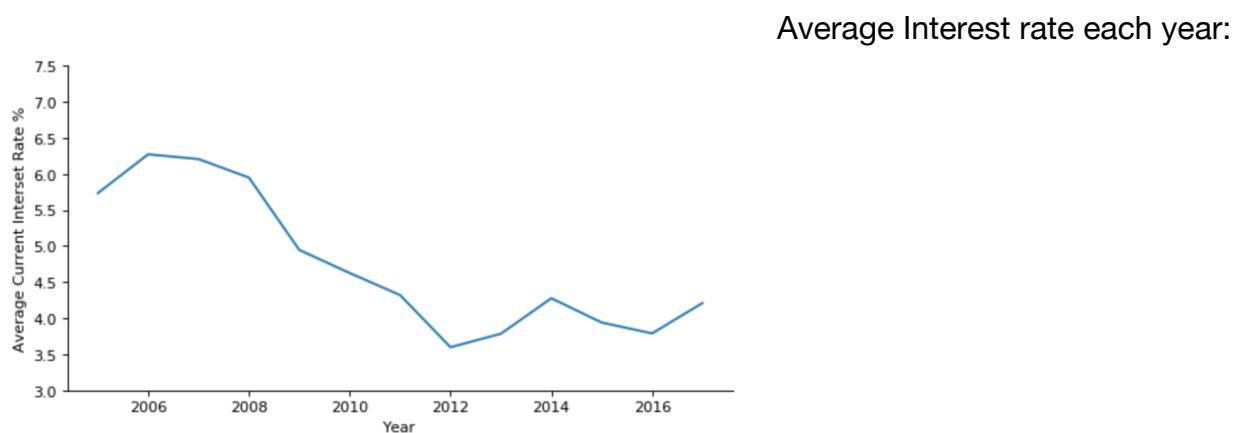
Performance Summary:



The longest loan month is decreasing because the more recent the more less

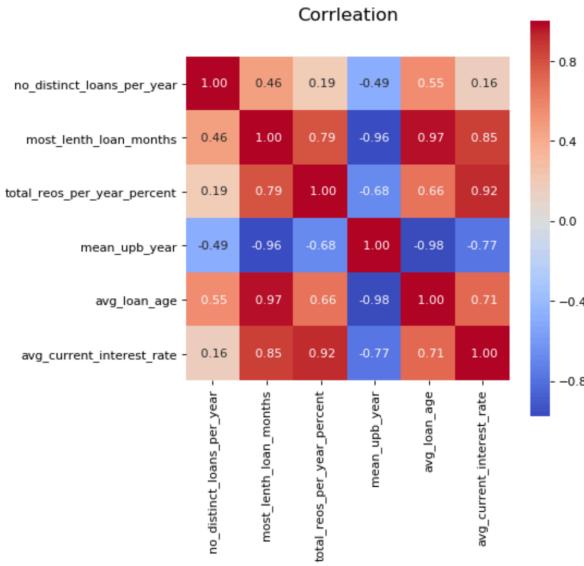


The total REO reach the peak value at 2007 and after 2008 it decrease significantly



The average interest rate is totally decreasing by year

Correlation:



We can see that some attribute have strong relation such as mean upb and loan age

Quarterly Analysis 2007:

```
orig_2007.head()
```

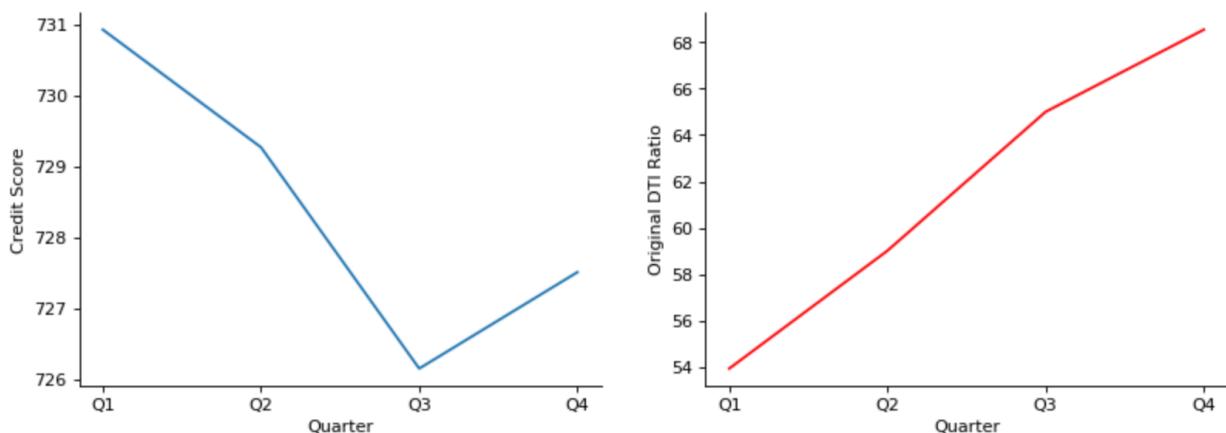
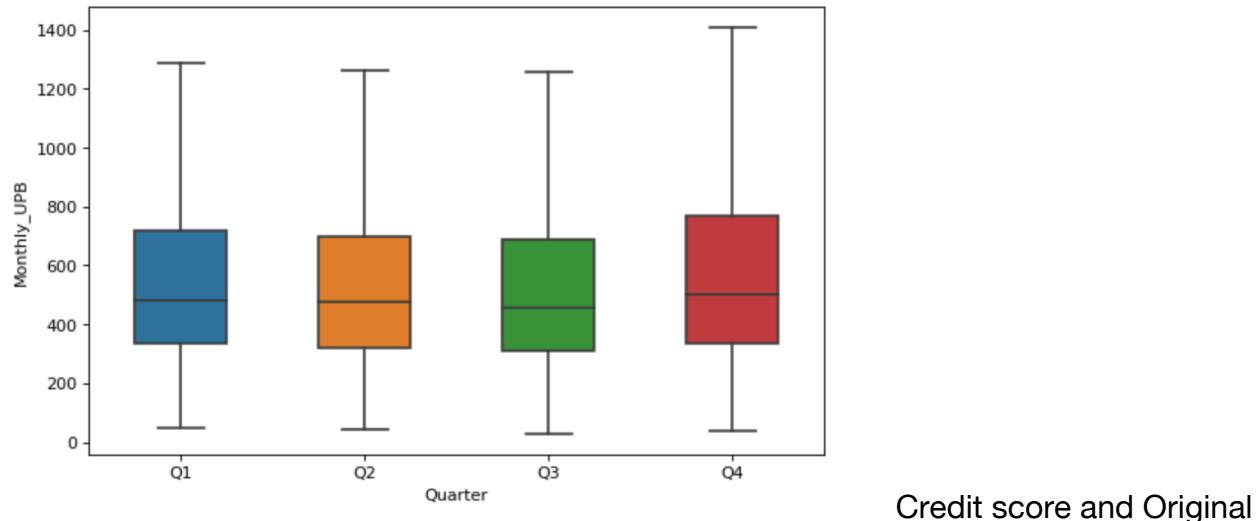
	credit_score	first_pay_date	first_time_homebuyer	maturity_date	msa	mi_percentage	no_of_units
0	606.0	200703	N	203702	0.0	0.0	1.0
1	712.0	200703	N	203702	0.0	0.0	1.0
2	698.0	200704	N	203703	24020.0	0.0	1.0
3	757.0	200703	N	203702	0.0	0.0	1.0
4	761.0	200703	N	203702	24660.0	0.0	1.0

Origination
file:

orig_2007.csv

2007 UPB Quarterly Analysis

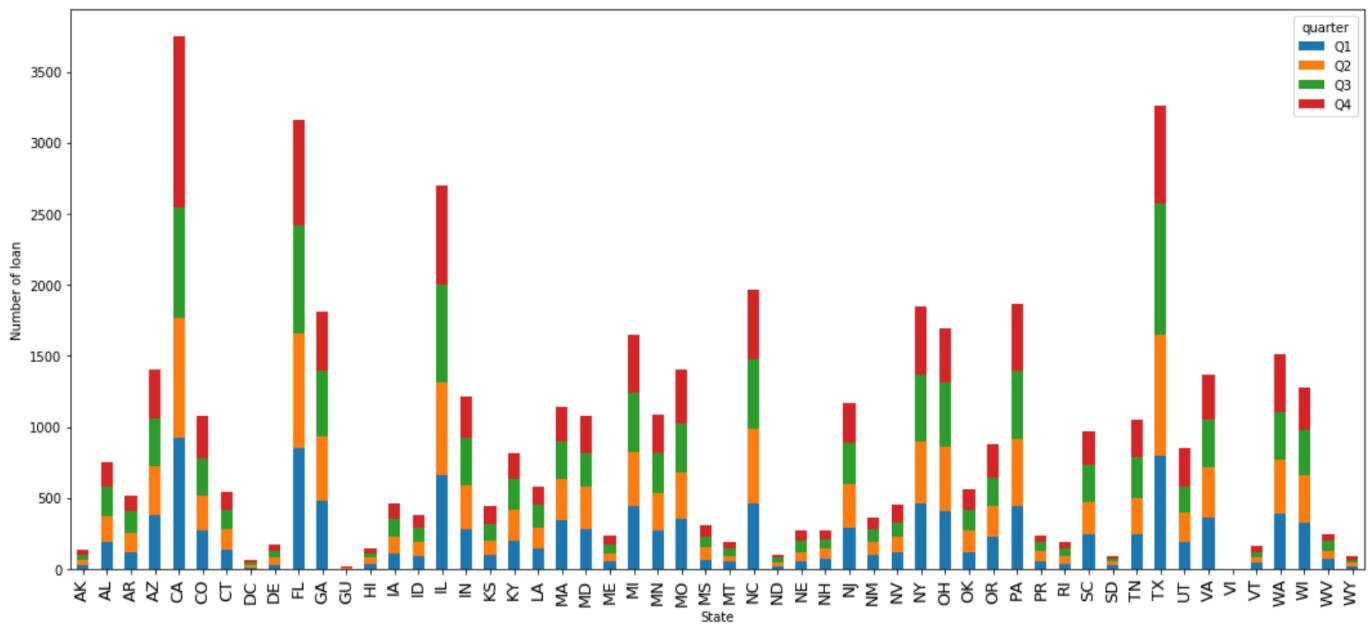
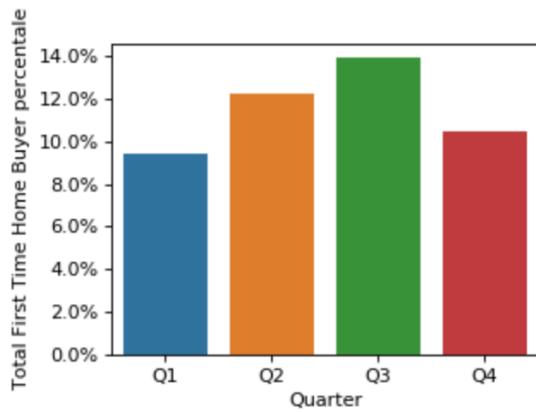
UPB Quarterly distribution: pretty the same



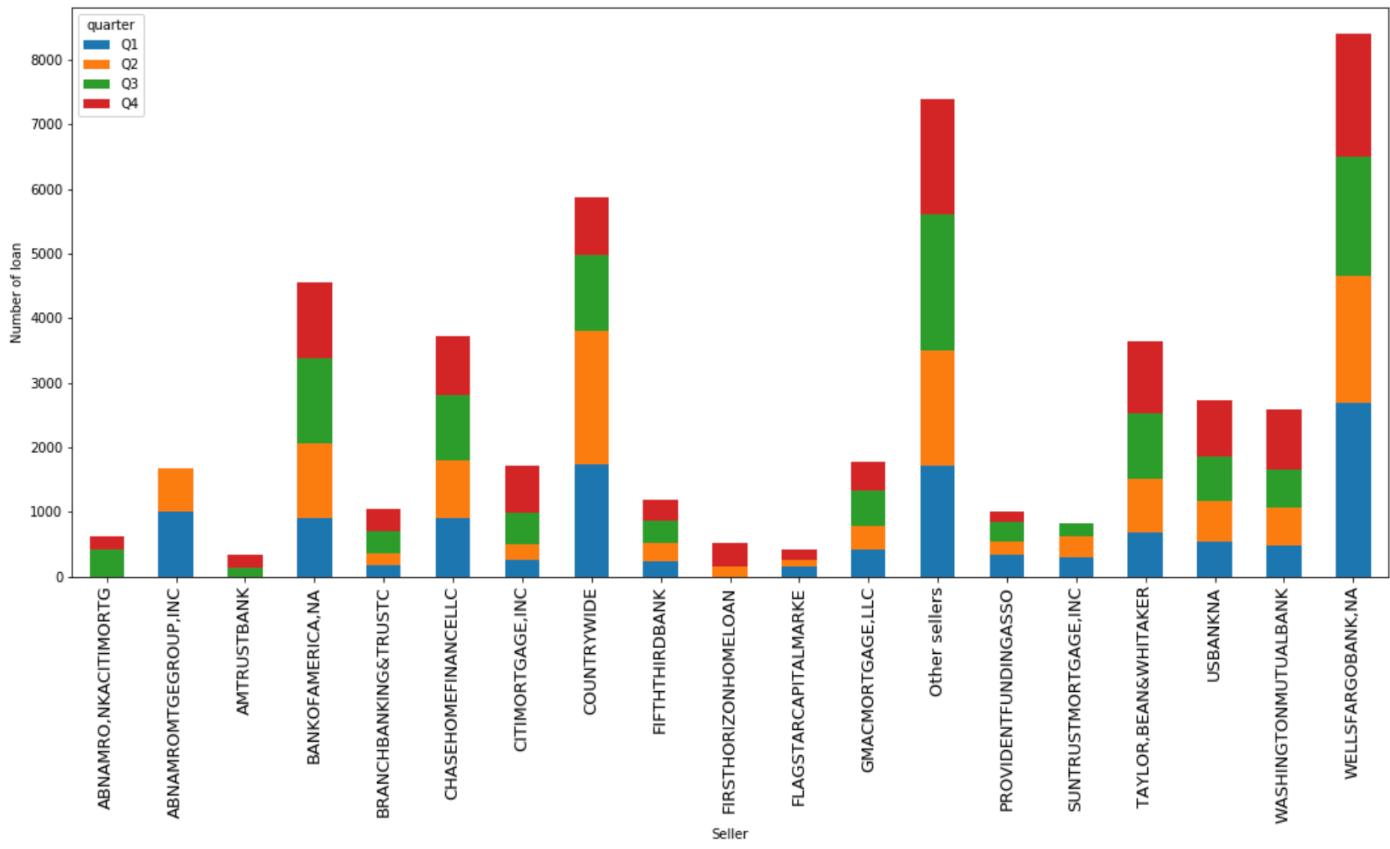
DTI Ratio:

The credit score decreased while DTI ratio is increased

First time home buyer percentage each quarter: Q3 have the most first time buyers



Number of loans in each state quarterly: totally CA no.1 but Q3 TX beaten CA

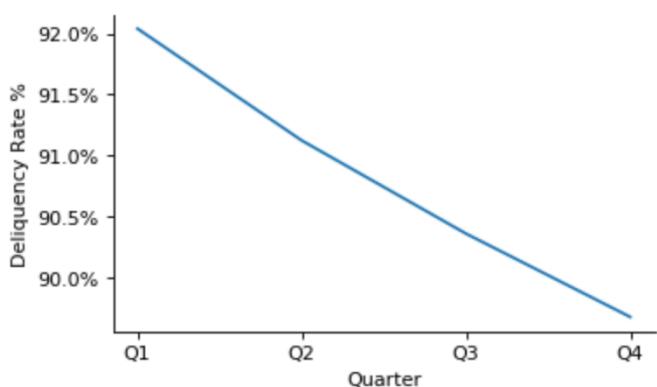


Number of loans from different seller quarterly: Some seller may sell 0 loan in a quarter

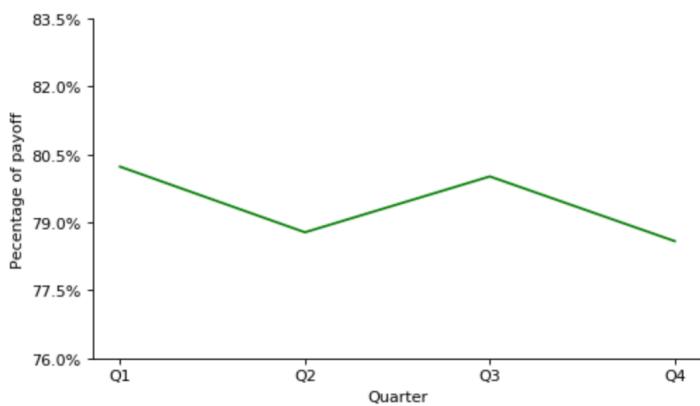
`svcg_2007.head()`

	loan_sequence_no	monthly_reporting_period	current_actual_upb	current_loan_delinquency_status
0	F107Q1000009	200702	110000.0	0
1	F107Q1000009	200703	110000.0	0
2	F107Q1000009	200704	110000.0	0
3	F107Q1000009	200705	110000.0	0
4	F107Q1000009	200706	110000.0	0

Performance file: `svcg_2007.csv`

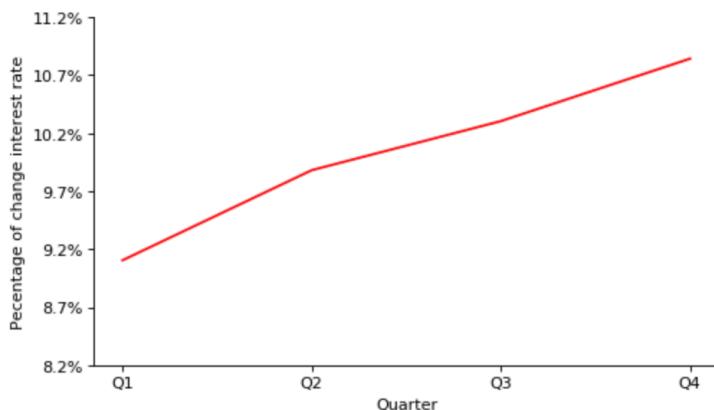


Percentage of delinquency status:
decreasing but not significantly



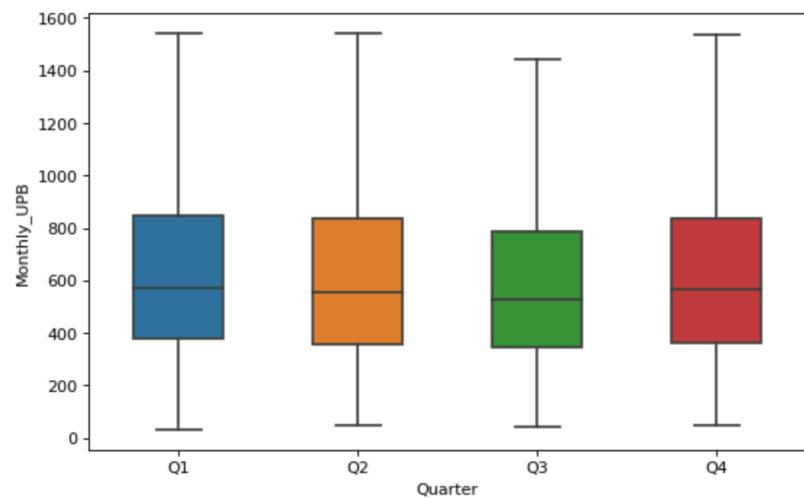
Percentage of loans which is paid
off : stable between 79-81 per-
cent

Percentage of loan which interest rate has changed: increasing

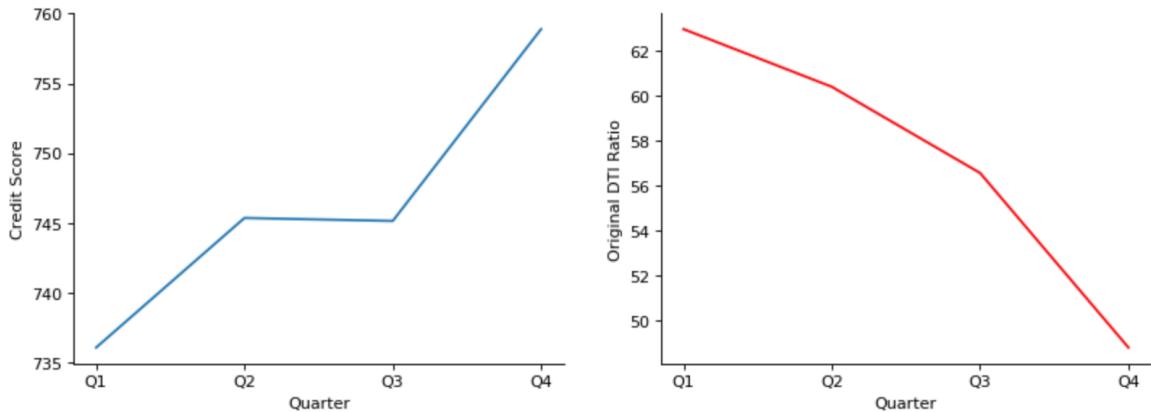


Quarterly Analysis 2008:

2008 UPB Quarterly Analysis

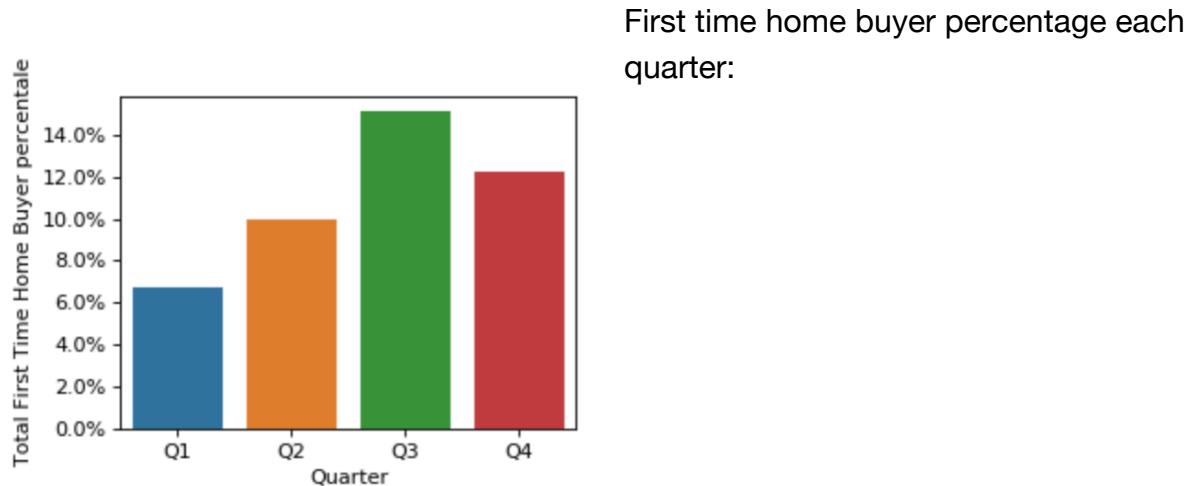


The un paid balance is stable for each quarter

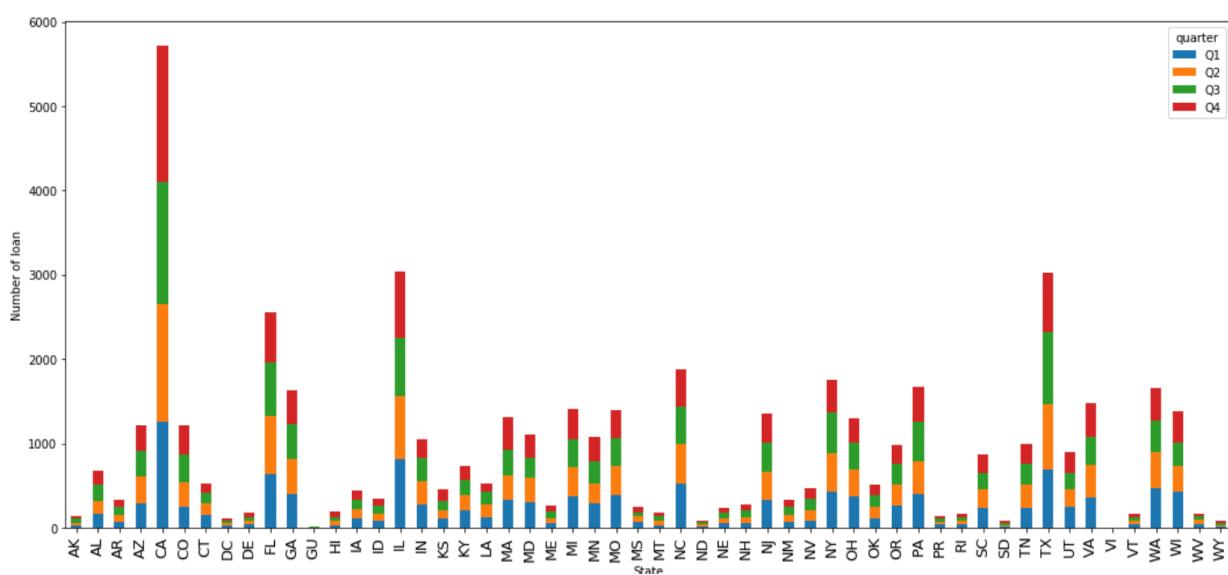


Credit score and Original DTI Ratio:

Different from 2007 the credit score in 2008 is increasing quarterly and the DTI ratio is decreasing so these two attribute may have negative correlation.

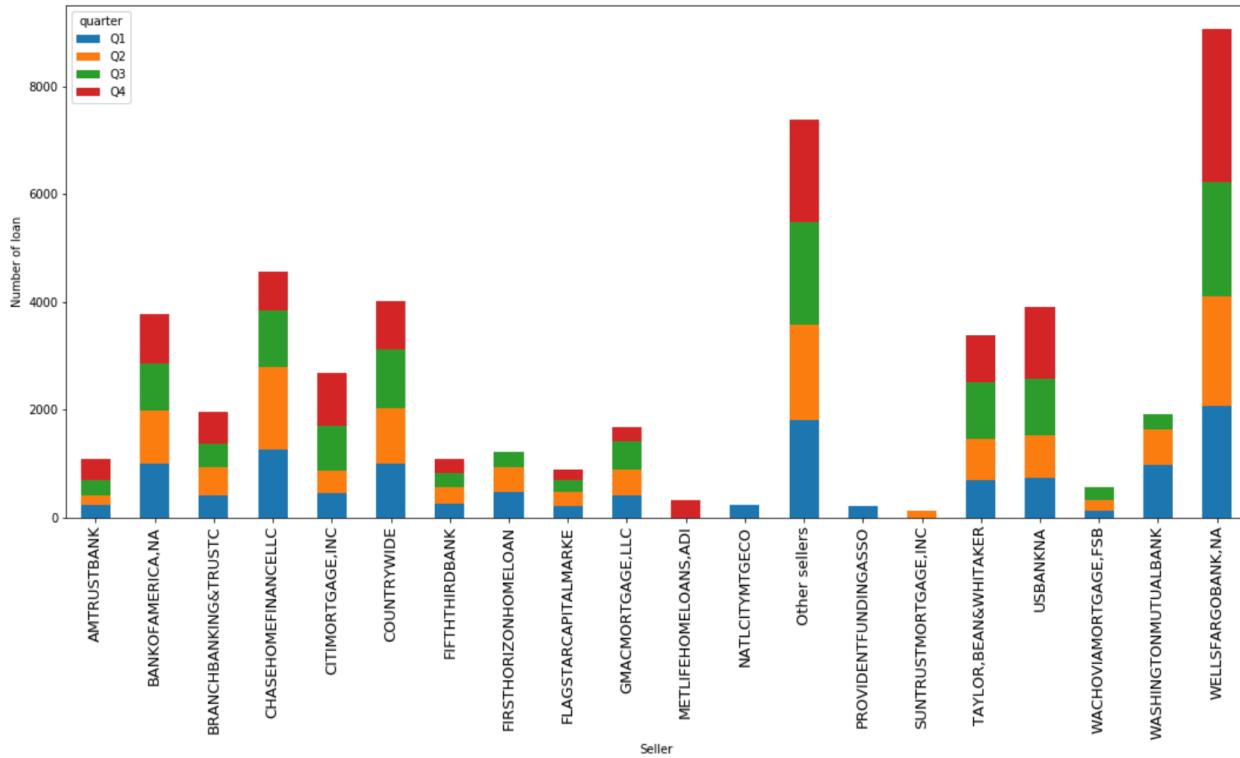


Same with 2007 the Q3 also has the most first time home buyers



Number of loans in each state quarterly:

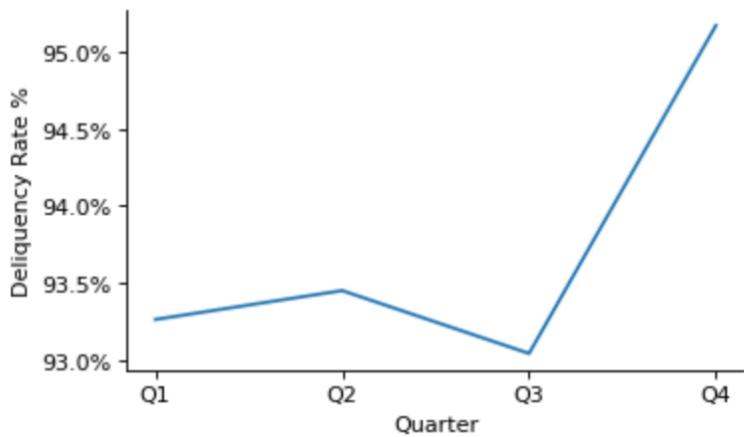
CA also the NO.1 state but in 2008, every quarter the CA have the largest loan number.



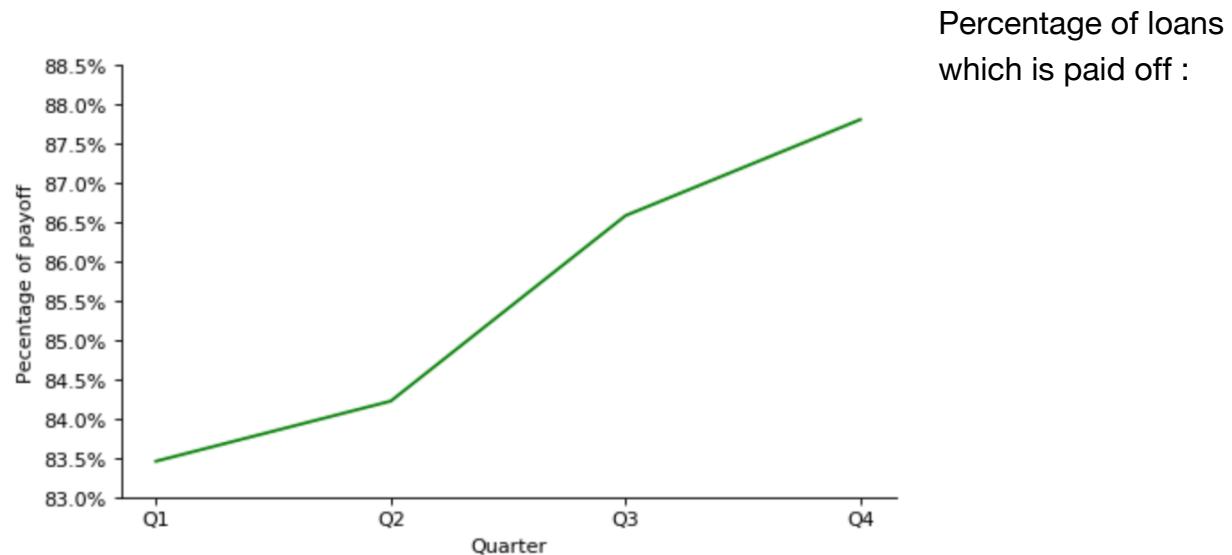
Number of loans from different seller quarterly:

We can see that some seller have 0 sales in a quarter and Wells Fargo also like 2007 is the largest seller

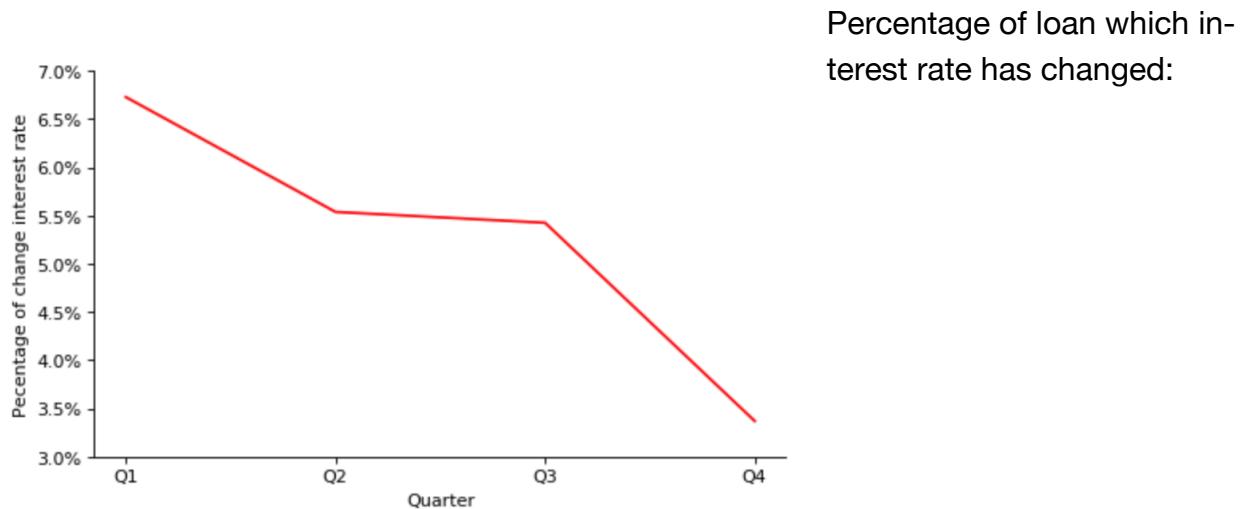
Percentage of delinquency status:



While Q3 has the lowest delinquency rate but it's still higher than the highest delinquency rate Q1 in 2007

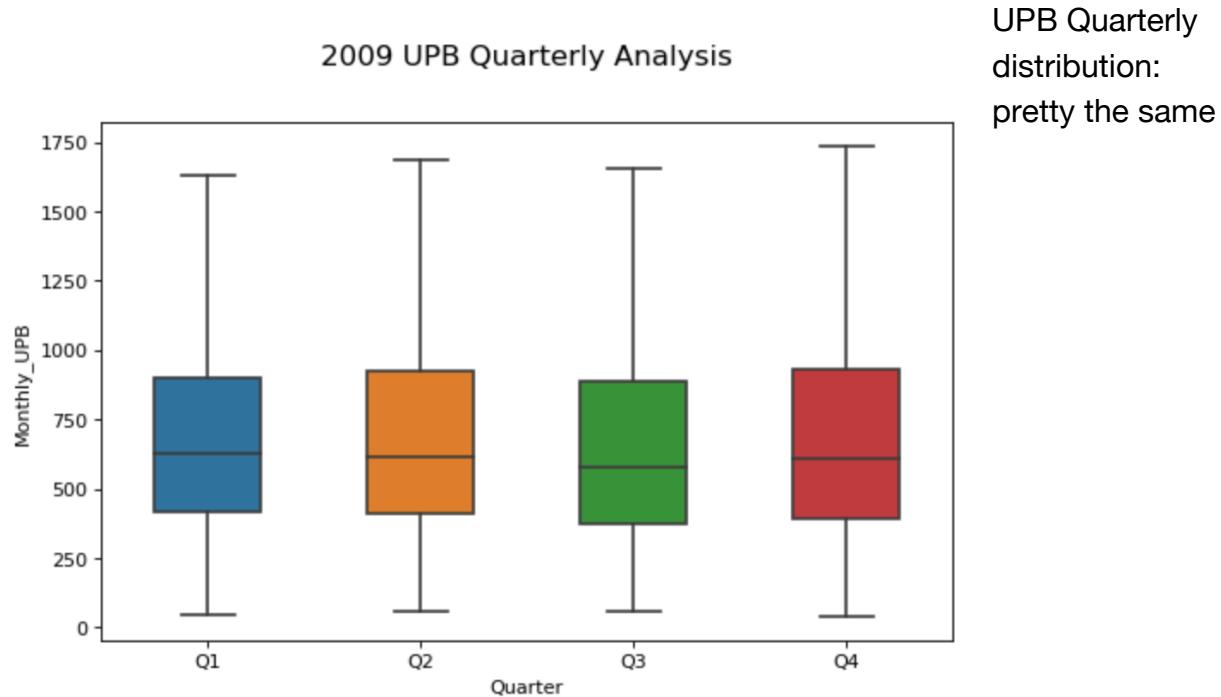


As we can see the percentage of payoff is increasing and higher than 2007

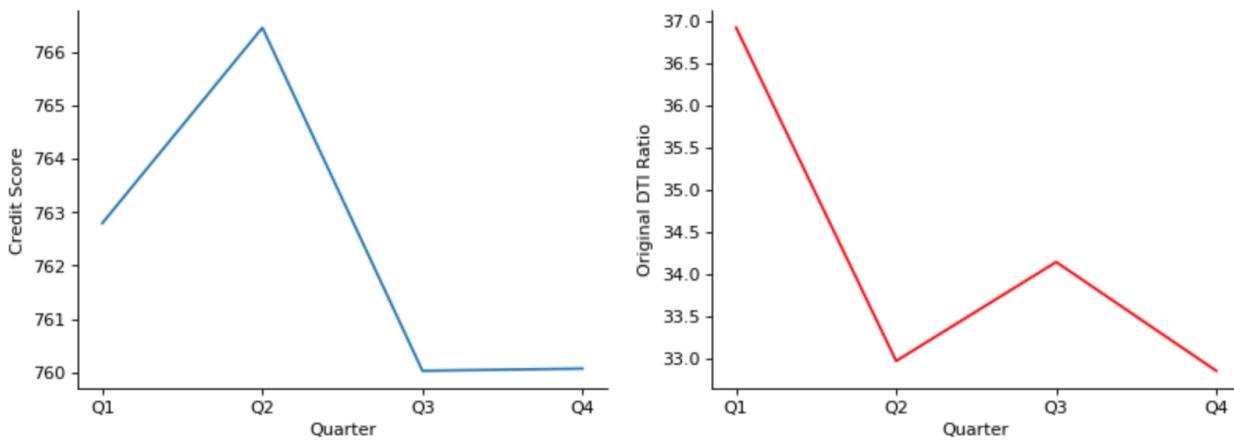


As we can see the changing interest rate number is decreasing from Q1-Q4 which is very different from 2007 and the peak value Q1 is lower than the bottom value Q1 in 2007.

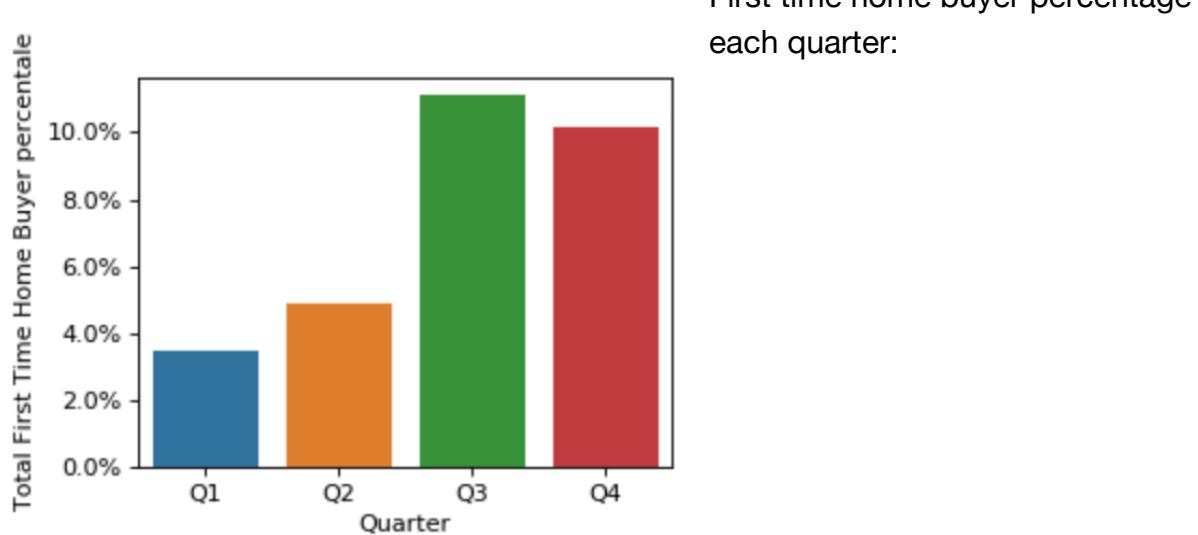
Quarterly Analysis 2009:



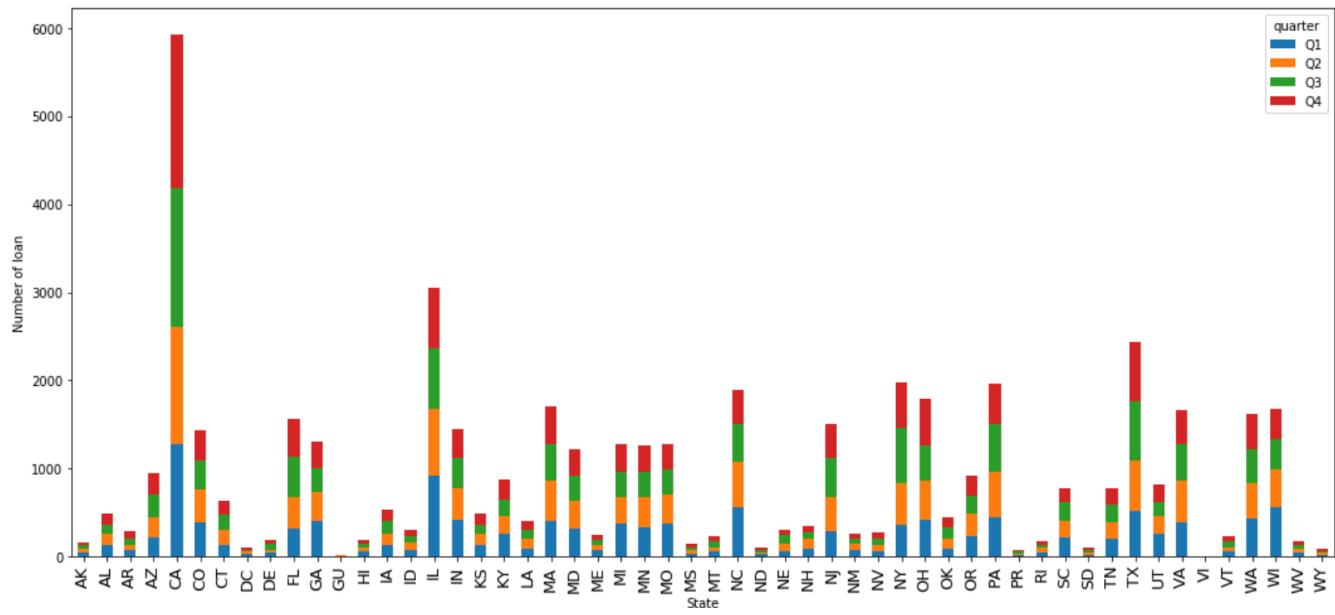
Credit score and Original DTI Ratio:



The credit score reached at the peak value at Q2 while the DTI reached the bottom value at Q2, and when the credit score increase the DTI ratio decrease. While Q1-Q2 the DTI ratio decreased significantly and the credit score not increased that significantly.

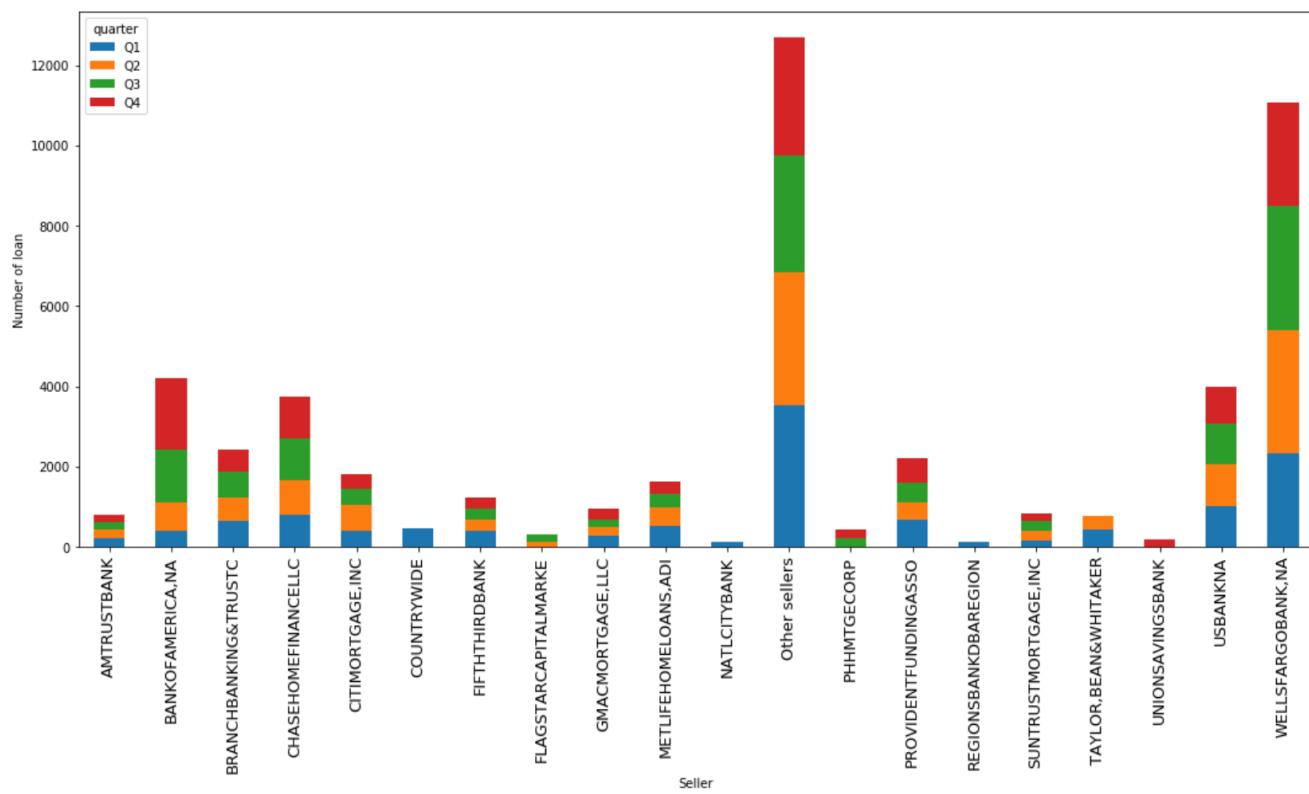


Same with 2007 and 2008, Q3 also has the most first time home buyers.



Number of loans in each state quarterly:

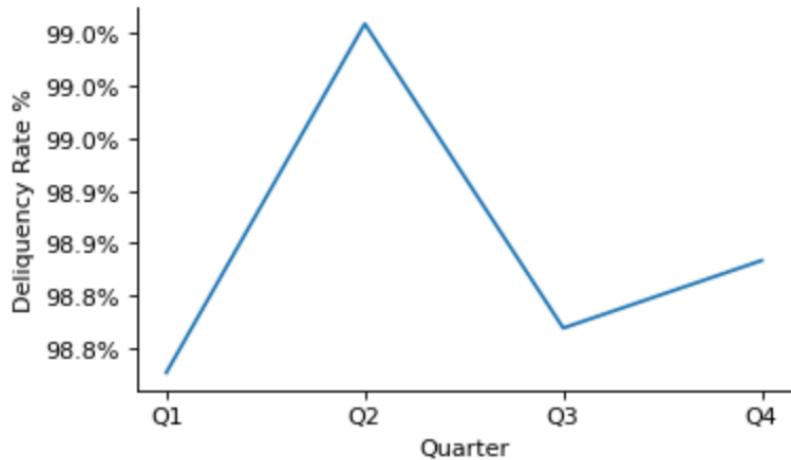
CA also the NO.1 state and every quarter the CA have the largest loan number.



Number of loans from different seller quarterly:

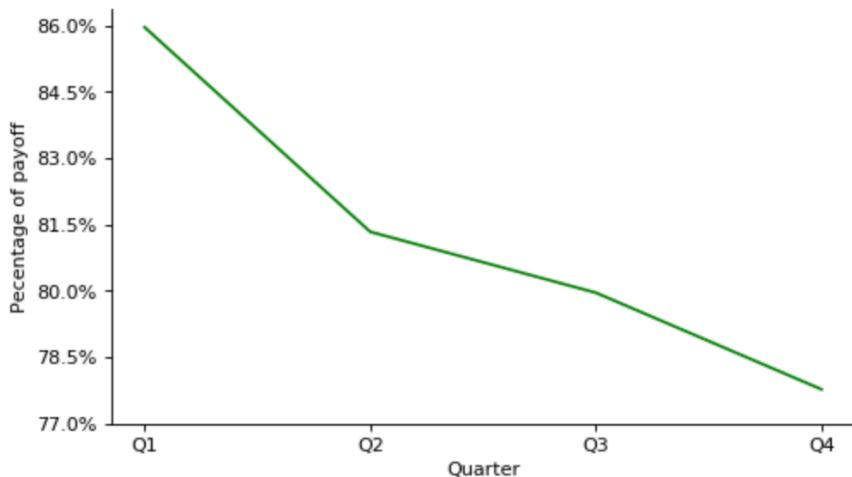
Different from 2007 and 2008 the largest seller is “other sellers” and for each quarter the other seller have the largest sales.

Percentage of delinquency status:

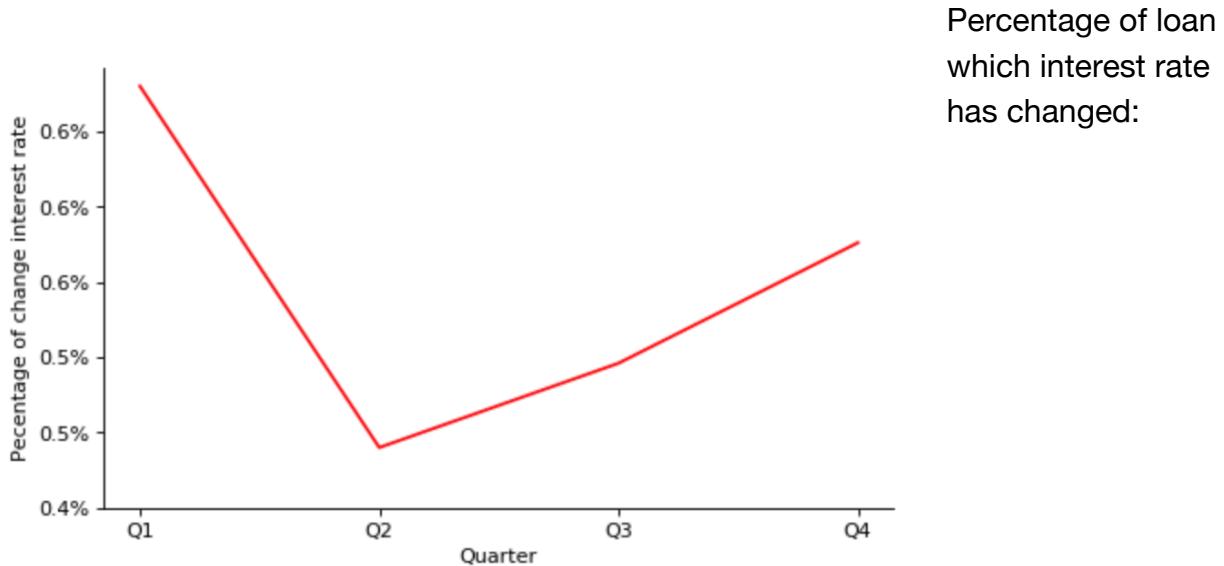


The delinquency rate is above all the 2008 and it change not very significantly only between 98-99 percent.

Percentage of loans which is paid off :



The payoff percentage is different form 2008. it's decreasing about 10 percent which is huge compared to 2007,2008.



The percentage of loan changes is pretty low compared to 2007 and 2008. only 0.4-0.6 percent.

Part II: Prediction

1. Download & pre-process the data :

Run a python script, in which includes the following major functions:

a. Data download. Parameters may passed to the function to download whatever file

```
def get_data_from_url(quarter):
    print('downloading...')
    urllib.request.urlretrieve('https://freddiemac.embs.com/FLoan/Data/historical_data1_' + str(quarter) + '.zip',
                                DIR_NAME + str(quarter) + '.zip')
    # unzip_files(year)
    try:
        zip_ref = zipfile.ZipFile(DIR_NAME + str(quarter) + '.zip', 'r')
        zip_ref.extractall(DIR_NAME)
        zip_ref.close()

        write_into_csv(quarter)
    except zipfile.BadZipfile:
        print (zipfile.BadZipfile)
```

requested.

```
def data_cleaning():
    fileList = fileList = glob.glob('data/*.csv')
    for file in fileList:
        df = pd.read_csv(file,error_bad_lines=False)
        print('categorical cleaning...')
        for col in ['C','G','H','M','N','O','P','Q','T','U','V','W','X']:
            mode = pd.DataFrame(df.groupby(col).size().rename('cnt')).idxmax()[0]
            df[col] = df[col].fillna(mode)
        print('numerical cleaning...')
        for col in ['E','F','I','J','K','L','R']:
            dfmean = df[(df[col] != 999)|(df[col] != None)]
            mean = int(dfmean[col].mean(axis=0))
            df[col] = df[col].fillna(mean)

    filename = "%sclean.csv"%file[:-4]
    df.to_csv(filename)

    print('cleaning finished.')
```

b. Data cleaning. The cleansed data is written to a csv file with renamed columns.

c. Data processing. The columns are processed as dummy values to be used against prediction.

d. Main function as the entrance point. In this function the program logs in and and

```
def start_execution():
    with requests.Session() as sess:
        sess.get(login_page_url);
        php_session_cookie = sess.cookies['PHPSESSID']
        login_payload = {'username' : USERNAME, 'password' : PASSWORD, 'cookie':php_session_cookie}
        sess.post(login_page_url, data = login_payload)
        download_page_payload = {'accept': 'Yes', 'action': 'acceptTandC', 'acceptSubmit': 'Continue', 'cookie': php_session_cookie}
        sess.post(download_page_url, data=download_page_payload)
        create_directory(DIR_NAME)
        # create_csv()

        get_data_from_url(START)
        get_data_from_url("Q12003")
        get_data_from_url("Q12005")
        get_data_from_url("Q12007")
        get_data_from_url("Q12009")
        get_data_from_url(END)

        get_data_from_url("Q22005")
        get_data_from_url("Q32005")
        get_data_from_url("Q42005")

        get_data_from_url("Q22007")
        get_data_from_url("Q32007")
        get_data_from_url("Q42007")

        data_cleaning()
        preprocessing()
```

uses the defined function.

Please note that all output files are put into a data/ directory.

2. Local Analysis

```
selector = SelectKBest(f_regression, k=20).fit(X_train,y_train)
k_best_features = X_train.columns.values[selector.get_support()]
```

According to:

```
k_best_features = [ 'D', 'F', 'I', 'K', 'L', 'U', 'C_Y', 'H_I', 'H_P', 'Q_MH', 'T_P',
'W_BANKOFAMERICA_NA', 'W_COUNTRYWIDE', 'W_GMACMTGECORP', 'W_NATLCITYMTGECO',
'W_WASHINGTONMUTUALBANK', 'X_BANKOFAMERICA_NA', 'X_COUNTRYWIDE',
'X_GMACMORTGAGE_LLC', 'X_NATLCITYMTGECO' ]
```

The program can select the best features for further analysis:

```

#Linear, random forestk, and neural network models
from sklearn.linear_model import LinearRegression

linear = LinearRegression()
linear.fit(X_train,y_train)
linearPredict = linear.predict(X_test)
printPerformance(linearPredict)

from sklearn.ensemble import RandomForestRegressor
randomForest = RandomForestRegressor(max_depth= 18, n_estimators = 16, random_state=2)
randomForest.fit(X_train,y_train)
randomForestPredict = randomForest.predict(X_test)
printPerformance(randomForestPredict)

from sklearn.neural_network import MLPRegressor
neuralNetwork = MLPRegressor()
neuralNetwork.fit(X_train,y_train)
neuralNetworkPredict = neuralNetwork.predict(X_test)
printPerformance(neuralNetworkPredict)

# Build RF classifier to use in feature selection
from sklearn.ensemble import RandomForestClassifier
from mlxtend.feature_selection import SequentialFeatureSelector as sfs
clf = RandomForestClassifier(n_estimators=100, n_jobs=-1)

```

These features are used against the following required algorithms and functions. For each algorithm, the program can also provide stepwise and exhaustive feature selection approaches.

```

#Forward Selection
fs = sfs(clf,
          k_features=9,
          forward=False,
          floating=False,
          n_jobs=-1,
          verbose=2,
          scoring='neg_mean_absolute_error',
          cv=10)
fs.fit(X_train,y_train)
print('Best MAE score: %.2f' % fs.best_score_* (-1))
print('Best subset:', fs.best_feature_names_)

```

```
#Backward Selection
bs = sfs(clf,
          k_features=9,
          forward=False,
          floating=False,
          n_jobs=-1,
          verbose=2,
          scoring='neg_mean_absolute_error',
          cv=10)
bs.fit(X_train,y_train)
```

```
#Exhaustive Selection
from mlxtend.feature_selection import ExhaustiveFeatureSelector as efs
es = efs(clf,
          min_features=8,
          max_features=11,
          scoring='neg_mean_absolute_error',
          n_jobs=-1,
          print_progress=True,
          cv=8)
es.fit(X_train,y_train_encoded)
print('Best MAE score: %.2f' % es.best_score_ * (-1))
print('Best subset:', es.best_feature_names_)
```

```
#Validation
from sklearn.model_selection import cross_val_score
scores = cross_val_score(linear, X_test, y_test, cv=3)
print(scores)
scores = cross_val_score(randomForest, X_test, y_test, cv=3)
print(scores)
scores = cross_val_score(neuralNetwork, X_test, y_test, cv=3)
print(scores)
```

...and then validates against each algorithm.

3. AutoML Usages

```
#tpot
from tpot import TPOTRegressor

pipeline_optimizer = TPOTRegressor(generations=5, population_size=20, cv=5,
                                   random_state=42, verbosity=2)
pipeline_optimizer.fit(X_train, y_train)
print(pipeline_optimizer.score(X_test, y_test))
pipeline_optimizer.export('tpot_exported_pipeline.py')
```

a.Run a TPOT python script like:

```
#Pipeline Exported
from sklearn.linear_model import LassoLarsCV
from sklearn.model_selection import train_test_split

# NOTE: Make sure that the class is labeled 'target' in the data file
tpot_data = df
features = tpot_data.drop('0', axis=1).values
training_features, testing_features, training_target, testing_target = \
    train_test_split(features, tpot_data['0'].values, random_state=42)

# Average CV score on the training set was:-6.254833533625527e-26
exported_pipeline = LassoLarsCV(normalize=True)

exported_pipeline.fit(training_features, training_target)
results = exported_pipeline.predict(testing_features)
```

...and it gives the following pipeline:

The program checks its performance by running the following function.

```

from sklearn.metrics import mean_absolute_error, mean_squared_error
from math import sqrt
def mean_absolute_percentage_error(y_true, y_pred):
    y_true, y_pred = np.array(y_true), np.array(y_pred)
    return np.mean(np.abs((y_true - y_pred) / y_true)) * 100

def printPerformance(pred):
    print(pred)
    print("RMSE: %.2f"
          % sqrt(mean_squared_error(y_test, pred)))
    print("MAPE: %.2f"
          % mean_absolute_percentage_error(y_test, pred)+'%')
    print("MAE: %.2f"
          % mean_absolute_error(y_test, pred))

```

b. Run a H2o jupyter notebook. Please note that this step requires your credentials after mounting driverless AI in docker.

```

import h2o
from h2o.automl import H2OAutoML
h2o.init()

h2oq12005 = h2o.import_file("data/Q12005.csv")
h2oq22005 = h2o.import_file("data/Q22005.csv")

X = h2oq22005.columns.remove("0")
y = "0"
h2oML = H2OAutoML(max_models = 30, max_runtime_secs=300, seed = 1)
h2oML.train(x= X, y = y, training_frame= h2oq12005, leaderboard_frame= h2oq22005)

aml.leaderboard.as_data_frame()

```

less AI in docker.

4. What-If analysis:

For this part, read different data frames as training sets and see the performance of each.

```
#What if there is a financial crisis...

df_train = pd.concat([df7,df72,df73,df74],axis = 0).sample(frac = 0.1)
df_test = pd.concat([df8,df72,df73,df74],axis = 0).sample(frac = 0.1)

#Pipeline random forest
from sklearn.ensemble import RandomForestRegressor
randomForest = RandomForestRegressor(max_depth= 18, n_estimators = 16, random_state=2)

X_train = df_train[k_best_features]
y_train = df_train['0']
X_test = df_test[k_best_features]
y_test = df_test['0']

randomForest.fit(X_train,y_train)
randomForest.predict(X_test)

#Two Years Later..
X_test = df9[k_best_features]
y_test = df9['0']
randomForest.fit(X_train,y_train)
printPerformance(randomForest.predict(X_test))
```

```
#What if there is a economy boom...
```

```
df_train = pd.concat([df0,df2,df4,df6,df8,df10,df12],axis = 0).sample(frac = 0.1)

X_train = df_train[k_best_features].head(shape)
y_train = df_train['0'].tail(df_test.shape)
X_test = df12
randomForest.fit(X_train,y_train)
printPerformance(randomForest.predict(X_test))
```

```
#What if there is a regime change from election
```

```
X_train = df6[k_best_features].head(shape)
y_train = df6['0'].tail(df_test.shape)
X_test = df12[k_best_features]
y_test = df12['0']

randomForest.fit(X_train,y_train)
printPerformance(randomForest.predict(X_test))
```

Analysis

Financial crisis

To analyze the global financial crisis of 2007 when there was an increase in perceived credit risk which meant that interest rates were at 45 years low at that point. Firstly, during our initial exploratory analysis, when we compared the average interest rates across all the years we observed that the data also suggest the same.

Hence, when predicting the interest rates on rolling four quarters of 2007 we wanted to predict the same trend and after applying the prediction algorithms we could achieve that.

This was a demonstration of a real-world prediction.

Results for year 2007 as input:

Q22007

Model	mae_test	mae_train	mape_test	mape_train	rms_test	rms_train
Regression	0.2629358023	0.2394470804	4.148025817	3.856505762	0.1203958029	0.10135417
RandomForest	0.269505549	0.08782097146	4.281543343	1.414535052	0.1233593407	0.01455494527
KNN	0.4159484609	0.4579470444	6.810792461	7.603761251	0.2484157089	0.2864749811
Nueral Network	408330.8397	408120.9876	6511743.454	6596818.878	167216364108	167060363261

Q32007

Model	mae_test	mae_train	mape_test	mape_train	rms_test	rms_train
Regression	0.4029513883	0.2627737642	5.91652137	4.168241375	0.2356904515	0.1186056597
RandomForest	0.5446795226	0.09148472963	8.072026386	1.451607726	0.4575119344	0.01566344685
KNN	0.6964585352	0.3861340897	10.2017319	5.919359159	0.6106584708	0.2498163864
Nueral Network	736705.2408	740009.5522	11103899.1	11815073.91	547035511716	551789417547

Q42007

Model	mae_test	mae_train	mape_test	mape_train	rms_test	rms_train
Regression	0.8454947093	0.2500554699	13.63193104	3.775162685	0.8470664749	0.1098105011
RandomForest	0.3943925252	0.09064244839	6.409776764	1.366924046	0.2344095519	0.01545857682
KNN	0.3792996396	0.3107574368	6.121126537	4.615202193	0.2285713076	0.1643232092
Nueral Network	266922.7753	267925.4544	4219615.518	4035007.816	71722454845	72196114218

Q12008

Model	mae_test	mae_train	mape_test	mape_train	rms_test	rms_train
Regression	0.5671321563	0.2687541438	10.07624573	4.227693005	0.4272141412	0.1275405334
RandomForest	0.6450613239	0.09872491777	11.51035578	1.554432517	0.5609488595	0.01823641337
KNN	0.4269412771	0.4216644719	7.459555023	6.408515619	0.2836636697	0.2915530914
Nueral Network	380155.8173	379119.4973	6515190.601	5994583.754	145356081225	144565666644

H2O

```
In [23]: aml.train(x=X,y=Y,training_frame=train_df)
AutoML progress: |██████████| 100%
```

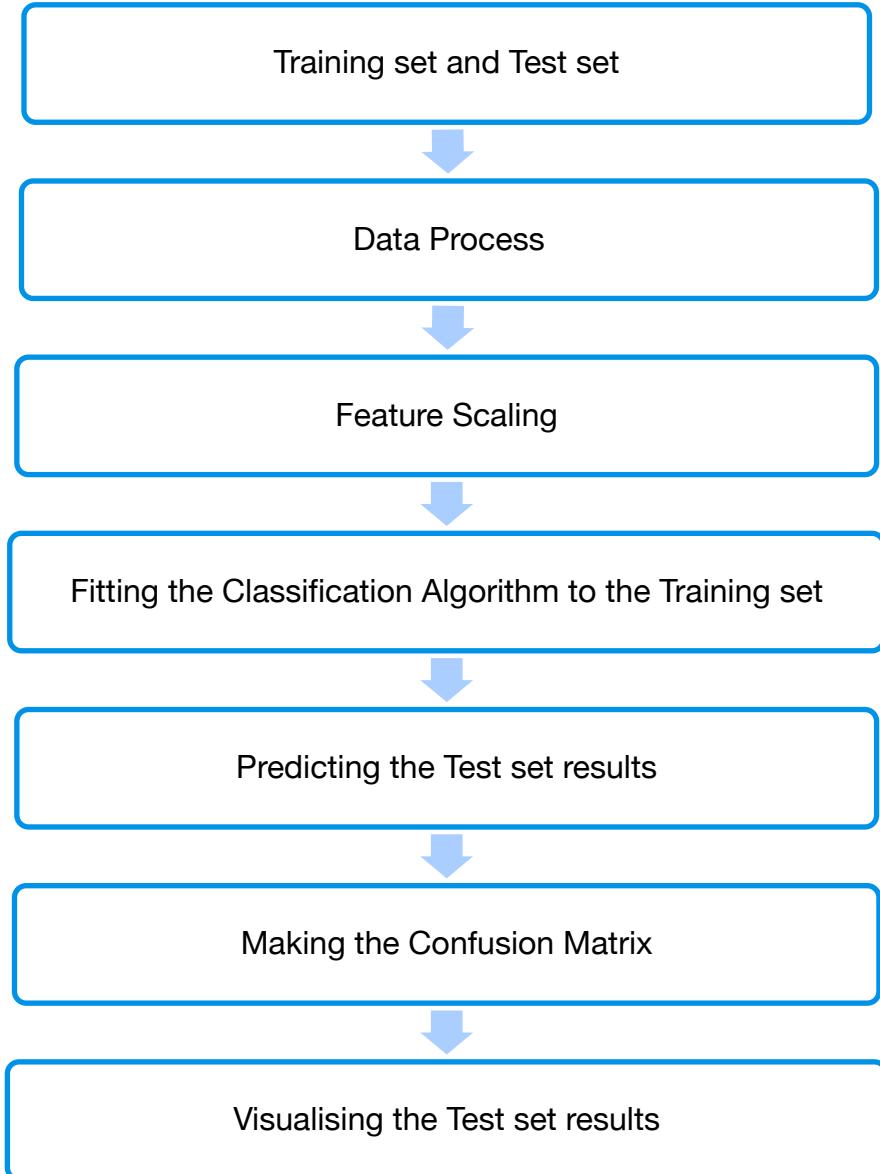
```
In [24]: aml_leaderboard_df=aml.leaderboard.as_data_frame()
aml_leaderboard_df
```

```
Out[24]:
```

	model_id	mean_residual_deviance	rmse	mse	mae	rmsle
0	StackedEnsemble_BestOfFamily_AutoML_20181207_1...	0.061811	0.248619	0.061811	0.182843	0.037073
1	StackedEnsemble_AllModels_AutoML_20181207_173712	0.061811	0.248619	0.061811	0.182843	0.037073
2	DRF_1_AutoML_20181207_173712	0.061938	0.248873	0.061938	0.182765	0.037102
3	XGBoost_1_AutoML_20181207_191701	0.061943	0.248884	0.061943	0.184050	0.037127
4	DRF_1_AutoML_20181207_191701	0.062034	0.249066	0.062034	0.182841	0.037129
5	XRT_1_AutoML_20181207_173712	0.065085	0.255118	0.065085	0.188245	0.038045
6	XRT_1_AutoML_20181207_191701	0.065114	0.255174	0.065114	0.188238	0.038047
7	GLM_grid_1_AutoML_20181207_173712_model_1	0.079482	0.281926	0.079482	0.207550	0.042004
8	GLM_grid_1_AutoML_20181207_191701_model_1	0.079482	0.281926	0.079482	0.207550	0.042004

```
In [25]: model_set=aml_leaderboard_df['model_id']
mod_best=h2o.get_model(model_set[0])
```

Part II: Classification



1. Programmatically downloads Q12005 and Q22005 origination data and pre-processes it.

```

#####
# Create Session #####
#####

USERNAME = 'wangying0327@hotmail.com'
PASSWORD = 'Y3ZUK;\D'

payload = {
    "username": USERNAME,
    "password": PASSWORD
}

session_requests = requests.session()

login_url = "https://freddiemac.embs.com/FLoan/secure/auth.php"

result = session_requests.post(
    login_url,
    data=payload,
    headers=dict(referer=login_url)
)

url = 'https://freddiemac.embs.com/FLoan/Data/download.php'
agreement_payload = {
    "accept": "Yes",
    "action": "acceptTandC",
    "acceptSubmit": "Continue"
}
result = session_requests.post(
    url,
    agreement_payload,
    headers=dict(referer=url)
)

```

2. Builds a Logistic regression model for the CURRENTLOANDELIN- QUENCYSTATUS

```

'''-----Logistic Regression-----'''

# Feature Scaling
from sklearn.preprocessing import StandardScaler
sc_X = StandardScaler()
X_train = sc_X.fit_transform(X_train)
X_test = sc_X.transform(X_test)

# Fitting the Logistic Regression to the Training set
from sklearn.linear_model import LogisticRegression
classifier = LogisticRegression(random_state = 0, solver='lbfgs')
classifier.fit(X_train, y_train)

# Predicting the Test set results
y_pred = classifier.predict(X_test)

# Making the Confusion Matrix
from sklearn.metrics import confusion_matrix
cm = confusion_matrix(y_test, y_pred)

```

3. using Q12005 data as training data (col 4). Note anytime col 4 is > 0, add a new variable as Delinquent. Use this variable as your “Y” variable. IGNORE COL 4 AND DON’T USE IT IN YOUR MODEL.

```

def f(row):
    if row['current_loan_delinquency_status'] > 0:
        val = 1
    else:
        val = 0
    return val
# Create dummy variables
df1_dummies = pd.get_dummies(df1[['repurchase_flag', 'modification_flag']])
df2_dummies = pd.get_dummies(df2[['repurchase_flag', 'modification_flag']])

df1_d = df1.drop(['loan_sequence_no', 'repurchase_flag', 'modification_flag'], axis=1)
df2_d = df2.drop(['loan_sequence_no', 'repurchase_flag', 'modification_flag'], axis=1)

global df1_final
global df2_final
df1_final = pd.concat([df1_d, df1_dummies], axis=1)
df2_final = pd.concat([df2_d, df2_dummies], axis=1)

# create target variable
df1_final['Delinquent'] = df1_final.apply(f, axis=1)
df2_final['Delinquent'] = df2_final.apply(f, axis=1)

# Create training and testing set

X_train = df1_final.drop(['current_loan_delinquency_status', 'Delinquent'], axis=1)
y_train = df1_final['Delinquent']

X_test = df2_final.drop(['current_loan_delinquency_status', 'Delinquent'], axis=1)
y_test = df2_final['Delinquent']

X_train = preprocessing.minmax_scale(np.array(X_train).astype(float)) # scale between 0 and 1
X_test = preprocessing.minmax_scale(np.array(X_test).astype(float))

return X_train, y_train, X_test, y_test

```

4. Validates against Q22005 data and selects the best Classification model Copyright QuantUniversity LLC. 2017 – Cannot be re-used/quoted without written permission
Random Forest

```

# random Forest
def rf(X_train, y_train, X_test, y_test, return_dict_rf):
    rf = RandomForestClassifier(n_estimators=100)
    rf.fit(X_train, y_train)

    y_train_predicted = rf.predict(X_train)
    y_test_predicted = rf.predict(X_test)

    conf_mat_rf = metrics.confusion_matrix(y_test, y_test_predicted)
    print(conf_mat_rf)

    # Compute ROC curve and AUC (Area under the Curve)
    false_positive_rate, true_positive_rate, thresholds = metrics.roc_curve(y_test, y_test_predicted)

    roc_auc_rf = metrics.auc(false_positive_rate, true_positive_rate)
    return_dict_rf['roc_auc_rf'] = roc_auc_rf
    return_dict_rf['conf_mat_rf'] = conf_mat_rf

    # Plot ROC Curve
    plt.title("Random Forest")
    plt.plot(false_positive_rate, true_positive_rate, 'b',
             label='AUC = %0.2f' % roc_auc_rf)
    plt.legend(loc='lower right')
    plt.plot([0, 1], [0, 1], 'r--')
    plt.xlim([-0.1, 1.2])
    plt.ylim([-0.1, 1.2])
    plt.ylabel('True Positive Rate')
    plt.xlabel('False Positive Rate')
    plt.show()

```

Neural Network

```

def nn(X_train, y_train, X_test, y_test, return_dict_nn):
    nn = MLPClassifier()

    nn.fit(X_train, y_train)

    y_train_predicted = nn.predict(X_train)
    y_test_predicted = nn.predict(X_test)

    conf_mat_nn = metrics.confusion_matrix(y_test, y_test_predicted)
    print(conf_mat_nn)

    # Compute ROC curve and AUC (Area under the Curve)
    false_positive_rate, true_positive_rate, thresholds = metrics.roc_curve(y_test, y_test_predicted)

    roc_auc_nn = metrics.auc(false_positive_rate, true_positive_rate)
    return_dict_nn['roc_auc_nn'] = roc_auc_nn
    return_dict_nn['conf_mat_nn'] = conf_mat_nn

    # Plot ROC Curve
    plt.title("Neural Network")
    plt.plot(false_positive_rate, true_positive_rate, 'b',
             label='AUC = %0.2f' % roc_auc_nn)
    plt.legend(loc='lower right')
    plt.plot([0, 1], [0, 1], 'r--')
    plt.xlim([-0.1, 1.2])
    plt.ylim([-0.1, 1.2])
    plt.ylabel('True Positive Rate')
    plt.xlabel('False Positive Rate')
    plt.show()

```

Assignment 3

Dec

```
: import h2o
h2o.init(strict_version_check=False, port=9957)

Checking whether there is an H2O instance running at http://localhost:9957..... not found.
Attempting to start a local H2O server...
Java Version: openjdk version "1.8.0_121"; OpenJDK Runtime Environment (Zulu 8.20.0.5-macosx) (build 1.8.0_121-b1
5); OpenJDK 64-Bit Server VM (Zulu 8.20.0.5-macosx) (build 25.121-b15, mixed mode)
Starting server from /anaconda3/lib/python3.6/site-packages/h2o/backend/bin/h2o.jar
Ice root: /var/folders/hf/c05h_6kd5wlds2ysv1yvpy1c0000gn/T/tmpdo3zi7w5
JVM stdout: /var/folders/hf/c05h_6kd5wlds2ysv1yvpy1c0000gn/T/tmpdo3zi7w5/h2o_wangying_started_from_python.out
JVM stderr: /var/folders/hf/c05h_6kd5wlds2ysv1yvpy1c0000gn/T/tmpdo3zi7w5/h2o_wangying_started_from_python.err
Server is running at http://127.0.0.1:9957
Connecting to H2O server at http://127.0.0.1:9957... successful.

H2O cluster uptime: 01 secs
H2O cluster timezone: America/New_York
H2O data parsing timezone: UTC
H2O cluster version: 3.22.0.2
H2O cluster version age: 15 days
H2O cluster name: H2O_from_python_wangying_kbgw7l
H2O cluster total nodes: 1
H2O cluster free memory: 3.556 Gb
H2O cluster total cores: 8
H2O cluster allowed cores: 8
H2O cluster status: accepting new members, healthy
H2O connection url: http://127.0.0.1:9957
H2O connection proxy: None
H2O internal security: False
H2O API Extensions: XGBoost, Algos, AutoML, Core V3, Core V4
Python version: 3.6.6 final
```

```
Out[27]:
model_id      auc  logloss  mean_per_class_error    rmse    mse
0  StackedEnsemble_AllModels_AutoML_20181207_133501  0.911731  0.139201  0.210976  0.187376  0.035110
1  StackedEnsemble_BestOfFamily_AutoML_20181207_1...  0.911731  0.139201  0.210976  0.187376  0.035110
2          DRF_1_AutoML_20181207_133501  0.907181  0.155591  0.219043  0.203407  0.041374
3          XRT_1_AutoML_20181207_133501  0.906353  0.157899  0.215183  0.205562  0.042256
4        GLM_grid_1_AutoML_20181207_133501_model_1  0.745118  0.221634  0.329522  0.241570  0.058356
```

```
In [28]: model_set=aml_leaderboard_df['model_id']
mod_best=h2o.get_model(model_set[0])
```

Best Model:

StackedEnsemble_AllModels_AutoML_20181207_133501

6. Parameterize the input (example it should take Q12005) and modify the code so that it outputs the 5 parameters listed in the matrix below.

matrix_classification

Quarter	No_of_actual_delq	No_of_pred_delq	No_of_records	No_of_delq_properly_classified	No_of_nonDelq_improperly_classified_as_delq
0	Q12005	8500	3761	200000	741
0	Q12005	8500	3761	200000	741

7. Write another script that calls the above classification script from Q11999-Q42016 and computes the following matrix.

```

arg_len = len(sys.argv)
quarters = []
end = ''
if arg_len == 3:
    startQuarter = sys.argv[1]
    endQuarter = sys.argv[2]
    end = endQuarter
    while(startQuarter != endQuarter):
        print(startQuarter)
        quarters.append(startQuarter)
        startQuarter = get_next_quarter(startQuarter)
        quarters.append(endQuarter)
elif arg_len == 2:
    startQuarter = sys.argv[1]
    print("----" + startQuarter)
    quarters.append(startQuarter)
    end = startQuarter
else:
    print("running for default Q12005")
    quarter = 'Q12005'
    quarters.append(quarter)
    end = quarter
download

```

8. • Extra credit: Note that each invocation of is independent. There is scope to parallelize it. Parallelize this so that you can run two or more models concurrently (task parallel). Review R and Python packages that can be used to parallelize it to generate the below matrix as a dataframe and export it to a csvfile.

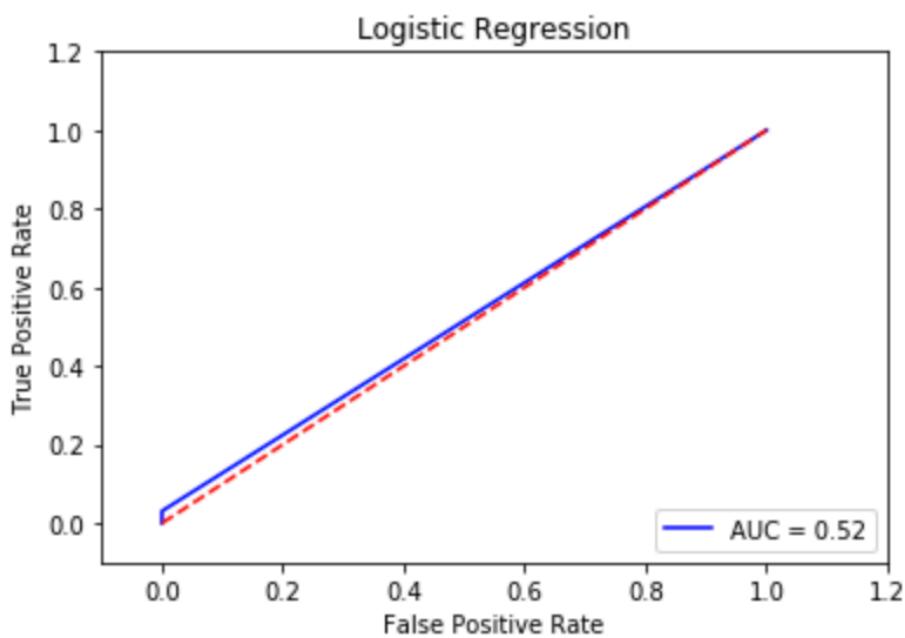
```
for q in quarters:
    nextQuarter = get_next_quarter(q)
    load_data_df(q, nextQuarter)
    remove_nan(df1)
    remove_nan(df2)
    X_train, y_train, X_test, y_test = process_data()
    print("process data executed")

    p1 = Process(target=logred,args=(X_train, y_train,X_test,y_test,return_dict_logred))
    p2 = Process(target=rf,args=(X_train, y_train,X_test,y_test,return_dict_rf))
    p3 = Process(target=nn,args=(X_train, y_train,X_test,y_test,return_dict_nn))
    p1.start()
    p2.start()
    p3.start()
    p1.join()
    p2.join()
    p3.join()
    matrix=genMatrix(q,return_dict_logred['roc_auc_logred'],return_dict_rf['roc_auc_rf'],return_dict_nn['roc_auc_nn'],return_dict_logred['conf_mat_logred'],return_dict_rf['conf_mat_rf'],return_dict_nn['conf_mat_nn'],matrix)

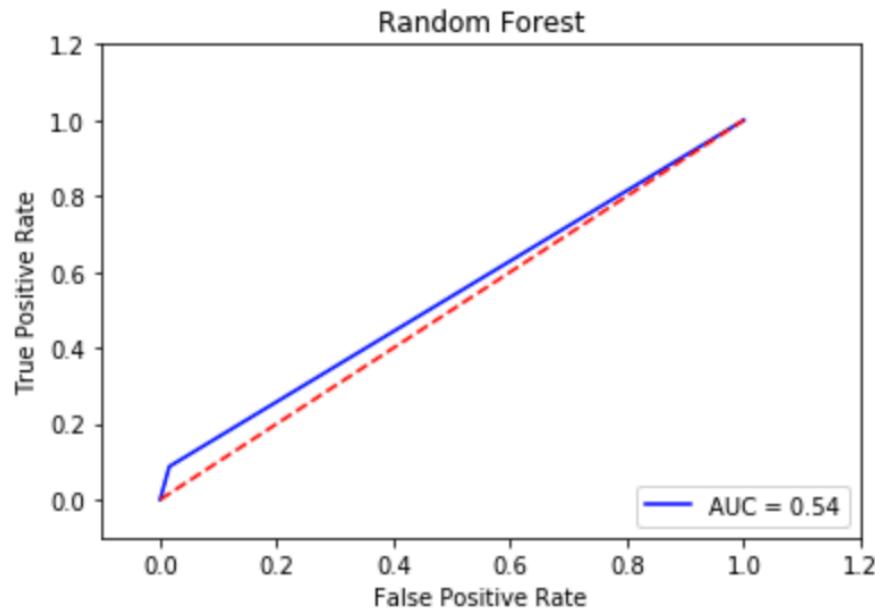
    print(matrix)
    matrix.to_csv('matrix_classification.csv')
```

Graph :

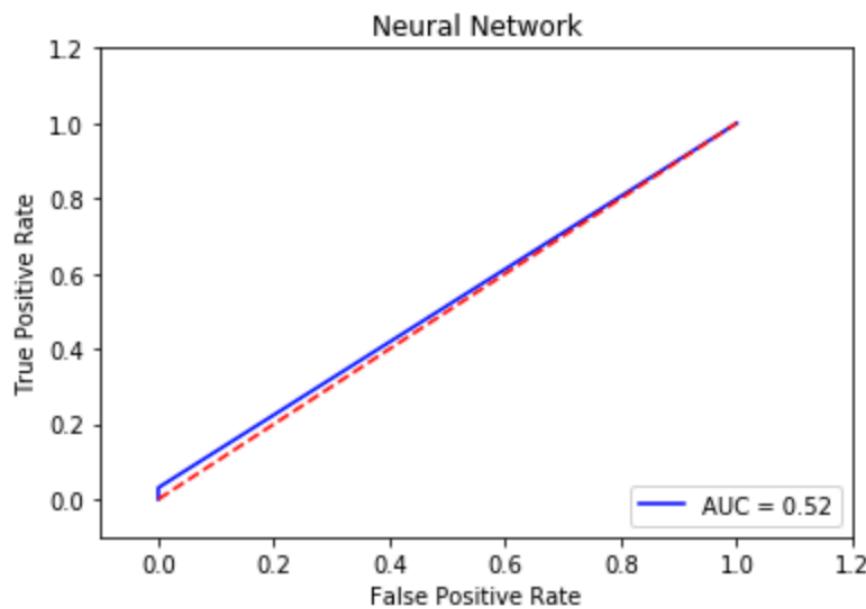
```
[[191498      2]
 [ 8239     261]]
```



```
[[188480    3020]
 [ 7759     741]]
```



```
[[191498      2]
 [ 8239     261]]
```



```
Quarter  No_of_actual_delq  No_of_pred_delq  No_of_records  \
0  Q12005          8500.0          3761.0        200000.0
0  Q12005          8500.0          3761.0        200000.0

No_of_delq_properly_classified  No_of_nonDelq_improperly_classified_as_delq
0                           741.0                           3020.0
0                           741.0                           3020.0
```

```
]:
```