

Foreclosure prediction using machine learning

Multiple Linear Regression/GradientBoostingRegressor/Ridge
Regression

Mika Breyfogle Dec 2020

GitHub: find my source code in [here](#).

Problem Statement

The goal of the analysis is to help us understanding the relation between unemployment rate/ income/ foreclosure and how zip code are used to determine it. I hope that creating a model will help predict foreclosure in near future. I also create recommended zip code for the home seeker in his early thirties who likes to go to coffee shops and like outdoor sports.

Data acquisition and cleaning

1. [City of los angles foreclosure data by zip code form 2015 to 2020](#)
2. [City of los angles unemployment data by zip 2020](#)
3. [Total unemployment per month from 1990 to 2020.](#)
4. [Unemployment LA County 1990](#)
5. [Income and Liability data 2018 by Zip code](#)
6. [Average Condo price by zip. 2015 ~ 2020](#)
7. Use Foreclosure to find coffee shops and sports venues

Updating missing data using data allocations. Creating new data frame for data analysis.

Training machine learning algorithms

Multiple Linear Regression, gradient boosting regression and Ridge Regression.

Step1 : Exploratory Data Analysis (EDA)

Import the data and check to see what kind of data we are dealing with:

- 1) Import required libraries.
- 2) Import Data to the data frame.
- 3) Merge datas into one Dataframe.
- 3) See the head of the data to know how the data looks like, also some data need to clean.

ML PROJECT

4) Use describe function to see the percentiles's

	zip code	unemployment_rate	yearmonth	foreclosure	income	condprice	lat	lon
0	90013	32.353884	2015-01-01	1.0	96694.701625	500461.0	34.04461	-118.241
1	90021	25.984959	2015-01-01	0.0	111552.909449	513140.0	34.02941	-118.240
2	90059	12.254380	2015-01-01	44.0	32530.141714	328985.0	33.92691	-118.248
3	90014	12.099669	2015-01-01	1.0	80076.552268	566374.0	34.04301	-118.251
4	90002	11.925620	2015-01-01	63.0	31013.428093	326200.0	33.94901	-118.247

Describe data to find %.

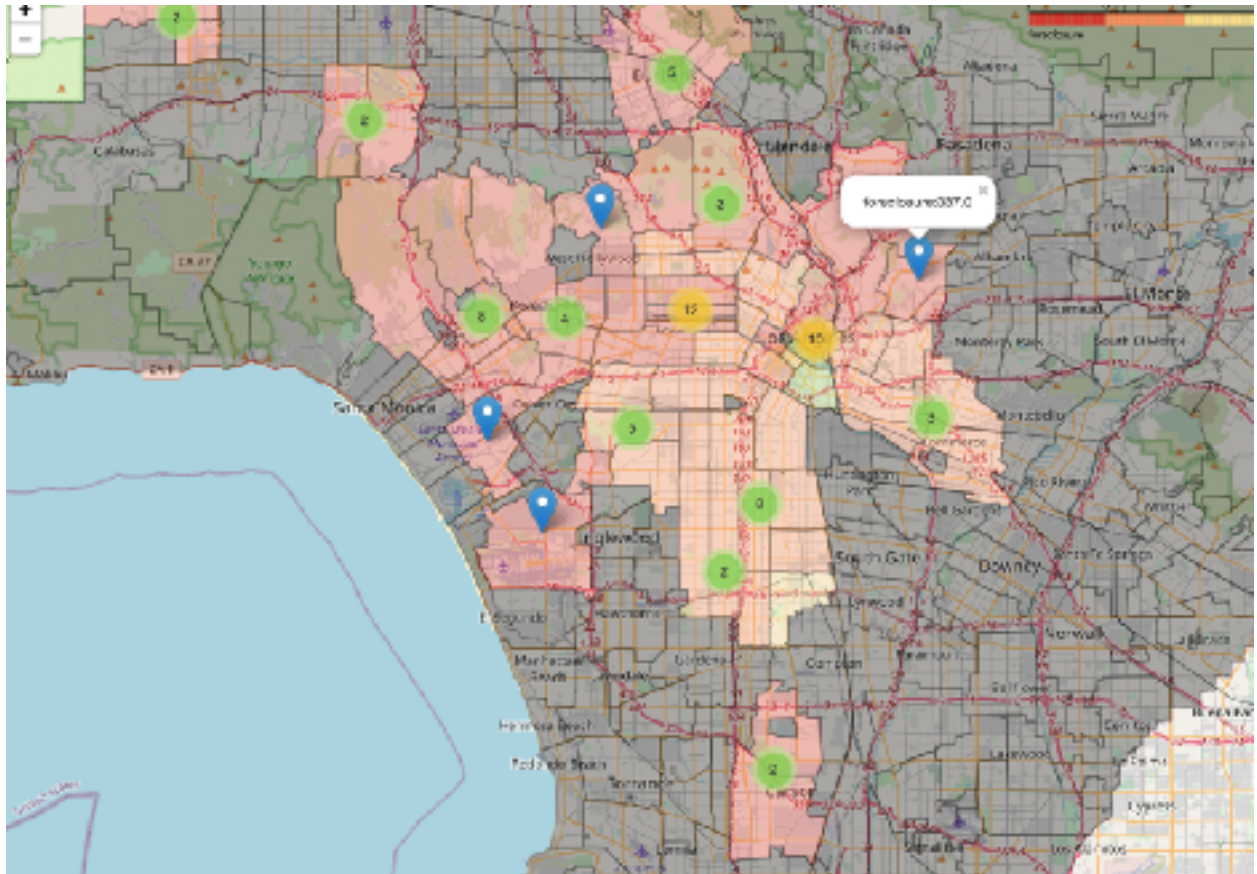
```
mydata.describe()
```

	unemployment_rate	foreclosure	income	condprice	lat	lon
count	4380.000000	4380.000000	4380.000000	4380.000000	4380.000000	4380.000000
mean	4.564478	2.772831	115556.737334	8.048009e+05	34.055348	-118.319938
std	3.147048	7.661656	141966.963846	2.348755e+05	0.071496	0.093424
min	0.981818	0.000000	29493.563483	2.595070e+05	33.824414	-118.624412
25%	2.563664	0.000000	41543.924521	4.303970e+05	34.023995	-118.373334
50%	4.025967	1.000000	63404.085511	5.594850e+05	34.057308	-118.303095
75%	5.605455	2.000000	110858.131093	7.255750e+05	34.097448	-118.258872
max	32.353884	138.000000	721863.508345	1.609523e+06	34.220394	-118.150703

```
mydata.head()
```

	zip code	unemployment_rate	yearmonth	foreclosure	income	condprice	lat	lon
0	90013	32.353884	2015-01-01	1.0	96694.701625	500461.0	34.044610	-118.241599
1	90021	25.984959	2015-01-01	0.0	111552.909449	513140.0	34.029411	-118.240375
2	90059	12.254380	2015-01-01	44.0	32530.141714	328985.0	33.926901	-118.240170
3	90014	12.099669	2015-01-01	1.0	80076.552268	566374.0	34.043081	-118.251755
4	90002	11.925620	2015-01-01	63.0	31013.428093	326200.0	33.949079	-118.247877

Now I created map to visualize the location of the foreclosure/unemployment rate based on latitude and longitude. I am going to see the common location and how the foreclosure are placed.

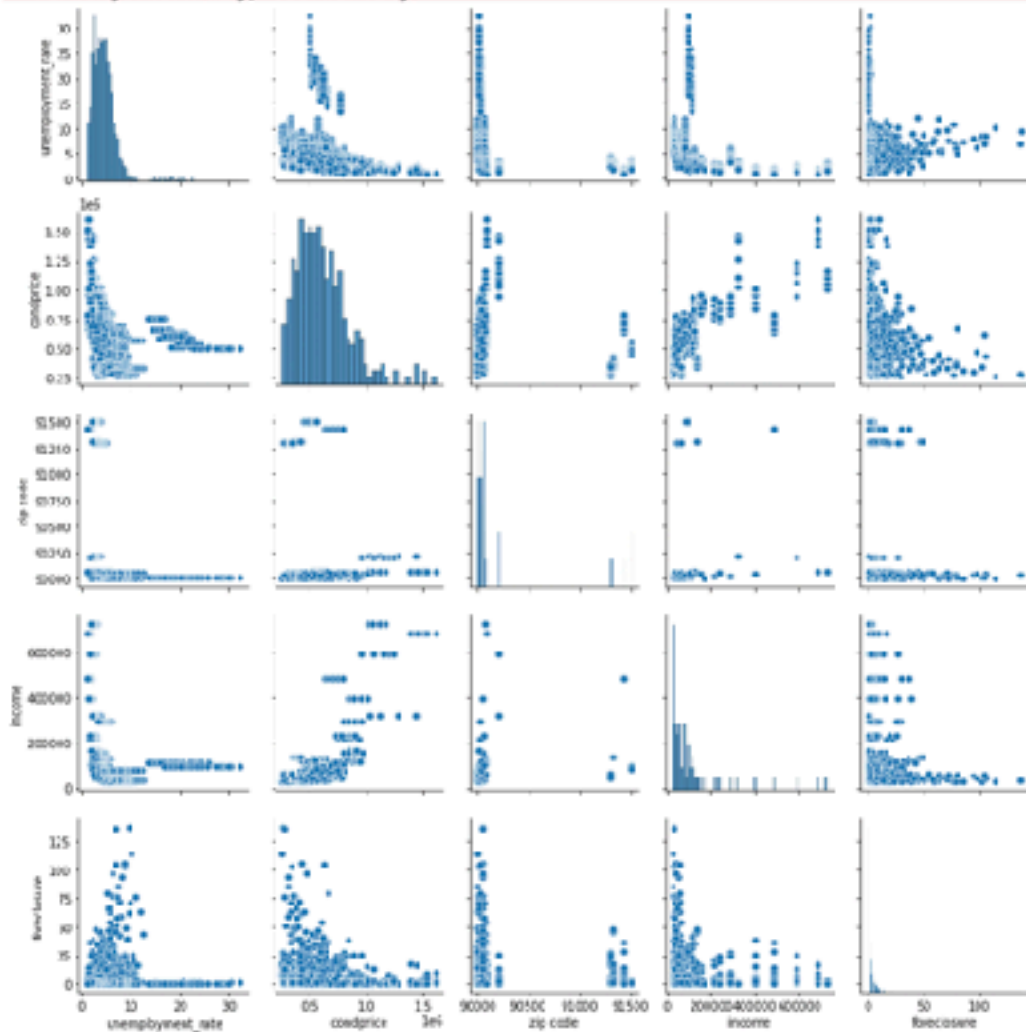


ML PROJECT

Let's plot couple of feature to get a better feel of the data.

```
# calculate and show pairplot
z=mydata[['unemployment_rate', 'condprice', 'zip_code', 'income', 'foreclosure']]
sns.pairplot(z, size=2.5)
plt.tight_layout()
```

/opt/anaconda3/lib/python3.8/site-packages/seaborn/axisgrid.py:1912: UserWarning: The `size` parameter has been re
warnings.warn(msg, UserWarning)

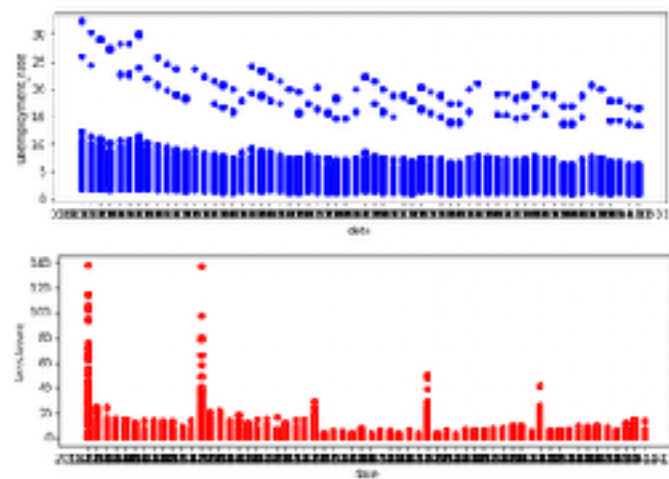


```

#simple regression zip vs price
msk=np.random.rand(len(mydatas))< 0.8
train=mydatas[msk]
test=mydatas[~msk]
fig = plt.figure(figsize=(10,7))
fig.add_subplot(2,1,1)
plt.scatter(train['yearmonth'],train.unemployment_rate,color='blue')
plt.xlabel("date")
plt.ylabel("unemployment_rate")
plt.show
#simple regression zip vs price
fig = plt.figure(figsize=(10,7))
fig.add_subplot(2,1,1)
plt.scatter(train['yearmonth'],train.foreclosure,color='red')
plt.xlabel("date")
plt.ylabel("foreclosure")
plt.show

```

```
<function matplotlib.pyplot.show(close=None, block=None)>
```



It looks like foreclosure and unemployment has some relationship.

ML PROJECT

Now I will find Coffee and sports venues using foursquare per each zip code

Here is a sample of code and Dataframe which I created.

```

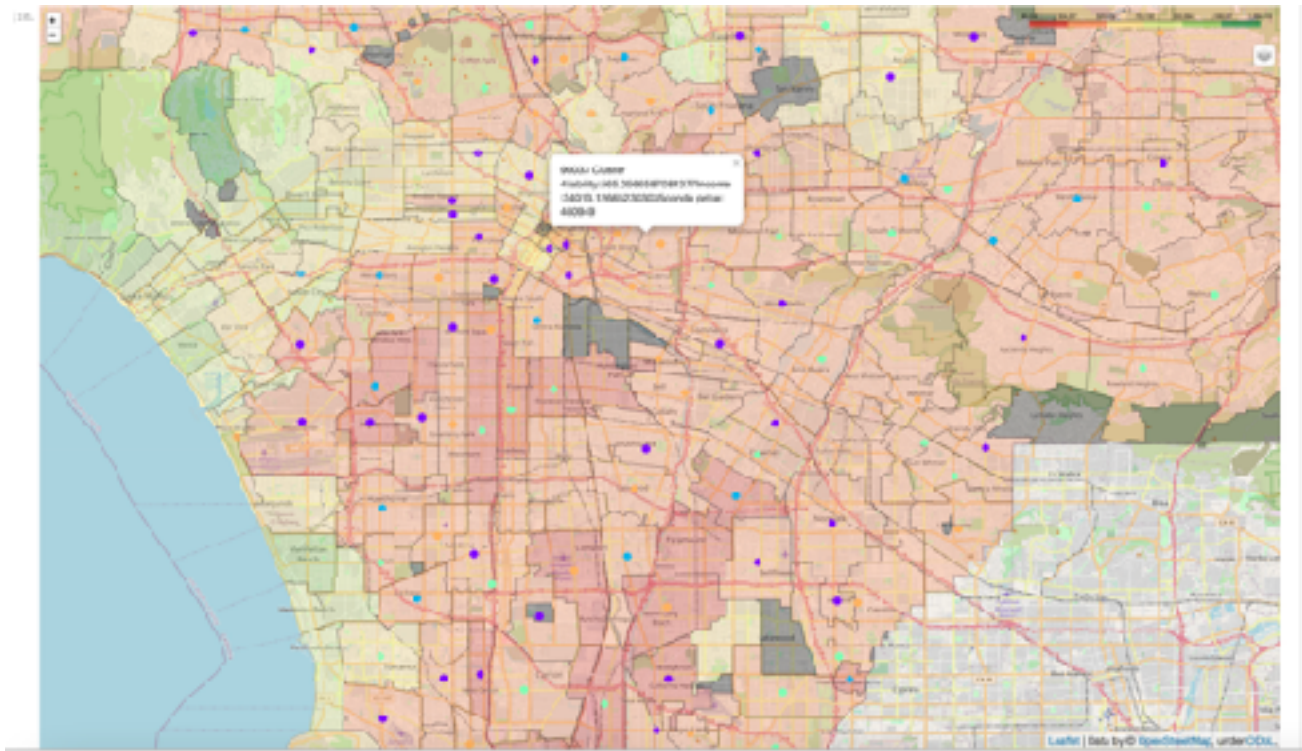
# Use the following coordinates (lat=50.8503386, lon=10.9941775, alt=508.036) or (lat=50.8503386, lon=10.9941775, alt=508.036) or (lat=50.8503386, lon=10.9941775, alt=508.036)
name=vers2
radius = 500
for name, lat, lon in zip(vers2['name'], vers2['lat'], vers2['lon']):
    # Create the GET request
    url = 'http://api.foursquare.com/v2/venues/explore?lat={lat}&lon={lon}&radius={radius}&category={category}'.format(
        CLINT_ID,
        CLINT_SECRET,
        VERSION,
        lat,
        lon,
        radius,
        2015,
        cat)
    # Make the GET request
    results = requests.get(url)
    # Return only relevant information for each nearby venue
    # Return only relevant information for each nearby venue
    venue_list.append([
        name,
        lat,
        lon,
        e['venue']['name'],
        e['venue']['location']['lat'],
        e['venue']['location']['lon'],
        e['venue']['categories'][0]['name'] for e in results])

nearby_venues = pd.DataFrame.from_records(venue_list)
nearby_venues.columns = ['Neighborhood',
                        'Neighborhood Latitude',
                        'Neighborhood Longitude',
                        'Venue',
                        'Venue Latitude',
                        'Venue Longitude',
                        'Venue Category']

```

year	zip code	city	community	median price	city	median	geo_location	income	quality	cluster_label	1st Most Common Value	2nd Most Common Value	3rd Most Common Value
0	2000	90000	Los Angeles (South Los Angeles), Florence-Graham	420000	Los Angeles	20700	["type": "Point", "coordinates": [-118.249540..., 34.0516332]]	314.7266	3	outdoor	Coffee/tea		
1	2010	90000	Los Angeles (South Los Angeles), Florence-Graham	484500	Los Angeles	20700	["type": "Point", "coordinates": [-118.249540..., 34.0516332]]	314.7266	3	outdoor	Coffee/tea		
2	2015	90000	Los Angeles (South Los Angeles), Florence-Graham	470000	Los Angeles	20700	["type": "Point", "coordinates": [-118.249540..., 34.0516332]]	314.7266	3	outdoor	Coffee/tea		
3	2015	90000	Los Angeles (South Los Angeles), Florence-Graham	441000	Los Angeles	20700	["type": "Point", "coordinates": [-118.249540..., 34.0516332]]	314.7266	3	outdoor	Coffee/tea		
4	2016	90000	Los Angeles (South Los Angeles), Florence-Graham	420000	Los Angeles	20700	["type": "Point", "coordinates": [-118.249540..., 34.0516332]]	314.7266	3	outdoor	Coffee/tea		

Now I created map to visualize the location of the venues based on latitude and longitude. I am going to see the common location and how the venues are placed. I cluster are labeled.



Now that I am familiar with all these representation and I can tell my own story .

Now I create a model to which would predict the foreclosure based upon the other factors such as income , zip code. I use train data and test data , train data to train the machine and test data to see if it has learnt the data well or not.

Step 2: Dataset Preparation (Splitting and Scaling)

Data is divided into the Train set and Test set.

Step 3: Model creation

Here is a code for Multiple Liner model coding

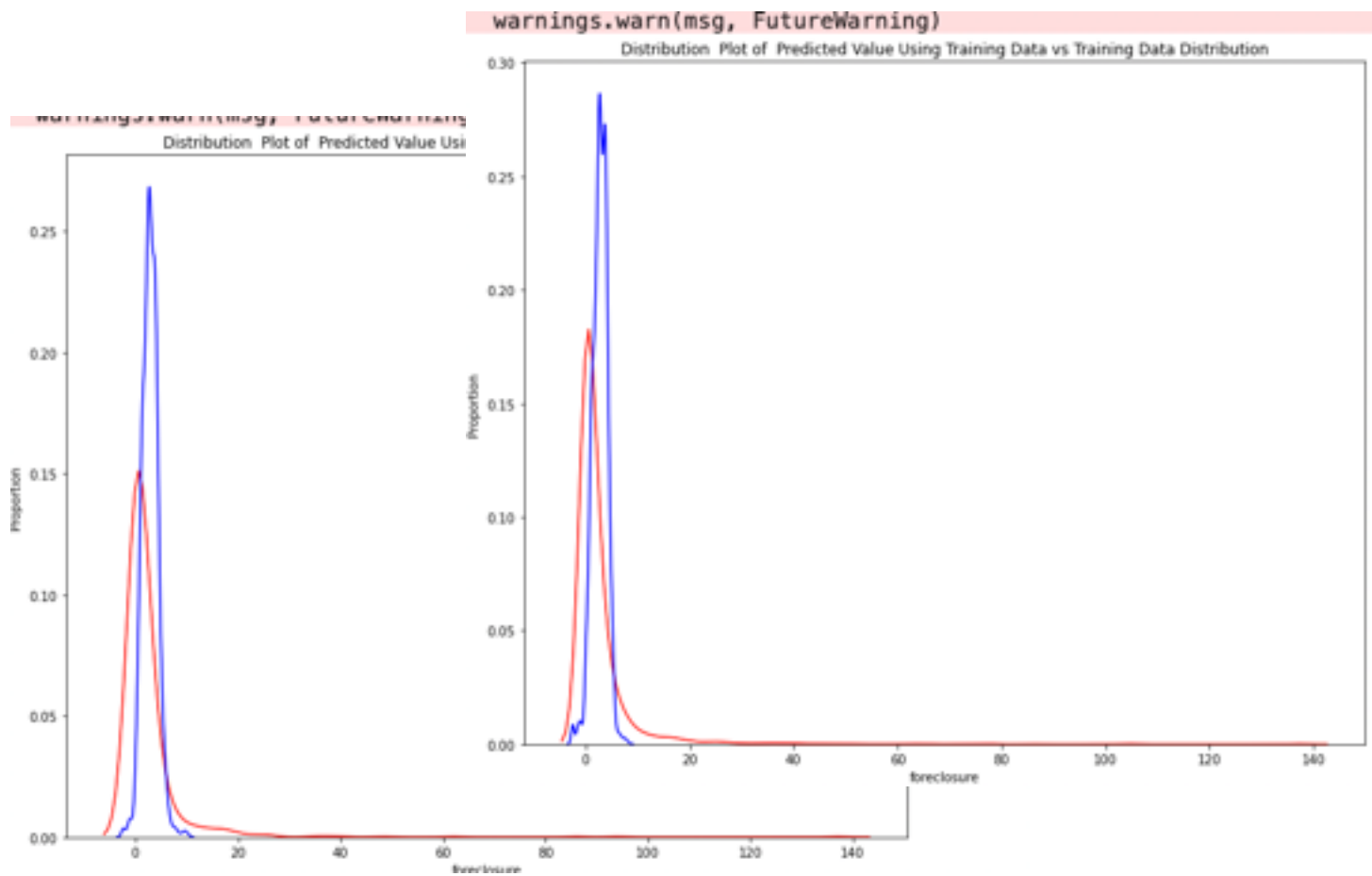
```
Type to enter text#Create mulitiple regression with train & test
mydataz=mydatas
lr1 = LinearRegression()
label=mydataz['foreclosure']
train1 = mydataz.drop(['foreclosure','lat','lon','yearmonth','zip code'],axis=1)
#Creating train / test set
x_train, x_test, y_train, y_test = train_test_split(train1, label, test_size=0.20, random_state=4)

lr1.fit(x_train, y_train)
print("number of test samples :", x_test.shape[0])
print("number of training samples:",x_train.shape[0])
yhat_train = lr1.predict(x_train)
yhat_test = lr1.predict(x_test)

Title = 'Distribution Plot of Predicted Value Using Training Data vs Training Data Distribution'
DistributionPlot(y_train, yhat_train, "Actual Values (Train)", "Predicted Values (Train)", Title)
Title='Distribution Plot of Predicted Value Using Test Data vs Data Distribution of Test Data'
DistributionPlot(y_test,yhat_test,"Actual Values (Test)","Predicted Values (Test)",Title)

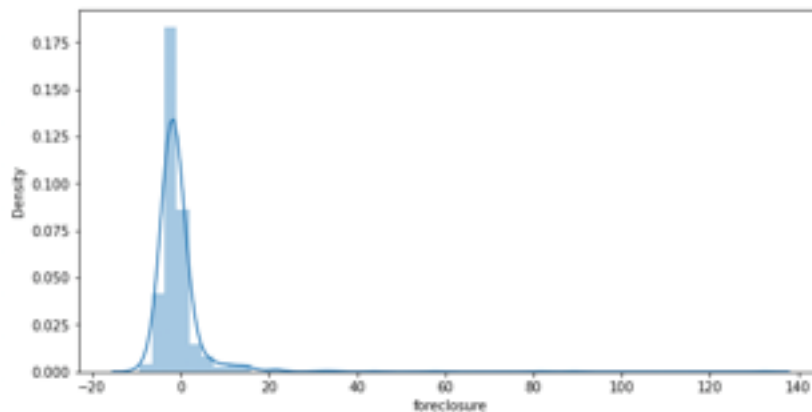
# visualizing residuals
fig = plt.figure(figsize=(10,5))
residuals = (y_test- yhat_test)
sns.distplot(residuals)

print("The R-square Linear is: ', lr1.score(train1, mydataz['foreclosure']))
# evaluate the performance of the algorithm (MAE - MSE - RMSE)
print('MAE:', metrics.mean_absolute_error(y_test, yhat_test))
print('MSE:', metrics.mean_squared_error(y_test, yhat_test))
print('RMSE:', np.sqrt(metrics.mean_squared_error(y_test, yhat_test)))
print('VarScore:',metrics.explained_variance_score(y_test,yhat_test))
```

warnings.warn(msg, FutureWarning)

The R-square Linear is: 0.0342734584356702
 MAE: 3.2251486345606897
 MSE: 60.809271268367986
 RMSE: 7.7980299607252075
 VarScore: 0.03576978059669422



As you can see both training and test is equally bad. R-square is .03. Error is high too.

Let train another model. GradientBoostingRegressor

```

params1 = {'n_estimators': 500,
           'max_depth': 4,
           'min_samples_split': 4,
           'learning_rate': 0.01,
           'loss': 'quantile',
           'alpha': .70}

label=mydataz['foreclosure']
train1 = mydataz.drop(['foreclosure','lat','lon','yearmonth','zip code'],axis=1)
x_train, x_test, y_train, y_test = train_test_split(train1, label, test_size=0.20, random_state=5)
clf = ensemble.GradientBoostingRegressor(**params1)
clf.fit(x_train,y_train)
clf.score(x_test,y_test)
yhat_test2=clf.predict(x_test)
mse = mean_squared_error(y_test, clf.predict(x_test))
print("The GRADIENT BOOST R-square is: ", clf.score(x_train,y_train))
print("MAE:", metrics.mean_absolute_error(y_test, yhat_test2))
print("MSE:", metrics.mean_squared_error(y_test, yhat_test2))
print("RMSE:", np.sqrt(metrics.mean_squared_error(y_test, yhat_test2)))
print("VarScore:",metrics.explained_variance_score(y_test,yhat_test2))

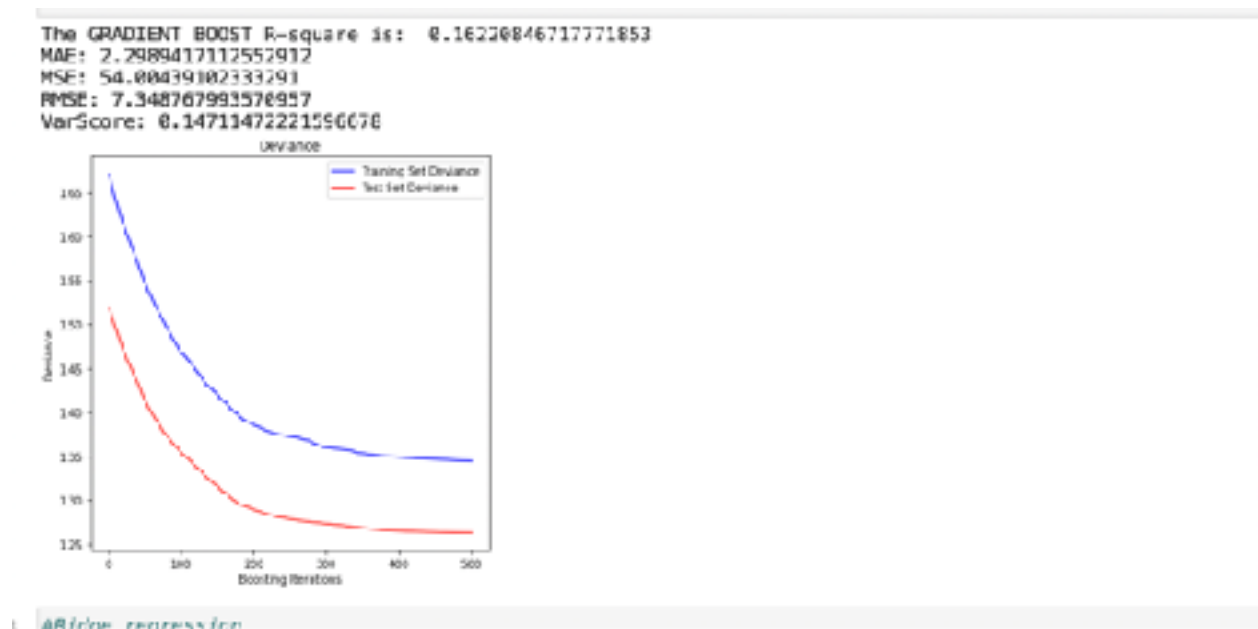
test_score = np.zeros((params['n_estimators'],), dtype=np.float64)
for i, y_pred in enumerate(clf.staged_predict(x_test)):
    test_score[i] = clf.loss_(y_test, y_pred)

fig = plt.figure(figsize=(6, 6))
plt.subplot(1, 1, 1)
plt.title('Deviance')
plt.plot(np.arange(params['n_estimators']) + 1, clf.train_score_, 'b-',
         label='Training Set Deviance')
plt.plot(np.arange(params['n_estimators']) + 1, test_score, 'r-',
         label='Test Set Deviance')
plt.legend(loc='upper right')
plt.xlabel('Boosting Iterations')
plt.ylabel('Deviance')
fig.tight_layout()
plt.show()

```

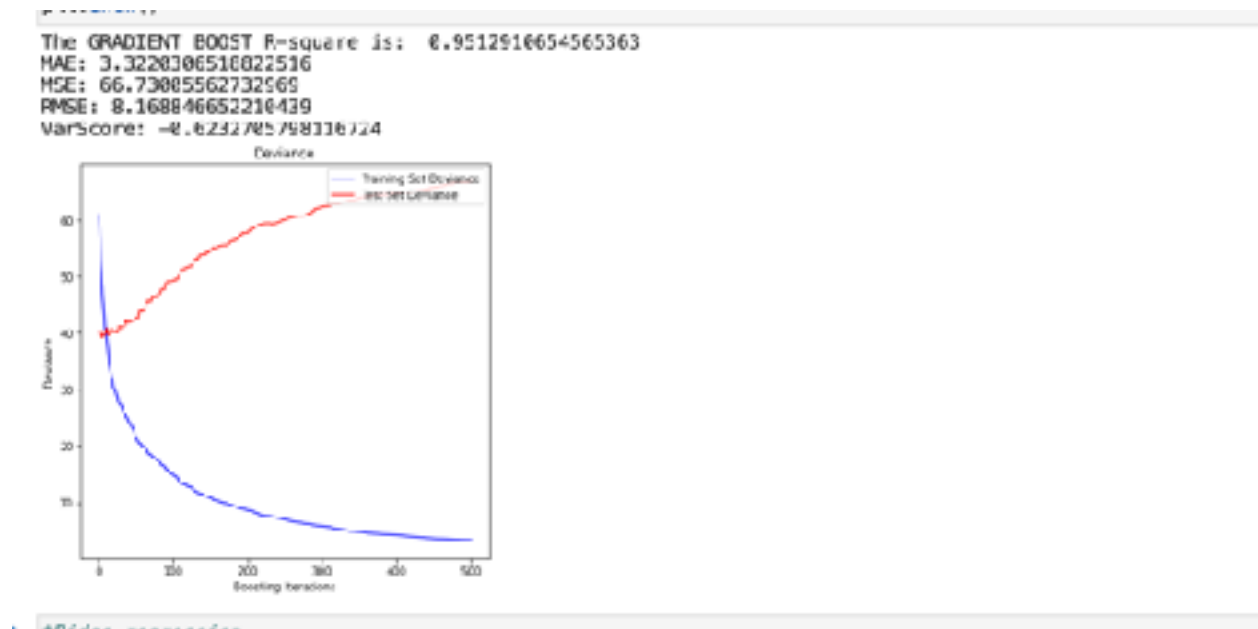
ML PROJECT

Here is a prediction looks like



Still R-square is too low

I will try another parameters and see what happen



R-square is near 1.00. this is good model, however MSE is still high

Now I am going to create Redge Model.

First I will find best alpha using GridSearch

Here is a code

```

#Ridge regression
pr=PolynomialFeatures(degree=2)
x_train_pr=pr.fit_transform(x_train)
x_test_pr=pr.fit_transform(x_test)

from sklearn.linear_model import Ridge

#para2= {'alpha': [0.001,0.1,1, 10, 100, 1000, 10000, 100000, 1000000], 'normalize':[True,False]}
para2= {'alpha': [0.001,0.1,1, 10, 100, 1000, 6000,10000, 50000,100000, 1000000]}
RR=Ridge()

Grid1 = GridSearchCV(RR, para2,cv=4)

Grid1.fit(train1, label)

GridSearchCV(cv=4, estimator=Ridge(),
              param_grid=[{'alpha': [0.001, 0.1, 1, 10, 100, 1000, 6000, 10000,
                                     50000, 100000, 1000000]}])

BestRR=Grid1.best_estimator_
BestRR

Ridge(alpha=10000)

```

Best Alpha is 10000, let create model

```

RidgeModel=Ridge(alpha=10000)
RidgeModel.fit(x_train_pr, y_train)
yhat2 = RidgeModel.predict(x_test_pr)
Rsqu_test = []
Rsqu_train = []
dummy1 = []
Alpha = 10 * np.array(range(0,1000))
for alpha in Alpha:
    RidgeModel = Ridge(alpha=alpha)
    RidgeModel.fit(x_train_pr, y_train)
    Rsqu_test.append(RidgeModel.score(x_test_pr, y_test))
    Rsqu_train.append(RidgeModel.score(x_train_pr, y_train))

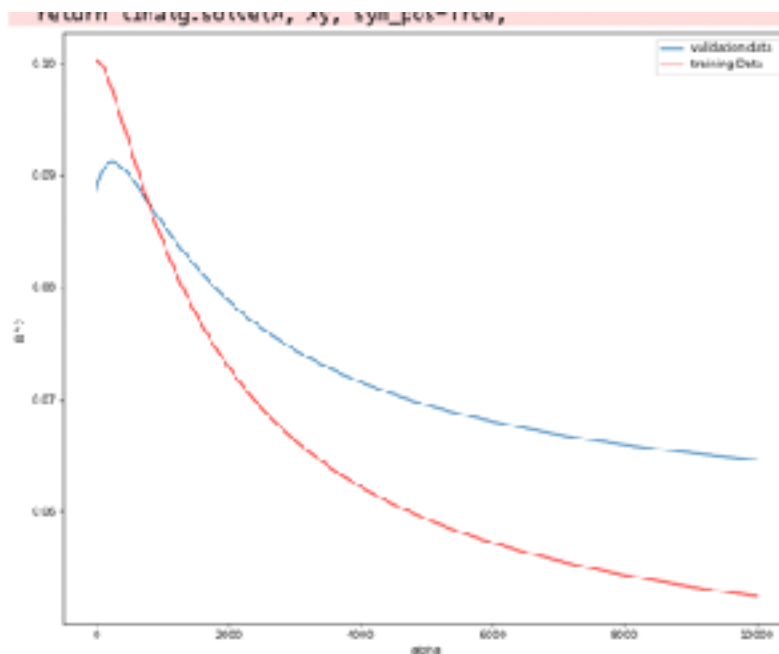
width = 12
height = 10
plt.figure(figsize=(width, height))

plt.plot(Alpha, Rsqu_test, label='validation data ')
plt.plot(Alpha, Rsqu_train, 'r', label='training Data ')
plt.xlabel('alpha')
plt.ylabel('R^2')
plt.legend()
print('R-score train:',RidgeModel.score(x_train_pr,y_train))
print('R-score test:',RidgeModel.score(x_test_pr,y_test))
print('The GRADIENT BOOST R-square is: ', RidgeModel.score(x_train_pr,y_train))
print('MAE:', metrics.mean_absolute_error(y_test, yhat2))
print('MSE:', metrics.mean_squared_error(y_test, yhat2))
print('RMSE:', np.sqrt(metrics.mean_squared_error(y_test, yhat2)))
print('VarScore:',metrics.explained_variance_score(y_test,yhat2))

```

ML PROJECT

Here is a prediction look like.



```
R-score train: 0.10027274079088344
R-score test: 0.08889618448379006
The GRADIENT BOOST R-square is: 0.10027274079088344
MAE: 3.109488590871713
MSE: 58.87722462607468
RMSE: 7.673149585710517
VarScore: 0.06556781107360133
```

```
/root/anaconda3/lib/python3.8/site-packages/sklearn/linear_model/_ridge.py:147: LinAlgWarning: Ill-condition
```

R-score is still low.

I will evaluate above models using actual number.

```
In [check result]

house1 = train1.loc[(mydata['zip code'] == 90005) & (mydata['unemployment_rate'] > 0.9) & (mydata['unemployment_rate'] < 4.5)]
house2 = gr.fat_transform(house1)
print('Linear :', lin.predict(house1).mean())
print('Gradient :', clf.predict(house2).mean())
print('Ridge :', RidgeModel.predict(house2).mean())
house3 = mydata.loc[(mydata['zip code'] == 90023) & (mydata['unemployment_rate'] > 0.9) & (mydata['unemployment_rate'] < 4.5)]
print(house3['foreclosure'].mean())
```

```
Linear : 2.8615970602007
Gradient : 0.41512522079152
Ridge : 2.8979124658712
0.586600020005866
```

```
In [check

house1 = train1.loc[(mydata['zip code'] == 90005) & (mydata['unemployment_rate'] > 0.9) & (mydata['unemployment_rate'] < 4.5)]
house2 = gr.fat_transform(house1)
print('Linear :', lin.predict(house1).mean())
print('Gradient :', clf.predict(house2).mean())
print('Ridge :', RidgeModel.predict(house2).mean())
house3 = mydata.loc[(mydata['zip code'] == 90001) & (mydata['unemployment_rate'] > 0.9) & (mydata['unemployment_rate'] < 4.5)]
print(house3['foreclosure'].mean())
```

```
Linear : 2.7018107100016945
Gradient : 2.0915017067110205
Ridge : 2.5137488648062527
2.9555555555555554
```

```
In [house1 = train1.loc[(mydata['zip code'] == 90004) & (mydata['unemployment_rate'] > 0.9) & (mydata['unemployment_rate'] > 4.5)]
house2 = gr.fat_transform(house1)
print('Linear :', lin.predict(house1).mean())
print('Gradient :', clf.predict(house2).mean())
print('Ridge :', RidgeModel.predict(house2).mean())
house3 = mydata.loc[(mydata['zip code'] == 90009) & (mydata['unemployment_rate'] > 0.9) & (mydata['unemployment_rate'] > 4.5)]
print(house3['foreclosure'].mean())
```

```
Linear : 3.5389186150807085
Gradient : 6.514467955070803
Ridge : 3.66293204131106
7.714205714205714
```

Conclusion:

Out of three model GradientBoostingRegressor create best result.

Most of bank foreclosed house on beginning of year. Maybe because of holiday back log.

When unemployment rate is real high (example 15% or more), foreclosure will not impact unemployment_rate. Probably, most of unemployed person do not own the house.