

Step: 1

Import pandas as pd

Import csv.

house = pd.read\_csv("Ames-House-sales-Cropped.csv")

house

house.head()

house.shape

df = pd.read\_csv("Ames-House-sales-Cropped.csv")

df

d = df.columns

d

df.info()

df.dtypes.value\_counts

Step: 2

house.drop(["Bldg Type"], axis=1)

df.pop("Central Air")

## Lab4. House Price Prediction using LR with Regularization

### Objectives

In this lab, you will learn how to build Regression Models to predict the Sales Price of houses using Ames, Iowa dataset.

### Learning Outcomes

After completing this lab, you will be able to

- Understand data and build baseline LR model
- Apply One Hot Encoding to categorical features and rebuild LR model
- Apply scaling using StandardScaler and MinMaxScaler and rebuild LR model
- Compare the performance of LR model with SGD Regressor, RidgeCV and LassoCV
- Compute RMSE for all models.

### Business Use Case

Ames house sales price dataset contains past sales price details of 1379 homes from Ames, USA. This dataset has 79 feature columns (independent variables) and the dependent variable, sale price ("SalePrice"). There are three different types: integers (int64), floats (float64), and strings (object, categoricals).

In this lab, you will use a cropped dataset where we have removed all object columns, except two object columns – building type ("BldgType") and centralized air conditioning ("CentralAir"). Therefore, it contains only 38 independent variables.

With the help of the regression model you will build on this reduced dataset, we can predict the sale price of a house given the values for all independent variables.

**Step1. [Import dataset].** Using Pandas, import "Ames\_House\_Sales\_Cropped.csv" file and print properties such as head, shape, columns, dtype, info and value\_counts.

**Step2. [Predict Sale Price without Categorical features].**

- Drop both categorical features – BldgType and CentralAir (USE drop() and pop() methods)
- Prepare X matrix (36 feature columns) and y vector (ie., SalePrice column)
- Split dataset for training and testing as X\_train, X\_test, y\_train, y\_test (use 25% test size).
- Create LinearRegression model, fit on training set and predict on test set
- Compute Mean Squared Error (MSE) on actual values and predicted values (you will get output as 1474827326.0).

**Step3. [Create Scatter Plot].** Plot Scatterplot between y\_test and y\_pred.

**Step4. [Encode Categorical columns].** Using get\_dummies() method, perform one hot encoding on the two categorical columns, BldgType and CentralAir. Now, you will get 5 columns for BldgType variable and 2 columns for CentralAir column. So, now you have 43 independent variables and 1 dependent variable.

**Step5. [Predict Sale Price with Categorical features]**

- Prepare X matrix (43 feature columns) and y vector (ie., SalePrice column)
- Split dataset for training and testing
- Create LinearRegression model, fit on training set and predict on test set
- Compute Mean Squared Error on actual values and predicted values (you will get output as 1461036570.0).

**Step6. [Normalize using StandardScaler and Predict Sale Price]**

- Using StandardScaler, perform fit\_transform() on X\_train and transform() on X\_test matrix that you already splitted.

```
new-df = df.drop(['Bldg Type'], axis=1)
```

```
new-df.
```

```
y = new-df['sales price']
```

```
y
```

```
inp = ['fst flrst', 'end flrst', '3.5th porch',
```

```
       'bed room brgr', 'BSM flrst']
```

```
x = new-df[inp]
```

```
x
```

```
from sklearn.model_selection import train-  
test-split.
```

```
x_train, x_test, y_train, y_test = train-test-  
split(x, y, train-size)
```

```
print(x_train.shape)
```

```
print(x_test.shape)
```

```
print(y_train.shape)
```

```
print(y_test.shape)
```

```
from sklearn.linear_model import Linear  
Regression
```

- Create a new Linear Regression model, fit on *scaled X\_train* and *y\_train* and predict on *scaled X\_test*.
- Compute Mean Squared Error (MSE) on actual values and predicted values (you will get output as 1461036570.0).

**Step7. [Normalize using MinMaxScaler and Predict Sale Price]**

- Repeat Step6 using MinMaxScaler
- Mean Squared Error will be: 1461036570.0

**Step8. [Predict using SGD Regressor]**

- Use scaled *X\_train* and *X\_test* using StandardScaler that you computed before
- Create SGDRegressor, fit and predict
- Compute MSE on *y\_test* and *y\_pred* this time. You will get output as 1592430104.0.

**Step8. [Predict using Ridge Regression]**

- Use scaled *X\_train* and *X\_test* using StandardScaler that you computed before
- Create RidgeCV, fit and predict
- Compute MSE on *y\_test* and *y\_pred* this time. You will get output as 1442196000.3367693.

**Step8. [Predict using Lasso Regression]**

- Use scaled *X\_train* and *X\_test* using StandardScaler that you computed before
- Create LassoCV, fit and predict
- Compute MSE on *y\_test* and *y\_pred* this time. You will get output as 1409368613.5329669.

**Step9.[RMSE].** Print Root Mean Squared Error values (use `numpy.sqrt()` method) as below and compare error values.

RMSE without one hot encoding: 38403.0

RMSE with One hot encoding: 38224.0

RMSE with OHE and Standard Scaling: 38224.0

RMSE with OHE and MinMax Scaling: 38224.0

RMSE of SGDRegressor with OHE and Standard Scaler: 38528.0

RMSE of RidgeCV with OHE and Standard Scaler: 37976.0

RMSE of LassoCV with OHE and Standard Scaler: 37542.0

```
model = LinearRegression()
```

```
model.fit(X_train, y_train)
```

```
X_pred = model.predict(X_test)
```

```
Y_pred
```

```
from sklearn.metrics import mean_squared_error
```

```
mean_squared_error(y_test, y_pred)
```

Step 3:-

```
plt.scatter(y_test, y_pred)
```

Step 4:-

```
en_df = pd.get_dummies(df, columns = ['Central Air',  
                                         'Bldg Type'])
```

```
en = df.head()
```

```
en = df.shape
```

```
en = df.columns
```

Step 5:-

```
e_y = en_df['Sales Price']
```

```
e_y
```

```
e_inp = ['1st Floor', '2nd Floor', '3rd Floor',
```

```
         'Bed room Area', 'Bsm Area', 'Enclosed porch',
```

```
         'Fireplaces', 'Full Bath', 'Garage Area',
```

```
         'Lot Area', 'Lot Frontage', 'LowQual Finsh',
```

```
         'HSS Uclass', 'Overall Qual', 'Pool Area', 'Screen  
        porch', 'TotRmsAbvGrd', 'Central Air - N']
```

NOTE

'Central Air - Y', 'BldgType - 1 Fam', 'BldgType - 2 Fam']

e.X = en\_df[e\_Pnp]

e.X\_train, ex\_test, ey\_train, ey\_test = train-test-split(e.X, e.Y)

model 1 = LinearRegression()

model 1.fit(ex\_train, ey\_train)

ey\_pred = model 1.predict(ex\_test)

ey\_pred

Step 6:

from sklearn.preprocessing import StandardScaler.

scaler = StandardScaler()

ss = scaler.fit\_transform(ex\_train)

ss

ss1 = scaler.transform(ex\_test)

ss1

model 2 = LinearRegression()

model 2.fit(ss, ey\_train)

se\_X\_pred = model 2.predict(ss1)

se\_X\_pred

mean-squared-error(ey\_test, se-y-pred)

Step 7:

from sklearn.preprocessing import MinMaxScaler

m\_scaler = MinMaxScaler()

m\_ss = m\_scaler.fit\_transform(ex\_train)

m\_ss

m\_ss1 = m\_scaler.transform(ex\_test)

m\_ss1

model 3 = LinearRegression()

model 3.fit(m\_ss, ey\_train)

me-y-pred = model 3.predict(m\_ss1)

me-x-pred

mean-squared-error(ey-test, me-y-pred)

Step 8:

from sklearn.linear\_model import SGDRegressor

from sklearn.pipeline import make\_pipeline

sgd = make\_pipeline(standard\_scaler(), SGDRegressor

(max\_iter=1000, t

SGD.fit(X, y)

y-y-pred = sgd.predict(X)

y-y-pred

mean-squared-error(ey\_train, ridge-y-pred)

Step 9:

from math import sqrt

print("RMSE without One hot encoding", sqrt(mean-squared-error))

NOTE

```
print("RMSE with one and standard scaling: ", sqrt(mse))  
print("RMSE with one and Min Max scaling: ", sqrt(mse -  
mm))  
print("RMSE of SGD Regressor with one and standard scaler: ",  
sqrt(mse - sgd))  
print("RMSE of Ridge with one and standard scaler: ",  
sqrt(mse - rid))  
print("RMSE of LassoCV with one and  
standard scaler: ", sqrt(mse - las))
```