

Use Of Linear and Logistics Regression Coefficients with Lasso (L1) and Ridge (L2) Regularization for Feature Selection In Machine Learning.

Linear Regression Equation.

```
In [170]: from IPython.display import Image
Image(filename='C:/Users/Microsoft/Desktop/pandas/linear-regression-equation-explained.png',height=400,width=400)
```

$$\hat{y} = \beta_0 + \beta_1 X$$

Predicted value Intercept Slope Predictor

```
In [171]: from IPython.display import Image
Image(filename='C:/Users/Microsoft/Desktop/pandas/47H1Yk.jpg',height=80,width=800)
```

Linear Regression: Single Variable

$$\hat{y} = \beta_0 + \beta_1 x + \epsilon$$

Predicted output Coefficients Input Error

Linear Regression: Multiple Variables

$$\hat{y} = \beta_0 + \beta_1 x_1 + \dots + \beta_p x_p + \epsilon$$

- Linear regression is a straight forward approach for predicting the quantitative response via on the basis on the different predictive variables of X1 and X2 etc.
- It assumes that there is Linear Relationship between the x and y.
- That means the by addition and subtraction y can be predicted based on the input features those input are x1,x2 we can say this is a vector.
- β_0 and β_1 are the coefficient of linear regression.
- ϵ (Epsilon) is the error:- This error having many types so we can say it is it's an error. There are different types of the error which get introduces to the model we train test our model.

Why we use Linear Regression ?

- While working on the linear regression we have to make some basic assumptions those basic assumptions are like.
- We can predict by using the linear combinations of the input vectors and feature space x should have **gaussian distribution**, and features **are not correlated with each other that means there should not be any multi-collinearity**.
- That means the feature space won't be correlated between and themselves.
- Features should be into the same scale x1,xn that means we need to do **Normalization or Standardization** bring this feature into a same scale.
- Since in this lesson we are talking how can get the feature selection based on the coefficient of the linear and logistics regression.
- This coefficient of regression we will calculate by the **LASSO and Ridge Regularization**

What is the Regularization ?

- Regularization is the technique to **discourage the complexity of the model** it does by the **penalizing the loss function**. It helps to **solve the overfitting problem into our machine learning model** and will improve the accuracy and reduce the complexity of the model and will increase the interpretability of the model.
- It adds the penalty on the different parameter of model to reduce the freedom of the model that means to reduce the complexity of model.
- Hence the model will be the less likely to fit the noise of the training data and it will improve the generalization ability of the model.

- L1 Regularization** (also called Lasso)
- L2 Regularization** (also called Ridge)
- L1/L2 Regularization** (also called Elastic Net)

- Lasso is help to shrink the coefficient which are less important to the zero. That means with the help of Lasso regularization we can remove a few features.

- Ridge regression it doesn't reduce the coefficient to zero but it reduce the regression coefficient and with the reduction of the coefficient we can identify which feature have important in final prediction.

For an example:-

- Here, β_1 is the linear or logistic regression coefficient and x1 and xp are the independent features.
- If we have the β_1 is more than the β_0 that time x1 is more important than the xp.
- A regression model that uses L1 regularization Techniques is called as LASSO regression.
- A regression model that uses L2 regularization Techniques is called as Ridge regression.

What is Lasso Regularization ?

- There are three sources of an error.
- That's means we having three types of errors by which the error introduce by testing and training model.
 - Noise
 - Bias
 - Variance
- We can't do anything about the noise but somehow we can do about the Bias and Variance.

Trade Off between the Bias-Variance.

We need to find such sweet spot on that spot we will minimum error and optimal complexity and minimum variance and baise there will found the optimal complexity.

Logic And Mathematics Behind The Lasso Regression.

- Lasso regularization adds penalty equivalent to absolute sum of regression coefficients values.
- In the equations:-
 - Error Function(RSS):- Loss function
 - λ :- Regularization Coefficient or Optimization Coefficient.

- If there are we are doing we add some penalty into a loss function.
- Now we have to find the Required sweet spot where the generation error should be the less and overall ingeneral model perform better.
- By using the L1 over the regression coefficient few parameter will become zero. Whichever the feature doesn't much information to identify the output.

- RSS** :- The residual sum of squares (RSS) is a statistical technique used to measure the amount of variance in a data set that is not explained by a regression model itself. Instead, it estimates the variance in the residuals, or error term.

- Equation = **$RSS(w) + \lambda ||w||$**
 - RSS** :- the loss function is the error function or residual sum square.
 - λ is the Regularization Coefficient or Optimization Coefficient or penalty that Lasso regression add.
 - w is absolute sum of the regression coefficients.

```
In [172]: from IPython.display import Image
Image(filename='C:/Users/Microsoft/Desktop/pandas/regularization-11example.png',height=400,width=400)
```

Weight Matrix:

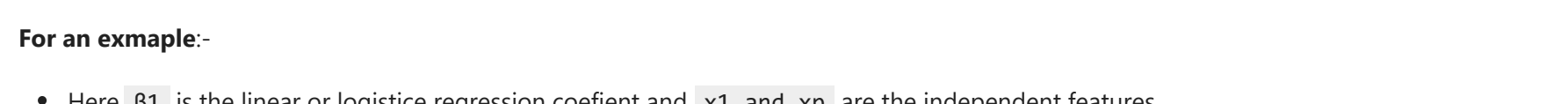
$$\begin{bmatrix} -1 & 5 \\ 3 & -9 \end{bmatrix}$$

L1 Penalty for $\lambda = 0.5$:
absolute value of weight matrix:

$$\begin{bmatrix} 1 & 5 \\ 3 & 9 \end{bmatrix}$$

L1 penalty:
 $\lambda(1 + 5 + 3 + 9) = \lambda(18) = (0.5)(18) = 9$

```
In [173]: from IPython.display import Image
Image(filename='C:/Users/Microsoft/Desktop/pandas/Feature-selection-using-LASSO-regression-model-A-Tuning-parameter.png',height=400,width=400)
```



How we can choose Best Lambda (λ) ?

- Here λ is nothing but **L1 Regularization Coefficient**.
- By doing the **Cross-Validation set** over the data we will get the best lambda value.

- Below we having the steps.**
 - Suppose We have sufficient amount the data.
 - We can split our data into the training set validations and testing set.
 - Then we test model performance on the validation set
 - If there is something which is not going give us good accuracy or something wrong in there so in order to detect that we have to validate the our model validate on the validation set. There we will find something which is good so go ahead then we changes the parameter goes back to the training set and do the optimization on the validation set.
 - Finally our features will do the generalised testing on the test set and this is how we choose the best and optimized lambda during the our regularizations.

What Is Ridge Regularization ?

- Ridge regularization is also know as **L2 regularization** which adds penalty is equal to the **sum of squared value of the coefficients of regression**.
- We add the penalties on the **sum of Squared regression coefficient**.
- Here the λ i.e. L2 Regularization Coefficient works as Tuning Parameter.

Comparison with RSS(Sum of Squared),Lambda AND Weight.

Equation = $RSS(w) + \text{square}(\lambda ||w||)$

Lambda is High :-

- L2 Regularization is high that is Lambda and that time the result we will get the biased somewhere else.
- That means in this case we will get the high bias low variance.
- Under the assumption if lambda is increased upto the infinity that time the w become so small almost zero i.e.0.

Lambda is Low :-

- L2 regularization is low that is Lambda and that time the result face the with number of features because of the L2 Regularization coefficient is low and it adds the less impact over the data.
- Due to this it will end up the overfitting of the due to high dimensional i.e. feature training.
- That will confuse the model.
- Model will train the add point so there will be the high variance low bias.

- λ is the tuning parameter which is performing perfect in our case.
- We can say the lambda is the strength of the regularization.
- L2 parameter we will force the parameter to be relatively small the bigger penalization (and the more robust) the coefficient are.

Difference Between the L1 and L2 Regularization

L1 Regularization:-

- L1 penalizes the sum of absolute values of weights.
- L1 has sparse solutions.
- L1 has multiple solutions.
- L1 has built in feature selection.
- L1 is robust to outliers.
- L1 generates a model that are simple and interpretable but can not learn the complex patterns.

L2 Regularization:-

- L2 Regularizations penalizes the sum of square weights.
- L2 has non sparse solutions.
- L2 has one solution.
- L2 has no feature selections.
- L2 does not robust to outliers.
- L2 gives better predictions even output variable is a function of all input features.
- L2 regularization is able to learn complex data patterns.

Start With Machine Learning.

```
In [13]: import numpy as np
import matplotlib.pyplot as plt
import pandas as pd
import seaborn as sns
import matplotlib inline
```

```
import warnings
warnings.filterwarnings('ignore')
```

```
In [12]: from sklearn.model_selection import train_test_split
from sklearn.feature_selection import SelectKBest
from sklearn.ensemble import RandomForestClassifier
from sklearn.preprocessing import StandardScaler
from sklearn.metrics import accuracy_score,classification_report,confusion_matrix
```

```
In [3]: from sklearn.linear_model import LogisticRegression,LogisticRegressionCV
from sklearn.feature_selection import SelectFromModel
```

```
In [4]: data = pd.read_csv('Data_train_reduced.csv')
```

```
In [6]: pd.options.display_max_columns = None
```

```
In [7]: data.head(1)
```

```
Out[7]:
```

Respondent.ID	Product.ID	Product	Instant.Liking	q1_1.personal.opinion.of.this.Deodorant	q2_all.words	q3_1.strength.of.the.Deodorant	q4
0	3800	121 Deodorant	B	1	4	1	4
1	3801	121 Deodorant	B	0	5	1	4
2	3802	121 Deodorant	B	0	6	1	3
3	3803	121 Deodorant	B	1	4	0	4
4	3804	121 Deodorant	B	1	4	1	2

```
In [8]: data.columns
```

```
Out[8]: Index(['Respondent.ID', 'Product.ID', 'Product', 'Instant.Liking', 'q1_1.personal.opinion.of.this.Deodorant', 'q2_all.words', 'q3_1.strength.of.the.Deodorant', 'q4_all.words', 'q1_2.artificial.chemical', 'q2_1.attractive', 'q3_3.bold', 'q4_3.casual', 'q1_5.masculine', 'q2_2.extractive', 'q3_4.boring', 'q4_4.easy.to.wear', 'q1_6.memorable', 'q2_3.sophisticated', 'q3_5.overpowering', 'q4_5.well.rounded', 'q1_7.natural', 'q2_4.well.rounded', 'q3_6.elegant', 'q4_6.cheap', 'q1_8.old.fashioned', 'q2_5.overpowering', 'q3_7.sharp', 'q4_7.easy.to.wear', 'q1_9.feminine', 'q2_6.addictive', 'q3_8.how.likely.would.you.be.to.purchase.this.Deodorant', 'q4_8.elegant', 'q1_10.feminine', 'q2_7.natural', 'q3_9.how.likely.would.you.be.to.purchase.this.Deodorant', 'q4_9.how.likely.would.you.be.to.purchase.this.Deodorant', 'q1_11.for.someone.like.me', 'q2_8.ethnic.background', 'q3_10.how.likely.would.you.be.to.purchase.this.Deodorant', 'q4_10.feminine', 'q1_12.heavy', 'q2_9.ethnic.background', 'q3_11.for.someone.like.me', 'q4_11.for.someone.like.me', 'q1_13.high.quality', 'q2_10.how.likely.would.you.be.to.purchase.this.Deodorant', 'q4_12.heavy', 'q1_14.long.lasting', 'q2_11.marital.status', 'q3_12.heavy', 'q4_12.heavy', 'q1_15.masculine', 'q2_12.sophisticated', 'q3_13.how.likely.would.you.be.to.purchase.this.Deodorant', 'q4_13.how.likely.would.you.be.to.purchase.this.Deodorant', 'q1_16.memorable', 'q2_13.sophisticated', 'q3_14.how.likely.would.you.be.to.purchase.this.Deodorant', 'q4_14.how.likely.would.you.be.to.purchase.this.Deodorant', 'q1_17.natural', 'q2_14.well.rounded', 'q3_15.masculine', 'q4_15.masculine', 'q1_18.old.fashioned', 'q2_15.masculine', 'q3_16.memorable', 'q4_16.memorable', 'q1_19.ordinary', 'q2_16.addictive', 'q3_17.natural', 'q4_17.natural', 'q1_20.overpowering', 'q2_17.natural', 'q3_18.how.likely.would.you.be.to.purchase.this.Deodorant', 'q4_18.how.likely.would.you.be.to.purchase.this.Deodorant', 'q1_21.sharp', 'q2_18.how.likely.would.you.be.to.purchase.this.Deodorant', 'q3_19.ordinary', 'q4_19.ordinary', 'q1_22.sophisticated', 'q2_19.ordinary', 'q3_20.overpowering', 'q4_20.overpowering', 'q1_23.sharp', 'q2_20.overpowering', 'q3_21.sharp', 'q4_21.sharp', 'q1_24.well.rounded', 'q2_21.marital.status', 'q3_22.sophisticated', 'q4_22.sophisticated', 'q1_25.masculine', 'q2_22.extractive', 'q3_23.sophisticated', 'q4_23.sophisticated', 'q1_26.memorable', 'q2_23.sophisticated', 'q3_24.well.rounded', 'q4_24.well.rounded', 'q1_27.natural', 'q2_24.well.rounded', 'q3_25.masculine', 'q4_25.masculine', 'q1_28.old.fashioned', 'q2_25.masculine', 'q3_26.addictive', 'q4_26.addictive', 'q1_29.ordinary', 'q2_26.addictive', 'q3_27.natural', 'q4_27.natural', 'q1_30.overpowering', 'q2_27.natural', 'q3_28.how.likely.would.you.be.to.purchase.this.Deodorant', 'q4_28.how.likely.would.you.be.to.purchase.this.Deodorant', 'q1_31.marital.status', 'q2_28.how.likely.would.you.be.to.purchase.this.Deodorant', 'q3_29.ordinary', 'q4_29.ordinary', 'q1_32.heavy', 'q2_29.ethnic.background', 'q3_30.how.likely.would.you.be.to.purchase.this.Deodorant', 'q4_30.how.likely.would.you.be.to.purchase.this.Deodorant', 'q1_33.bottles.of.Deodorant.do.you.currently.own', 'q2_30.ethnic.background', 'q3_31.bottles.of.Deodorant.do.you.currently.own', 'q4_31.bottles.of.Deodorant.do.you.currently.own', 'dtype='object']
```

```
In [9]: data.info()
```

```
Out[9]: <class 'pandas.core.frame.DataFrame'>
```

```
RangeIndex: 2500 entries, 0 to 2499
```

```
Data columns (total 64 columns):
```

```
Respondent.ID      2500 non-null int64
```

```
Product.ID         2500 non-null int64
```

```
Product            2500 non-null object
```

```
Instant.Liking     2500 non-null int64
```

```
q1_1.personal.opinion.of.this.Deodorant 2500 non-null int64
```

```
q2_all.words       2500 non-null object
```

```
q3_1.strength.of.the.Deodorant          2500 non-null int64
```

```
q4_all.words       2500 non-null object
```

```
q1_2.artificial.chemical                 2500 non-null int64
```

```
q2_1.attractive    2500 non-null int64
```

```
q3_3.bold          2500 non-null int64
```

```
q4_3.casual        2500 non-null int64
```

```
q1_5.masculine     2500 non-null int64
```

```
q2_2.extractive    2500 non-null int64
```

```
q3_4.boring        2500 non-null int64
```

```
q4_4.easy.to.wear  2500 non-null int64
```

```
q1_6.memorable     2500 non-null int64
```

```
q2_3.sophisticated 2500 non-null int64
```

```
q3_5.overpowering  2500 non-null int64
```

```
q4_5.well.rounded  2500 non-null int64
```

```
q1_7.natural        2500 non-null int64
```

```
q2_4.well.rounded  2500 non-null int64
```

```
q3_6.elegant        2500 non-null int64
```

```
q4_6.cheap          2500 non-null int64
```

```
q1_8.old.fashioned 2500 non-null int64
```

```
q2_5.overpowering  2500 non-null int64
```

```
q3_7.sharp          2500 non-null int64
```

```
q4_7.easy.to.wear  2500 non-null int64
```

```
q1_9.feminine       2500 non-null int64
```

```
q2_6.addictive      2500 non-null int64
```

```
q3_8.how.likely.would.you.be.to.purchase.this.Deodorant 2500 non-null int64
```

```
q4_8.elegant        2500 non-null int64
```

```
q1_10.feminine      2500 non-null int64
```

```
q2_7.natural         2500 non-null int64
```

```
q3_9.how.likely.would.you.be.to.purchase.this.Deodorant 2500 non-null int64
```

```
q4_9.how.likely.would.you.be.to.purchase.this.Deodorant 2500 non-null int64
```

```
q1_11.for.someone.like.me 2500 non-null int64
```

```
q2_8.ethnic.background 2500 non-null int64
```

```
q3_10.how.likely.would.you.be.to.purchase.this.Deodorant 2500 non-null int64
```

```
q4_10.feminine      2500 non-null int64
```

```
q1_12.heavy         2500 non-null int64
```

```
q2_9.ethnic.background 2500 non-null int64
```

```
q3_11.for.someone.like.me 2500 non-null int64
```

```
q4_11.for.someone.like.me 2500 non-null int64
```

```
q1_13.high.quality   2500 non-null int64
```

```
q2_10.how.likely.would.you.be.to.purchase.this.Deodorant 2500 non-null int64
```

```
q4_12.heavy         2500 non-null int64
```

```
q1_14.long.lasting   2500 non-null int64
```

```
q2_11.marital.status 2500 non-null int64
```

```
q3_12.heavy         2500 non-null int64
```

```
q4_12.heavy         2500 non-null int64
```

```
q1_15.masculine      2500 non-null int64
```

```
q2_12.sophisticated 2500 non-null int64
```

```
q3_13.how.likely.would.you.be.to.purchase.this.Deodorant 2500 non-null int64
```

```
q4_13.how.likely.would.you.be.to.purchase.this.Deodorant 2500 non-null int64
```

```
q1_16.memorable      2500 non-null int64
```

```
q2_13.sophisticated 2500 non-null int64
```

```
q3_14.how.likely.would.you.be.to.purchase.this.Deodorant 2500 non-null int64
```

```
q4_14.how.likely.would.you.be.to.purchase.this.Deodorant 2500 non-null int64
```

```
q1_17.natural         2500 non-null int64
```

```
q2_14.well.rounded   2500 non-null int64
```

```
q3_15.masculine      2500 non-null int64
```

```
q4_15.masculine      2500 non-null int64
```

```
q1_18.old.fashioned 2500 non-null int64
```

```
q2_15.masculine      2500 non-null int64
```

```
q3_16.memorable      2500 non-null int64
```

```
q4_16.memorable      2500 non-null int64
```

```
q1_19.ordinary        2500 non-null int64
```

```
q2_16.addictive      2500 non-null int64
```

```
q3_17.natural         2500 non-null int64
```

```
q4_17.natural         2500 non-null int64
```

```
q1_20.overpowering    2500 non-null int64
```

```
q2_17.natural         2500 non-null int64
```

```
q3_18.how.likely.would.you.be.to.purchase.this.Deodorant 2500 non-null int64
```

```
q4_18.how.likely.would.you.be.to.purchase.this.Deodorant 2500 non-null int64
```

```
q1_21.sharp           2500 non-null int64
```

```
q2_18.how.likely.would.you.be.to.purchase.this.Deodorant 2500 non-null int64
```

```
q3_19.ordinary        2500 non-null int64
```

```
q4_19.ordinary        2500 non-null int64
```

```
q1_22.sophisticated   2500 non-null int64
```

```
q2_19.ordinary        2500 non-null int64
```

```
q3_20.overpowering    2500 non-null int64
```

```
q4_20.overpowering    2500 non-null int64
```

```
q1_23.sharp           2500 non-null int64
```

```
q2_21.marital.status 2500 non-null int64
```

```
q3_22.sophisticated   2500 non-null int64
```

```
q4_22.sophisticated   2500 non-null int64
```

```
q1_24.well.rounded    2500 non-null int64
```

```
q2_22.extractive      2500 non-null int64
```

```
q3_23.sophisticated   2500 non-null int64
```

```
q4_23.sophisticated   2500 non-null int64
```

```
q1_25.masculine       2500 non-null int64
```

```
q2_23.sophisticated   2500 non-null int64
```

```
q3_24.well.rounded    2500 non-null int64
```

```
q4_24.well.rounded    2500 non-null int64
```

```
q1_26.memorable       2500 non-null int64
```

```
q2_24.well.rounded    2500 non-null int64
```

```
q3_25.masculine       2500 non-null int64
```

```
q4_25.masculine       2500 non-null int64
```

```
q1_27.natural         2500 non-null int64
```

```
q2_25.masculine       2500 non-null int64
```

```
q3_26.addictive       2500 non-null int64
```

```
q4_26.addictive       2500 non-null int64
```

```
q1_28.old.fashioned 2500 non-null int64
```

```
q2_26.addictive       2500 non-null int64
```

```
q3_27.natural         2500 non-null int64
```

```
q4_27.natural         2500 non-null int64
```

```
q1_29.ordinary        2500 non-null int64
```

```
q2_27.natural         2500 non-null int64
```

```
q3_28.how
```



```
Index(['q1_1.personal.opinion.of.this.Deodorant', 'q2.all.words',
      'q3_1.artificial.chemical', 'q3_2.attractive', 'q3_3.bold',
      'q3_4.boring', 'q3_6.cheap', 'q3_8.easy.to.wear', 'q3_9.elegant',
      'q3_11.for.someone.like.me', 'q3_13.high.quality', 'q3_14.long.lasting',
      'q3_15.muscular', 'q3_16.memorable', 'q3_18.old.fashioned',
      'q3_19.ordinary', 'q3_20.overpowering', 'q3_21.sharp', 'q3_23.upscale',
      'q3_24.well.rounded', 'q3', 'q3.5', 'q3.6', 'q3.12', 'q3.13', 'q3.19',
      'q3_1.time.of.day.would.this.Deodorant.be.appropriate',
      'q3_2.which.occasions.would.this.Deodorant.be.appropriate',
      'q3_3_linking.after.30.minutes',
      'q3_4.Deodorant.oversell.on.a.scale.from.1.to.10', 'Val8eqb',
      's3.ethnic.background', 's10.income', 's11.marital.status',
      's13a.h.most.often', 's13b.bottles.of.Deodorant.do.you.currently.own'],
      dtype='object')
```

```
In [113]: from sklearn.ensemble import RandomForestClassifier

In [114]: classifier = RandomForestClassifier(n_estimators=100, criterion='entropy')
classifier.fit(X_train,y_train_res)
Y_pred = classifier.predict(X_test)
print('Accuracy : ',accuracy_score(y_test,Y_pred))

Accuracy : 1.0

In [115]: classifier = RandomForestClassifier(n_estimators=100, criterion='gini')
classifier.fit(X_train,y_train_res)
Y_pred = classifier.predict(X_test)
print('Accuracy : ',accuracy_score(y_test,Y_pred))

Accuracy : 1.0

In [116]: print('Model score in training data :-',classifier.score(X_train,y_train_res))
print('Model score in testing data :-',classifier.score(X_test,y_test))

Model score in training data :- 1.0
Model score in testing data :- 1.0
```

L2 Regularization.

```
In [146]: sel = SelectFromModel(LogisticRegression(C=0.05, class_weight=None,
sel.fit(X_train_res,y_train_res)

Out[146]: SelectFromModel(estimator=LogisticRegression(C=0.05, class_weight=None,
duid=False, fit_intercept=True,
intercept_scaling=1, l1_ratio=None,
max_iter=100, multi_class='auto',
n_jobs=None, penalty='l2',
random_state=None,
solver='liblinear', tol=0.0001,
verbose=0, warm_start=False),
max_features=None, norm_order=1,
prefit=False, threshold=None)

In [151]: X_train_l2 = sel.transform(X_train_res)
X_test_l2 = sel.transform(X_test_res)

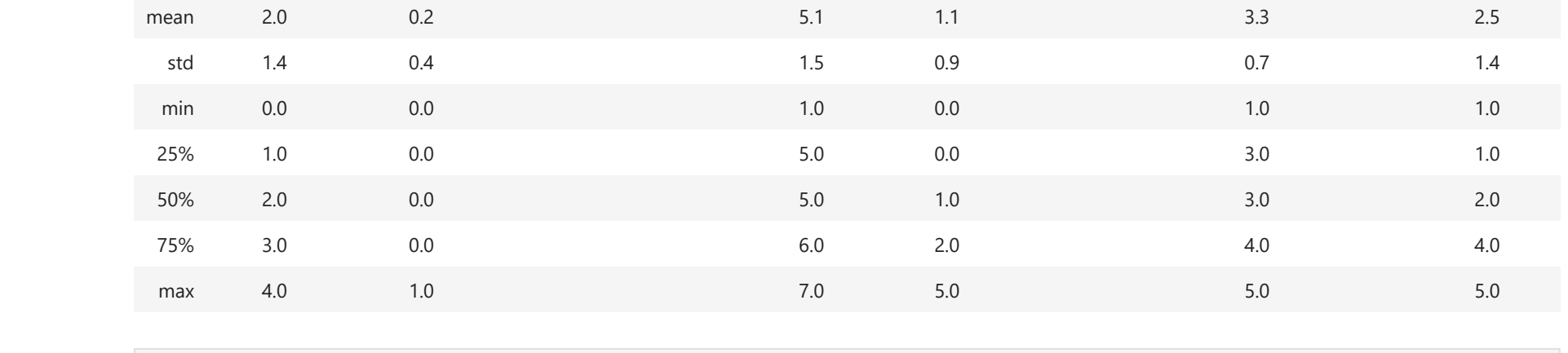
In [153]: %time
classifier = RandomForestClassifier(n_estimators=100, criterion='gini')
classifier.fit(X_train_l2,y_train_res)
Y_pred = classifier.predict(X_test_l2)
print('Accuracy : ',accuracy_score(y_test,Y_pred))

Accuracy : 1.0
Wall time: 632 ms

In [154]: mat = confusion_matrix(y_test,Y_pred)
mat

Out[154]: array([[385, 0],
[ 0, 115]], dtype=int64)

In [155]: sns.heatmap(mat,annot=True,fmt='d')
plt.show()
```



```
print(classification_report(y_test,Y_pred))
```

	precision	recall	f1-score	support
0	1.00	1.00	1.00	385
1	1.00	1.00	1.00	115
accuracy	1.00	1.00	1.00	500
macro avg	1.00	1.00	1.00	500
weighted avg	1.00	1.00	1.00	500

Linear Regression Regularization.

```
In [126]: data1 = pd.read_csv('k5_house_data.csv')

In [127]: data1.head()
```

```
Out[127]:
```

	id	date	price	bedrooms	bathrooms	sqft_living	sqft_lot	floors	waterfront	view	condition	grade	sqft_above
0	7129300520	20141013T000000	2219000	3	1.00	1180	5650	1.0	0	0	3	7	1180
1	6414100192	20141209T000000	5380000	3	2.25	2570	7242	2.0	0	0	3	7	2170
2	5631500400	20150225T000000	1800000	2	1.00	770	10000	1.0	0	0	3	6	770
3	2487200875	20141209T000000	6040000	4	3.00	1960	5000	1.0	0	0	5	7	1050
4	1954400510	20150218T000000	5100000	3	2.00	1680	8080	1.0	0	0	3	8	1680

```
In [128]: data1.tail()
```

```
Out[128]:
```

	id	date	price	bedrooms	bathrooms	sqft_living	sqft_lot	floors	waterfront	view	condition	grade	sqft_ab
21608	263000018	20140521T000000	3600000	3	2.50	1530	1131	3.0	0	0	3	8	1
21609	6600600120	20150223T000000	4000000	4	2.50	2310	5813	2.0	0	0	3	8	2
21610	1523300141	20140623T000000	4021010	2	0.75	1020	1350	2.0	0	0	3	7	1
21611	2913101000	20150116T000000	4000000	3	2.50	1600	2388	2.0	0	0	3	8	1
21612	1523300157	20141015T000000	3250000	2	0.75	1020	1076	2.0	0	0	3	7	1

```
In [129]: data1.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 21613 entries, 0 to 21612
Data columns (total 21 columns):
id                21613 non-null object
price             21613 non-null float64
bedrooms         21613 non-null int64
bathrooms        21613 non-null float64
sqft_living       21613 non-null int64
sqft_lot         21613 non-null float64
floors           21613 non-null float64
waterfront       21613 non-null int64
condition        21613 non-null int64
grade            21613 non-null int64
sqft_above       21613 non-null int64
sqft_basement    21613 non-null int64
yr_built         21613 non-null int64
yr_renovated     21613 non-null int64
zipcode         21613 non-null float64
long            21613 non-null float64
sqft_living15    21613 non-null int64
sqft_lot15      21613 non-null int64
dtype: object
memory usage: 3.4+ MB
```

```
In [130]: data1.shape

Out[130]: (2500, 55)
```

```
In [131]: round(data1.describe(),1)
```

```
Out[131]:
```

	Product	Instant_Living	q1_1.personal.opinion.of.this.Deodorant	q2.all.words	q3_1.strength.of.the.Deodorant	q4_1.artificial.chemical	q4_2
count	2500.0	2500.0		2500.0	2500.0	2500.0	2500.0
mean	2.0	0.2		5.1	1.1	3.3	2.5
std	1.4	0.4		1.5	0.9	0.7	1.4
min	0.0	0.0		1.0	0.0	1.0	1.0
25%	1.0	0.0		5.0	0.0	3.0	1.0
50%	2.0	0.0		5.0	1.0	3.0	2.0
75%	3.0	0.0		6.0	2.0	4.0	4.0
max	4.0	1.0		7.0	5.0	5.0	5.0

```
In [132]: data1.isnull().sum()
```

```
Out[132]: id                0
date                0
price              0
bedrooms          0
bathrooms         0
sqft_living       0
sqft_lot          0
floors            0
waterfront       0
view             0
condition         0
grade            0
sqft_above       0
sqft_basement    0
yr_built         0
yr_renovated     0
zipcode          0
lat              0
long            0
sqft_living15    0
sqft_lot15      0
dtype: int64
```

```
In [133]: len(data1.columns)

Out[133]: 21
```

```
In [134]: data1.drop(columns=['id','date','zipcode'],inplace=True)

Out[134]: data1['yr_renovated'].value_counts().sort_values(ascending=False).head()
```

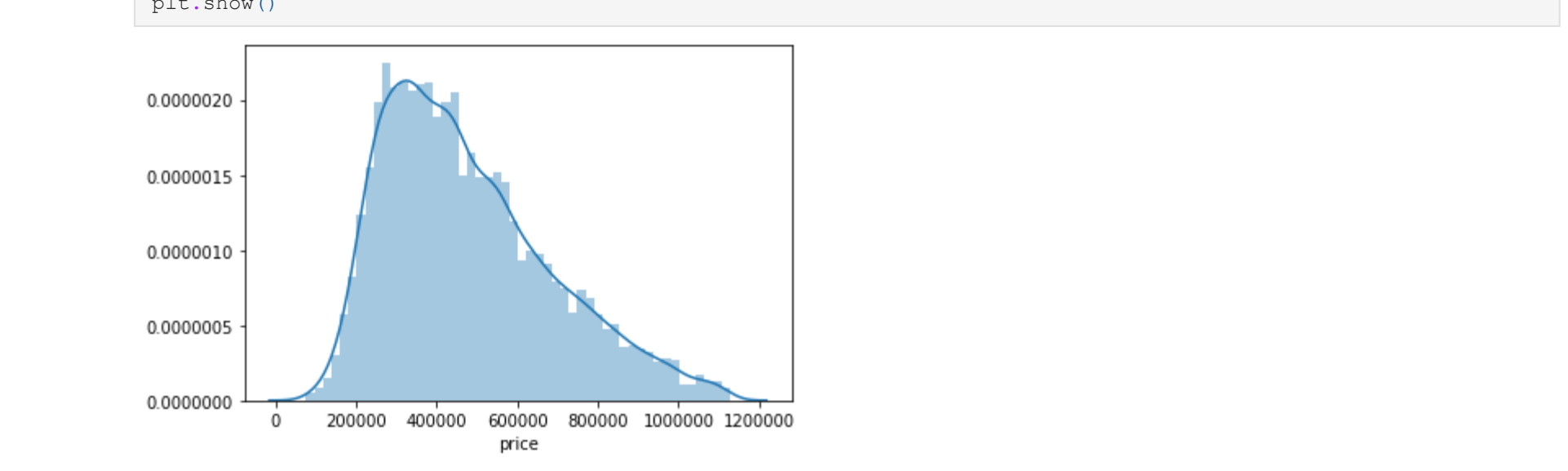
```
20699
2014
2013    37
2003    36
2000    35
Name: yr_renovated, dtype: int64
```

```
In [136]: data1['yr_built'].value_counts().sort_values(ascending=False).head()
```

```
Out[136]: 2014    559
2006    454
2005    450
2004    433
2003    422
Name: yr_built, dtype: int64
```

```
In [137]: data1['price'].hist(figsize=(15,8))

Out[137]: <matplotlib.axes._subplots.AxesSubplot at 0x7f3d170>
```



```
In [138]: x = data1.drop(columns=['price'])
y = data1['price']
```

```
In [139]: X_train,X_test,Y_train,Y_test =train_test_split(x,y,random_state=0,train_size=0.80)
```

```
In [233]: X_train.shape,X_test.shape

Out[233]: ((17290, 17), (4323, 17))
```

```
In [234]: linear= LinearRegression()
```

```
In [235]: sel = SelectFromModel(LinearRegression())
sel.fit(X_train,Y_train)
```

```
Out[235]: SelectFromModel(estimator=LinearRegression(copy_X=True, fit_intercept=True,
n_jobs=None, normalize=False),
max_features=None, norm_order=1,
prefit=False, threshold=None)
```

```
In [236]: sel.estimator_.coef_

array([-3.31909409e+04, 4.03002875e+04, 1.11182207e+02, 1.71360381e-01,
1.80903018e+03, 6.09402810e+05, 4.94058679e+04, 3.07240636e+04,
9.47870908e+04, 7.27869231e+01, 3.83951041e+01, -2.45235095e+03,
2.18772868e+01, 5.63493177e+05, -1.26112606e+05, 2.97732684e+01,
-4.76002806e-01])
```

```
In [237]: X_train_sel = sel.transform(X_train)
X_test_sel = sel.transform(X_test)
```

```
In [238]: from sklearn.metrics import r2_score
from sklearn.tree import DecisionTreeRegressor
```

```
In [239]: model = LinearRegression()
model.fit(X_train_sel,Y_train)
Y_pred = model.predict(X_test_sel)
print('R2 Score :',r2_score(Y_test,Y_pred))
R2 Score : 0.5554708605617198
```

```
In [240]: corr =data1.corr()['price']
corr
```

```
Out[240]: price            1.000000
bedrooms          0.308350
bathrooms         0.223138
sqft_living       0.702035
sqft_lot          0.089661
floors            0.257784
waterfront       0.266369
view             0.397293
condition         0.036362
grade            0.677434
sqft_above       0.603467
sqft_basement    0.323816
yr_built         0.054012
yr_renovated     0.126434
zipcode          0.307003
lat              0.021626
long            0.585779
sqft_living15    0.082447
sqft_lot15      0.082447
Name: price, dtype: float64
```

```
In [241]: good_corr = corr[corr>0.3]
len(good_corr)

Out[241]: 10
```

```
In [242]: corrdata = data1[good_corr.index]
```

```
In [243]: corrdata.head()
```

```
Out[243]:
```

	price	bedrooms	bathrooms	sqft_living	view	grade	sqft_above	sqft_basement	lat	sqft_living15
0	2219000.0	3	1.00	1180	0	7	1180	0	47.5112	1340
1	5380000.0	3	2.25	2570	0	7	2170	400	47.7210	1690
2	1800000.0	2	1.00	770	0	6	770	0	47.7379	2720
3	6040000.0	4	3.00	1960	0	7	1050	910	47.5208	1360
4	5100000.0	3	2.00	1680	0	8	1680	0	47.6168	1800

```
In [244]: plt.figure(figsize=(20,8))
sns.heatmap(corrdata.corr(),annot=True)
plt.show()
```



```
In [245]: q1= corrdata['price'].quantile(0.25)
q3 = corrdata['price'].quantile(0.75)
IQR = q3-q1
```

```
Out[245]: 323050.0
```

```
In [246]: dataaig= corrdata[~(corrdata['price']<q1-1.5*IQR) | (corrdata['price']>q3+1.5*IQR)]
```

```
In [247]: sns.distplot(dataaig['price'])
plt.show()
```



```
In [248]: dataaig['price'].kurt()
```

```
Out[248]: -0.006910700414262116
```

```
In [249]: dataaig['price'].skew()
```

```
Out[249]: 0.7501833824182833
```

```
In [250]: x =dataaig.drop(columns='price')
y = dataaig['price']
```

```
In [251]: X_train,X_test,y_train,y_test = train_test_split(x,y,random_state=42,train_size=0.80)
```

```
Out[251]: X_train.shape,y_train.shape

Out[252]: ((16373, 9), (16373,))
```

```
In [253]: model = LinearRegression()
model.fit(X_train,y_train)
y_pred = model.predict(X_test)
print('R2 score :',r2_score(y_test,y_pred))
```

R2 score : 0.628973577101415

Here the below graph showing the coefficients are getting zero in the case of Lasso Regression whereas in the case of Ridge regression we will not get exactly zero but we will values in the lower values to some values range and from there we can consider which feature is most important for us.

