

Lecture-03

→ Practicals [Linear, Ridge & Lasso, Logistic]

→ Naive Bayes's Intuition.

→ KNN algorithms.

Practicals: [Linear Regression & Ridge and Lasso]

Program:

```
from sklearn.datasets import load_boston
```

```
import numpy as np
```

```
import pandas as pd
```

```
import seaborn as sns
```

```
import matplotlib.pyplot as plt
```

```
% Matplotlib inline
```

```
df = load_boston()
```

```
type(df)
```

```
dataset = pd.DataFrame(df.data)
```

```
dataset.columns = df.feature_names
```

```
dataset.head()
```

Create dependent feature

dataset['price'] = df.target

dataset.head()

Dividing the dataset into independent and dependent features

X = dataset.iloc[:, :-1] # independent features

Y = dataset.iloc[:, -1] # dependent features

X.head()

Y.head()

Linear Regression

from sklearn.linear_model import LinearRegression

from sklearn.model_selection import cross_val_score

lin_reg = LinearRegression()

MSE = cross_val_score(lin_reg, X, Y, scoring = 'neg_mean_squared_error',

CV=5)

MSE # return 5 value

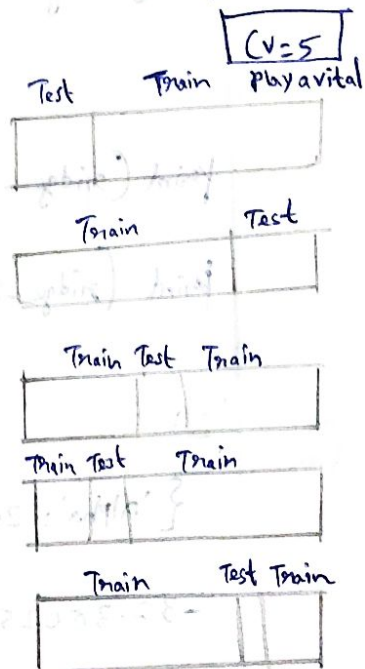
Output:

[-12.46030, -26.0486, -33.074, -80.762,
-33.31360]

Mean_mse = np.mean(MSE)

print(mean_mse)

Output: -37.1318074



To predict with test data.

`lin_reg.predict()` # To pass the test data to make prediction -

Ridge Regression

`from sklearn.linear_model import Ridge`

`from sklearn.model_selection import GridSearchCV`

`GridSearchCV` play among all (α) parameters to give best (α) value for our model.

`ridge = Ridge()`

`Params = {'alpha': [1e-15, 1e-10, 1e-8, 1e-3, 1e-2, 1, 5, 10, 20]}`

`ridge_regressor = GridSearchCV(ridge, Params, Scoring = 'neg-Mean-Squared-Error', CV=5)`

`ridge_regressor.fit(x, y)`

`print(ridge_regressor.best_params_)` # return best α value

`print(ridge_regressor.best_score_)` # return best mse

Output:

`{'alpha': 20}`

`-32.380250`

MSE for both regression,

Linear Reg

(-37)



This is best
because it is low.

Ridge Reg

(-32)

But, Ridge also tries
to reduce overfitting.

In this scenario, we get linear regression because
it is good.

Lasso Regression [Same code for ridge but name change to Lasso]

Also play a n number of α parameters like 30, 35, 40, to give
better accuracy.

We are already know cross validation splits the data
and we also have test-train split.

Code

from sklearn.model_selection import train-test-split

X_train, X_test, y_train, y_test = train-test-split

(X, y, test-size = 0.33, random-state = 42)

Then pass the code like,

mse = cross_val_score (lin-reg, X_train, y_train, Scoring)

Like this ↗

What is difference between Cross validation and Train Test Split ?

Some points,

→ Train-Test Split only splits based on random state and test-size. It gives different accuracy when changing the random state.

→ But, Cross Validation performs correctly based on the CV Count.

Suppose, 10k datapoints

$$CV = \frac{10}{10} = 1 \text{ [for test]} \quad 10\%$$

$$10 - 1 = 9 \text{ [for train]} \quad 90\%$$

It gives n number of accuracy based on CV Count and take better accuracy with the help of mean.

Remember:

The performance metrics is evaluated by the test data.

MSE, MAE, RMSE \rightarrow It is used for mainly to

evaluate the prediction error rates. Lower MSE, MAE,

RMSE the closer is to [i.e. cost function = 0]. It is better fit.

R^2 and adjusted R^2 :

The R^2 and adjusted R^2 is used to evaluate model performance like 80%, 95% and so on.

Both are regression evaluation metrics based on the usage we can use. Mostly, MSE, MAE, RMSE is used only for evaluate prediction rates. R^2 and adjusted R^2 are used to model performance.

Code:

```
# Find Accuracy of the Model
```

```
y_pred = ridge_regressor.predict(x_test)
```

```
from sklearn.metrics import r2_score
```

```
r2_score1 = r2_score(y_pred, y_test)
```

```
Print (r2_score1).
```

Output: 0.650955

65% Accuracy

Logistic Regression:

Code:

```
from sklearn.linear_model import LogisticRegression
```

```
from sklearn.datasets import load_breast_cancer
```

```
df = load_breast_cancer()
```

Independent feature

```
X = pd.DataFrame(df['data'], columns=df['feature_names'])
```

```
X.head()
```

Dependent feature

```
Y = pd.DataFrame(df['target'], columns=['Target'])
```

Y

Check dataset is balanced or imbalanced

```
Y['Target'].value_counts()
```

```
out: 1    357
```

```
0    212
```

The dataset is balanced One . . .

Train Test Split

```
from sklearn.model_selection import train_test_split
```

```
X_train, X_test, y_train, y_test = train_test_split
```

```
(X, y, test_size = 0.33, random_state = 42)
```

```
Params = [{ 'C': [1, 5, 10] }, { 'max_iter': [100, 150] }]
```

```
Model1 = LogisticRegression (C=100, max_iter=100)
```

```
Model = GridSearchCV (Model1, param_grid = Params, Scoring = 'f1',  
(V=5)
```

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

```
# Check best Params
```

```
Model.best_params_
```

```
Out: { 'max_iter': 150 }
```

```
Model.best_score
```

```
Out: 0.95751136
```

```
X_test y_pred = Model.predict (X_test)
```

```
y_pred
```


Confusion Matrix

Confusion_matrix (y-test, y-pred)

Out: array ([[63, 4],

[3, 118]], dtype = int64)

Evaluate Metrics

from sklearn.metrics import Confusion_matrix, Classification_report, accuracy_score

accuracy_score (y-test, y-pred)

Classification_report (y-test, y-pred) # This show all the Precision, Recall and F1-Score

Remember:

Max_iter default value is 100, why set the max_iter.

Because of Some time leads to infinite loop and

May be not converge.

C → Inverse of regularization strength

↳ weak Regularization

↳ Strong Regularization

Weak Regularization:

- The weak regularization, the penalty term added to the loss function is relatively small.
- It means the model's coefficient (or) slope can take larger values, allowing to fit the training data more closely.
- It is suitable for large and complex datasets and there's a need to capture intricate patterns in data.
- However, when applied, the model still overfits and performs poorly on new data.

Strong Regularization:

- The larger penalty term in the loss function.
- The penalty discourages the model coefficients from taking large values, reducing their impact on prediction.
- It is useful for small and noisy datasets and also prevents overfitting.
- It encourages the model to focus on the most important features and prevents overfitting.

Naive Bayes Intuition: [Classification]

This algorithm is used for classification

Problem.

This algorithm is worked with the help of Baye's theorem.

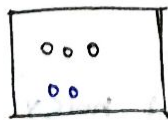
Rolling a Dice, [Independent events]

$\{1, 2, 3, 4, 5, 6\}$

$$P(1) = \frac{1}{6}, P(2) = \frac{1}{6}, P(3) = \frac{1}{6}, P(4) = \frac{1}{6}$$

So, this event is called as the independent event.

Dependent events:



○ → Blue Ball

● → Black Ball

$$P(\text{Black Ball}) = \frac{3}{5}$$

after that draw the,

$$P(\text{Blue Ball}) = \frac{2}{4} = \frac{1}{2}$$

This is dependent events.

Suppose we take,

$$P(\text{Black and Blue}) = P(\text{Black}) * P(\text{Blue} | \text{Black})$$

→ It is something called
Conditional probability

So, let's write

$$P(A \text{ and } B) = P(A) * P(B|A)$$

Let's derive something,

$$P(A \text{ and } B) = P(B \text{ and } A) \quad \left[\text{Both are equal, so just take this} \right]$$

$$P(A) * P(B|A) = P(B) * P(A|B)$$

$$P(B|A) = \frac{P(B) * P(A|B)}{P(A)}$$

→ Baye's theorem
So, This CRUX behind
Naive Bayes

How to apply in real time?

Consider,

$x_1 \quad x_2 \quad x_3 \quad x_4 \quad x_5 \quad \dots \quad x_n$

↳ I/P features

y

dependant or O/P

So, let's write

$$P(y | x_1 x_2 x_3 x_4 \dots x_n) = \frac{P(y) * P(x_1 x_2 x_3 \dots x_n | y)}{P(x_1 x_2 x_3 \dots x_n)}$$

$$= \frac{P(y) * P(x_1/y_1) * P(x_2/y_2) * P(x_3/y_3) \dots P(x_n/y_n)}{P(x_1) * P(x_2) * P(x_3) * \dots * P(x_n)}$$

Suppose the dataset contains,

x_1	x_2	x_3	x_4	y
Independent O/P				dependent O/P
-	-	-	-	Yes
-	-	-	-	No

$$P(y=Yes | x_i) = \frac{P(Yes) * P(x_1/Yes) * P(x_2/Yes) * P(x_3/Yes) * P(x_4/Yes)}{P(x_1) * P(x_2) * P(x_3) * P(x_4)}$$

↳ Fixed [constant]

$$P(y=No | x_i) = \frac{P(No) * P(x_1/No) * P(x_2/No) * P(x_3/No) * P(x_4/No)}{P(x_1) * P(x_2) * P(x_3) * P(x_4)}$$

↳ Fixed [constant]

x_i $\begin{cases} \rightarrow \text{Yes} \\ \rightarrow \text{No} \end{cases}$

Suppose, we have

$$P(Yes | x_i) = 0.13$$

$$P(No | x_i) = 0.05$$

In binary classification, we decide

$$\geq 0.5 \Rightarrow 1$$

$$< 0.5 \Rightarrow 0$$

[0.5 is the threshold for

Comparing binary Classification]

So, let's apply normalization

$$P(\text{Yes} | x_i) = \frac{0.13}{0.13 + 0.05} = 0.72 \Rightarrow 72\%$$

$$P(\text{No} | x_i) = 1 - 0.72 = 0.28 \Rightarrow 28\%$$

This is known as the intuition of Naive Bayes.

Real world Example:-

Consider the DATASET,

DAY	OUTLOOK	TEMPERATURE	HUMIDITY	WIND	PLAY TENNIS
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

So, let's take Outlook feature

Outlook:

	Yes	No	$P(y)$	$P(N)$
Sunny	2	3	$\frac{2}{9}$	$\frac{3}{5}$
Overcast	4	0	$\frac{4}{9}$	$\frac{0}{5}$
Rain	3	2	$\frac{3}{9}$	$\frac{2}{5}$
Total	9	5		

$P(y)$ → Probability of yes for given outlook values.

$P(N)$ → Probability of No for given outlook values.

So, let's take the temperature feature

Temperature:

	Yes	No	$P(y)$	$P(N)$
Hot	2	2	$\frac{2}{9}$	$\frac{2}{5}$
Mild	4	2	$\frac{4}{9}$	$\frac{2}{5}$
Cold	3	1	$\frac{3}{9}$	$\frac{1}{5}$
Total	9	5		

Then also take a PLAY feature

↳ dependent feature

So, Conclude that

Play:

Yes 9

$$P(\text{Yes}) = 9/14$$

No 5

$$P(\text{No}) = 5/14$$

Total 14

Suppose the new data will be Come,

Test

→ (Sunny, Hot) → what is the O/P?

So, let's find out

$$P(\text{Yes} | (\text{Sunny}, \text{Hot})) = \frac{P(\text{Yes}) * P(\text{Sunny} | \text{Yes}) * P(\text{Hot} | \text{Yes})}{P(\text{Sunny}) * P(\text{Hot})}$$

$$P(\text{Sunny}) * P(\text{Hot})$$

$$= \frac{9/14 * 2/9 * 2/9}{5/14 * 4/14}$$

$$= \frac{2/9}{5/14 * 4/14}$$

[That's why Just ignore because it is fixed]

$$= \frac{2/9}{5/14 * 4/14}$$

$$= \frac{2}{63} \Rightarrow 0.031$$

$$P(\text{No/Sunny, hot}) = \frac{P(\text{No}) * P(\text{Sunny/No}) * P(\text{hot/No})}{P(\text{Sunny}) * P(\text{hot})}$$

$$P(\text{Sunny}) * P(\text{hot}) \rightarrow \text{Just ignore [It's fixed]}$$

$$= \frac{5}{14} * \frac{3}{8} * \frac{1}{5}$$

$$= \frac{3}{35} = 0.085$$

So,

$$P(\text{yes/Sunny, hot}) = 0.031$$

$$P(\text{No/Sunny, hot}) = 0.085$$

The final thing is Normalize, (joint, priors) ←

$$\text{Normalized } P(\text{yes/Sunny, hot}) = \frac{P(\text{yes/Sunny, hot})}{P(\text{yes/Sunny, hot}) + P(\text{No/Sunny, hot})}$$

$$\frac{(0.031) + (0.085)}{0.031 + 0.085} = \frac{0.031}{0.116}$$

$$= 0.267241$$

$$= 0.267241 \Rightarrow 27\%$$

$$\Rightarrow 27\%$$

$$\text{Normalized } P(\text{No/Sunny, hot}) = \frac{P(\text{No/Sunny, hot})}{P(\text{yes/Sunny, hot}) + P(\text{No/Sunny, hot})}$$

$$= \frac{0.085}{0.116} = 0.732759 \Rightarrow 73\%$$

Just apply formula (or) Subtract whole from 1.

$$= 1 - 0.267241$$

$$= 0.732 \Rightarrow 73\%$$

So,

$$P(\text{Yes} / \text{Sunny, hot}) = 27\%$$

$$P(\text{No} / \text{Sunny, hot}) = 73\%$$

So, the threshold value is 0.5,

$$P(\text{Yes} / \text{Sunny, hot}) < 0.5$$

So, obviously take,

$$P(\text{No} / \text{Sunny, hot}) = 73\%$$

Therefore,

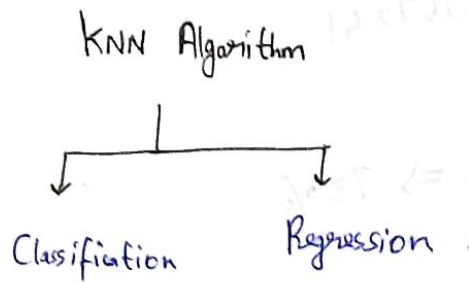
If a person play with (Sunny, hot) \rightarrow O/P?

The probability of No is high.

So, The Answer is No.

This is all about the Naive Bayes for binary Classification.

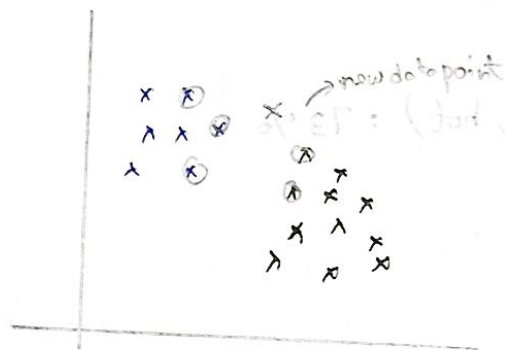
KNN Algorithm:



→ KNN means k-nearest neighbours.

→ It also support to solve both regression and Classification problem.

Classification:



Suppose, $k=5$

$k=5$

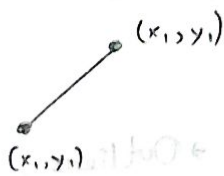
Take 5 nearest points and calculate the distance.


In this case, maximum point taking from blue group. So, the new point consider as a blue group.

Blue point - 3 (obviously this is less distance)

Black point - 2

The Calculation is done by two distance formula,

→ Euclidean Distance . 

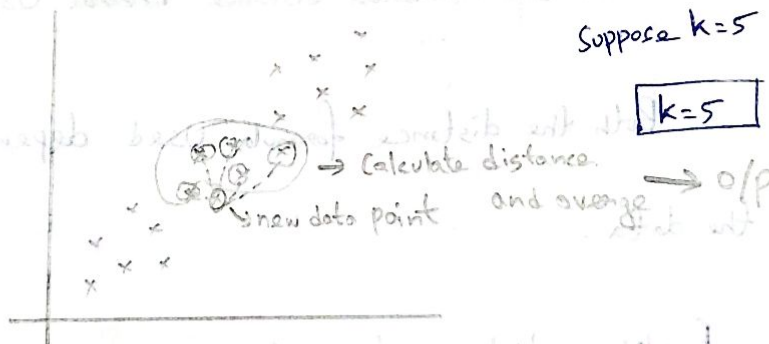
→ Manhattan Distance . 

How Categorical feature plotted in KNN?

Before applying KNN, The various feature encoding techniques will be used. This will be used to Convert Categorical feature into Numerical values.

Then plot KNN .

Regression:



Take nearest five points distance and Calculate the average . This will be prediction.

k is the hyperparameter .

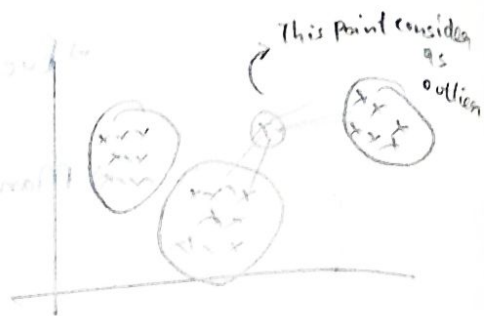
Apply the k value 1 to 50 and check the error rate .

If error rate is low, take this model for analysis .

K-nearest is working bad in two cases are,

→ Outliers.

→ Imbalanced dataset.



Formula,

$$\text{Euclidean distance} = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$$

$$\text{Manhattan distance} = |(x_2 - x_1) + (y_2 - y_1)| \quad (\text{or})$$

$$|x_1 - x_2| + |y_1 - y_2|$$

When Euclidean and Manhattan distance ~~work~~ Use?

Both the distance formula used depends nature of the data.

Euclidean distance tends to work well when the data is distributed in a more Circular Manner.

While, The Manhattan distance can be more

Suitable for Grid-like (or) lattice like structure.

(It means where values are arranged in rows and columns).