

4/11/22
Friday

Naïve Bayes

Date: _____
Page: _____

Topic 17 Naïve Bayes — Naïve Bayes classifier is a machine learning algorithm based on Supervised learning, that can be used for solving classification problem.

It is a probabilistic classifier, which means it predicts on the basis of the probability of an object.

i) Comprises of two Names Naïve Bayes?

- Naïve — It is called Naïve because it makes the assumption that all attributes are independent of each other.
- Bayes — It's called Bayes because it depends on the principle of 'Bay's' theorem.

ii) Where Naïve Bayes is used? —

- Face Recognition Software
- Weather Forecast
- News categorization (Eg - political news, ^{Film} news, ^{Sport} news)
- Real-time prediction

(iii) Independence of Naive Bayes

$$\begin{aligned} \Rightarrow P(A \& B) &= P(A) * P(B/A) \\ &= P(B \& A) = P(B) * P(A/B) \\ \Rightarrow P(A) * P(B/A) &= P(B) * P(A/B) \end{aligned}$$

$$\frac{P(B/A)}{P(A)} = \frac{P(B) * P(A/B)}{P(A)} \Rightarrow \text{Naive Bayes} \Rightarrow \text{Bayes's Theorem}$$

$P(B/A)$ = B given A = Conditional Probability

(iv) Types of Naive Bayes classifier -

- Multinomial Naive Bayes
- Bernoulli Naive Bayes
- Gaussian Naive Bayes - Normal Distribution

Code of Naives Bayes -

Code no ①

```
>>> import numpy as np
>>> import pandas as pd
```

Code no ②

```
>>> from sklearn.datasets import load_breast_cancer
>>> data = load_breast_cancer()
```

Code no ③

```
>>> data.data
>>> data.feature_names
>>> data.target
>>> data.target_names
```

Code no ④

∴ Our dataset is in array. So we have to change it into Dataframe as we directly import it from sklearn.

```
>>> df = pd.DataFrame(np.c_[data.data, data.target],
                      columns = list(data.feature_names) + ['target'])
```

Code no ⑤

```
>>> X = df.iloc[:, 0:-1]    (# independence variable)
>>> y = df.iloc[:, -1]     (# dependence variable)
```

Code no ⑥

```
>>> from sklearn.model_selection import train_test_split
>>> x_train, x_test, y_train, y_test = train_test_split(
    X, y, test_size=0.2, random_state=2020)
```


import naive bayes & train the data —

```

code ⑦ >>> from sklearn.naive_bayes import GaussianNB
>>> classifier = GaussianNB()
>>> classifier.fit(X_train, y_train)

code ⑧ >>> classifier.score(X_test, y_test)
>>> classifier.predict(X_test)
>>> y_test

```

(to compare)

Note :- we can also check ✓
check other naive bayes
classifier like MultinomialNB or
BernoulliNB.

Naive Bayes

$$P(A|B) = \frac{P(A|B) \times P(B)}{P(A)} = 1 \text{ for Categorical dataset}$$

$$P(A|B) = \frac{P(B|A) \times P(A)}{P(B)} = \text{for Categorical dataset}$$

KNN

1. KNN \rightarrow K-Nearest Neighbor is Supervised ML algorithm that can perform both classification and regression tasks using number (K) of neighbors (Instance).

KNN - Classifier

- K value should be odd.

- Steps of prediction of New instance -

- Get Data

- Define K Neighbors

- Calculate the Neighbors distance.

- Assign New instances to Majority of Neighbors.

To calculate the distance -

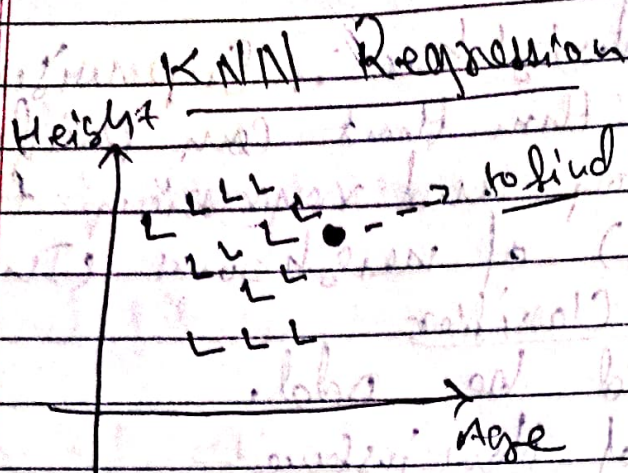
$$\text{Euclidean Distance} = d = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$$

$$\text{Manhattan dist} = d = |x_2 - x_1| + |y_2 - y_1|$$

1. \rightarrow Mode

```
>>> from sklearn.neighbors import KNeighborsClassifier  
>>> classifier = KNeighborsClassifier(n_neighbors=5)  
>>> classifier.fit(x_train, y_train)
```

Note If there is so much Outlier in data set then KNN will not work



If $K = 5$ then we will find average of 5 point near the predid value.

```
>>> from sklearn.neighbors import KNeighborsRegressor
>>> regression = KNeighborsRegressor(n_neighbors=5)
>>> regression.fit(x_train, y_train)
```

In KNN Regression the output will be the Average of the $(K-5)$ & K nearest value.

Limitation of KNN

- we can't use in large dataset
- Outlier & Sensitivity
- Sensitive to missing value

Parameter in KNN

Weight parameter in KNN:- If distance is less weight may be high if distance is more then weight will be less

p = take a distance = $\cdot 1$ = for ^{manhattan distance} ~~Manhattan~~
 2 = for euclidean distance

• We can use Hyper parameter tuning after use KNN with Grid Search CV & Random Search CV.