



Complete Data Science Course

\$147
for just \$93

[View Course](#)

Python

Decorators in Python –
How to enhance functions
without changing the
code?

Generators in Python –
How to lazily return values
only when needed and
save memory?

Iterators in Python – What
are Iterators and Iterables?

Python Module – What are
modules and packages in
python?

Object Oriented
Programming (OOPS) in
Python

How to create conda
virtual environment

How to use Numpy
Random Function in
Python

How Naive Bayes Algorithm Works? (with example and full code)

- by Selva Prabhakaran

Naive Bayes is a probabilistic machine learning algorithm based on the Bayes Theorem, used in a wide variety of classification tasks.

In this post, you will gain a clear and complete understanding of the Naive Bayes algorithm and all necessary concepts so that there is no room for doubts or gap in understanding.

Contents

-
1. Introduction
 2. What is Conditional Probability?

cProfile – How to profile your python code

Dask Tutorial – How to handle big data in Python

Numpy Reshape – How to reshape arrays and what does -1 mean?

Modin – How to speedup pandas

What does Python Global Interpreter Lock – (GIL) do?

Python Yield – What does the yield keyword do?

Lambda Function in Python – How and When to use?

Investor's Portfolio Optimization with Python

datetime in Python – Simplified Guide with Clear Examples

Python Collections – Complete Guide

pdb – How to use Python debugger

Python JSON – Guide

How to use tf.function to speed up Python code in Tensorflow

List Comprehensions in Python – My Simplified Guide

Mahalonobis Distance – Understanding the math with examples (python)

Parallel Processing in Python – A Practical Guide with Examples

Python @Property Explained – How to Use and When? (Full Examples)

Python Logging – Simplest Guide with Full Code and

3. The Bayes Rule
4. The Naive Bayes
5. Naive Bayes Example by Hand
6. What is Laplace Correction?
7. What is Gaussian Naive Bayes?
8. Building a Naive Bayes Classifier in R
9. Building Naive Bayes Classifier in Python
10. Practice Exercise: Predict Human Activity Recognition (HAR)
11. Tips to improve the model

1. Introduction

Naive Bayes is a probabilistic machine learning algorithm that can be used in a wide variety of classification tasks.

Typical applications include filtering spam, classifying documents, sentiment prediction etc. It is based on the works of Rev. Thomas Bayes (1702) and hence the name.

But why is it called 'Naive'?

The name `naive` is used because it assumes the features that go into the model is independent of each other. That is changing the value of one feature, does not directly influence or change the value of any of the other features used in the algorithm.

Alright. By the sounds of it, Naive Bayes does seem to be a simple yet powerful algorithm. But why is it so popular?

That's because there is a significant advantage with NB. Since it is a probabilistic model, the algorithm can be coded up easily and the predictions made real quick. Real-time quick.

Because of this, it is easily scalable and is traditionally the algorithm of choice for real-world applications

Examples

Python Regular

Expressions Tutorial and Examples: A Simplified Guide

Requests in Python

Tutorial – How to send HTTP requests in Python?

Simulated Annealing Algorithm Explained from Scratch

Setup Python environment for ML

Numpy.median() – How to compute median in Python

add Python to PATH – How to add Python to the PATH environment variable in Windows?

Install pip mac – How to install pip in MacOS?: A Comprehensive Guide

Install opencv python – A Comprehensive Guide to Installing "OpenCV-Python"

(apps) that are required to respond to user's requests instantaneously.

But before you go into Naive Bayes, you need to understand what '**Conditional Probability**' is and what is the '**Bayes Rule**'.

And by the end of this tutorial, you will know:

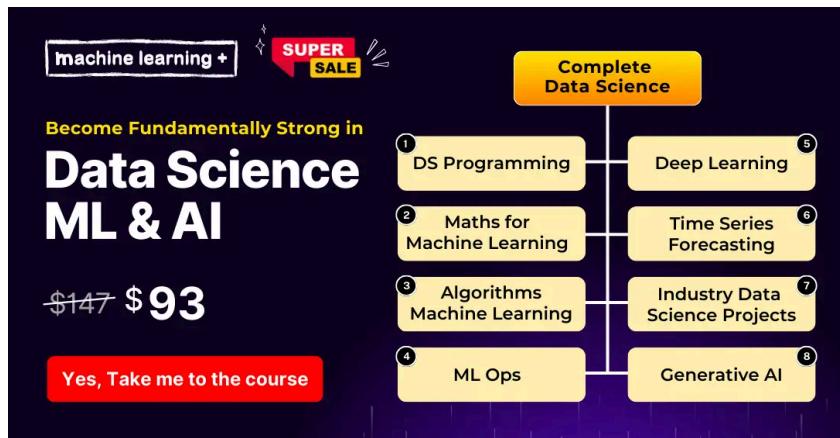
- How exactly Naive Bayes Classifier works step-by-step
- What is Gaussian Naive Bayes, when is it used and how it works?
- How to code it up in R and Python
- How to improve your Naive Bayes models?

Cool? Let's begin.

Also: You might enjoy our **Industrial project course based on a real world problem**. It comes with a Full Hands-On Walk-through of multiple ML solution strategies: [Microsoft Malware Detection](#).

2. What is Conditional Probability?

Lets start from the basics by understanding conditional probability.



Coin Toss and Fair Dice Example When you flip a fair coin, there is an equal chance of getting either heads or tails. So you can say the probability of getting heads is 50%. Similarly what would be the probability of getting

a 1 when you roll a dice with 6 faces? Assuming the dice is fair, the probability of $1/6 = 0.166$.

Alright, one final example with playing cards.

Playing Cards Example If you pick a card from the deck, can you guess the probability of getting a queen given the card is a spade? Well, I have already set a condition that the card is a spade.

So, the denominator (eligible population) is 13 and not 52. And since there is only one queen in spades, the probability it is a queen given the card is a spade is $1/13 = 0.077$

This is a classic example of conditional probability.

So, when you say the conditional probability of A given B, it denotes the probability of A occurring given that B has already occurred.

Mathematically, Conditional probability of A given B can be computed as: $P(A|B) = P(A \text{ AND } B) / P(B)$ **School Example**

Let's see a slightly complicated example.

Consider a school with a total population of 100 persons. These 100 persons can be seen either as 'Students' and 'Teachers' or as a population of 'Males' and 'Females'.

With below tabulation of the 100 people, what is the conditional probability that a certain member of the school is a 'Teacher' given that he is a 'Man'?

	Female	Male	Total
Teacher	8	12	20
Student	32	48	80
Total	40	60	100

To calculate this, you may intuitively filter the sub-population of 60 males and focus on the 12 (male) teachers. So the required conditional probability

$$P(\text{Teacher} \mid \text{Male}) = 12 / 60 = 0.2.$$

$$P(\text{Teacher} \mid \text{Male}) = \frac{P(\text{Teacher} \cap \text{Male})}{P(\text{Male})} = 12/60 = 0.2$$

This can be represented as the intersection of Teacher (A) and Male (B) divided by Male (B). Likewise, the conditional probability of B given A can be computed. The Bayes Rule that we use for Naive Bayes, can be derived from these two notations.

$$P(A \mid B) = \frac{P(A \cap B)}{P(B)} \quad (1)$$

$$P(B \mid A) = \frac{P(A \cap B)}{P(A)} \quad (2)$$

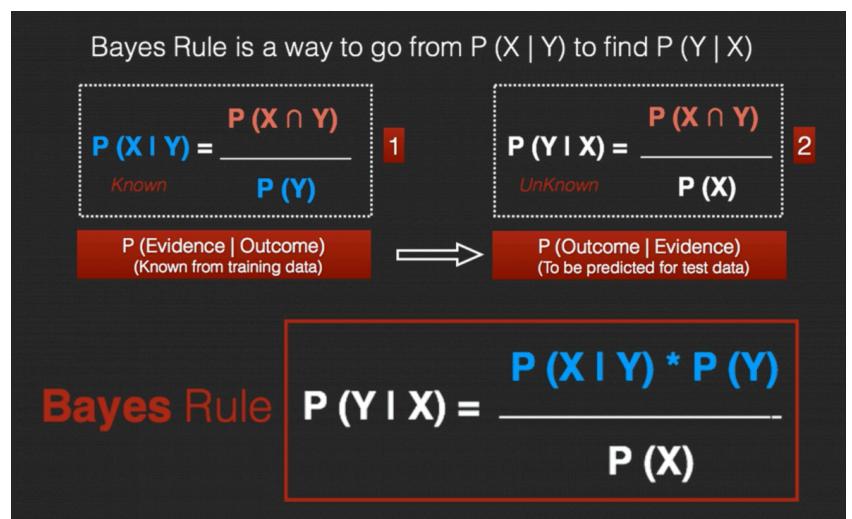
3. The Bayes Rule

The Bayes Rule is a way of going from $P(X|Y)$, known from the training dataset, to find $P(Y|X)$.

To do this, we replace A and B in the above formula, with the feature X and response Y.

For observations in test or scoring data, the X would be known while Y is unknown. And for each row of the test dataset, you want to compute the probability of Y given the X has already happened..

What happens if Y has more than 2 categories? we compute the probability of each class of Y and let the highest win.



4. The Naive Bayes

The Bayes Rule provides the formula for the probability of Y given X.

But, in real-world problems, you typically have multiple X variables.

When the features are independent, we can extend the Bayes Rule to what is called Naive Bayes.

It is called 'Naive' because of the naive assumption that the X's are independent of each other.

Regardless of its name, it's a powerful formula.

When there are multiple X variables, we simplify it by assuming the X's are independent, so the **Bayes** rule

$$P(Y=k | X) = \frac{P(X | Y=k) * P(Y=k)}{P(X)}$$

where, k is a class of Y

becomes, Naive **Bayes**

$$P(Y=k | X_1..X_n) = \frac{P(X_1 | Y=k) * P(X_2 | Y=k) ... * P(X_n | Y=k) * P(Y=k)}{P(X_1) * P(X_2) ... * P(X_n)}$$

$$P(Y=k | X_1..X_n) = \frac{P(X_1 | Y=k) * P(X_2 | Y=k) ... * P(X_n | Y=k) * P(Y=k)}{P(X_1) * P(X_2) ... * P(X_n)}$$

can be understood as ..

$$\text{Probability of Outcome | Evidence} \quad (\text{Posterior Probability}) = \frac{\text{Probability of Likelihood of evidence} * \text{Prior}}{\text{Probability of Evidence}}$$

Probability of Evidence is same for all classes of Y

In technical jargon, the left-hand-side (LHS) of the equation is understood as the **posterior probability** or simply the **posterior**.

The RHS has 2 terms in the numerator. The first term is called the '**Likelihood of Evidence**'. It is nothing but the conditional probability of each X's given Y is of particular class 'c'. Since all the X's are assumed to be independent of each other, you can just multiply the

'likelihoods' of all the X's and called it the '**Probability of likelihood of evidence**'.

This is known from the training dataset by filtering records where $Y=c$. The second term is called the prior which is the overall probability of $Y=c$, where c is a class of Y. In simpler terms, `Prior = count(Y=c) / n_Records`.

An example is better than an hour of theory. So let's see one.

5. Naive Bayes Example by Hand

Say you have 1000 fruits which could be either 'banana', 'orange' or 'other'. These are the 3 possible classes of the Y variable. We have data for the following X variables, all of which are binary (1 or 0).

- Long
- Sweet
- Yellow

The first few rows of the training dataset look like this:

Fruit	Long (x1)	Sweet (x2)	Yellow (x3)
Orange	0	1	0
Banana	1	0	1
Banana	1	1	1
Other	1	1	0
..

For the sake of computing the probabilities, let's aggregate the training data to form a counts table like

this.

Type	Long	Not Long	Sweet	Not Sweet	Yellow	Not Yellow	Total
Banana	400	100	350	150	450	50	500
Orange	0	300	150	150	300	0	300
Other	100	100	150	50	50	150	200
Total	500	500	650	350	800	200	1000

So the objective of the classifier is to predict if a given fruit is a 'Banana' or 'Orange' or 'Other' when only the 3 features (long, sweet and yellow) are known.

Let's say you are given a fruit that is: Long, Sweet and Yellow, can you predict what fruit it is?

This is the same of predicting the Y when only the X variables in testing data are known.

Let's solve it by hand using Naive Bayes. The idea is to compute the 3 probabilities, that is the probability of the fruit being a banana, orange or other. Whichever fruit type gets the highest probability wins.

All the information to calculate these probabilities is present in the above tabulation.

Step 1: Compute the 'Prior' probabilities for each of the class of fruits. That is, the proportion of each fruit class out of all the fruits from the population.

You can provide the 'Priors' from prior information about the population. Otherwise, it can be computed from the training data. For this case, let's compute from the training data. Out of 1000 records in training data, you have 500 Bananas, 300 Oranges and 200 Others.

So the respective priors are 0.5, 0.3 and 0.2.

$$P(Y=\text{Banana}) = 500 / 1000 = 0.50 \quad P(Y=\text{Orange}) = 300 / 1000 = 0.30 \quad P(Y=\text{Other}) = 200 / 1000 = 0.20$$

Step 2: Compute the probability of evidence that goes in the denominator. This is nothing but the product of P of Xs for all X. This is an optional step because the denominator is the same for all the classes and so will not affect the probabilities. $P(x_1=\text{Long}) =$

$$500 / 1000 = 0.50 \quad P(x_2=\text{Sweet}) = 650 / 1000 = 0.65$$

$$P(x_3=\text{Yellow}) = 800 / 1000 = 0.80$$

Step 3: Compute the probability of likelihood of evidences that goes in the numerator. It is the product of conditional probabilities of the 3 features. If you refer back to the formula, it says $P(X_1 | Y=k)$.

Here X_1 is 'Long' and k is 'Banana'.

That means the probability the fruit is 'Long' given that it is a Banana. In the above table, you have 500 Bananas. Out of that 400 is long.

So, $P(\text{Long} | \text{Banana}) = 400/500 = 0.8$. Here, I have done it for Banana alone.

Probability of Likelihood for Banana $P(x_1=\text{Long} | Y=\text{Banana}) = 400 / 500 = 0.80$ $P(x_2=\text{Sweet} | Y=\text{Banana}) = 350 / 500 = 0.70$ $P(x_3=\text{Yellow} | Y=\text{Banana}) = 450 / 500 = 0.90$.

So, the overall probability of Likelihood of evidence for Banana = $0.8 * 0.7 * 0.9 = 0.504$

Step 4: Substitute all the 3 equations into the Naive Bayes formula, to get the probability that it is a banana.

Step 4: If a fruit is 'Long', 'Sweet' and 'Yellow', what fruit is it?

$$P(\text{Banana} | \text{Long, Sweet and Yellow}) = \frac{P(\text{Long} | \text{Banana}) * P(\text{Sweet} | \text{Banana}) * P(\text{Yellow} | \text{Banana}) * P(\text{banana})}{P(\text{Long}) * P(\text{Sweet}) * P(\text{Yellow})}$$

$$= \frac{0.8 * 0.7 * 0.9 * 0.5}{P(\text{Evidence})} = 0.252 / P(\text{Evidence})$$

$$P(\text{Orange} | \text{Long, Sweet and Yellow}) = 0, \text{ because } P(\text{Long} | \text{Orange}) = 0$$

$$P(\text{Other Fruit} | \text{Long, Sweet and Yellow}) = 0.01875 / P(\text{Evidence})$$

Answer: Banana - Since it has highest probability amongst the 3 classes

Similarly, you can compute the probabilities for 'Orange' and 'Other fruit'. The denominator is the same for all 3 cases, so it's optional to compute. Clearly, Banana gets the highest probability, so that will be our predicted class.

6. What is Laplace Correction?

The value of $P(\text{Orange} \mid \text{Long, Sweet and Yellow})$ was zero in the above example, because, $P(\text{Long} \mid \text{Orange})$ was zero.

That is, there were no 'Long' oranges in the training data.

It makes sense, but when you have a model with many features, the entire probability will become zero because one of the feature's value was zero. To avoid this, we increase the count of the variable with zero to a small value (usually 1) in the numerator, so that the overall probability doesn't become zero. This approach is called '**Laplace Correction**'.

Most Naive Bayes model implementations accept this or an equivalent form of correction as a parameter.

7. What is Gaussian Naive Bayes?

So far we've seen the computations when the X's are categorical.

But how to compute the probabilities when X is a continuous variable?

If we assume that the X follows a particular distribution, then you can plug in the probability density function of that distribution to compute the probability of likelihoods.

If you assume the X's follow a Normal (aka Gaussian) Distribution, which is fairly common, we substitute the corresponding probability density of a Normal distribution and call it the Gaussian Naive Bayes.

You need just the mean and variance of the X to compute this formula.

$$P(X|Y = c) = \frac{1}{\sqrt{2\pi\sigma_c^2}} e^{-\frac{(x-\mu_c)^2}{2\sigma_c^2}}$$

where mu and sigma are the mean and variance of the continuous X computed for a given class 'c' (of Y).

To make the features more Gaussian like, you might consider transforming the variable using something like the [Box-Cox](#) to achieve this.

That's it. Now, let's build a Naive Bayes classifier.

8. Building a Naive Bayes Classifier in R

Understanding Naive Bayes was the (slightly) tricky part. Implementing it is fairly straightforward.

In R, Naive Bayes classifier is implemented in packages such as `e1071` , `klaR` and `bnlearn` .

In Python, it is implemented in `scikit learn`, `h2o` etc.

For sake of demonstration, let's use the standard `iris` dataset to predict the `Species` of flower using 4 different features: `Sepal.Length` , `Sepal.Width` , `Petal.Length` , `Petal.Width`

```
# Import Data
training <- read.csv('https://raw.githubusercontent.com/justmarkham/DAT8/master/data/iris.csv')
test <- read.csv('https://raw.githubusercontent.com/justmarkham/DAT8/master/data/iris-test.csv')
```

The training data is now contained in `training` and test data in `test` dataframe. Lets load the `KlaR` package and build the naive bayes model.

```
# Using klaR for Naive Bayes
library(klaR)
nb_mod <- NaiveBayes(Species ~ ., data=training)
pred <- predict(nb_mod, test)
```

Lets see the confusion matrix.

```
# Confusion Matrix
tab <- table(pred$class, test$Species)
caret::confusionMatrix(tab)
```

```
#> Confusion Matrix and Statistics
```

```
#>           setosa versicolor virginica
#> setosa          15            0            0
#> versicolor       0            11            0
#> virginica        0            4            15
```

```
#> Overall Statistics
```

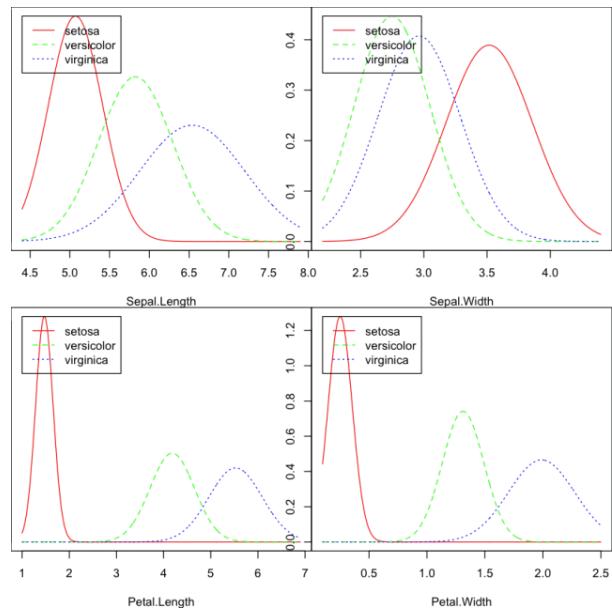
```
#> Accuracy : 0.9111
#> 95% CI : (0.7878, 0.9752
#> No Information Rate : 0.3333
#> P-Value [Acc > NIR] : 8.467e-16
```

```
#> Kappa : 0.8667
#> Mcnemar's Test P-Value : NA
```

```
#> Statistics by Class:
```

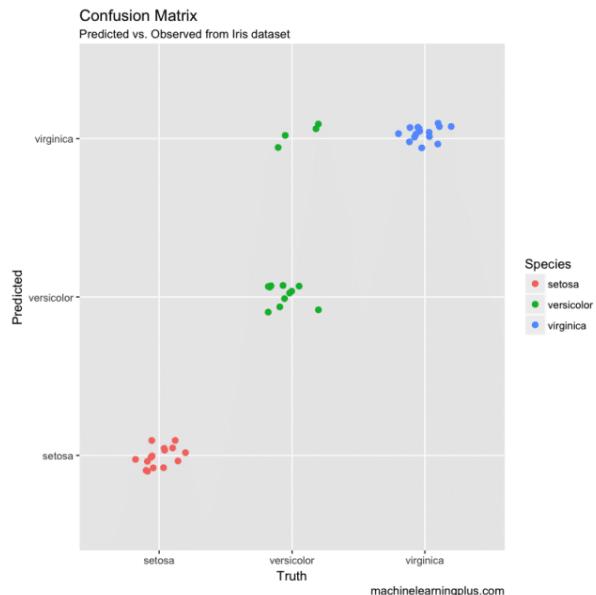
```
#>           Class: setosa Class:
#> Sensitivity          1.0000
#> Specificity          1.0000
#> Pos Pred Value       1.0000
#> Neg Pred Value       1.0000
#> Prevalence           0.3333
#> Detection Rate       0.3333
#> Detection Prevalence 0.3333
#> Balanced Accuracy    1.0000
```

```
# Plot density of each feature using nb_mod
opar = par(mfrow=c(2, 2), mar=c(4,0,0,0))
plot(nb_mod, main="")
par(opar)
```



```
# Plot the Confusion Matrix
```

```
library(ggplot2)
test$pred <- pred$class
ggplot(test, aes(Species, pred, color = Species))
  geom_jitter(width = 0.2, height = 0.1, size = 1)
  labs(title="Confusion Matrix",
       subtitle="Predicted vs. Observed from Iris dataset",
       y="Predicted",
       x="Truth",
       caption="machinelearningplus.com")
```



9. Building Naive Bayes Classifier in Python

```

# Import packages
from sklearn.naive_bayes import GaussianNB
from sklearn.model_selection import train_test_split
from sklearn.metrics import confusion_matrix
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns; sns.set()

# Import data
training = pd.read_csv('https://raw.githubusercontent.com/justmarkham/DAT8/master/data/breast_cancer.csv')
test = pd.read_csv('https://raw.githubusercontent.com/justmarkham/DAT8/master/data/breast_cancer.csv')

# Create the X, Y, Training and Test
xtrain = training.drop('Species', axis=1)
ytrain = training.loc[:, 'Species']
xtest = test.drop('Species', axis=1)
ytest = test.loc[:, 'Species']

# Init the Gaussian Classifier
model = GaussianNB()

# Train the model
model.fit(xtrain, ytrain)

# Predict Output
pred = model.predict(xtest)

# Plot Confusion Matrix
mat = confusion_matrix(pred, ytest)
names = np.unique(pred)
sns.heatmap(mat, square=True, annot=True, fmt='d',
            xticklabels=names, yticklabels=names)
plt.xlabel('Truth')
plt.ylabel('Predicted')

```

	setosa	versicolor	virginica
setosa	15	0	0
versicolor	0	11	0
virginica	0	4	15

10. Practice Exercise: Predict Human Activity Recognition (HAR)

The objective of this practice exercise is to predict current human activity based on physiological activity measurements from 53 different features based in the [HAR dataset](#).

The [training](#) and [test](#) datasets are provided.

Build a Naive Bayes model, predict on the test dataset and compute the confusion matrix. [Show R Solution >](#)

[Show Python Solution >](#)

11. Tips to improve the model

1. Try transforming the variables using transformations like BoxCox or YeoJohnson to make the features near Normal.
2. Try applying Laplace correction to handle records with zeros values in X variables.
3. Check for correlated features and try removing the highly correlated ones. Naive Bayes is based on the assumption that the features are independent.
4. Feature engineering. Combining features (a product) to form new ones that makes intuitive sense might help.

5. Try providing more realistic prior probabilities to the algorithm based on knowledge from business, instead of letting the algo calculate the priors based on the training sample.

For this case, ensemble methods like bagging, boosting will help a lot by reducing the variance.

Recommended: Industrial project course (Full Hands-On Walk-through): [Microsoft Malware Detection](#).

More Articles

K-Means Clustering Algorithm from Scratch

Apr 26, 2020

How Naive Bayes Algorithm Works? (with example and full code)

Nov 04, 2018