# Introduction to Naive Bayes Classification

**Naive Bayes** is a simple, yet effective and commonly-used, machine learning classifier. It is a probabilistic classifier that makes classifications using the Maximum A Posteriori decision rule in a Bayesian setting. It can also be represented using a very simple Bayesian network. Naive Bayes classifiers have been especially popular for text classification, and are a traditional solution for problems such as spam detection.

## The Model

The goal of any probabilistic classifier is, with features $x\_0$ through $x\_n$ and classes $c\_0$ through $c\_k$, to determine the probability of the features occurring in each class, and to return the most likely class. Therefore, for each class, we want to be able to calculate $P(c\_i \mid x\_0, ..., x\_n)$.

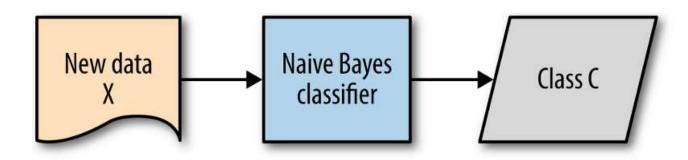In order to do this, we use **Bayes rule**. Recall that Bayes rule is the following:

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

In the context of classification, you can replace A with a class, $c\_i$, and B with our set of features, $x\_0$ through $x\_n$. Since $P(B)$ serves as normalization, and we are usually unable to calculate $P(x\_0, ..., x\_n)$, we can simply ignore that term, and instead just state that $P(c\_i \mid x\_0, ..., x\_n) \propto P(x\_0, ..., x\_n \mid c\_i) * P(c\_i)$, where $\propto$ means "is proportional to". $P(c\_i)$ is simple to calculate; it is just the proportion of the data-set that falls in class i. $P(x\_0, ..., x\_n \mid c\_i)$ is more difficult to compute. In order to simplify its computation, we make the assumption that $x\_0$ through $x\_n$ are **conditionally independent** given $c\_i$, which allows us to say that $P(x\_0, ..., x\_n \mid c\_i) = P(x\_0 \mid c\_i) * P(x\_1 \mid c\_i) * ... * P(x\_n \mid c\_i)$. This assumption is most likely not true—hence the name *naive* Bayes classifier, but the classifier nonetheless performs well in most situations. Therefore, our final representation of class probability is the following:

$$P(c_i|x_0, \ldots, x_n) \propto P(x_0, \ldots, x_n|c_i)P(c_i)$$

$$\propto P(c_i) \prod_{j=1}^{n} P(x_j|c_i)$$

Calculating the individual P(x_j | c_i) terms will depend on what distribution your features follow. In the context of text classification, where features may be word counts, features may follow a **multinomial distribution**. In other cases, where features are continuous, they may follow a **Gaussian distribution**.

Note that there is very little explicit training in Naive Bayes compared to other common classification methods. The only work that must be done before prediction is finding the parameters for the features' individual probability distributions, which can typically be done quickly and deterministically. This means that Naive Bayes classifiers can perform well even with high-dimensional data points and/or a large number of data points.



## Classification

Now that we have a way to estimate the probability of a given data point falling in a certain class, we need to be able to use this to produce classifications. Naive Bayes handles this in a very simple manner; simply pick the $c_i$ that has the largest probability given the data point's features.

$$y = \underset{c_i}{argmax}\ P(c_i) \prod_{j=1}^{n} P(x_j|c_i)$$

This is referred to as the **Maximum A Posteriori** decision rule. This is because, referring back to our formulation of Bayes rule, we only use the P(B|A) and P(A) terms, which are the likelihood and prior terms, respectively. If we only used P(B|A), the likelihood, we would be using a **Maximum Likelihood** decision rule.