

Derivation of Naive Bayes Classifier

The Naive Bayes classifier decomposes a single high dimensional probability estimation problem into many 1-dimensional estimation problems by assuming the features are conditionally independent given the class label. We want to estimate $\hat{P}(y|\mathbf{x})$ with the Naive Bayes classifier. To do so, we will make use of **Bayes' Rule and the Naive Bayes assumption**.

Breaking Down Naive Bayes

You can estimate $\hat{P}(y|\mathbf{x})$ by calculating $P(y)$ and $P(\mathbf{x}|y)$, since by Bayes' Rule:

$$P(y|\mathbf{x}) = \frac{P(\mathbf{x}|y)P(y)}{P(\mathbf{x})}$$

Estimating $\mathbf{P}(\mathbf{y} = \mathbf{c})$ (the probability that \mathbf{y} takes on some value \mathbf{c}) is simple - it is just the fraction of data points where $\mathbf{y} = \mathbf{c}$.

We can write this formally using an indicator variable I , which is 1 when the argument holds and 0 when it doesn't. If y takes on discrete binary values, such as 0 or 1, it simply counts how many times we observe each outcome:

$$P(y = c) = \frac{\sum_{i=1}^n I(y_i = c)}{n} = \hat{\pi}_c$$

Estimating $P(\mathbf{x}|y)$, however, is not so simple. To do so, you will make use of the Naive Bayes assumption from earlier:

$$P(\mathbf{x}|y) = \prod_{\alpha=1}^d P(x_{\alpha}|y)$$

where $\mathbf{x}_\alpha = [\mathbf{x}]_\alpha$ is the value for feature α .

Below, you'll see an illustration of how the Naive Bayes algorithm decomposes a multi-dimensional probability estimation problem (first image) into many 1-dimensional estimation problems. First, you estimate $\mathbf{P}(\mathbf{x}_\alpha|\mathbf{y})$ independently in each dimension (center two images) and then obtain an estimate of the full data distribution by assuming conditional independence $\mathbf{P}(\mathbf{x}|\mathbf{y}) = \prod_\alpha \mathbf{P}(\mathbf{x}_\alpha|\mathbf{y})$ (rightmost image).

Original	Estimation of first dimension	Estimation of second dimension	Resulting data distribution
<p>Scatter plot showing two classes of data points: blue circles ($y = 1$) and red squares ($y = 2$).</p>	<p>Plot showing the estimated first dimension. The data points are projected onto the horizontal axis. The conditional distributions $P(\mathbf{x}_1 y=1)$ and $P(\mathbf{x}_1 y=2)$ are shown as shaded regions.</p>	<p>Plot showing the estimated second dimension. The data points are projected onto the vertical axis. The conditional distributions $P(\mathbf{x}_2 y=1)$ and $P(\mathbf{x}_2 y=2)$ are shown as shaded regions.</p>	<p>Plot showing the resulting data distribution. The data points are projected onto the 2D space. The conditional distributions $P(\mathbf{x}_\alpha y=1)$ and $P(\mathbf{x}_\alpha y=2)$ are shown as shaded regions.</p>

Derivation of Naive Bayes Classifier

Now, we can begin our derivation of the Naive Bayes classifier as follows:

$$\begin{aligned} h(\mathbf{x}) &= \arg \max_y P(y|\mathbf{x}) \\ &= \arg \max_y \frac{P(\mathbf{x}|y)P(y)}{P(\mathbf{x})} && \text{(Bayes rule)} \\ &= \arg \max_y P(\mathbf{x}|y)P(y) && \text{(P(\mathbf{x}) does not depend on y)} \\ &= \arg \max_y \prod_{\alpha=1}^d P(x_{\alpha}|y)P(y) && \text{(by the Naive Bayes assumption)} \\ &= \arg \max_y \sum_{\alpha=1}^d \log(P(x_{\alpha}|y)) + \log(P(y)) && \text{(as log is a monotonic function)} \end{aligned}$$

Estimating $\mathbf{P}(x_\alpha | \mathbf{y})$ is easy since we only need to consider one single dimension, so MLE should work well. Estimating $\mathbf{P}(\mathbf{y})$ is not affected by the assumption and can similarly done with straight-forward MLE.

★ Key Points

First you need to calculate the probabilities of the labels given the observed data, using Bayes' Rule.

Then, you choose y such that this probability is maximized, considering each dimension of the feature vector independently.