

Our training consists of the set

$$D = \{(x_1, y_1), \dots, (x_n, y_n)\}$$

drawn from some unknown distribution $P(X, Y)$.

Since all pairs are sampled iid, we obtain

$$P(D) = P((x_1, y_1), \dots, (x_n, y_n)) = \prod_{\alpha=1}^n P(x_\alpha, y_\alpha)$$

If we have enough data, we could estimate $P(X, Y)$ similar to the coin example in lecture 4, where we imagine a gigantic die that has one side for each possible value of (\vec{x}, y) .

We can estimate the probability that one specific side comes up through counting:

$$\hat{P}(\vec{x}, y) = \frac{\sum_{i=1}^n I(x_i = x \wedge y_i = y)}{n}$$

where

$$I(x_i = x \wedge y_i = y) = \begin{cases} 1 & \text{if } x_i = x \wedge y_i = y \\ 0 & \text{o/w} \end{cases}$$

Of course we are primarily interested in predicting the label y from the features \vec{x} , we may estimate $P(Y|X)$ directly instead of $P(X, Y)$.

We can then use the Bayes Optimal Classifier for a specific $\hat{P}(y|\vec{x})$ to make predictions.

How do we estimate $\hat{P}(y|\vec{x})$?

Well, know

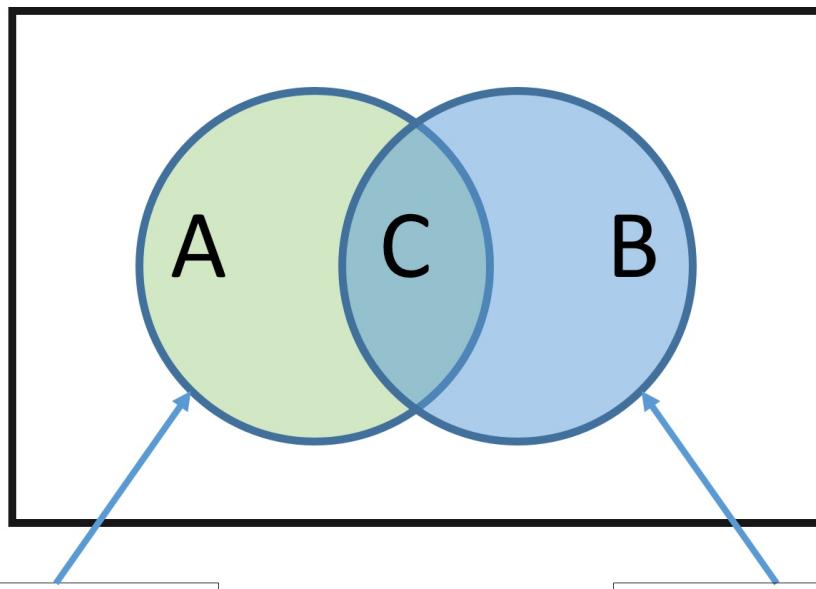
$$\hat{P}(\vec{x}) = \frac{\sum_{i=1}^n I(x_i=x)}{n} ; \hat{P}(y) = \frac{\sum_{i=1}^n I(y_i=y)}{n}$$

and

$$\hat{P}(\vec{x}, y) = \frac{\sum_{i=1}^n I(x_i=x \wedge y_i=y)}{n}$$

Thus,

$$\hat{P}(y|\vec{x}) = \frac{\hat{P}(y, \vec{x})}{\hat{P}(\vec{x})} = \frac{\frac{\sum_{i=1}^n I(x_i=x \wedge y_i=y)}{n}}{\frac{\sum_{i=1}^n I(x_i=x)}{n}} = \frac{\sum_{i=1}^n I(x_i=x \wedge y_i=y)}{\sum_{i=1}^n I(x_i=x)}$$



The Venn diagram illustrates that the MLE method estimates $\hat{P}(y|\vec{x})$ as $\hat{P}(y|\vec{x}) = \frac{|C|}{|B|}$

PROBLEM!

The MLE estimate is only good if there are many training vectors w/ the same identical features as \vec{x} !
In high dimensional spaces (or w/ continuous \vec{x}), this NEVER happens! So $|B| \rightarrow 0$ and $|C| \rightarrow 0$.

Naive Bayes

We can approach this dilemma w/ a simple trick, and an additional assumption.

The trick: estimate $P(y)$ and $P(\vec{x}|y)$ instead of $P(\vec{x}, y)$
generative learning

Thus,

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

Estimating $P(y)$ is easy.

For example, if Y takes on discrete binary values estimating $P(Y)$ reduces to coin tossing. Just need to count how many times we observe each outcome

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

Estimating $P(\vec{x}|y)$ however is NOT easy!

Need to make an additional assumption, the Naive Bayes assumption

Naive Bayes Assumption:

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

where $x_\alpha = [x]_\alpha$ is the value for feature α

i.e. feature values are independent given the label! This is a **BOLD** assumption, Sameer

For example, a setting where Naive Bayes is often used is spam filtering.

Here the data is emails, and the label is spam or not spam

The Naive Bayes assumption \Rightarrow implies that the words in an email are conditionally independent, given that you know whether an email is spam or not

Clearly, this is NOT true Sameer!

Neither the words of spam or not spam email are drawn independently at random.

However, the resulting classifiers can work well in practice even if this assumption is violated.

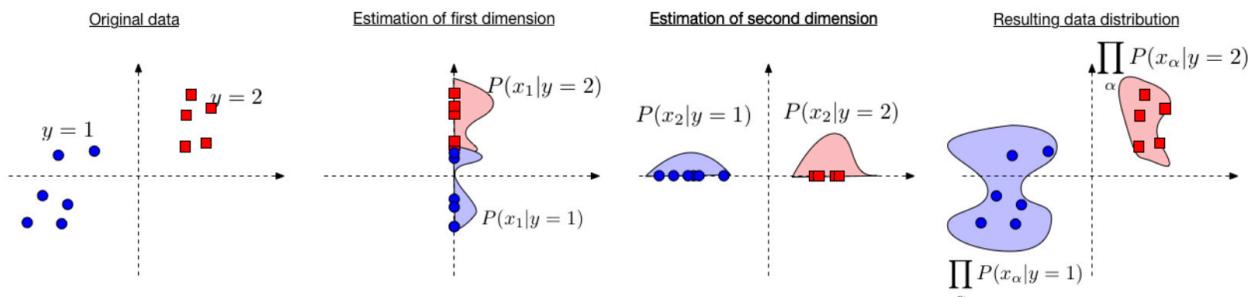


Illustration behind the Naive Bayes algorithm. We estimate $P(x_\alpha|y)$ independently in each dimension (middle two images) and then obtain an estimate of the full data distribution by assuming conditional independence $P(\mathbf{x}|y) = \prod_\alpha P(x_\alpha|y)$ (very right image).

So, for now assume the Naive Bayes assumption holds. Then the Bayes classifier can be defined as

$$\begin{aligned}
 h(\vec{x}) &= \operatorname{argmax}_y P(y|\vec{x}) \\
 &= \operatorname{argmax}_y \frac{P(\vec{x}|y) P(y)}{P(\vec{x})} \\
 &= \operatorname{argmax}_y P(\vec{x}|y) P(y) \quad P(x) \text{ has no } y \text{ dependence} \\
 &= \operatorname{argmax}_y \prod_{\alpha=1}^d P(x_\alpha|y) P(y) \quad \text{by Naive Bayes Assumption} \\
 &= \operatorname{argmax}_y \sum_{\alpha=1}^d \log(P(x_\alpha|y)) + \log(P(y)) \quad \log \text{ is monotonic}
 \end{aligned}$$

Estimating $P(x_\alpha|y)$ is easy as we only need to consider one dimension. And estimating $P(y)$ is NOT impacted by assumption.

Estimating $P([\vec{x}]_\alpha|y)$

Now that we know how we can use our assumption to make the definition of $P(y|\vec{x})$ tractable.

There are 3 notable cases we can use our Naive Bayes classifier.

Case 1: Categorical Features



Features:

$$[\hat{x}]_\alpha \in \{f_1, f_2, \dots, f_{K_\alpha}\}$$

Each feature α falls into one of K_α categories. (Binary, $K_\alpha = 2$). An example of such a setting could be medical data where one feature could be gender or marital status.

Illustration of categorical NB. For d dimensional data, there exist d independent dice for each class. Each feature has one die per class. We assume training samples were generated by rolling one die after another. The value in dimension i corresponds to the outcome that was rolled with the i^{th} die.

Model $P(x_\alpha | y)$:

$$P(x_\alpha = j | y=c) = [\theta_{j,c}]_\alpha$$

and

$$\sum_{j=1}^{K_\alpha} [\theta_{j,c}]_\alpha = 1 \quad \text{with } P(x_\alpha = j | y=c)$$

where $[\theta_{j,c}]_\alpha$ is the probability of feature α having value j , given that the label is c . And the constraint indicates that x_α must have one of the categories $\{1, \dots, K_\alpha\}$

Parameter Estimation:

$$[\hat{\theta}_{j,c}]_\alpha = \frac{\sum_{i=1}^n I(y_i=c) I(x_{i\alpha} = j) + l}{\sum_{i=1}^n I(y_i=c) + l K_\alpha}$$

where $x_{i\alpha} = [\hat{x}_i]_\alpha$ and l = smoothing parameter.

$l = 0 \rightarrow$ MLE estimator

$l > 0 \rightarrow$ MAP

$l = +1 \rightarrow$ Laplace smoothing

In words (w/o hallucinated 1 samples) this means

$$\frac{\text{# of sample w/ label } c \text{ that have feature } \alpha \text{ w/ value } j}{\text{# of samples w/ label } c}$$

Stay w/ me Sameer, this next part gets a little confusing

- Essentially the categorical feature model associates a special coin w/ each feature and label.
- The generative model we are assuming is that the data was generated by first choosing the label.
 - that label comes w/ a set of d dice; one for each dimension
- The generator picks each die, tosses it, and fills in the feature value w/ the outcome of the coin toss
- So for C possible labels and d dimensions we are estimating $d \times C$ "dice" from the data.
- **HOWEVER**, per data pt only d dice are tossed (1 per dimension)
- Die α (for any label) has K_α possible "sides"
- Of course, this is NOT how data is generated in reality - but it is a modeling assumption we make
- We then learn these models from the data and during test time see which model is more likely given the sample

Prediction:

$$\operatorname{argmax}_y P(y=c | \vec{x}) \propto \operatorname{argmax}_y \hat{\pi}_c \prod_{\alpha=1}^d [\hat{\theta}_{j,c}]_\alpha$$

Case #2: Multinomial Features

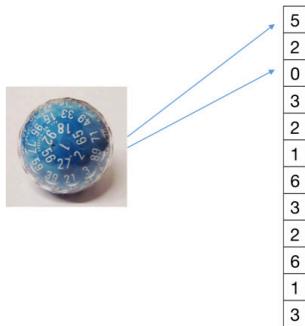


Illustration of multinomial NB. There are only as many dice as classes. Each die has d sides. The value of the i^{th} feature shows how many times this particular side was rolled.

If feature values don't represent categories (e.g. male/female) but counts we need to use a different model. E.g. in the text document categorization, feature value $x_\alpha = j$ means that in this particular document \mathbf{x} the α^{th} word in my dictionary appears j times. Let us consider the example of spam filtering. Imagine the α^{th} word is indicative towards ``spam''. Then if $x_\alpha = 10$ means that this email is likely spam (as word α appears 10 times in it). And another email with $x'_\alpha = 20$ should be even more likely to be spam (as the spammy word appears twice as often). With categorical features this is not guaranteed. It could be that the training set does not contain any email that contain word α exactly 20 times. In this case you would simply get the hallucinated smoothing values for both spam and not-spam - and the signal is lost. We need a model that incorporates our knowledge that features are counts - this will help us during estimation (you don't have to see a training email with exactly the same number of word occurrences) and during inference/testing (as you will obtain these monotonicities that one might expect). The multinomial distribution does exactly that.

Features:

$$x_\alpha \in \{0, 1, 2, \dots, m\} \text{ and } m = \sum_{\alpha=1}^d x_\alpha$$

Each feature α represents a count and m is the length of the sequence. An example of this could be the count of a specific word α in a document of length m and d is the size of the vocabulary.

Model $P(\vec{x}|y)$: Use multinomial distribution

$$P(\vec{x}|m, y=c) = \frac{m!}{x_1! x_2! \cdots x_d!} \prod_{\alpha=1}^d (\theta_{ac})^{x_\alpha}$$

ways to reshuffle \vec{x} given $y=c$
 probability of getting $\vec{x}|y=c$
 amount of times observed

where θ_{ac} is the probability of selecting x_α and $\sum_{\alpha=1}^d \theta_{ac} = 1$

So, we can use this to generate a spam email!

i.e. a document \vec{x} of class $y=\text{spam}$ by picking m words independently at random from the vocabulary of d words using $P(\vec{x}|y=\text{spam})$

Parameter Estimation:

$$\hat{\theta}_{ac} = \frac{\sum_{i=1}^n I(y_i=c) x_{ia} + l}{\sum_{i=1}^n I(y_i=c) m_i + l \cdot d}$$

all counts for features x_a
sums up all counts
of all features across
ALL data points

where

$$m_i = \sum_{\beta=1}^d x_{i\beta}$$

denotes the number of words in document i .

Once again, l is smoothing parameter.

Prediction:

$$\arg \max_c P(y=c | \vec{x}) \propto \arg \max_c \pi_c \prod_{a=1}^d \hat{\theta}_{ac}^{x_a}$$

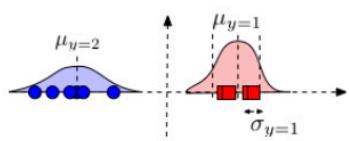
$$P(Y=1 | \vec{x} = x_{test}) = \frac{P(\vec{x} = x_{test} | Y=1) P(Y=1)}{P(\vec{x} = x_{test})}$$

$$P(Y=-1 | \vec{x} = x_{test}) = \frac{P(\vec{x} = x_{test} | Y=-1) P(Y=-1)}{P(\vec{x} = x_{test})}$$

$$= \frac{P(\vec{x} = x_{test} | Y=1) P(Y=1)}{P(\vec{x} = x_{test} | Y=-1) P(Y=-1)}$$

Case 3: Continuous Features (Gaussian Naive Bayes)

Features:



$x_\alpha \in \mathbb{R}$ (each feature takes on real value)

Illustration of Gaussian NB. Each class conditional feature distribution $P(x_\alpha | y)$ is assumed to originate from an independent Gaussian distribution with its own mean $\mu_{\alpha,y}$ and variance $\sigma_{\alpha,y}^2$.

Model:

$$P(x_\alpha | y=c) = N(\mu_{ac}, \sigma_{ac}^2) = \frac{1}{\sqrt{2\pi} \sigma_{ac}} e^{-\frac{(x-\mu_{ac})^2}{2\sigma_{ac}^2}}$$

Note that the model is based on our assumption about the data - that each feature α comes from a class-conditional Gaussian distribution.

The full distribution

$$P(\vec{x}|y) \sim N(\mu_y, \Sigma_y) \quad \text{diagonal covariance matrix w/ } (\Sigma_y)_{x,y} = \sigma_{xy}^2$$

Parameter Estimation:

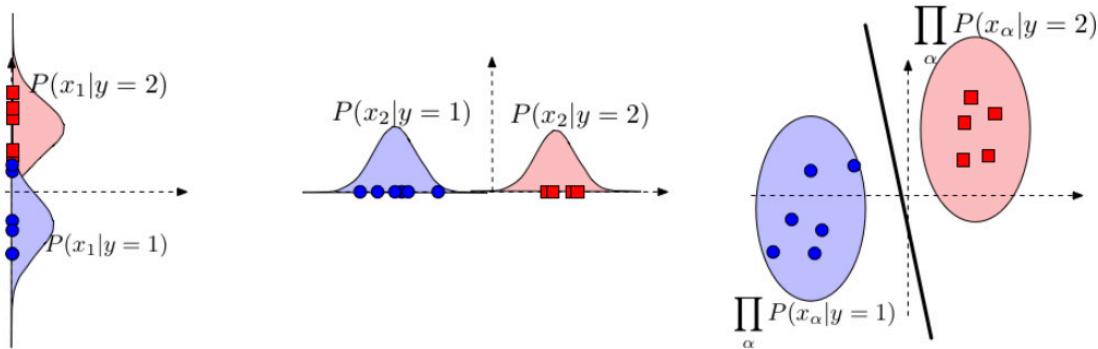
We estimate the parameters of the distributions for each dimension and class independently. Gaussian distributions only have 2 parameters.

$$\mu_{ac} \leftarrow \frac{1}{n_c} \sum_{i=1}^n I(y_i=c) x_{ia} \quad \begin{matrix} \text{avg feature value of dimension } \alpha \\ \text{from all samples w/ label } y \end{matrix}$$

$$\text{where } n_c = \sum_{i=1}^n I(y_i=c)$$

$$\sigma_{ac}^2 \leftarrow \frac{1}{n_c} \sum_{i=1}^n I(y_i=c) (x_{ia} - \mu_{ac})^2 \quad \text{variance of estimate } \mu_{ac}$$

Naive Bayes is a Linear Classifier



Naive Bayes leads to a linear decision boundary in many common cases. Illustrated here is the case where $P(x_\alpha|y)$ is Gaussian and where $\sigma_{\alpha,c}$ is identical for all c (but can differ across dimensions α).

The boundary of the ellipsoids indicate regions of equal probabilities $P(\mathbf{x}|y)$. The red decision line indicates the decision boundary where $P(y=1|\mathbf{x}) = P(y=2|\mathbf{x})$.

1. Suppose $y_i \in \{-1, +1\}$ and features are multinomial

We can show that

$$h(\vec{x}) = \operatorname{argmax}_y P(y) \prod_{\alpha=1}^d P(x_\alpha|y) = \operatorname{sign}(\vec{w}^T \vec{x} + b)$$

That is,

$$\vec{w}^T \vec{x} + b > 0 \Leftrightarrow h(\vec{x}) = +1$$

As before, we define

$$P(x_\alpha|y=+1) \propto e^{-\frac{x_\alpha - \mu_+}{\sigma^2}} \quad \text{and} \quad P(Y=+1) = \pi_+ :$$

$$|\vec{w}|_\alpha = \log(\theta_{\alpha+}) - \log(\theta_{\alpha-})$$

$$b = \log(\pi_+) - \log(\pi_-)$$

If we use the above to do classification, we can compute
 $\vec{w}^T \vec{x} + b$

Simplifying further leads to

$$\vec{w}^T \vec{x} + b > 0 \Leftrightarrow \sum_{\alpha=1}^d [x]_\alpha (\log(\theta_{\alpha+}) - \log(\theta_{\alpha-})) + \log(\pi_+) - \log(\pi_-) > 0$$

exponential both sides
 $\log a = \frac{a}{e}$
 $\exp(a+b) = \frac{\exp(a)\exp(b)}{\exp(b)}$

$$\Leftrightarrow \exp \left(\sum_{\alpha=1}^d [x]_\alpha (\log(\theta_{\alpha+}) - \log(\theta_{\alpha-})) + \log(\pi_+) - \log(\pi_-) \right) > 1$$

$$\Leftrightarrow \prod_{\alpha=1}^d \frac{\exp(\log(\theta_{\alpha+}^{[x]_\alpha}) + \log(\pi_+))}{\exp(\log(\theta_{\alpha-}^{[x]_\alpha}) + \log(\pi_-))} > 1$$

$$\Leftrightarrow \prod_{\alpha=1}^d \frac{\theta_{\alpha+}^{[x]_\alpha} \pi_+}{\theta_{\alpha-}^{[x]_\alpha} \pi_-} > 1$$

$$\Leftrightarrow \frac{\prod_{\alpha=1}^d P([\vec{x}]_\alpha | Y=+1) \pi_+}{\prod_{\alpha=1}^d P([\vec{x}]_\alpha | Y=-1) \pi_-} > 1$$

$$\Leftrightarrow \frac{P(\vec{x} | Y=+1) \pi_+}{P(\vec{x} | Y=-1) \pi_-} > 1$$

$$\Leftrightarrow \frac{P(Y=+1 | \vec{x})}{P(Y=-1 | \vec{x})} > 1$$

$$\Leftrightarrow P(Y=+1 | \vec{x}) > P(Y=-1 | \vec{x})$$

$$\Leftrightarrow \operatorname{argmax}_y P(Y=y | \vec{x}) = +1 \quad \text{i.e. the point } \vec{x} \text{ lies on positive side of hyperplane iff Naive Bayes predicts 1}$$

2. In the case of Gaussian continuous features we can show that

$$P(y | \vec{x}) = \frac{1}{1 + e^{-y(\vec{w}^T \vec{x} + b)}} \quad \text{also called logistic regression}$$

NB & LR produce asymptotically the same model if naive Bayes assumption holds.