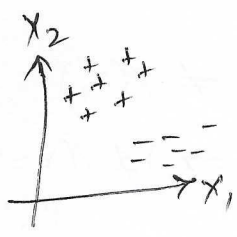


# Naive Bayes

(our first probabilistic classifier)

I/p features:  $\vec{x} \in \mathbb{R}^d$   
O/p:  $y \in \{1 \dots k\}$   
or  $y \in \{0, 1\}$   
or  $y \in \{-1, +1\}$



$$\hat{y} = \arg \max_y p(Y=y | X)$$

Generative approach

Discriminative approach

$$p(Y=y | X=x) = \frac{p(x|y)p(y)}{p(x)}$$

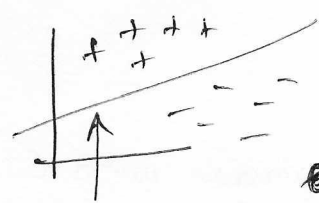
& estimate:  $p(X=x | Y=y)$   
&  $p(Y=y)$



- Try to "explain" your data for two classes
- Can create new data by sampling

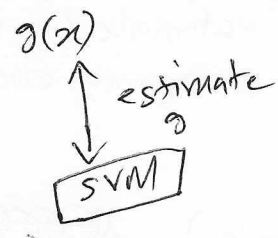
Naive Bayes

Directly estimate the posterior:  
 $p(Y=y | X=x)$   
or, estimate (non-probabilistic approach)  
Directly  $g: X \rightarrow Y$



No "explanation" of your data

Logistic regression



Consider binary features (for simplicity):  
 $x_1 \dots x_d \in \{0, 1\}$

$y \in \{1 \dots k\}$

We need to estimate  $p(Y=y)$  &  $p(X_1=x_1 \dots X_d=x_d | Y=y)$

probability distribution of length  $k$  that sums to 1

Massive table  
# of parameters:  $(2^d - 1)k$

If  $d=100$ , then  $2^{100} = 10^{30} \gg$  available data

Huge number of parameters  $\rightarrow$  sparse table  $\rightarrow$  hard to solve

The rescue is  $\rightarrow$  assume independence.

$\Downarrow$   
Naive Bayes

Assumption:

$$p(x_1 = x_1, \dots, x_d = x_d | y = y) = \prod_{i=1}^d p(x_i = x_i | y = y)$$

features are conditionally independent given the class.

$$\Rightarrow p(x_1, x_2) \neq p(x_1)$$

$$\& p(x_1 = x_1, x_2 = x_2) \neq p(x_1 = x_1)$$

$$p(x_2 = x_2)$$

(features are not independent)

NB  
 $\Downarrow$   
True Bayes

Axis-aligned features.

Number of parameters in NB?

- for binary features, for each feature  $x_i$  & each class  $y$ , we estimate 1 parameter (e.g.  $p(x_i = 1 | y)$ ), as  $p(x_i = 0 | y) = 1 - p(x_i = 1 | y)$ . So for  $d$  binary features &  $k$  classes, total # of parameters =  $\boxed{dk}$

$\uparrow$   
much less than the previous case.

- for categorical features with  $c$  categories, for each feature  $x_i$  & each class  $y$ , we estimate  $(c-1)$  parameters. So for  $n$  categorical features with  $c$  categories each, &  $k$  classes, total # of parameters =  $nk(c-1)$

$$p(x_i = x_i | y = c)$$

$$\downarrow$$
$$d \cdot (2^1 - 1) \cdot k = dk$$

(16)

$$\prod_{j=1}^d P(X_j = x_j | Y = y)$$

each of these is a table.

& we have to estimate prior probability

$$P(Y)$$

a vector where each entry is just a count:

$$\hat{P}_y = \frac{\text{count}(Y=y)}{N}$$

↑  
# of training samples with class  $Y=y$ , divided by total # of samples.

we want to estimate this table

$X_j$	1	...	K
$ X_j $			

Every column is also a categorical distribution.

basically an MLE.

we can do the same

If I tell you  $Y=y$ , then  $X_j = x_j$  is a categorical distribution since it sums to 1.

Shorthand notation:  $\theta_a^{jc} = P(X_j = a | Y = c)$

← for  $j$ -th feature,  $c$ -th class, of a particular value 'a'

$$\text{MLE } \theta_a^{jc} = \frac{\text{count}(X_j = a, Y = c)}{\text{count}(Y = c)}$$

Prediction in NB will be then,

$$\hat{y}_{\text{MAP}} = \max_y P(Y=y | X=x) = \arg \max_y \frac{\prod_{j=1}^d P(X_j = x_j | Y=y) P(Y=y)}{P(X=x)} \leftarrow \text{ignore this.}$$

$$\Rightarrow \hat{y}_{\text{MAP}} = \arg \max_y \left[ \log P(Y=y) + \sum_{j=1}^d \log P(X_j = x_j | Y=y) \right]$$



## Laplacian smoothing

$$P(X_i = a | Y = c) \Rightarrow \frac{\text{Count}(X_i = a, Y = c)}{\text{Count}(Y = c)}$$

If a feature value never appears with a class in the training data, the numerator gets zero. So,

$$P(X_i = a | Y = c) = 0$$

Since NB multiplies probabilities,  $P(x|y) = \prod_i P(x_i | y)$

One zero term kills the entire product, making  $P(y|x) = 0$

Even if all other evidence strongly supports the class.

This is called zero-frequency problem.

↓ solution - Laplacian smoothing.

Pretend that we have seen every possible feature value at least once for every class. So instead of trusting raw counts completely add a small constant to every count.

$$P(X_i = a | Y = c) = \frac{\text{Count}(X_i = a, Y = c) + 1}{\text{Count}(Y = c) + |X_i|}$$

↖ total no. of possible values  $X_i$  can have

↗ numerator: add 1 to every count

denominator: add  $|X_i|$  to keep probabilities summing to 1.