

# Multinomial Naive Bayes

Input  $x^{(i)}$  that is a list of words with length  $d_i$ .

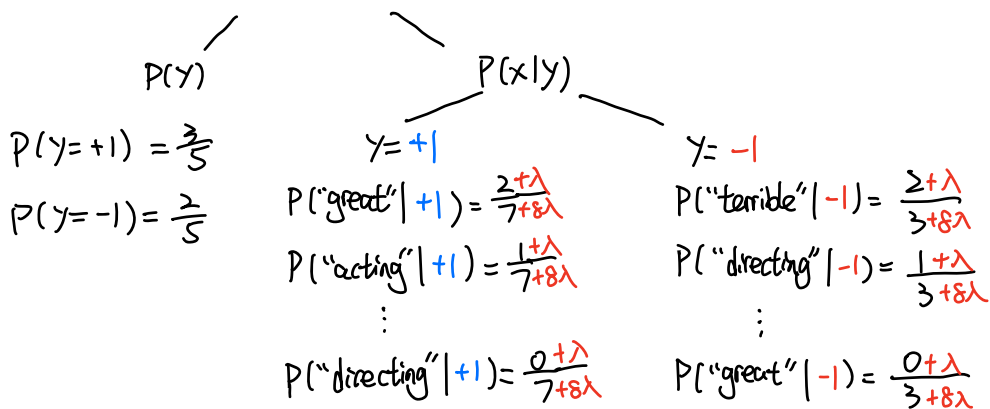
y	x
+1	great acting and score
-1	terrible directing
+1	great execution
-1	terrible
+1	amazing

NB assumption:  $P(x^{(i)}|y^{(i)}) = \prod_{j=1}^{x_i} P(x_j^{(i)}|y^{(i)})$

e.g.  $P(\text{"great directing"}|+1)$   
 $= P(\text{"great"}|+1) P(\text{"directing"}|+1)$

Additional assumption: word position doesn't matter

Step 1: Parameter Estimation (a.k.a. training)



Step 2: Inference

Given input  $x = \text{"great directing"}$ , compute  $P(y|x = \text{"great directing"})$

$y=+1: \frac{3}{5} \cdot \frac{3}{15} \cdot \frac{1}{15} = 0.008$

$\xrightarrow{\text{red arrow}} P(y=+1)$    
  $\xrightarrow{\text{red arrow}} P(\text{"great"}|+1)$    
  $\xrightarrow{\text{red arrow}} P(\text{"directing"}|+1)$

$= \frac{P(y) P(x = \text{"g d"}|y)}{P(x = \text{"g d"})}$   
 normalizing constant

$$y = -1: \frac{2}{5} \cdot \frac{1}{11} \cdot \frac{2}{11} = 0.0066$$

$$P(y = +1 | "g d") = \frac{0.008}{0.008 + 0.0066} = 0.55$$

If documents are long, you have to multiply many small probabilities together.

Numerical underflow i.e. on computer, everything  $\rightarrow 0$

Solution: work in log space

$$y = +1: \text{compute } \log(P(y = -1) \cdot P("g d" | +1)) \\ = \log\left(\frac{2}{5}\right) + \log\left(\frac{3}{15}\right) + \log\left(\frac{1}{15}\right) \approx -2.1$$

$$y = -1: \log\left(\frac{2}{5}\right) + \log\left(\frac{1}{11}\right) + \log\left(\frac{2}{11}\right) \approx -2.2$$

To get probabilities, compute max log score = -2.1

subtract that from everything  $y = +1 = 0$

$$y = -1 = -0.1$$

then exponentiate

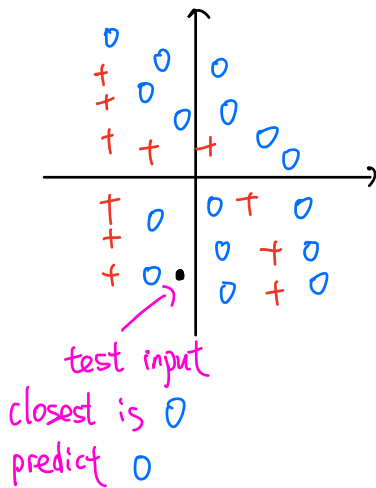
$$y = +1 = 1$$

$$y = -1 = e^{-0.1} \approx 0.9$$

finally normalize

$$P(y = +1) = \frac{1}{1 + 0.9}$$

	model $P(y x)$ Discriminative	$P(y) P(x y)$ Generative
parametric learn some params, afterward can throw away training data	logistic regression ( softmax regression [ $w$ $w^{(1)}, \dots, w^{(c)}$ ]	naive bayes / $P(y)$ a.k.a. $\pi$ \ $P(x y)$ $\tau$ Puk's
non-parametric use training set when making decisions	K-Nearest Neighbors	



1-NN: predict the same class as closest training example

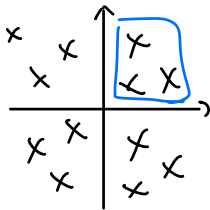
e.g. Euclidean distance

Generalisation: K-NN

- 1) find K-nearest neighbors of test input  $x$
- 2) predict label that's most common among neighbors  
"majority vote"

### Curse of Dimensionality

In high dimensions, you very rarely have nearby neighbors.



← prob in the same quadrant =  $\frac{1}{4}$

If  $d = 1000$ ,  $P = (\frac{1}{2})^{1000}$

No close neighbors to test data  $\Rightarrow$  using nearest neighbors to predict may not be accurate