

use to learn y

$$[x^1, x^2, \dots, x^D] \rightarrow y \text{ discrete}$$

given features, how to find label

features: ex. a word (word cup) to guess what kind of document (in document, or count in document)

$p(y | [x^1 \dots x^D])$ probability of y given the features
pick the one with the maximum probability

features:	"world cup"	"Germany"	"Brazil"
binary:	1	0	1

← work not in document

$$p(y | \underbrace{[x^1 \dots x^D]}_x) = \frac{P(x | y) P(y)}{P(x)} \leftarrow \text{Bayes rule}$$

$$= \frac{P([x^1, x^2, \dots, x^D] | y) P(y)}{P([x^1, x^2, \dots, x^D])}$$

features are independent given the labels
feature: words
label: "sports document"

features are independent given the document

given this assumption of independence,

$$= \frac{\prod_{i=1}^D p(x^i | y) P(y)}{P(x)}$$

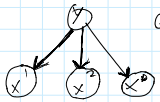
probability of word in document

how can we calculate $p(y)$?
feature \rightarrow class
estimate of the probability of y
ratio; count compared to total
$$Y=1 = \frac{\#}{\text{total } \#}$$

(Y) document sports?
mention "worldcup winner" — probability: highest "germany" (since they won)

(Y) tourism
"visit" — probability: any nation

so, given



Given y , these features are independent

or do dependencies on previous word $p(x^i | x^j)$, etc

avoid $p=0$ with y having no occurrences?

prior (our beliefs about the documents)

then adjust by ratio from training data

add certain value, normalize.

5.2 Bayes rule for medical diagnosis

After your yearly checkup, the doctor has bad news and good news. The bad news is that you tested positive for a serious disease, and that the test is 99% accurate (i.e., the probability of testing positive given that you have the disease is 0.99, as is the probability of testing negative given that you don't have the disease). The good news is that this is a rare disease, striking only one in 10,000 people. What are the chances that you actually have the disease?

$$P(\text{test} = \text{true} | \text{Disease}) = 0.99$$

$$P(\text{test} = \text{true} | \text{No Disease}) = 0.01$$

$$P(\text{Dis}) = 10^{-4}$$

$$P(\text{Dis} | \text{test} = \text{pos}) = \frac{P(\text{test} = \text{true}) \cdot P(\text{Dis})}{P(\text{test} = \text{pos})}$$

← $\begin{cases} \text{pos, dis} \\ \text{pos, no dis} \end{cases}$