

CLASSIFYING DOCUMENTS

USING A

NAIVE BAYES CLASSIFIER

Probability

I have a coin which is biased towards heads.

Probability of heads = 0.75

① QUESTION: What's the probability of getting 9 heads when the coin is tossed 10 times?



This is an easy question

Statistical Inference

I have a coin which I've tossed 10 times and got 9 heads.

② QUESTION: What is the bias of the coin ?



This is a hard question!

And it's critical to all of
Science.

Let's answer this hard question

Hypothesis The coin's bias

towards heads is 0.9.

Evidence / Data Toss the coin

10 times and observe 9 heads
in total.

But how to calculate this?

Answer:

$\text{Prob}(\underbrace{\text{bias} = 0.9_0}_{\text{H}} \mid \underbrace{9 \text{ heads in}}_{\text{E}} \underbrace{10 \text{ tosses}}_{\text{in}})$

Symbol for "given that"

The structural difference between

$P_{\text{rob}}(H | E)$ ← What we
want
(but hard to get)
vs.

$P_{\text{rob}}(E | H)$ ← what we
can usually
get

In general,

$$\text{Prob}(H|E) \neq \text{Prob}(E|H)$$

Simple Examples

$$\text{Prob}(\text{animal} / \text{human}) = 1$$

$$\text{Prob}(\text{human} / \text{animal}) \neq 1$$

$$\text{Prob}(\text{dog barks} / \text{burglar enters}) = A$$

$$\text{Prob}(\text{burglar enters} / \text{dog barks}) = B$$

For many dogs, $A \gg B$

\nearrow
S

symbol for
"much greater than"

$$P(\text{Win the lottery} \mid \text{have \# 1234567})$$

$$= \frac{1}{\text{\# tickets issued by the lottery commission}}$$

$$P(\text{have \# 1234567} \mid \text{Win the lottery})$$

$$= ? \text{ if's either 0 or 1}$$

Usually, it is easy to calculate

 Probability calculation

$P(E|H)$ but difficult to

calculate $P(H|E)$.

And

 Statistical
information gathering

Usually, it's $P(H|E)$ is what

we'd most like to know in

Science.

Bayes' Rule

$$P(H|E) = \frac{P(E|H) P(H)}{P(E)}$$

From
The Cookie Problem (Allen B. Downey
in Think Bayes)

Vanilla = 30

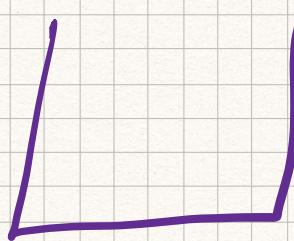
Chocolate = 10



Bowl 1

vanilla = 20

chocolate = 20



Bowl 2

Without looking, pick a cookie from a bowl. Suppose you get a vanilla cookie.

QUESTION: What is the probability that the cookie came from Bowl 1?

$\text{Prob}(\text{Bowl 1} | \text{vanilla}) = ?$

try

$\text{Prob}(\text{vanilla} | \text{Bowl 1}) \leftarrow$ ^{this is} easy

$$= \frac{30}{40} = 0.75$$

Use Bayes' Theorem

$$\begin{aligned}\text{Prob}(\text{Bowl 1} | \text{vanilla}) &= \\ \underbrace{\text{Prob}(\text{Vanilla} | \text{Bowl 1}) * \overbrace{\text{Prob}(\text{Bowl 1})}^{0.5}}_{\text{Prob}(\text{vanilla})} &\rightarrow \frac{50}{80} = 5/8\end{aligned}$$

OK, what does this have to do with spam?

	features			Target
	a_1	a_2	a_3	y
row 1	0	1	1	0
row 2	1	1	0	1
:	:			:
row n	0	1	0	1

usual logistic regression problem

all features are binary 0 or 1.

Goal: Use this data to learn a
classification.

How?

Usual Approach

- ① Traditional logistic regression model:

$$w_0 + a_1 \cdot w_1 + a_2 \cdot w_2 + a_3 \cdot w_3 \\ = \hat{y}^1$$

- ② Apply a penalty for each row
- ③ Find the values of w_0, w_1, w_2, w_3 that minimize the cost (total penalty over the entire dataset).

Another way to approach the problem

	features			Target
	a_1	a_2	a_3	y
row 1	0	1	1	0
row 2	1	1	0	1
:	:			:
row m	0	1	0	1

Question

$$\text{Prob}(y=0 \mid a_1=0 \text{ and } a_2=1 \text{ and } a_3=1) = ?$$

↓ take advantage of Bayes' theorem

$$\begin{aligned} & \text{Prob}(a_1=0, a_2=1, a_3=1 \mid y=0) * \underline{\text{Prob}(y=0)} \\ &= \underline{\text{Prob}(a_1=0, a_2=1, a_3=1)} \end{aligned}$$

Independent Events

$$P(A|B) = P(A) \text{ and } P(B|A) = P(B)$$

In words: The probability of A
is not affected by the occurrence
or non-occurrence of B and
vice versa.

$$\begin{aligned}
 & \text{Prob}(a_1=0, a_2=1, a_3=1 | y=0) * \text{Prob}(y=0) \\
 & \quad \text{Independence} \quad \text{Prob}(a_1=0, a_2=1, a_3=1) \\
 & \quad \textcircled{2} \quad \text{Prob}(a_1=0) * \text{Prob}(a_2=0) * \text{Prob}(a_3=0) \\
 & \textcircled{1} \quad \text{Prob}(a_1=0 | y=0) * \text{Prob}(a_2=1 | y=0) * \text{Prob}(a_3=1 | y=0)
 \end{aligned}$$

①, ②, and ③ can be readily
 calculated from the training data set.

	features			Target
	a_1	a_2	a_3	y
row 1	0	1	1	0
row 2	1	1	0	1
:				:
row n	0	1	0	1

$$P(y=0) = \frac{\# \text{ of } y=0}{\# \text{ rows of data}}$$

similar
expressions
for $y=1$

$$\text{Prob}(a_1=0) = \frac{\# \text{ of } a_1=0}{\# \text{ rows of data}}$$

$$\text{Prob}(a_1=0 | y=0) = \begin{matrix} \text{Pick all rows} \\ \text{where } y=0 \end{matrix}$$

$$\frac{k}{s}$$

Count the number of
these rows ($= s$)

In these rows, Count
the # of $a_1=0$. ($= k$)

Now we have all we need to use
Bayes' rule and calculate the probability
of spam or ham.

Where do the $a_1=0$, $a_2=1$, etc.

Come from?

They are actually ways to represent

words.

bag of words

Word₁, word₂, ... word_n → document₁

: A set of
: documents
is a Corpus

Word₁, word₃, ... word_k → document_m

Dictionary

The unique list of words that occur
in the complete set of documents

A simple way to represent a document
= Bag of words

Document 1

"The quick fox jumped over the hen"



Document 1
bag of
words

fox
hen
jumped
over
quick
the
the
over

Structure of
the sentence
is lost.
Word order
is not preserved.

Document 2

"Free! Don't Miss Out!!"



Document 2
bag of words

don't } Punctuation
free } ignored.
miss } Capitalization
out } ignored.

Dictionary = { don't, fox, free, hen, ... ,
out, over }

There are 11 unique words in the dictionary. These are the attributes of the dataset.

	<i>a₁</i>	<i>a₂</i>	...	<i>a₁₁</i>	<i>over target</i>
<i>document 1</i>	0	0	...	1	0
⋮	⋮	⋮	⋮	⋮	⋮
<i>document n</i>	1	1	...	0	1

↑ ↑ ↗
 The number of times
 a word occurs in the
 document.

The Naive Bayes model is trained on this dataset of documents that are "vectorized" or embedded using the bag of words representation.

What does the Naive Bayes model
"learn"?

- $P(a_1=0|y=0), P(a_1=1|y=0)$ } for
 $P(a_1=0|y=1), P(a_1=1|y=1)$ } all attributes
 a_1, a_2, \dots, a_{11}

- $P(y=0), P(y=1)$

- $P(a_1=0), P(a_1=1), \dots$

- $P(a_{11}=0), P(a_{11}=1)$

If turns out that

We don't need to learn these
at all.

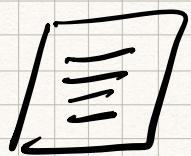
while these can
be learned from
the dataset,
they are usually
just set to

0.5 each.

Putting it all together

New Document

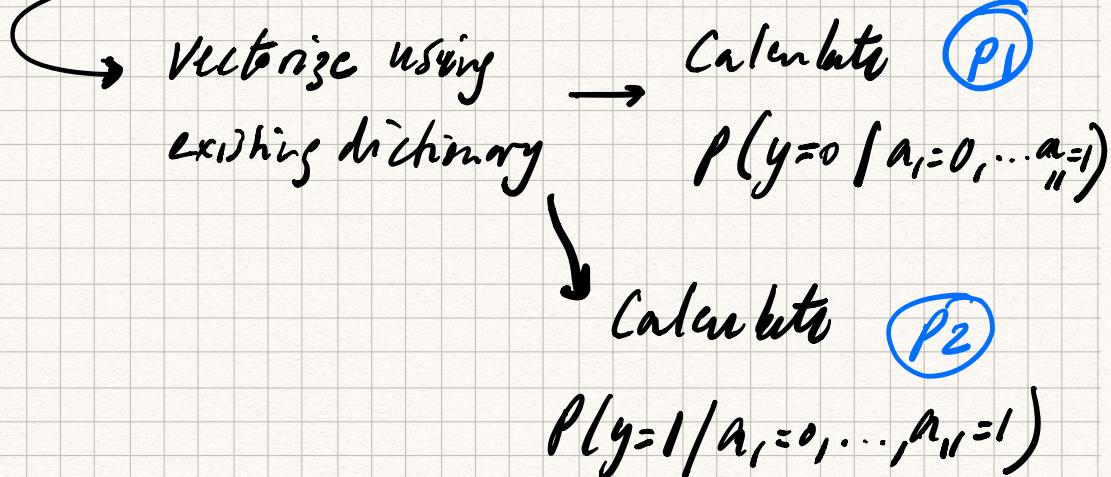
Convert to bag of words



→ tokenize → clean

Split into
words

remove
punctuation,
capitalization, ...



$$\text{Prediction} = \text{Max} (P_1, P_2)$$

Summary

- The difference between probability and statistical inference.

$$P(H|E) \neq P(E|H)$$

- When the probability you want is hard to calculate, try Bayesian inversion!

- Text can be represented by lists of numbers - the standard machine learning techniques apply.

- Naive Bayes is a surprisingly good classifier.

but is it machine learning?