

CS440/ECE448

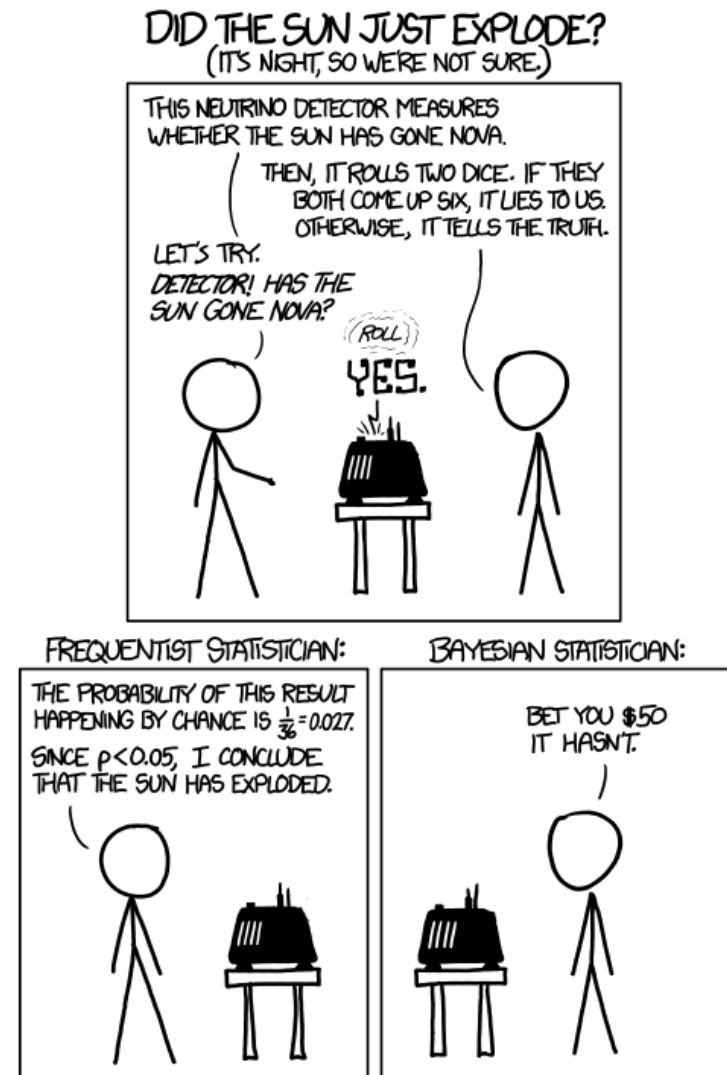
Lecture 4: Naïve Bayes

Mark Hasegawa-Johnson

Lecture slides: CC0



Some images may have other license terms.



© <https://www.xkcd.com/1132/>

Naïve Bayes

- naïve Bayes
- unigrams and bigrams
- parameter estimation

MPE classifier using Bayes' rule

$$f(x) = \operatorname{argmax}_y P(Y = y|X = x)$$

$$= \operatorname{argmax}_y \frac{P(Y = y)P(X = x|Y = y)}{P(X = x)}$$

$$= \operatorname{argmax}_y P(Y = y)P(X = x|Y = y)$$

The problem with likelihood: Too many words

What does it mean to say that the words, x , have a particular probability?

Suppose our training corpus contains two sample emails:

Email1: $Y = \text{spam}$, $X = \text{"Hi there man – feel the vitality! Nice meeting you..."}$

Email2: $Y = \text{ham}$, $X = \text{"This needs to be in production by early afternoon..."}$

Our test corpus is just one email:

Email1: $X = \text{"Hi! You can receive within days an approved prescription for increased vitality and stamina"}$

How can we estimate $P(X = \text{"Hi! You can receive within days an approved prescription for increased vitality and stamina"} | Y = \text{spam})$?

Unigram Naïve Bayes: the “Bag-of-words” model

We can estimate the likelihood of an e-mail by pretending that the e-mail is just a bag of words (order doesn't matter).

With only a few thousand spam e-mails, we can get a pretty good estimate of these things:

- $P(W = \text{“hi”} | Y = \text{spam}), P(W = \text{“hi”} | Y = \text{ham})$
- $P(W = \text{“vitality”} | Y = \text{spam}), P(W = \text{“vitality”} | Y = \text{ham})$
- $P(W = \text{“production”} | Y = \text{spam}), P(W = \text{“production”} | Y = \text{ham})$

Then we can approximate $P(X|Y)$ by assuming that the words, W , are **conditionally independent of one another given the category label**:

$$P(X = x | Y = y) \approx \prod_{i=1}^n P(W = w_i | Y = y)$$



The unigram Naïve Bayes classifier

$$f(x) = \operatorname{argmax}_y P(Y = y) \prod_{i=1}^n P(W = w_i | Y = y)$$



Dear Sir.

First, I must solicit your confidence in this transaction, this is by virtue of its nature as being utterly confidential and top secret. ...



TO BE REMOVED FROM FUTURE MAILINGS, SIMPLY REPLY TO THIS MESSAGE AND PUT "REMOVE" IN THE SUBJECT.

99 MILLION EMAIL ADDRESSES
FOR ONLY \$99



Ok, I know this is blatantly OT but I'm beginning to go insane. Had an old Dell Dimension XPS sitting in the corner and decided to put it to use, I know it was working pre being stuck in the corner, but when I plugged it in, hit the power nothing happened.

Floating-point underflow

$$f(x) = \operatorname{argmax}_y P(Y = y) \prod_{i=1}^n P(W = w_i | Y = y)$$

- That equation has a computational issue. Suppose that the probability of any given word is roughly $P(W = w_i | Y = y) \approx 10^{-3}$, and suppose that there are 103 words in an email. Then $\prod_{i=1}^n P(W = w_i | Y = y) = 10^{-309}$, which gets rounded off to zero. This phenomenon is called “floating-point underflow.”
- To avoid floating-point underflow, we can take the logarithm of the equation above:

$$f(x) = \operatorname{argmax}_y \left(\ln P(Y = y) + \sum_{i=1}^n \ln P(W = w_i | Y = y) \right)$$

Naïve Bayes

- naïve Bayes
- unigrams and bigrams
- parameter estimation

Reducing the naivety of naïve Bayes

Unigram naïve Bayes is unable to represent this fact:

True Statement:

$$P(X = \text{for you} | Y = \text{Spam}) > P(W = \text{for} | Y = \text{Spam})P(W = \text{you} | Y = \text{Spam})$$

We can modify naïve Bayes model to give it this power, using bigrams.

N-Grams

Claude Shannon, in his 1948 book *A Mathematical Theory of Communication*, proposed that the probability of a sequence of words could be modeled using N-grams: sequences of N consecutive words.

- **Unigram**: a unigram (1-gram) is an isolated word, e.g., “you”
- **Bigram**: a bigram (2-gram) is a pair of words, e.g., “for you”
- **Trigram**: a trigram (3-gram) is a triplet of words, e.g., “prescription for you”
- **4-gram**: a 4-gram is a 4-tuple of words, e.g., “approved prescription for you”

Bigram naïve Bayes

A bigram naïve Bayes model approximates the bigrams as conditionally independent, instead of the unigrams. For example,

$$P(X = \text{"approved prescription for you"} | Y = \text{Spam}) \approx$$

$$\begin{aligned} &P(B = \text{"approved prescription"} | Y = \text{Spam}) \times \\ &P(B = \text{"prescription for"} | Y = \text{Spam}) \times \\ &P(B = \text{"for you"} | Y = \text{Spam}) \end{aligned}$$

Advantages and disadvantages of bigram models relative to unigram models

- Advantage:
 - Slightly more accurate
- Disadvantage:
 - Far more parameters to store
 - Slightly more computational complexity
 - Far more possibility that we learn useless quirks of the training data, instead of true facts about spam and ham emails. This is called “overtraining the model.”

Naïve Bayes

- naïve Bayes
- unigrams and bigrams
- parameter estimation

Parameter estimation

Model parameters: feature likelihoods $P(\text{Word} \mid \text{Class})$ and priors $P(\text{Class})$

- How do we obtain the values of these parameters?

prior

spam:	0.33
¬spam:	0.67

$P(\text{word} \mid \text{spam})$

the :	0.0156
to :	0.0153
and :	0.0115
of :	0.0095
you :	0.0093
a :	0.0086
with:	0.0080
from:	0.0075
...	

$P(\text{word} \mid \text{ham})$

the :	0.0210
to :	0.0133
of :	0.0119
2002:	0.0110
with:	0.0108
from:	0.0107
and :	0.0105
a :	0.0100
...	

Maximum likelihood parameter estimation

The likelihood, $P(W = w_i | Y = y)$, can be estimated by counting.

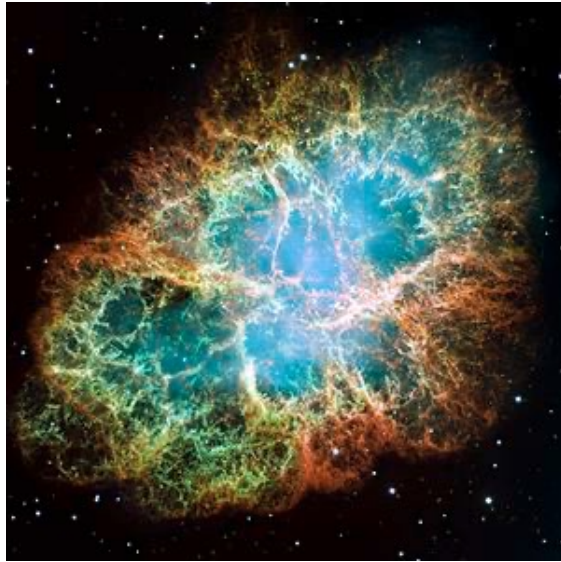
The “maximum likelihood estimate of the parameter” is the most intuitively obvious estimate:

$$P(W = w_i | Y = \text{Spam}) = \frac{\text{Count}(W = w_i, Y = \text{Spam})}{\text{Count}(Y = \text{Spam})}$$

where “ $\text{Count}(W = w_i, Y = \text{Spam})$ ” means the number of times that the word w_i occurs in the Spam portion of the training corpus, and “ $\text{Count}(Y = \text{Spam})$ ” is the total number of words in the Spam portion.

What is the probability that the sun will fail to rise tomorrow?

- # times we have observed the sun to rise = 1,825,000
- # times we have observed the sun not to rise = 0
- Estimated probability the sun will not rise = $\frac{0}{0+1,825,000} = 0$



Oops....

Laplace Smoothing

- The basic idea: add k “unobserved observations” to every possible event
- # times the sun has risen or might have ever risen = $1,825,000+k$
- # times the sun has failed to rise or might have ever failed to rise = $0+k$
- Estimated probability the sun will rise tomorrow = $\frac{1,825,000+k}{1,825,000+2k}$
- Estimated probability the sun will not rise = $\frac{k}{1,825,000+2k}$
- Notice that, if you add these two probabilities together, you get 1.0.

Laplace smoothing with only in-vocabulary events

If the domain \mathcal{W} is known in advance, the Laplace smoothed estimates are:

$$P(W = w|Y = y) = \frac{k + \text{Count}(w, y)}{\sum_{v \in \mathcal{W}} (k + \text{Count}(v, y))}$$

...which satisfies $\sum_{w \in \mathcal{W}} P(W = w|Y = y) = 1$.

Laplace smoothing with out-of-vocabulary events

Suppose we don't know all the possible words, \mathcal{W} . Our “vocabulary” includes some of the words that might be used, but not all.

- In that case, we need to allocate some probability to the “out of vocabulary (OOV) words:” :

$$P(W = w|Y = y) = \begin{cases} \frac{k + \text{Count}(w, y)}{k + \sum_{v \in \mathcal{W}} (k + \text{Count}(v, y))} & \text{in - vocabulary} \\ \frac{k}{k + \sum_{v \in \mathcal{W}} (k + \text{Count}(v, y))} & \text{OOV} \end{cases}$$

...which satisfies $\sum_{w \in \{\mathcal{W}, \text{OOV}\}} P(W = w|Y = y) = 1$.

Types vs. Tokens

- A “type” is a word that is spelled differently from all other words.
 - The number of “types,” M , is the number of words in the dictionary
- A “token” is one instance of a word that appears on the page.
 - The number of “tokens,” N , is the number of words in the training dataset.

For example, this page has 67 tokens, but only 40 types.

Laplace smoothing with vs w/o OOVs

N = # tokens, M = # types.

Without OOVs:

$$P(W = w|Y = y) = \frac{k + \text{Count}(w, y)}{N + kM}$$

With OOVs:

$$P(W = w|Y = y) = \begin{cases} \frac{k + \text{Count}(w, y)}{N + k(M + 1)} & \text{in - vocabulary} \\ \frac{k}{N + k(M + 1)} & \text{OOV} \end{cases}$$

...which satisfies $\sum_{w \in \{w, \text{OOV}\}} P(W = w|Y = y) = 1$.

Quiz!

- Go to the PrairieLearn, try the quiz!

Conclusions

- Naïve Bayes classifier:

$$f(x) = \operatorname{argmax}(\log P(Y = y) + \log P(X = x|Y = y))$$

$$\log P(X = x|Y = y) \approx \sum_{i=1}^n \log P(W = w_i|Y = y)$$

- maximum likelihood parameter estimation:

$$P(W = w_i|Y = y) = \frac{\operatorname{Count}(w_i, y)}{N}$$

- Laplace Smoothing:

$$P(W = w|Y = y) = \begin{cases} \frac{k + \operatorname{Count}(w, y)}{N + k(M + 1)} & \text{in - vocabulary} \\ \frac{k}{N + k(M + 1)} & \text{OOV} \end{cases}$$