

Classification: Naive Bayes

APAM E4990
Modeling Social Data

Jake Hofman

Columbia University

March 29, 2019

Learning by example

Fwd: Yahoo! supercomputing cluster RFP - i have no idea. i have no idea. O
non urgent - whoops! yes that's what i meant, thanks for decoding my questi
SourceForge.net: variational bayes for network modularity - can i get admin
Byline - iPhone Apps, iPhone 3G apps and iPod touch Applications Gallery a
Laurence J. Peter: Facts are stubborn things, but statistics are more pliable.
Re: JAFOS 2008, Applied Math Session - yes. the listening post dude. On N
Access to over 5,000 Health Plan Choices! - Affordable health insurance. Ins
More effective - If you are having trouble viewing this email click here. Thurs
Special Offer! Cialis, Viagra, VicodinES! - Order all your Favorite Rx~Medica
Financial Aid Available: Find Funding for Your Education - Get the financial a
Find The Perfect School and Financial Aid for your College Degree - HI ! It h
PHARMA_viagra_PHARMA_cialis - Wanted: web store with remedies. N

Learning by example

Fwd: Yahoo! supercomputing cluster RFP - i have no idea. i have no idea. O
non urgent - whoops! yes that's what i meant, thanks for decoding my questi
SourceForge.net: variational bayes for network modularity - can i get admin
Byline - iPhone Apps, iPhone 3G apps and iPod touch Applications Gallery a
Laurence J. Peter: Facts are stubborn things, but statistics are more pliable.
Re: JAFOS 2008, Applied Math Session - yes. the listening post dude. On N
Access to over 5,000 Health Plan Choices! - Affordable health insurance. Ins
More effective - If you are having trouble viewing this email click here. Thurs
Special Offer! Cialis, Viagra, VicodinES! - Order all your Favorite Rx~Medica
Financial Aid Available: Find Funding for Your Education - Get the financial a
Find The Perfect School and Financial Aid for your College Degree - HI ! It h
PHARMA_viagra_PHARMA_cialis - Wanted: web store with remedies. N

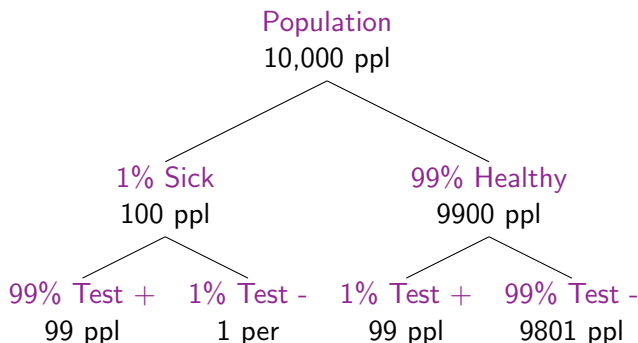
- How did you solve this problem?
- Can you make this process explicit (e.g. write code to do so)?

Diagnoses a la Bayes¹

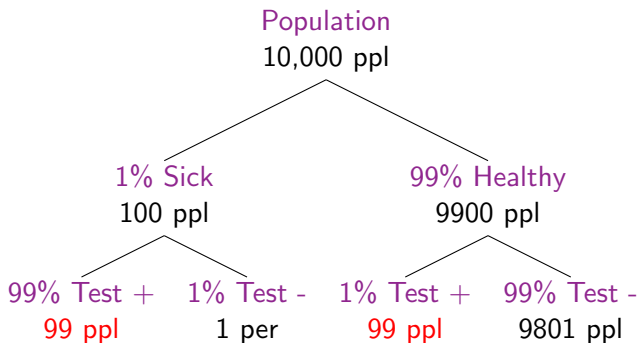
- You're testing for a rare disease:
 - 1% of the population is infected
- You have a highly sensitive and specific test:
 - 99% of sick patients test positive
 - 99% of healthy patients test negative
- Given that a patient tests positive, what is probability the patient is sick?

¹Wiggins, SciAm 2006

Diagnoses a la Bayes

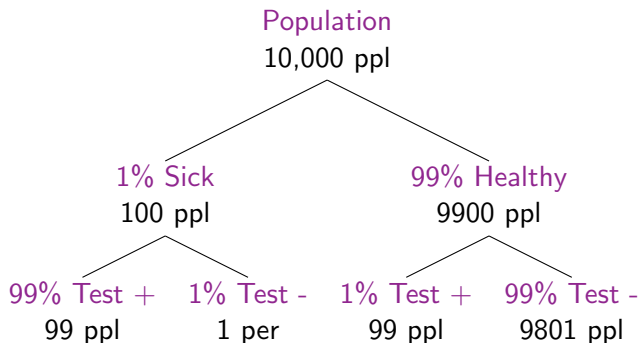


Diagnoses a la Bayes



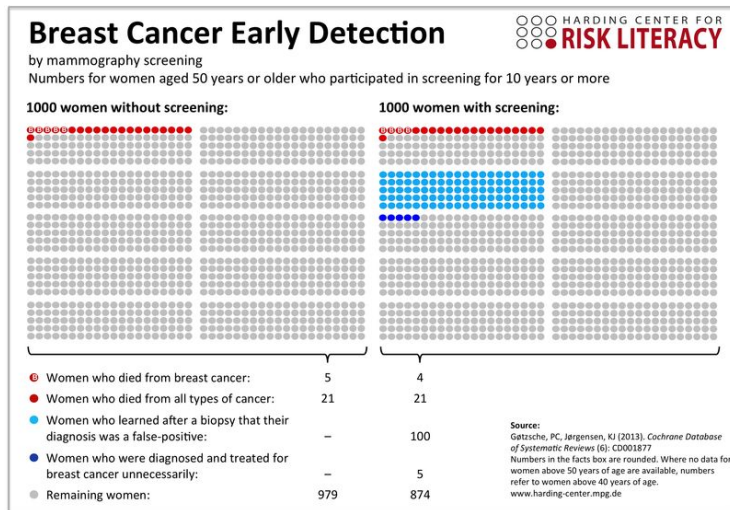
So given that a patient tests positive (198 ppl), there is a 50% chance the patient is sick (99 ppl)!

Diagnoses a la Bayes



The small error rate on the large healthy population produces many false positives.

Natural frequencies a la Gigerenzer²



²<http://bit.ly/ggbbc>

Inverting conditional probabilities

Bayes' Theorem

Equate the far right- and left-hand sides of product rule

$$p(y|x) p(x) = p(x, y) = p(x|y) p(y)$$

and divide to get the probability of y given x from the probability of x given y :

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

where $p(x) = \sum_{y \in \Omega_Y} p(x|y) p(y)$ is the normalization constant.

Diagnoses a la Bayes

Given that a patient tests positive, what is probability the patient is sick?

$$p(sick|+) = \frac{\overbrace{p(+|sick)}^{99/100} \overbrace{p(sick)}^{1/100}}{\underbrace{p(+)}_{99/100^2 + 99/100^2 = 198/100^2}} = \frac{99}{198} = \frac{1}{2}$$

where $p(+) = p(+|sick) p(sick) + p(+|healthy) p(healthy)$.

(Super) Naive Bayes

We can use Bayes' rule to build a one-word spam classifier:

$$p(\text{spam}|\text{word}) = \frac{p(\text{word}|\text{spam}) p(\text{spam})}{p(\text{word})}$$

where we estimate these probabilities with ratios of counts:

$$\hat{p}(\text{word}|\text{spam}) = \frac{\# \text{ spam docs containing word}}{\# \text{ spam docs}}$$

$$\hat{p}(\text{word}|\text{ham}) = \frac{\# \text{ ham docs containing word}}{\# \text{ ham docs}}$$

$$\hat{p}(\text{spam}) = \frac{\# \text{ spam docs}}{\# \text{ docs}}$$

$$\hat{p}(\text{ham}) = \frac{\# \text{ ham docs}}{\# \text{ docs}}$$

(Super) Naive Bayes

```
$ ./enron_naive_bayes.sh meeting  
1500 spam examples  
3672 ham examples  
16 spam examples containing meeting  
153 ham examples containing meeting
```

```
estimated P(spam) = .2900  
estimated P(ham) = .7100  
estimated P(meeting|spam) = .0106  
estimated P(meeting|ham) = .0416
```

```
P(spam|meeting) = .0923
```

(Super) Naive Bayes

```
$ ./enron_naive_bayes.sh money  
1500 spam examples  
3672 ham examples  
194 spam examples containing money  
50 ham examples containing money
```

```
estimated  $P(\text{spam}) = .2900$   
estimated  $P(\text{ham}) = .7100$   
estimated  $P(\text{money}|\text{spam}) = .1293$   
estimated  $P(\text{money}|\text{ham}) = .0136$ 
```

```
 $P(\text{spam}|\text{money}) = .7957$ 
```

(Super) Naive Bayes

```
$ ./enron_naive_bayes.sh enron
1500 spam examples
3672 ham examples
0 spam examples containing enron
1478 ham examples containing enron
```

```
estimated P(spam) = .2900
estimated P(ham) = .7100
estimated P(enron|spam) = 0
estimated P(enron|ham) = .4025
```

```
P(spam|enron) = 0
```

Naive Bayes

Represent each document by a binary vector \vec{x} where $x_j = 1$ if the j -th word appears in the document ($x_j = 0$ otherwise).

Modeling each word as an *independent* Bernoulli random variable, the probability of observing a document \vec{x} of class c is:

$$p(\vec{x}|c) = \prod_j \theta_{jc}^{x_j} (1 - \theta_{jc})^{1-x_j}$$

where θ_{jc} denotes the probability that the j -th word occurs in a document of class c .

Naive Bayes

Using this likelihood in Bayes' rule and taking a logarithm, we have:

$$\begin{aligned}\log p(c|\vec{x}) &= \log \frac{p(\vec{x}|c) p(c)}{p(\vec{x})} \\ &= \sum_j x_j \log \frac{\theta_{jc}}{1 - \theta_{jc}} + \sum_j \log(1 - \theta_{jc}) + \log \frac{\theta_c}{p(\vec{x})}\end{aligned}$$

where θ_c is the probability of observing a document of class c .

Naive Bayes

We can eliminate $p(\vec{x})$ by calculating the log-odds:

$$\log \frac{p(1|\vec{x})}{p(0|\vec{x})} = \sum_j x_j \underbrace{\log \frac{\theta_{j1}(1 - \theta_{j0})}{\theta_{j0}(1 - \theta_{j1})}}_{w_j} + \underbrace{\sum_j \log \frac{1 - \theta_{j1}}{1 - \theta_{j0}} + \log \frac{\theta_1}{\theta_0}}_{w_0}$$

which gives a linear classifier of the form $\vec{w} \cdot \vec{x} + w_0$

Naive Bayes

We train by counting words and documents within classes to estimate θ_{jc} and θ_c :

$$\begin{aligned}\hat{\theta}_{jc} &= \frac{n_{jc}}{n_c} \\ \hat{\theta}_c &= \frac{n_c}{n}\end{aligned}$$

and use these to calculate the weights \hat{w}_j and bias \hat{w}_0 :

$$\begin{aligned}\hat{w}_j &= \log \frac{\hat{\theta}_{j1}(1 - \hat{\theta}_{j0})}{\hat{\theta}_{j0}(1 - \hat{\theta}_{j1})} \\ \hat{w}_0 &= \sum_j \log \frac{1 - \hat{\theta}_{j1}}{1 - \hat{\theta}_{j0}} + \log \frac{\hat{\theta}_1}{\hat{\theta}_0}.\end{aligned}$$

We predict by simply adding the weights of the words that appear in the document to the bias term.

Naive Bayes

In practice, this works better than one might expect given its simplicity³

³<http://www.jstor.org/pss/1403452>

Naive Bayes

Training is computationally cheap and scalable, and the model is easy to update given new observations³

³<http://www.springerlink.com/content/wu3g458834583125/>

Naive Bayes

Performance varies with document representations and corresponding likelihood models³

³<http://ceas.cc/2006/15.pdf>

Naive Bayes

It's often important to smooth parameter estimates (e.g., by adding pseudocounts) to avoid overfitting

$$\hat{\theta}_{jc} = \frac{n_{jc} + \alpha}{n_c + \alpha + \beta}$$