

MITx: 6.008.1x Computational Probability and Inference

Heli

**Bookmarks** 

- **▶** Introduction
- Part 1: Probability and Inference
- Part 2: Inference in Graphical Models
- ▼ Part 3: Learning Probabilistic Models

Week 8: Introduction to Learning Probabilistic Models

Week 8: Introduction to
Parameter Learning Maximum Likelihood and
MAP Estimation

<u>Exercises due Nov 10, 2016 at 01:30 IST</u>

Week 8: Homework 6
Homework due Nov 10, 2016 at
01:30 IST

Part 3: Learning Probabilistic Models > Week 9: Parameter Learning - Naive Bayes Classification > The Naive Bayes Classifier: Training

## The Naive Bayes Classifier: Training

☐ Bookmark this page

THE NAIVE BAYES CLASSIFIER: TRAINING (PREFACE)

You should try to answer the following *before* watching the video below which presents the solution.

Note: "log" means natural log in these videos/notes.

**Practice problem:** Show that the log likelihood for our email spam detection setup can be written as

$$\log\Biggl(\prod_{i=1}^n p_{C,Y_1,\ldots,Y_J}(c^{(i)},y_1^{(i)},\ldots,y_J^{(i)}; heta)\Biggr) = f(s) + \sum_{j=1}^J g_j(p_j) + \sum_{j=1}^J h_j(q_j)$$

for some functions  $f, g_1, g_2, \ldots, g_J, h_1, h_2, \ldots, h_J$ . What are these functions? Note that what the above equation is saying is that the log likelihood decouples into functions where each function depends on just one of the parameters. Thus, when we want to maximize over  $\theta$ , what we can do is maximize over each parameter in  $\theta$  separately! For example, to find what the ML estimate for s is, we only need to look at f(s).

The Naive Bayes Classifier: Training | Week 9: Parameter Learning - Naive Bayes Classification | 6.008.1x Courseware | edX

Week 9: Parameter
Learning - Naive Bayes
Classification

Week 9: Mini-project on Email Spam Detection

Mini-projects due Nov 17, 2016 at 01:30 IST

*Hint:* To show the above equation, you may find it helpful that we can write:

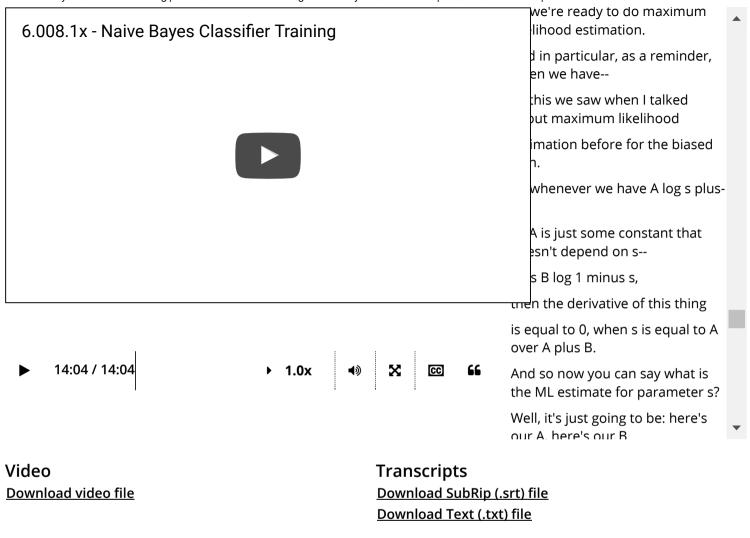
$$egin{aligned} p_C(c; heta) &= s^{1\{c= ext{spam}\}} (1-s)^{1-1\{c= ext{spam}\}} & ext{for } c \in \{ ext{spam}, ext{ham}\} \ p_{Y_j|C}(y_j \mid ext{ham}; heta) &= p_j^{y_j} (1-p_j)^{1-y_j} & ext{for } y_j \in \{0,1\} \ p_{Y_j|C}(y_j \mid ext{spam}; heta) &= q_j^{y_j} (1-q_j)^{1-y_j} & ext{for } y_j \in \{0,1\} \end{aligned}$$

**Practice problem:** After you figure out what the functions f,  $g_j$ 's, and  $h_j$ 's are, obtain the ML estimate for each of the parameters  $s, p_1, \ldots, p_J, q_1, \ldots, q_J$  by setting derivatives equal to 0.

*Hint:* You may find it helpful that for nonzero constants  $m{A}$  and  $m{B}$ ,

$$rac{d}{dt}\{A\log t + B\log(1-t)\} = 0 \qquad ext{when} \qquad t = rac{A}{A+B}.$$

The Naive Bayes Classifier: Training



These notes cover roughly the same content as the video:

THE NAIVE BAYES CLASSIFIER: TRAINING (COURSE NOTES)

## **Simplifying the log likelihood:** The log likelihood is given by

$$\begin{split} &\log\left(\prod_{i=1}^{n}p_{C,Y_{1},\ldots,Y_{J}}(c^{(i)},y_{1}^{(i)},\ldots,y_{J}^{(i)};\theta)\right) \\ &= \log\left(\prod_{i=1}^{n}\left[p_{C}(c^{(i)};\theta)\prod_{j=1}^{J}p_{Y_{j}|C}(y_{j}^{(i)}|c^{(i)};\theta)\right]\right) \\ &= \sum_{i=1}^{n}\left[\log p_{C}(c^{(i)};\theta) + \sum_{j=1}^{J}\log p_{Y_{j}|C}(y_{j}^{(i)}|c^{(i)};\theta)\right] \\ &= \underbrace{\sum_{i=1}^{n}\log p_{C}(c^{(i)};\theta)}_{(*)} + \underbrace{\sum_{i=1}^{n}\sum_{j=1}^{J}\log p_{Y_{j}|C}(y_{j}^{(i)}|c^{(i)};\theta)}_{(**)}. \end{split}$$

We next simplify the expressions (\*) and (\*\*).

First let's simplify term (\*):

$$egin{aligned} (*) &= \sum_{i=1}^n \log p_C(c^{(i)}; heta) \ &= \sum_{i=1}^n \left[ \mathbf{1}\{c = ext{``spam "}\} \log s + \mathbf{1}\{c = ext{``ham "}\} \log (1-s) 
ight] \ &= \left[ \sum_{i=1}^n \mathbf{1}\{c = ext{``spam "}\} 
ight] \log s + \left[ \sum_{i=1}^n \mathbf{1}\{c = ext{``ham "}\} 
ight] \log (1-s) \ & riangleq f(s). \end{aligned}$$

Next, we simplify (\*\*), splitting it up as to decouple  $p_j$  and  $q_j$ . To do this, we can split the summation over i into two sums, one accounting for all the ham emails and one accounting for all the spam emails:

$$\begin{split} & = \sum_{i=1}^{n} \sum_{j=1}^{J} \log p_{Y_{j}|C}(y_{j}^{(i)}|c^{(i)};\theta) \\ & = \sum_{i=1}^{n} \mathbf{1}\{c^{(i)} = \text{``ham''}\} \sum_{j=1}^{J} \log p_{Y_{j}|C}(y_{j}^{(i)}|c^{(i)};\theta) \\ & + \sum_{i=1}^{n} \mathbf{1}\{c^{(i)} = \text{``spam''}\} \sum_{j=1}^{J} \log p_{Y_{j}|C}(y_{j}^{(i)}|c^{(i)};\theta) \\ & = \sum_{i=1}^{n} \mathbf{1}\{c^{(i)} = \text{``ham''}\} \sum_{j=1}^{J} \left[y_{j}^{(i)} \log p_{j} + (1 - y_{j}^{(i)}) \log(1 - p_{j})\right] \\ & + \sum_{i=1}^{n} \mathbf{1}\{c^{(i)} = \text{``spam''}\} \sum_{j=1}^{J} \left[y_{j}^{(i)} \log q_{j} + (1 - y_{j}^{(i)}) \log(1 - q_{j})\right] \\ & = \sum_{j=1}^{J} \sum_{i=1}^{n} \mathbf{1}\{c^{(i)} = \text{``ham''}\} \left[y_{j}^{(i)} \log p_{j} + (1 - y_{j}^{(i)}) \log(1 - p_{j})\right] \\ & \stackrel{\triangleq g_{j}(p_{j})}{=} \\ & + \sum_{j=1}^{J} \sum_{i=1}^{n} \mathbf{1}\{c^{(i)} = \text{``spam''}\} \left[y_{j}^{(i)} \log q_{j} + (1 - y_{j}^{(i)}) \log(1 - q_{j})\right]. \end{split}$$

In summary:

$$f(s) = \left[\sum_{i=1}^n \mathbf{1}\{c^{(i)} = ext{``spam "}\}
ight] \log s + \left[\sum_{i=1}^n \mathbf{1}\{c^{(i)} = ext{``ham "}\}
ight] \log (1-s), \ g_j(p_j) = \left[\sum_{i=1}^n \mathbf{1}\{c^{(i)} = ext{``ham "}\}y_j^{(i)}
ight] \log p_j + \left[\sum_{i=1}^n \mathbf{1}\{c^{(i)} = ext{``ham "}\}(1-y_j^{(i)})
ight] \log (1-p_j), \ h_j(q_j) = \left[\sum_{i=1}^n \mathbf{1}\{c^{(i)} = ext{``spam "}\}y_j^{(i)}
ight] \log q_j + \left[\sum_{i=1}^n \mathbf{1}\{c^{(i)} = ext{``spam "}\}(1-y_j^{(i)})
ight] \log (1-q_j).$$

**Setting derivatives to 0:** The ML estimate for s is  $\hat{s}=rg\max_{s\in[0,1]}f(s)$ , which occurs when  $\frac{df}{ds}=0$ . Using the hint, we see that

$$f(s) = \underbrace{\left[\sum_{i=1}^n \mathbf{1}\{c^{(i)} = \text{``spam''}\}
ight]}_{A} \log s + \underbrace{\left[\sum_{i=1}^n \mathbf{1}\{c^{(i)} = \text{``ham''}\}
ight]}_{B} \log(1-s)$$

has derivative equal to 0 when

$$\hat{s} = \frac{A}{A+B} = \frac{\sum_{i=1}^{n} \mathbf{1}\{c^{(i)} = \text{``spam''}\}}{\sum_{i=1}^{n} \mathbf{1}\{c^{(i)} = \text{``spam''}\} + \sum_{i=1}^{n} \mathbf{1}\{c^{(i)} = \text{``ham''}\}} = \frac{\sum_{i=1}^{n} \mathbf{1}\{c^{(i)} = \text{``spam''}\}}{n}.$$

This result is intuitive — it's the number of emails labeled "spam" divided by the total number of emails.

The ML estimate for  $p_j$  is  $\hat{p}_j = \arg\max_{p_j \in [0,1]} g_j(p_j)$ , which occurs when  $\frac{dg_j}{dp_j} = 0$ . Again using the hint, we see that

$$g_j(p_j) = \underbrace{\left[\sum_{i=1}^n \mathbf{1}\{c^{(i)} = ext{``ham "}\}y_j^{(i)}
ight]}_A \log p_j + \underbrace{\left[\sum_{i=1}^n \mathbf{1}\{c^{(i)} = ext{``ham "}\}(1-y_j^{(i)})
ight]}_B \log (1-p_j)$$

has derivative equal to 0 when

$$egin{aligned} \hat{p}_j &= rac{A}{A+B} \ &= rac{\sum_{i=1}^n \mathbf{1}\{c^{(i)} = ext{``ham "}\}y_j^{(i)}}{\sum_{i=1}^n \mathbf{1}\{c^{(i)} = ext{``ham "}\}y_j^{(i)} + \sum_{i=1}^n \mathbf{1}\{c^{(i)} = ext{``ham "}\}(1-y_j^{(i)})} \ &= rac{\sum_{i=1}^n \mathbf{1}\{c^{(i)} = ext{``ham "}\}y_j^{(i)}}{\sum_{i=1}^n \mathbf{1}\{c^{(i)} = ext{``ham "}\}}. \end{aligned}$$

This result is also intuitive — it's the number of times word j occurred in an email labeled "ham" divided by the total number of emails labeled "ham".

Finally, by pattern-matching, the ML estimate for  $oldsymbol{q_j}$  is

$$\hat{q}_j = rac{\sum_{i=1}^n \mathbf{1}\{c^{(i)} = \text{``spam''}\}y_j^{(i)}}{\sum_{i=1}^n \mathbf{1}\{c^{(i)} = \text{``spam''}\}}.$$

Wonderful, now we can write up an algorithm that computes all those ML estimates above. Once we learn the parameters  $\theta$ , we can treat them as fixed and start doing prediction.

## Discussion

**Topic:** Parameter Learning - Naive Bayes Classification / The Naive Bayes

Classifier: Training

**Show Discussion** 

© All Rights Reserved



© 2016 edX Inc. All rights reserved except where noted. EdX, Open edX and the edX and Open EdX logos are registered trademarks or trademarks of edX Inc.















