

Naïve Bayes Model

Lecture 5: Sept 25, 2013

CS886-2 Natural Language Understanding
University of Waterloo

CS886 Lecture Slides (c) 2013 P. Poupart

1

Classification

- NLP tasks:
 - Spam filtering
 - Sentiment analysis
 - Text categorization
- Classification techniques
 - Decision trees
 - Support Vector Machines
 - Naïve Bayes

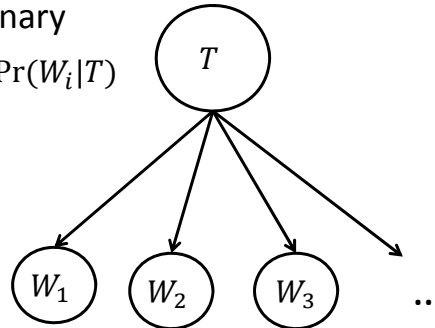
CS886 Lecture Slides (c) 2013 P. Poupart

2

Naïve Bayes Model

Generative model

- Sample topic from $\Pr(T)$
- For each word in dictionary
Sample $W_i = 1/0$ from $\Pr(W_i|T)$



CS886 Lecture Slides (c) 2013 P. Poupart

3

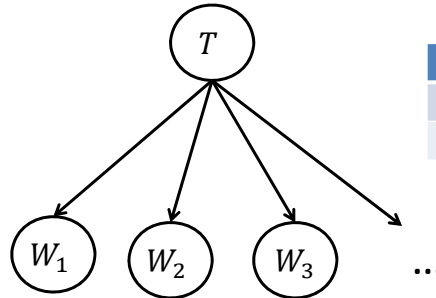
Simple Bayesian Network

- The name **Naïve Bayes** comes from the fact that it is a very **simple Bayesian network**.
- Bayesian network:
 - Graphical model:
 - Nodes: random variables
 - Edges: conditional dependencies
 - Conditional distributions: $\Pr(V|\text{parents}(V))$

CS886 Lecture Slides (c) 2013 P. Poupart

4

Model Specification



T	$\Pr(T)$
finance	0.2
sports	0.8

W_1	T	$\Pr(W_1 T)$	W_2	T	$\Pr(W_2 T)$	W_3	T	$\Pr(W_3 T)$
1	finance	0.1	1	finance	0.7	1	finance	0.75
0	finance	0.9	0	finance	0.3	0	finance	0.25
1	sports	0.8	1	sports	0.05	1	sports	0.2
0	sports	0.2	0	sports	0.95	0	sports	0.8

CS886 Lecture Slides (c) 2013 P. Poupart

5

Joint distribution

- Joint distribution:

$$\Pr(T, W_1, W_2, W_3, \dots) = \Pr(T) \prod_i \Pr(W_i|T)$$

- Joint assignment: document (bag of words):

$$d_1: T = \textit{sports}, W_1 = 1, W_2 = 1, W_3 = 0, \dots$$

$$d_2: T = \textit{finance}, W_1 = 0, W_2 = 1, W_3 = 1, \dots$$

...

CS886 Lecture Slides (c) 2013 P. Poupart

6

Classification

- Classification by probabilistic inference:

$$\Pr(T|W_1, W_2, W_3, \dots, W_n)$$

- Computation:

$$\begin{aligned} \Pr(T|W_1, W_2, W_3, \dots, W_n) &= \frac{\Pr(T, W_1, W_2, W_3, \dots, W_n)}{\Pr(W_1, W_2, W_3, \dots, W_n)} \\ &= \frac{\Pr(T) \prod_i \Pr(W_i|T)}{\Pr(W_1, W_2, W_3, \dots, W_n)} \\ &= k \Pr(T) \prod_i \Pr(W_i|T) \end{aligned}$$

CS886 Lecture Slides (c) 2013 P. Poupart

7

Example

- $d_1 : \Pr(T|W_1 = 1, W_2 = 1, W_3 = 0)$

$$T = \textit{sports}: k(0.8)(0.8)(0.05)(0.8) = \mathbf{0.88}$$

$$T = \textit{finance}: k(0.2)(0.1)(0.7)(0.25) = \mathbf{0.12}$$

- $d_2 : \Pr(T|W_1 = 0, W_2 = 1, W_3 = 1)$

$$T = \textit{sports}: k(0.8)(0.2)(0.05)(0.2) = \mathbf{0.17}$$

$$T = \textit{finance}: k(0.2)(0.9)(0.7)(0.75) = \mathbf{0.83}$$

CS886 Lecture Slides (c) 2013 P. Poupart

8

Parameter Estimation

- Idea: set conditional distributions to **relative frequency counts**

$$\Pr(T = \textit{finance}) = \frac{\#(T=\textit{finance})}{\#\textit{documents}}$$

$$\Pr(W_i = 1 | T = \textit{finance}) = \frac{\#(W_i=1 \wedge T=\textit{finance})}{\#(T=\textit{finance})}$$

Maximum Likelihood

- Parameter estimation by relative frequency follows from the principle of **maximum likelihood**
- Maximum likelihood:

$$\theta^* = \operatorname{argmax}_{\theta} \Pr(\textit{data} | \theta)$$

where θ is the set of parameters

Estimation of Topic Distribution

- Let $\theta = \Pr(T = \textit{sports})$
- Corpus:
 - $\#s$: # of documents about sports
 - $\#f$: # of documents about finance

CS886 Lecture Slides (c) 2013 P. Poupart

11

Estimation of Topic Distribution

- 1) Likelihood expression

$$\Pr(\textit{data}|\theta) = \theta^{\#s}(1 - \theta)^{\#f}$$
- 2) log likelihood

$$\log \Pr(\textit{data}|\theta) = \#s \log \theta + \#f \log(1 - \theta)$$
- 3) log likelihood derivative

$$\frac{\partial(\log \Pr(\textit{data}|\theta))}{\partial \theta} = \frac{\#s}{\theta} - \frac{\#f}{1 - \theta}$$
- 4) ML hypothesis

$$\frac{\#s}{\theta} - \frac{\#f}{1 - \theta} = 0 \rightarrow \theta = \frac{\#s}{\#s + \#f}$$

CS886 Lecture Slides (c) 2013 P. Poupart

12

Full Parameter Estimation

- Let $\theta = \Pr(T = \textit{sports})$
 $\theta_{si} = \Pr(W_i = 1 | T = \textit{sports})$
 $\theta_{fi} = \Pr(W_i = 1 | T = \textit{finance})$
- Corpus:
 - $\#s$: # of documents about sports
 - $\#f$: # of documents about finance
 - $\#s1i$: # of documents about sports with w_i
 - $\#s0i$: # of documents about sports without w_i
 - $\#f1i$: # of documents about finance with w_i
 - $\#f0i$: # of documents about finance without w_i

CS886 Lecture Slides (c) 2013 P. Poupart

13

Full Parameter Estimation

- 1) Likelihood expression
 $\Pr(\textit{data}|\theta)$

$$= \theta^{\#s} (1 - \theta)^{\#f} \prod_i \theta_{si}^{\#s1i} (1 - \theta_{si})^{\#s0i} \theta_{fi}^{\#f1i} (1 - \theta_{fi})^{\#f0i}$$
- ...
- 4) ML hypothesis

$$\frac{\#s}{\theta} - \frac{\#f}{1-\theta} = 0 \rightarrow \theta = \frac{\#s}{\#s+\#f}$$

$$\frac{\#s1i}{\theta_{si}} - \frac{\#s0i}{1-\theta_{si}} = 0 \rightarrow \theta_{si} = \frac{\#s1i}{\#s1i+\#s0i}$$

$$\frac{\#f1i}{\theta_{fi}} - \frac{\#f0i}{1-\theta_{fi}} = 0 \rightarrow \theta_{fi} = \frac{\#f1i}{\#f1i+\#f0i}$$

CS886 Lecture Slides (c) 2013 P. Poupart

14

Laplace smoothing

- An important case of overfitting happens when there is no sample for a certain outcome

- E.g. no document about sports contains w_i

$$\Pr(W_i = 1 | T = \text{sports}) = \theta_{s1i} = \frac{\#s1i}{\#s1i + \#s0i} = 0$$

- Zero probabilities are dangerous: they rule out outcomes

- Solution: Laplace (add-one) smoothing

- Add 1 to all frequencies

$$\Pr(W_i = 1 | T = \text{sports}) = \theta_{s1i} = \frac{\#s1i + 1}{\#s1i + \#s0i + 2}$$

- Much better results in practice