

Expectation Maximization

Lecture 6: Sept 27, 2013

CS886-2 Natural Language Understanding
University of Waterloo

CS886 Lecture Slides (c) 2013 P. Poupart

1

Text Categorization

- Algorithms
 - Naïve Bayes model
 - Support vector machines
 - Decision Trees
- Assumption: **labeled articles available for training**
 - **Hand labeling is expensive**
 - **Lots of unlabeled documents available**

CS886 Lecture Slides (c) 2013 P. Poupart

2

Text Categorization

- **Supervised learning: categorization**
 - Data: all training documents labeled with category
 - Alg: Relative frequency counts (maximum likelihood)
- **Unsupervised learning: clustering**
 - Data: no training document labeled with category
 - Alg: Expectation maximization (approx. max. likelihood)
- **Semi-supervised learning: mix of categorization and clustering**
 - Data: some labeled documents with some unlabeled documents
 - Alg: Expectation maximization (approx. max. likelihood)

CS886 Lecture Slides (c) 2013 P. Poupart

3

Supervised Maximum Likelihood

- **Notation:**
 - h : hypothesis
 - E : observable variables
 - e : evidence (data)
- **Supervised maximum likelihood**

$$h_{ML} = \operatorname{argmax}_h \Pr(e|h)$$
- **Supervised text categorization**

$$h_{ML} = \operatorname{argmax}_h \Pr(t, w_1, w_2, w_3, \dots |h)$$

CS886 Lecture Slides (c) 2013 P. Poupart

4

Unsupervised Maximum Likelihood

- Notation:
 - \mathbf{Z} : hidden (unobservable) variables
- Unsupervised Maximum likelihood

$$\begin{aligned} h_{ML} &= \operatorname{argmax}_h \Pr(\mathbf{e}|h) \\ &= \operatorname{argmax}_h \sum_{\mathbf{Z}} \Pr(\mathbf{Z}, \mathbf{e}|h) \end{aligned}$$

- Unsupervised Text Categorization

$$\begin{aligned} h_{ML} &= \operatorname{argmax}_h \Pr(w_1, w_2, w_3, \dots |h) \\ &= \operatorname{argmax}_h \sum_T \Pr(T, w_1, w_2, w_3, \dots |h) \end{aligned}$$

CS886 Lecture Slides (c) 2013 P. Poupart

5

“Direct” maximum likelihood

- Notation: V_1, V_2, \dots (all hidden and obs. variables)
 - E.g.: $V_1 = T, V_2 = W_1, V_3 = W_2, \dots$
- Directly maximizing likelihood in the unsupervised case is difficult
- $h_{ML} = \operatorname{argmax}_h \sum_{\mathbf{Z}} \Pr(\mathbf{Z}, \mathbf{e}|h)$

$$= \operatorname{argmax}_h \sum_{\mathbf{Z}} \prod_l \Pr(V_l | \text{parents}(V_l))$$

$$= \operatorname{argmax}_h \log \sum_{\mathbf{Z}} \prod_l \Pr(V_l | \text{parents}(V_l))$$

Problem: can't push log past sum to linearize product

CS886 Lecture Slides (c) 2013 P. Poupart

6

Expectation-Maximization (EM)

- EM algorithm
 - Intuition: if we knew the missing values, computing h_{ML} would be trivial
- Guess h_{ML}
- Iterate
 - **Expectation**: based on h_{ML} , compute expectation of the missing values
 - **Maximization**: based on expected missing values, compute new estimate of h_{ML}

CS886 Lecture Slides (c) 2013 P. Poupart

7

Expectation-Maximization (EM)

- More formally:
 - Approximate maximum likelihood
 - Iteratively compute:

$$h_{j+1} = \underset{h}{\operatorname{argmax}} \underbrace{\sum_{\mathbf{Z}} \Pr(\mathbf{Z} | h_j, \mathbf{e}) \log \Pr(\mathbf{e}, \mathbf{Z} | h)}_{\text{Expectation}}$$

$$\underbrace{\hspace{10em}}_{\text{Maximization}}$$

CS886 Lecture Slides (c) 2013 P. Poupart

8

Expectation-Maximization (EM)

- Derivation

$$\begin{aligned} \log \Pr(e|h) &= \log \left[\frac{\Pr(e, Z|h)}{\Pr(Z|e, h)} \right] \\ &= \log \Pr(e, Z|h) - \log \Pr(Z|e, h) \\ &= \sum_Z \Pr(Z|e, h) \log \Pr(e, Z|h) \\ &\quad - \sum_Z \Pr(Z|e, h) \log \Pr(Z|e, h) \\ &\geq \sum_Z \Pr(Z|e, h) \log \Pr(e, Z|h) \end{aligned}$$

- EM finds a **local maximum** of

$$\sum_Z \Pr(Z|e, h) \log \Pr(e, Z|h)$$

which is a **lower bound** of $\log P(\mathbf{e}|h)$

CS886 Lecture Slides (c) 2013 P. Poupart

9

Expectation-Maximization (EM)

- **Log inside sum can linearize product**

$$\begin{aligned} h_{j+1} &= \operatorname{argmax}_h \sum_Z \Pr(\mathbf{Z}|h_j, \mathbf{e}) \log \Pr(\mathbf{e}, \mathbf{Z}|h) \\ &= \operatorname{argmax}_h \sum_Z \Pr(\mathbf{Z}|h_j, \mathbf{e}) \log \prod_l \Pr(V_l | \text{parents}(V_l), h) \\ &= \operatorname{argmax}_h \sum_Z \Pr(\mathbf{Z}|h_j, \mathbf{e}) \sum_l \log \Pr(V_l | \text{parents}(V_l), h) \end{aligned}$$

- **Monotonic improvement of likelihood**

$$\Pr(\mathbf{e}|h_{j+1}) \geq \Pr(\mathbf{e}|h_j)$$

CS886 Lecture Slides (c) 2013 P. Poupart

10

Algorithm

- Iterative Optimization

$$h_{j+1} = \operatorname{argmax}_h \sum_Z \Pr(\mathbf{Z}|h_j, \mathbf{e}) \sum_l \log \Pr(V_l | \text{parents}(V_l), h)$$

Initialize h_0

Alternate between

1) Estimate Z based on h_j by computing
 $\Pr(\mathbf{Z}|h_j, \mathbf{e})$

2) Find h_{j+1} according to

$$h_{j+1} = \operatorname{argmax}_h \sum_Z \Pr(\mathbf{Z}|h_j, \mathbf{e}) \sum_l \log \Pr(V_l | \text{parents}(V_l), h)$$

Text Categorization Example

- Let $\theta = \Pr(T = \text{sports})$
 $\theta_{si} = \Pr(W_i = 1 | T = \text{sports})$
 $\theta_{fi} = \Pr(W_i = 1 | T = \text{finance})$

- 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

EM for Text Categorization

Algorithm summary

Initialize $\theta, \theta_{si}, \theta_{fi}$ (with any values)

Repeat

For each document d compute

$$\text{prob}(d, s) \leftarrow \Pr(T^d = \text{sports} | \theta, \theta_{si}, \theta_{fi}, w_i^d \forall i)$$

$$\text{prob}(d, f) \leftarrow \Pr(T^d = \text{finance} | \theta, \theta_{si}, \theta_{fi}, w_i^d \forall i)$$

$$E[\#s] \leftarrow \sum_d \text{prob}(d, s), \quad E[\#f] \leftarrow \sum_d \text{prob}(d, f)$$

$$E[\#s1i] \leftarrow \sum_d \delta(w_i^d = 1) \text{prob}(d, s), \quad E[\#s0i] \leftarrow \sum_d \delta(w_i^d = 0) \text{prob}(d, s) \quad \forall i$$

$$E[\#f1i] \leftarrow \sum_d \delta(w_i^d = 1) \text{prob}(d, f), \quad E[\#f0i] \leftarrow \sum_d \delta(w_i^d = 0) \text{prob}(d, f) \quad \forall i$$

$$\theta \leftarrow \frac{E[\#s]}{E[\#s] + E[\#f]}, \quad \theta_{si} \leftarrow \frac{E[\#s1i]}{E[\#s1i] + E[\#s0i]}, \quad \theta_{fi} \leftarrow \frac{E[\#f1i]}{E[\#f1i] + E[\#f0i]} \quad \forall i$$

Until convergence

CS886 Lecture Slides (c) 2013 P. Poupart

13

Data

- Data: 1000 documents
- Frequency: # of documents

W_1	W_2	W_3	Frequency
1	1	1	273
1	1	0	93
1	0	1	104
1	0	0	90
0	1	1	79
0	1	0	100
0	0	1	94
0	0	0	167

CS886 Lecture Slides (c) 2013 P. Poupart

14

Initialization

- Set the parameters to some initial values
 - Make educated guess
 - Since EM converges to local optimum, initial values influence the final values that it will converge to
- Example:

$$\theta = 0.5, \theta_{si} = 0.8 \forall i, \theta_{fi} = 0.3 \forall i$$

CS886 Lecture Slides (c) 2013 P. Poupart

15

Estimate T

- 8 types of documents
- Estimate T for each type of document

w_1	w_2	w_3	$prob(d, s)$	$prob(d, f)$
1	1	1	$k(0.6)(0.8)(0.8)(0.8) = 0.966$	$k(0.4)(0.3)(0.3)(0.3) = 0.034$
1	1	0	$k(0.6)(0.8)(0.8)(0.2) = 0.753$	$k(0.4)(0.3)(0.3)(0.7) = 0.247$
1	0	1	$k(0.6)(0.8)(0.2)(0.8) = 0.753$	$k(0.4)(0.3)(0.7)(0.3) = 0.247$
1	0	0	$k(0.6)(0.8)(0.2)(0.2) = 0.246$	$k(0.4)(0.3)(0.7)(0.7) = 0.754$
0	1	1	$k(0.6)(0.2)(0.8)(0.8) = 0.753$	$k(0.4)(0.7)(0.3)(0.3) = 0.247$
0	1	0	$k(0.6)(0.2)(0.8)(0.2) = 0.246$	$k(0.4)(0.7)(0.3)(0.7) = 0.754$
0	0	1	$k(0.6)(0.2)(0.2)(0.8) = 0.246$	$k(0.4)(0.7)(0.7)(0.3) = 0.754$
0	0	0	$k(0.6)(0.2)(0.2)(0.2) = 0.034$	$k(0.4)(0.7)(0.7)(0.7) = 0.966$

CS886 Lecture Slides (c) 2013 P. Poupart

16

Expected Frequencies

$$E[\#s] = 0.966(273) + 0.753(93 + 104 + 79) \\ + 0.246(90 + 100 + 94) + 0.034(167) = 547$$

$$E[\#f] = 1000 - 547 = 453$$

$$E[\#s_{11}] = 0.966(273) + 0.753(93 + 104) + 0.246(90) = 434$$

$$E[\#s_{12}] = 0.966(273) + 0.753(93 + 79) + 0.246(100) = 418$$

$$E[\#s_{13}] = 0.966(273) + 0.753(104 + 79) + 0.246(94) = 425$$

$$E[\#f_{11}] = 0.034(273) + 0.247(93 + 104) + 0.754(90) = 126$$

$$E[\#f_{12}] = 0.034(273) + 0.247(93 + 79) + 0.754(100) = 127$$

$$E[\#f_{13}] = 0.034(273) + 0.247(104 + 79) + 0.754(94) = 125$$

CS886 Lecture Slides (c) 2013 P. Poupart

17

Revised Parameters After One Iteration

$$\theta = \frac{547}{1000} = 0.547$$

$$\theta_{s1} = \frac{434}{547} = 0.793, \quad \theta_{f1} = \frac{126}{453} = 0.278$$

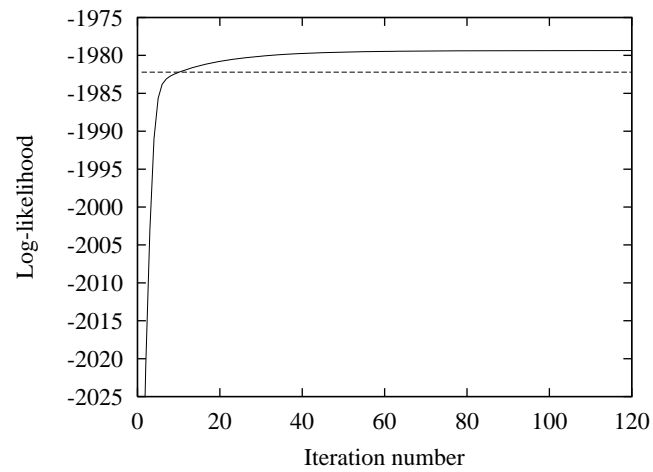
$$\theta_{s2} = \frac{418}{547} = 0.764, \quad \theta_{f2} = \frac{127}{453} = 0.280$$

$$\theta_{s3} = \frac{425}{547} = 0.777, \quad \theta_{f3} = \frac{125}{453} = 0.276$$

CS886 Lecture Slides (c) 2013 P. Poupart

18

Typical Log-Likelihood Curve



CS886 Lecture Slides (c) 2013 P. Poupart

19