



**Samueli**  
Computer Science



# CS145 Discussion: Week 10

## Naive Bayes, Topic Modeling, Final Review

Junheng Hao  
Friday, 12/11/2020

<b>Homework #6</b>	Due on 11:59 PM, <b>Dec 14</b> (Monday, final week)
<b>Final Exam (100 minutes via CCLE)</b>	Morning Session: Around 8:00-9:40 AM, <b>Dec 16</b> (Wednesday, final Week) Evening Session: Around 6:00-7:40 PM, <b>Dec 16</b> (Wednesday, final Week)
<b>Project: Final Report</b>	Due on 11:59 PM, <b>Dec 18</b> (Friday, final week)
<b>All regrade request</b>	Due on 11:59 PM, <b>Dec 21</b> (Monday, the week after final week)

**Note:** Practice exam for final exam setup has been added on CCLE, under the tab of “Week 11”.



- Basics
  - Text Data
  - MLE v.s. MAP
- Naïve Bayes
- pLSA
- Final Exam Q & A

# How to represent text data?

- Most common way: Bag-of-Words
  - Ignore the order of words
  - keep the count

	c1	c2	c3	c4	c5	m1	m2	m3	m4
c1: <i>Human machine interface for Lab ABC computer applications</i>	1	0	0	1	0	0	0	0	0
c2: <i>A survey of user opinion of computer system response time</i>	1	0	1	0	0	0	0	0	0
c3: <i>The EPS user interface management system</i>	1	1	0	0	0	0	0	0	0
c4: <i>System and human system engineering testing of EPS</i>	0	1	1	0	1	0	0	0	0
c5: <i>Relation of user-perceived response time to error measurement</i>	0	1	1	2	0	0	0	0	0
m1: <i>The generation of random, binary, unordered trees</i>	0	1	0	0	1	0	0	0	0
m2: <i>The intersection graph of paths in trees</i>	0	1	0	0	1	0	0	0	0
m3: <i>Graph minors IV: Widths of trees and well-quasi-ordering</i>	0	0	1	1	0	0	0	0	0
m4: <i>Graph minors: A survey</i>	0	1	0	0	0	0	0	0	1
	0	0	0	0	0	1	1	1	0
	0	0	0	0	0	0	1	1	1
	0	0	0	0	0	0	0	1	1

For document  $d$ ,  $\mathbf{x}_d = (x_{d1}, x_{d2}, \dots, x_{dN})$ , where  $x_{dn}$  is the number of words for  $n$ th word in the vocabulary

**Vector space model**

- *Maximum Likelihood* (maximize the likelihood):

$$h_{ML} = \arg \max_{h \in H} P(X | h)$$

- *Maximum a posteriori* (maximize the posterior):
  - Useful observation: it does not depend on the denominator  $P(X)$

$$h_{MAP} = \arg \max_{h \in H} P(h | X) = \arg \max_{h \in H} P(X | h)P(h)$$

- *Maximum Likelihood* (maximize the likelihood):

h is a parameter

$$h_{ML} = \arg \max_{h \in H} P(X | h)$$

- *Maximum a posteriori* (maximize the posterior):

- Useful observation: it does not depend on the denominator  $P(X)$

h is a random variable

$$h_{MAP} = \arg \max_{h \in H} P(h | X) = \arg \max_{h \in H} P(X | h)P(h)$$

# Naïve Bayes: Modeling



- A document is represented as a bag of words
  - $\mathbf{x}_d = (x_{d1}, x_{d2}, \dots, x_{dN})$ , where  $x_{dn}$  is the number of words for  $n$ th word in the vocabulary
- Model  $p(\mathbf{x}_d | y)$  for class  $y$ 
  - Follow multinomial distribution with parameter vector  $\boldsymbol{\beta}_y = (\beta_{y1}, \beta_{y2}, \dots, \beta_{yN})$ , i.e.,
    - $$p(\mathbf{x}_d | y) = \frac{(\sum_n x_{dn})!}{x_{d1}! x_{d2}! \dots x_{dN}!} \prod_n \beta_{yn}^{x_{dn}}$$
- Model  $p(y = j)$ 
  - Follow categorical distribution with parameter vector  $\boldsymbol{\pi} = (\pi_1, \pi_2, \dots, \pi_m)$ , i.e.,
    - $p(y = j) = \pi_j$

# Naïve Bayes: Inference

- Find  $y$  that maximizes  $p(y|\mathbf{x}_d)$ , which is equivalently to maximize

$$\begin{aligned}
 y^* &= \underset{y}{\operatorname{argmax}} p(\mathbf{x}_d, y) \\
 &= \underset{y}{\operatorname{argmax}} p(\mathbf{x}_d|y)p(y) \\
 &= \underset{y}{\operatorname{argmax}} \frac{(\sum_n x_{dn})!}{x_{d1}! x_{d2}! \dots x_{dN}!} \prod_n \beta_{yn}^{x_{dn}} \times \pi_y
 \end{aligned}$$

Constant for every class,  
denoted as  $c_d$

$$\begin{aligned}
 &= \underset{y}{\operatorname{argmax}} \prod_n \beta_{yn}^{x_{dn}} \times \pi_y \\
 &= \underset{y}{\operatorname{argmax}} \sum_n x_{dn} \log \beta_{yn} + \log \pi_y
 \end{aligned}$$



# Naïve Bayes: Learning

- Given a corpus and labels for each document
  - $D = \{(\mathbf{x}_d, y_d)\}$
  - Find the MLE estimators for  $\Theta = (\boldsymbol{\beta}_1, \boldsymbol{\beta}_2, \dots, \boldsymbol{\beta}_m, \boldsymbol{\pi})$
- The log likelihood function for the training dataset

$$\begin{aligned} \log L &= \log \prod_d p(\mathbf{x}_d, y_d | \Theta) = \sum_d \log p(\mathbf{x}_d, y_d | \Theta) \\ &= \sum_d \log p(\mathbf{x}_d | y_d) p(y_d) = \sum_d x_{dn} \log \beta_{yn} + \log \pi_{y_d} + \boxed{\log c_d} \end{aligned}$$

- The optimization problem

$$\max_{\Theta} \log L$$

s. t.

$$\pi_j \geq 0 \text{ and } \sum_j \pi_j = 1$$

$$\beta_{jn} \geq 0 \text{ and } \sum_n \beta_{jn} = 1 \text{ for all } j$$

Does not involve  
parameters, can be  
dropped for optimization  
purpose

# Naïve Bayes: Learning



- $\hat{\beta}_{jn} = \frac{\sum_{d:y_d=j} x_{dn}}{\sum_{d:y_d=j} \sum_{n'} x_{dn'}}$ 
  - $\sum_{d:y_d=j} x_{dn}$ : total count of word  $n$  in class  $j$
  - $\sum_{d:y_d=j} \sum_{n'} x_{dn'}$ : total count of words in class  $j$
- $\hat{\pi}_j = \frac{\sum_d 1(y_d=j)}{|D|}$ 
  - $1(y_d = j)$  is the indicator function, which equals to 1 if  $y_d = j$  holds
  - $|D|$ : total number of documents

- What if some word  $n$  does not appear in some class  $j$  in training dataset?
  - $\hat{\beta}_{jn} = \frac{\sum_{d:y_d=j} x_{dn}}{\sum_{d:y_d=j} \sum_{n'} x_{dn'}} = 0$
  - $\Rightarrow p(\mathbf{x}_d | y = j) \propto \prod_n \beta_{yn}^{x_{dn}} = 0$
  - But other words may have a strong indication the document belongs to class  $j$
- Solution: add-1 smoothing or Laplacian smoothing
  - $\hat{\beta}_{jn} = \frac{\sum_{d:y_d=j} x_{dn} + 1}{\sum_{d:y_d=j} \sum_{n'} x_{dn'} + N}$
  - $N$ : total number of words in the vocabulary
  - Check:  $\sum_n \hat{\beta}_{jn} = 1$ ?

# Naïve Bayes: Example

Index	Word	# in class1	# in class2
1	I	1	9
2	like	2	1
3	data	3	9
4	mining	4	1

beta = ?

pi = ?

classification result of “I like data”?

# Naïve Bayes: Example

Index	Word	# in class1	# in class2
1	I	1	9
2	like	2	1
3	data	3	9
4	mining	4	1

$$\text{beta}_{11} = 1/10$$

$$\text{beta}_{12} = 2/10$$

$$\text{beta}_{13} = 3/10$$

$$\text{beta}_{14} = 4/10$$

$$\text{beta}_{21} = 9/20$$

$$\text{beta}_{22} = 1/20$$

$$\text{beta}_{23} = 9/20$$

$$\text{beta}_{24} = 1/20$$

$$P(\text{class1} \mid \text{"I like data"}) \propto 1/3 * 1/10 * 2/10 * 3/10 = 6/3000$$

$$P(\text{class2} \mid \text{"I like data"}) \propto 2/3 * 9/20 * 1/20 * 9/20 = 81/12000$$

$$\text{pi}_1 = 1/3$$

$$\text{pi}_2 = 2/3$$

# Generative vs Discriminative Models



- Training classifiers involve estimating  $f: X \rightarrow Y$ , or  $P(Y|X)$
  - Generative classifiers → “*distribution*”
    - Assume some functional form for  $P(Y)$ ,  $P(X|Y)$
    - Estimate parameters of  $P(X|Y)$ ,  $P(Y)$  directly from training data
    - Use Bayes rule to calculate  $P(Y|X)$
    - Actually learn the underlying structure of the data
  - Discriminative Classifiers → “*boundary*”
    - Assume some functional form for  $P(Y|X)$
    - Estimate parameters of  $P(Y|X)$  directly from training data.
    - Learn the mappings directly from the points to the classes
  - Generative classifiers
    - Naive Bayes
  - Discriminative classifiers
    - Logistic Regression
    - Neural Network
    - SVM
- Q:** With the aim of classification only, which type of models may be less expensive?

Credit:

<https://stats.stackexchange.com/questions/12421/generative-vs-discriminative>

# Topic Models

Each document can have more than one topic.

A more flexible model than Naïve Bayes

"Arts"	"Budgets"	"Children"	"Education"
NEW	MILLION	CHILDREN	SCHOOL
FILM	TAX	WOMEN	STUDENTS
SHOW	PROGRAM	PEOPLE	SCHOOLS
MUSIC	BUDGET	CHILD	EDUCATION
MOVIE	BILLION	YEARS	TEACHERS
PLAY	FEDERAL	FAMILIES	HIGH
MUSICAL	YEAR	WORK	PUBLIC
BEST	SPENDING	PARENTS	TEACHER
ACTOR	NEW	SAYS	BENNETT
FIRST	STATE	FAMILY	MANIGAT
YORK	PLAN	WELFARE	NAMPHY
OPERA	MONEY	MEN	STATE
THEATER	PROGRAMS	PERCENT	PRESIDENT
ACTRESS	GOVERNMENT	CARE	ELEMENTARY
LOVE	CONGRESS	LIFE	HAITI

The William Randolph Hearst Foundation will give \$1.25 million to Lincoln Center, Metropolitan Opera Co., New York Philharmonic and Juilliard School. "Our board felt that we had a real opportunity to make a mark on the future of the performing arts with these grants an act every bit as important as our traditional areas of support in health, medical research, education and the social services," Hearst Foundation President Randolph A. Hearst said Monday in announcing the grants. Lincoln Center's share will be \$200,000 for its new building, which will house young artists and provide new public facilities. The Metropolitan Opera Co. and New York Philharmonic will receive \$400,000 each. The Juilliard School, where music and the performing arts are taught, will get \$250,000. The Hearst Foundation, a leading supporter of the Lincoln Center Consolidated Corporate Fund, will make its usual annual \$100,000 donation, too.



- Word, document, topic
  - $w, d, z$
- Word count in document
  - $c(w, d)$
- Word distribution for each topic ( $\beta_z$ )
  - $\beta_{zw}: p(w|z)$
- Topic distribution for each document ( $\theta_d$ )
  - $\theta_{dz}: p(z|d)$  (Yes, soft clustering)

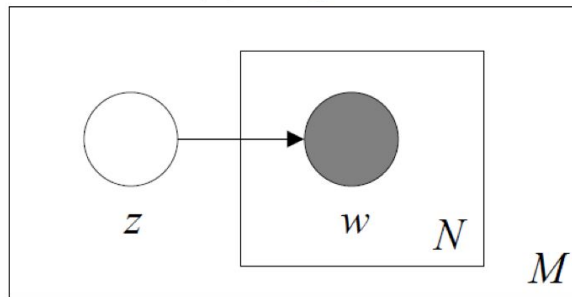


- For each position in  $d$ ,  $n = 1, \dots, N_d$ 
  - Generate the topic for the position as
$$z_n \sim \text{Multinoulli}(\boldsymbol{\theta}_d), \text{ i.e., } p(z_n = k) = \theta_{dk}$$
(Note, 1 trial multinomial, i.e., categorical distribution)
  - Generate the word for the position as
$$w_n \sim \text{Multinoulli}(\boldsymbol{\beta}_{z_n}), \text{ i.e., } p(w_n = w) = \beta_{z_n w}$$

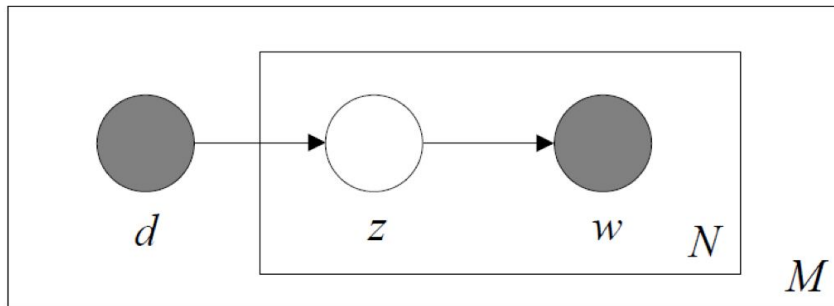
# pLSA: Graphical model



Naïve Bayes  
(Multinomial mixture model)



pLSA



- Probability of a word

$$p(w|d) = \sum_k p(w, z = k|d) = \sum_k p(w|z = k)p(z = k|d) = \sum_k \beta_{kw}\theta_{dk}$$

$$\begin{aligned} \max \log L &= \sum_{dw} c(w, d) \log \sum_z \theta_{dz} \beta_{zw} \\ \text{s. t. } \sum_z \theta_{dz} &= 1 \text{ and } \sum_w \beta_{zw} = 1 \end{aligned}$$

- Repeat until converge
  - E-step: for each word in each document, calculate its conditional probability belonging to each topic
 
$$p(z|w, d) \propto p(w|z, d)p(z|d) = \beta_{zw}\theta_{dz} \text{ (i. e., } p(z|w, d) = \frac{\beta_{zw}\theta_{dz}}{\sum_{z'} \beta_{z'w}\theta_{dz'}})$$
  - M-step: given the conditional distribution, find the parameters that can maximize the expected complete log-likelihood

$$\beta_{zw} \propto \sum_d p(z|w, d)c(w, d) \text{ (i. e., } \beta_{zw} = \frac{\sum_d p(z|w, d)c(w, d)}{\sum_{w', d} p(z|w', d)c(w', d)})$$

$$\theta_{dz} \propto \sum_w p(z|w, d)c(w, d) \text{ (i. e., } \theta_{dz} = \frac{\sum_w p(z|w, d)c(w, d)}{N_d})$$

# pLSA: Example

Two documents d1 and d2  
Two topics z1 and z2

Initialize  $\beta_{wz1} = 1/6$   
Initialize  $\beta_{wz2}$   
=  $1/12$  for  $i \leq 4$   
 $1/3$  for  $i > 4$

For all z and d,  
Initialize  $\theta_{zd} = 1/2$

E step?  
M step?

d1

Index	Word (w)	Count	$P(z1   w, d1) \propto$	$P(z2   w, d1) \propto$
1	frequent	4	?	?
2	pattern	3	?	?
3	data	2	?	?
5	cat	1	?	?

d2

Index	Word (w)	Count	$P(z1   w, d2) \propto$	$P(z2   w, d2) \propto$
3	data	1	?	?
4	mining	2	?	?
5	cat	3	?	?
6	dog	4	?	?

# pLSA: Example

Two documents d1 and d2  
Two topics z1 and z2

Initialize  $\beta_{wz1} = 1/6$   
Initialize  $\beta_{wz2}$   
=  $1/12$  for  $i \leq 4$   
 $1/3$  for  $i > 4$

For all z and d,  
Initialize  $\theta_{zd} = 1/2$

d1

Index	Word (w)	Count	$P(z1   w, d1) \propto$	$P(z2   w, d1) \propto$
1	frequent	4	1/12	1/24
2	pattern	3	1/12	1/24
3	data	2	1/12	1/24
5	cat	1	1/12	1/6

d2

Index	Word (w)	Count	$P(z1   w, d2) \propto$	$P(z2   w, d2) \propto$
3	data	1	1/12	1/24
4	mining	2	1/12	1/24
5	cat	3	1/12	1/6
6	dog	4	1/12	1/6

E step result

# pLSA: Example

Two documents d1 and d2  
Two topics z1 and z2

Initialize  $\beta_{wz1} = 1/6$   
Initialize  $\beta_{wz2}$   
=  $1/12$  for  $i \leq 4$   
 $1/3$  for  $i > 4$

For all z and d,  
Initialize  $\theta_{zd} = 1/2$

M step result:

$\beta_{w1z1} \propto 4/12$        $\beta_{w1z2} \propto 4/24$

$\beta_{w2z1} \propto 3/12$        $\beta_{w2z2} \propto 3/24$

$\beta_{w3z1} \propto 3/12$        $\beta_{w3z2} \propto 3/24$

$\beta_{w4z1} \propto$        $\beta_{w4z2} \propto$

d1

Index	Word (w)	Count	$P(z1   w, d1) \propto$	$P(z2   w, d1) \propto$
1	frequent	4	1/12	1/24
2	pattern	3	1/12	1/24
3	data	2	1/12	1/24
5	cat	1	1/12	1/6

d2

Index	Word (w)	Count	$P(z1   w, d2) \propto$	$P(z2   w, d2) \propto$
3	data	1	1/12	1/24
4	mining	2	1/12	1/24
5	cat	3	1/12	1/6
6	dog	4	1/12	1/6

$\theta_{z1d1} \propto$

10/12

$\theta_{z2d1} \propto$

$\theta_{z1d2} \propto$

10/12

$\theta_{z2d2} \propto$



- Topics after the midterm (80%)
  - Frequent pattern mining
  - Sequential pattern mining
  - Time series data
  - Naïve Bayes
  - Topic modeling
- Topics before the midterm (20%, most likely only T/F)
  - Linear regression
  - Logistic regression
  - SVM
  - Decision Trees
  - Neural Network
  - KNN
  - Clustering

# Continuing Grad-level Courses

---



- CS247: Advanced Data Mining (*coming Spring 2021*)
- CS245: Graph Neural Network (*coming Winter 2021*)



**Samueli**  
Computer Science



# Thank you!

## Q & A

**Reminder:** Students have until **8:00 AM Wednesday, December 16** to log into MyUCLA to complete evaluations for COM SCI 145 section 1C.

---

*This is an intentionally blank page.*