

CISC 372

Text Analytic IIII

EM & LDA



Last lectures

- Expectation Maximization
 - Mode-based clustering
 - Iterative approach for optimization
 - Domain knowledge integration

Today

- EM
 - An example
- LDA & Gibb Sampling

Coin Flipping

H T H H H H H H T H H T H
H H H H H T H H H T H T H
H H H H H H H H H H T H H
H T H H T H H H H H H T H

Binominal Distribution: $B(n, p)$



Coin Flipping using Two coins

H T H H H H H H T H H T H
T T T T H H T T T T H T H
H H H H H T H H H T H T H
T T T T T T T H T H T H T
H T T T T T T T T H T H T
T T T T T T H T T T T T H
H H H H H H H H H H T H H
T H T H T T T T T T T T H
H T H H T H H H H H H T H



Flipping coins

- there are two coins
- And each of them has a different chance of observing head/tail

Binominal Distribution: $B(n_1, p_1)$ $B(n_2, p_2)$

What if we don't know which coin was used in which flip??

```
H T H H H H H H T H H T H
T T T T H H T T T H T H
H H H H H T H H H T H T H
T T T T T T T H T H T H T
H T T T T T T T T H T H T
T T T T T H T T T T T H
H H H H H H H H H T H H
T H T H T T T T T T T H
H T H H T H H H H H H T H
```



Flipping coins

- there are two coins
- And each of them has a different chance of observing head/tail

Binominal Distribution: $B(n_1, p_1)$ $B(n_2, p_2)$

What if we don't know which coin was used in which flip??

H T H H H H H H T H H T H
T T T T H H T T T T H T H
H H H H H T H H H T H T H
T T T T T T T H T H T H T
H T T T T T T T T H T H T
T T T T T T H T T T T T H
H H H H H H H H H H T H H
T H T H T T T T T T T T H
H T H H T H H H H H H T H

$$x = \{x_1, x_2, \dots, x_m\}$$

$$z_i \in \{0,1\}$$

Flipping coins

- there are two coins
- And each of them has a different chance of observing head/tail

Binominal Distribution: $B(n_1, p_1)$ $B(n_2, p_2)$

What if we don't know which coin was used in which flip??

$$x = \{x_1, x_2, \dots, x_m\}$$
$$z_i \in [0,1]$$

```
H T H H H H H H T H H T H
T T T T H H T T T T H T H
H H H H H T H H H T H T H
T T T T T T T H T H T H T
H T T T T T T T T H T H T
T T T T T T H T T T T T H
H H H H H H H H H H T H H
T H T H T T T T T T T T H
H T H H T H H H H H H T H
```

Flipping coins

- there are two coins
- And each of them has a different chance of observing head/tail

EM Algorithm:

1. Randomly assign each flip to a coin (random z)
2. Estimate the parameter for that coin w.r.t. n & p
3. Update the flip assignment Z based on the updated parameters for these two coins
4. Go back to 2, until converge

Binominal Distribution: $B(n_1, p_1)$ $B(n_2, p_2)$

LDA – Document Generation (no labels)

How would you generate (write)
a document??



LDA – Document Generation

1. What it is about??

- Pick a set (mixture of topics)

2. Start writing word-by-word

- Each word is generated one-by-one
- Which one to use?
- Depends on the what is the topic to be covered

LDA – Document Generation

- i denotes the index of a document

$$i = \{1, \dots, N_D\}$$

- v denotes the index of a word

$$v = \{1, \dots, N_w\}$$

- k denotes the index of a topic

$$k = \{1, \dots, N_k\}$$

LDA – Document Generation

- To generate a document, we need a mixture of topics:

$$\pi_i \sim \text{Dir}(\pi_i \mid \alpha) \quad i = \{1, \dots, N_D\}$$

- How about topic? Topic itself is a mixture of words.

$$\mathbf{b}_k \sim \text{Dir}(\mathbf{b}_k \mid \gamma) \quad k = \{1, \dots, N_k\}$$

- Given a document, to generate a word, we need to know what topics to be used

$$z_{iv} \sim \text{Cat}(z_{iv} \mid \pi_i) \quad v = \{1, \dots, N_w\}$$

- Then given the topic for this word under this document, we can generate the actual word (observation)

$$y_{iv} \sim \text{Cat}(y_{iv} \mid z_{iv} = k, \mathbf{B})$$

LDA – Document Generation

- Learning – Gibbs Sampler
 - Bayesian Inference
 - Iterative approach
 - Update the parameters/assignment for each variable at a time
 - A variants of generalized EM algorithm
 - Hyperparameter: number of topics
- Nice implementation:
 - Gensim
 - <https://radimrehurek.com/gensim/models/ldamodel.html>
- Example:
 - https://radimrehurek.com/gensim/auto_examples/tutorials/run_lda.html

LDA – After training

- A document is a distribution over topics:

$$\pi_i \sim \text{Dir}(\pi_i \mid \alpha) \quad i = \{1, \dots, N_D\}$$

- A topic is a distribution over words

$$\mathbf{b}_k \sim \text{Dir}(\mathbf{b}_k \mid \gamma) \quad k = \{1, \dots, N_k\}$$

LDA – Document Generation

- How about topic? Topic itself is a mixture of words.

$$\mathbf{b}_k \sim \text{Dir}(\mathbf{b}_k \mid \gamma) \quad k = \{1, \dots, N_k\}$$

Topic 1		Topic 2		Topic 3	
term	weight	term	weight	term	weight
game	0.014	space	0.021	drive	0.021
team	0.011	nasa	0.006	card	0.015
hockey	0.009	earth	0.006	system	0.013
play	0.008	henry	0.005	scsi	0.012
games	0.007	launch	0.004	hard	0.011