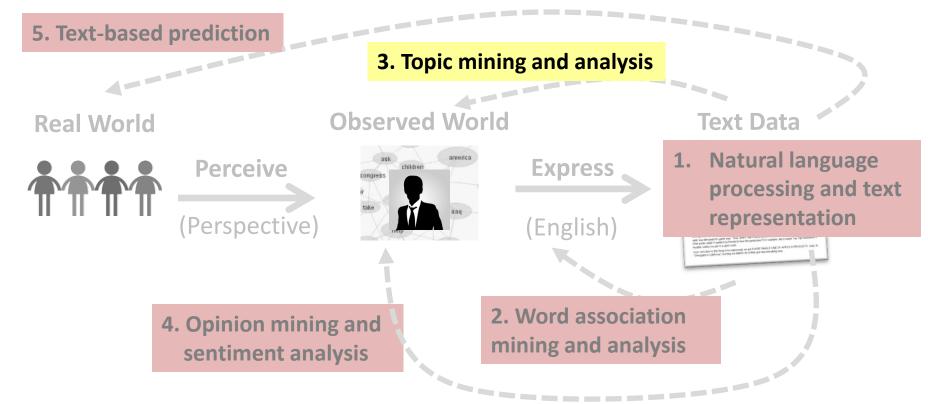# Topic Mining and Analysis: Probabilistic Topic Models

ChengXiang "Cheng" Zhai
Department of Computer Science
University of Illinois at Urbana-Champaign

1

# Topic Mining and Analysis:
# Probabilistic Topic Models

**5. Text-based prediction**

**3. Topic mining and analysis**

Real World

Observed World

Text Data

Perceive

(Perspective)

Express

(English)

1. Natural language processing and text representation

4. Opinion mining and sentiment analysis

2. Word association mining and analysis

# Problems with "Term as Topic"

- Lack of expressive power     ➔ **Topic = {Multiple Words}**
  - Can only represent simple/general topics
  - Can't represent complicated topics
- Incompleteness in vocabulary coverage     **+ weights on words**
  - Can't capture variations of vocabulary (e.g., related words)
- Word sense ambiguity     ➔ **Split an ambiguous word**
  - A topical term or related term can be ambiguous (e.g., basketball star vs. star in the sky)

**A probabilistic topic model can do all these!**

# Improved Idea: Topic = Word Distribution

$\theta_1$ "**Sports**"  $\theta_2$ "**Travel**"  $\bullet\bullet\bullet$  $\theta_k$ "**Science**"

$P(w|\theta_1)$  $P(w|\theta_2)$  $P(w|\theta_k)$

习题·

| | |
|---|---|
| sports | 0.02 |
| game | 0.01 |
| basketball | 0.005 |
| football | 0.004 |
| play | 0.003 |
| star | 0.003 |
| ... | |
| nba | 0.001 |
| ... | |
| travel | 0.0005 |
| ... | |

| | |
|---|---|
| travel | 0.05 |
| attraction | 0.03 |
| trip | 0.01 |
| flight | 0.004 |
| hotel | 0.003 |
| island | 0.003 |
| ... | |
| culture | 0.001 |
| ... | |
| play | 0.0002 |
| ... | |

所有词的
概率都
不为0.

| | |
|---|---|
| science | 0.04 |
| scientist | 0.03 |
| spaceship | 0.006 |
| telescope | 0.004 |
| genomics | 0.004 |
| star | 0.002 |
| ... | |
| genetics | 0.001 |
| ... | |
| travel | 0.00001 |
| ... | |

$$\sum_{w \in V} p(w|\theta_i) = 1$$

**Vocabulary Set: V={w1, w2,....}**

# Probabilistic Topic Mining and Analysis

- Input
  - A **collection** of **N** text documents **C={d$_1$, ..., d$_N$}**
  - **Vocabulary set: V={w$_1$, ..., w$_M$}**
  - **Number of topics**: **k**
- Output
  - **k topics, each a word distribution**: **{ θ$_1$, ..., θ$_k$ }**
  - **Coverage of topics in each d$_i$: { π$_{i1}$, ..., π$_{ik}$ }**
  - π$_{ij}$=prob. of d$_i$ covering topic θ$_j$

$$\sum_{w \in V} p(w \mid \theta_i) = 1$$

$$\sum_{j=1}^{k} \pi_{ij} = 1$$

# The Computation Task

**INPUT: C, k, V**

docs  topics  词汇表
词库

**Text Data**

**OUTPUT: $\{ \theta_1, ..., \theta_k \}, \{ \pi_{i1}, ..., \pi_{ik} \}$**

| | Doc 1 | Doc 2 | • • • | Doc N |
|---|---|---|---|---|

$\theta_1$

sports  0.02
game  0.01
basketball 0.005
football  0.004
...

**30%**  $\pi_{11}$     $\pi_{21}=0\%$     $\pi_{N1}=0\%$

$\theta_2$

travel  0.05
attraction  0.03
trip  0.01
...

**12%**  $\pi_{12}$     $\pi_{22}$     $\pi_{N2}$

• • •

$\theta_k$

science  0.04
scientist  0.03
spaceship 0.006
...

**8%**  $\pi_{1k}$     $\pi_{2k}$     $\pi_{Nk}$

# Generative Model for Text Mining

**Modeling of Data Generation: P(Data |Model, $\Lambda$)**
$\Lambda = (\{ \theta_1, ..., \theta_k \}, \{ \pi_{11}, ..., \pi_{1k} \}, ..., \{ \pi_{N1}, ..., \pi_{Nk} \})$

**How many parameters in total?**

**Parameter Estimation/ Inferences**
$\Lambda^* = \text{argmax}_{\Lambda} \, p(\text{Data}| \text{Model}, \Lambda)$

P(Data |Model, $\Lambda$)

$\Lambda^*$

$\Lambda$

Text Data

# Summary

- Topic represented as word distribution
  – Multiple words: allow for describing a complicated topic
  – Weights on words: model subtle semantic variations of a topic
- Task of topic mining and analysis
  – Input: collection C, number of topics k, vocabulary set V
  – Output: a set of topics, each a word distribution; coverage of all topics in each document

$$\Lambda=(\{\ \theta_1,\ ...,\ \theta_k\ \},\ \{\ \pi_{11},\ ...,\ \pi_{1k}\ \},\ ...,\ \{\ \pi_{N1},\ ...,\ \pi_{Nk}\ \})$$

$$\forall j \in [1,k], \sum_{w \in V} p(w \mid \theta_j) = 1$$

$$\forall i \in [1,N], \sum_{j=1}^{k} \pi_{ij} = 1$$

# Summary (cont.)

- **Generative model** for text mining
  - **Model data generation** with a prob. model: **P(Data |Model, $\Lambda$)**
  - **Infer the most likely parameter values $\Lambda^*$** given a particular data set: **$\Lambda^* = \text{argmax}_\Lambda$ p(Data| Model, $\Lambda$)**
  - **Take $\Lambda^*$ as the "knowledge"** to be mined for the text mining problem
  - **Adjust** the design of the model to discover different knowledge