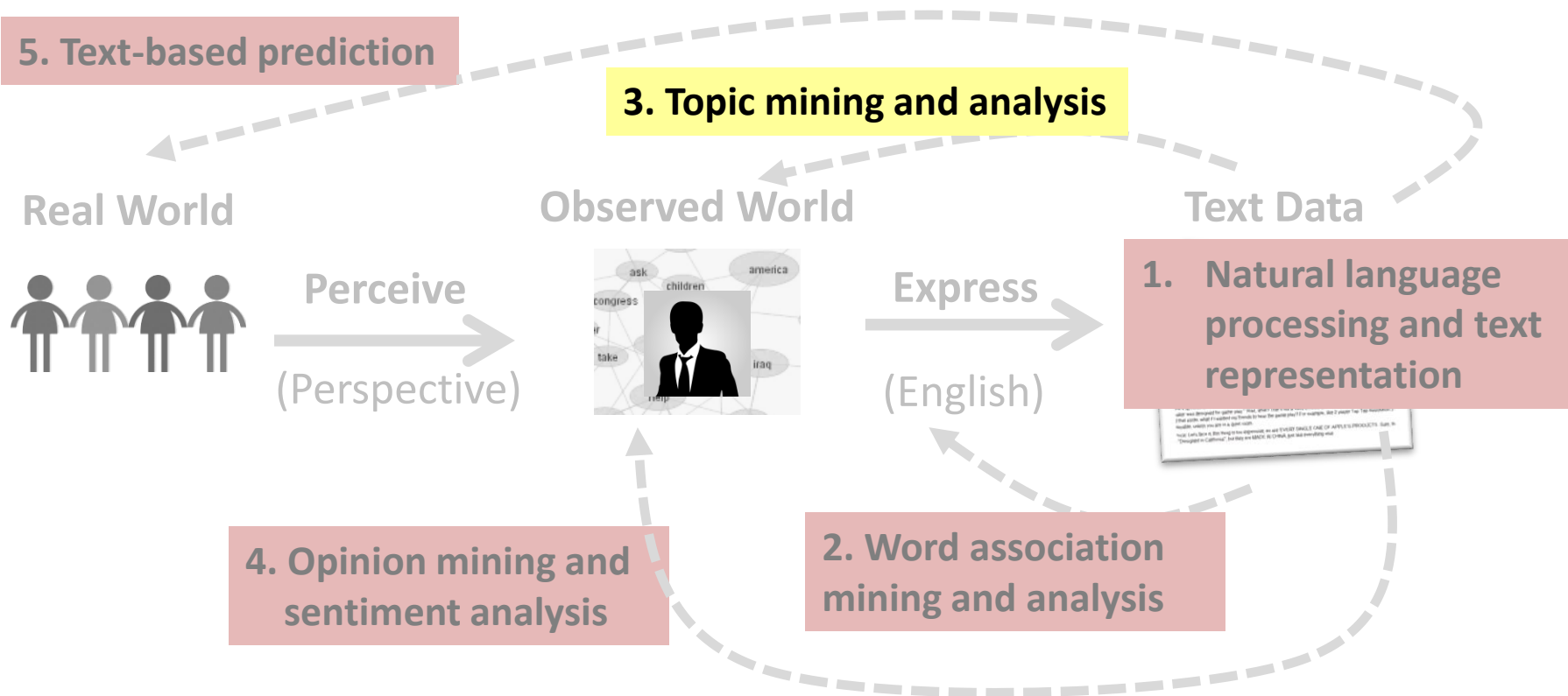# Text Clustering: Generative Probabilistic Models

Part 1

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

# Text Clustering: Generative Probabilistic Models (Part 1)

**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**

# Overview

- What is text clustering?
- Why text clustering?
- How to do text clustering?
  - **Generative probabilistic models** 生成概率模型
  - Similarity-based approaches 相似度为基础的方法.
- How to evaluate clustering results?

# Topic Mining Revisited

**INPUT: C, k, V**

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

**Text Data**

**Doc 1**   **Doc 2**   ● ● ●   **Doc N**

$\theta_1$
- sports  0.02
- game  0.01
- basketball 0.005
- football  0.004
- ...

$\theta_2$
- travel  0.05
- attraction  0.03
- trip  0.01
- ...

● ● ●

$\theta_k$
- science  0.04
- scientist  0.03
- spaceship 0.006
- ...

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

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

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

# One Topic(=cluster) Per Document

**INPUT: C, k, V**

**OUTPUT: $\{ \theta_1, ..., \theta_k \}$, $\{ c_1, ..., c_N \}$ $c_i \in [1,k]$**

**Text Data**

$\theta_1$
sports  0.02
game    0.01
basketball 0.005
football   0.004
...

$\theta_2$
travel  0.05
attraction  0.03
trip      0.01
...

$\theta_k$
science  0.04
scientist   0.03
spaceship 0.006
...

**Doc 1**   **Doc 2**   • • •   **Doc N**

$\pi_{11}=100\%$     $\pi_{21}=0\%$     $\pi_{N1}=100\%$

$\pi_{12}=0$     $\pi_{22}=100\%$     $\pi_{N2}=0$

$\pi_{1k}=0$     $\pi_{1k}=0$     $\pi_{Nk}=0$

5

# Mining One Topic Revisited

OUTPUT: { θ}

**Text Data**

$P(w|\theta)$

Doc d

100%

θ

text ?
mining ?
association ?
database ?

**(1 Doc, 1 Topic)**
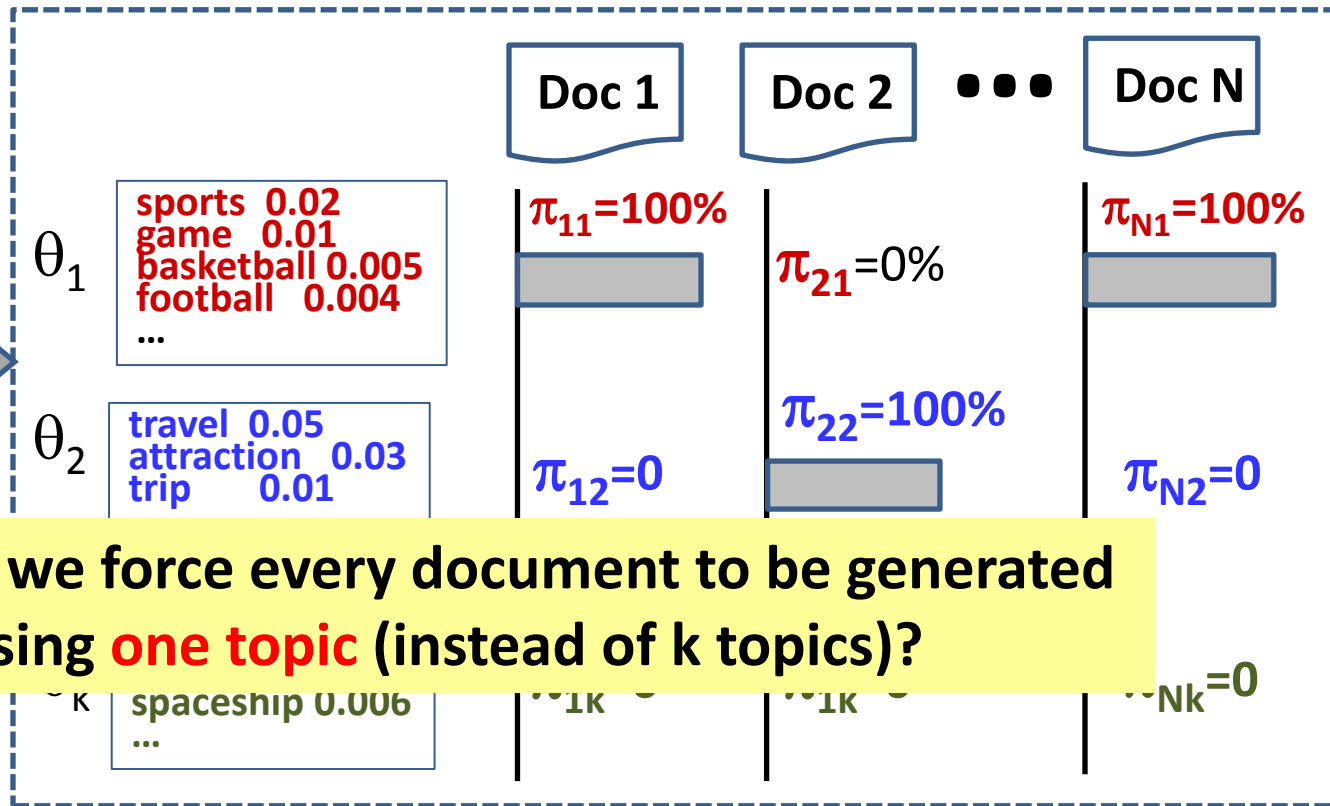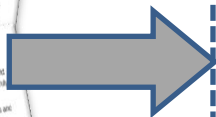➜ **(N Docs, N Topics)**     k<N    也时 docs 会物至
         分享 topics
       ➜ **(N Docs, k Underline{Shared} Topics)=Clustering!**

# What Generative Model Can Do Clustering?

**INPUT: C, k, V**

**OUTPUT: { $\theta_1$, …, $\theta_k$ },**  **{ $c_1$, …, $c_N$ } $c_i \in [1,k]$**

**Text Data**

Doc 1    Doc 2  ● ● ●  Doc N

$\theta_1$
sports  0.02
game   0.01
basketball 0.005
football   0.004
…

$\pi_{11}$=100%      $\pi_{N1}$=100%

$\pi_{21}$=0%

$\theta_2$
travel  0.05
attraction   0.03
trip      0.01

$\pi_{22}$=100%

$\pi_{12}$=0      $\pi_{N2}$=0

**How can we force every document to be generated using one topic (instead of k topics)?**

$\theta_k$  spaceship 0.006
…

$\pi_{1k}$  $\pi_{1k}$  $\pi_{Nk}$=0

# Generative Topic Model Revisited

Why can't this model be used for clustering?

d

word 可以从这个 全部
什或 所以不能间
聚类.

"the"?

"text"?

w

$P(w | \theta_1)$

$p(w | \theta_2)$

xt 0.04
mining 0.035
association 0.03
clustering 0.005
…
the 0.000001
$\theta_1$

the 0.03
a 0.02
is 0.015
we 0.01
food 0.003
…
text 0.000006
$\theta_2$

$p(\theta_1) + p(\theta_2) = 1$

$P(\theta_1) = 0.5$
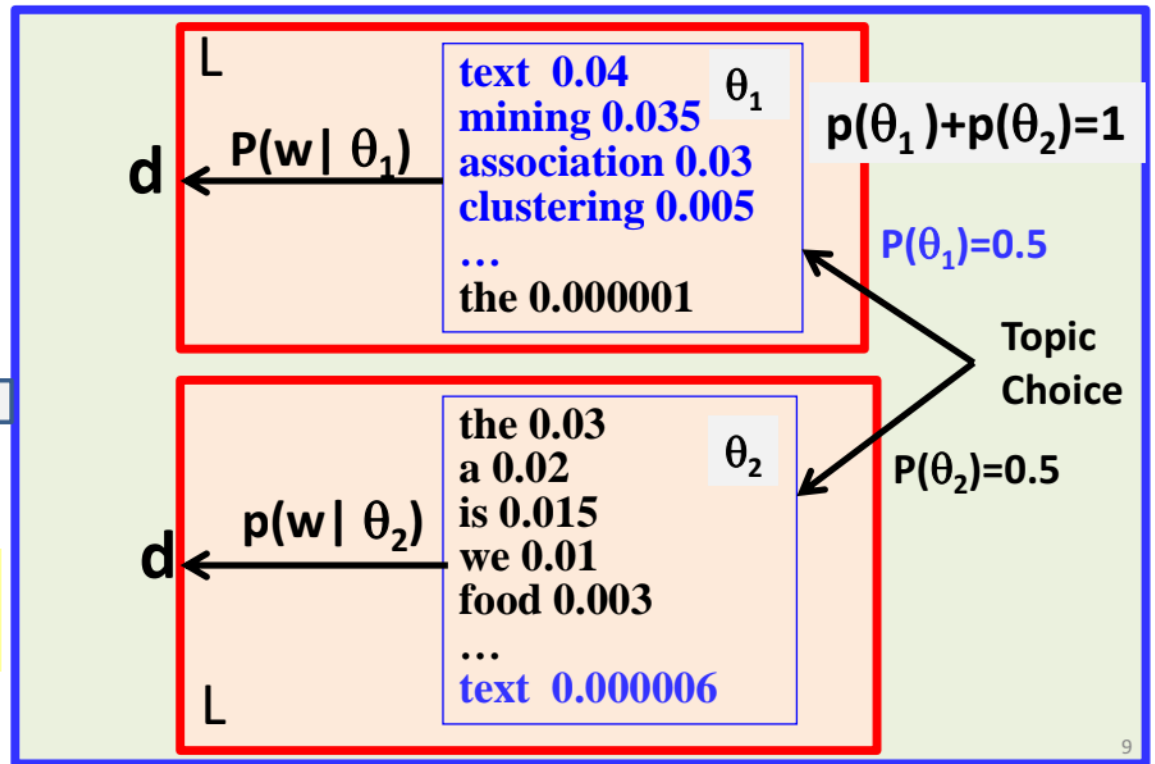
Topic
Choice

$P(\theta_2) = 0.5$

# Mixture Model for Document Clustering

Difference from topic model?

$d = x_1\ x_2 \ldots x_L$

What if $P(\theta_1)=1$ or $P(\theta_2)=1$?

就是一个分布来 式 word.

L

$P(w | \theta_1)$

text  0.04
mining 0.035
association 0.03
clustering 0.005
...
the 0.000001

$\theta_1$

$p(\theta_1)+p(\theta_2)=1$

$P(\theta_1)=0.5$

the 0.03
a 0.02
is 0.015
we 0.01
food 0.003
...
text  0.000006

$\theta_2$

$p(w | \theta_2)$

d

L

Topic Choice

$P(\theta_2)=0.5$

9

$$p(d) = p(\theta_1)p(d \mid \theta_1) + p(\theta_2)p(d \mid \theta_2)$$

$$= p(\theta_1)\prod_{i=1}^{L} p(x_i \mid \theta_1) + p(\theta_2)\prod_{i=1}^{L} p(x_i \mid \theta_2)$$

先选择一个分布，
后一直用个分布生成所有的词

$x_i$ 为某一个词.

$d = x_1\ x_2\ \ldots\ x_L$

**How is this different from a topic model?**

$P(\theta_1)=0.5$

Topic
Choice

the 0.03

food 0.003
…
text 0.000006

L

$$topic\ model: \quad p(d) = \prod_{i=1}^{L}[p(\theta_1)p(x_i \mid \theta_1) + p(\theta_2)p(x_i \mid \theta_2)]$$

每个词是独立生成的
考虑每一个词的时候，
再选择需要的分布.

( + ) × ( + )、

week3 的 quiz 里有提及.