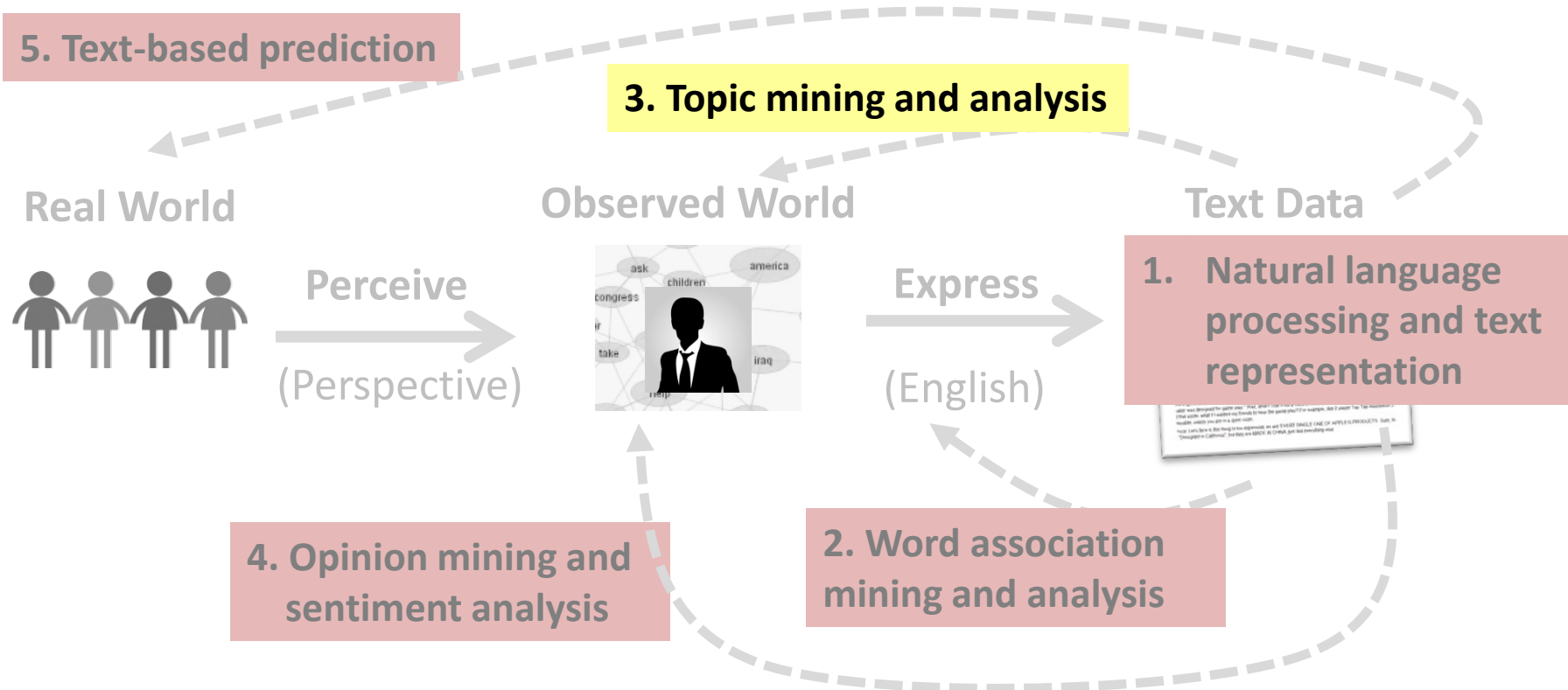# Probabilistic Topic Models: Mixture of Unigram Language Models

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

# Probabilistic Topic Models: Mixture of Unigram LMs

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

# Factoring out Background Words

**p(w| θ)**

d

Text mining paper

the **0.031**
a **0.018**

…
**text 0.04**
**mining 0.035**
**association 0.03**
**clustering 0.005**
**computer 0.0009**
…
food **0.000001**
…

How can we get rid of these common words?

# Generate d Using Two Word Distributions

**d**

Text mining paper

**Topic:** $\theta_d$

$P(w| \theta_d)$

text  0.04
mining 0.035
association 0.03
clustering 0.005
…
the 0.000001

$p(w| \theta_B)$

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

**Background (topic)** $\theta_B$

$p(\theta_d )+(\theta_B)=1$

$P(\theta_d)=0.5$

**Topic Choice**

$P(\theta_B)=0.5$

使用另外的 模型去
封 生成 common words

4

# What's the probability of observing a word w?

后模型下使式指定单词的概率.

选择模型的未联素

**Topic:** $\theta_d$

text 0.04
mining 0.035

$p(\theta_d)+(\theta_B)=1$

d

**P("the")=$p(\theta_d)p("the"|\theta_d) + p(\theta_B)p("the"|\theta_B)$**
**= 0.5\*0.000001+0.5\*0.03**

=0.5

the 0.000001

"the"?

Topic

"text"?

**P("text")=$p(\theta_d)p("text"|\theta_d) + p(\theta_B) p("text"|\theta_B)$** hoice
**= 0.5\*0.04+0.5\*0.000006**

.5

we 0.01
food 0.003
…
text 0.000006
…

**Background (topic) $\theta_B$**

5

# The Idea of a Mixture Model



**Mixture Model**

"the"?

"text"?

w

$P(w \mid \theta_d)$

$p(w \mid \theta_B)$

text  0.04
mining 0.035
association 0.03
clustering 0.005
…
the 0.000001

$\theta_d$

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

$\theta_B$

$p(\theta_d) + (\theta_B) = 1$

$P(\theta_d) = 0.5$

Topic Choice

$P(\theta_B) = 0.5$

# As a Generative Model…

text 0.04
mining 0.035   $\theta_d$
association 0.03
clustering 0.005

$p(\theta_d)+(\theta_B)=1$

**W**

**Formally defines the following generative model:**
$$p(w)=p(\theta_d)p(w|\theta_d) + p(\theta_B)p(w|\theta_B)$$

the 0.05
a 0.02   $\theta_B$

**Estimate of the model "discovers"**
**two topics + topic coverage**

...

**What if $p(\theta_d)=1$ or $p(\theta_B)=1$?** 只选择一种生成方法

# Mixture of Two Unigram Language Models

- **Data**: Document d
- Mixture **Model**: **parameters** $\Lambda=(\{p(w|\theta_d)\}, \{p(w|\theta_B)\}, p(\theta_B), p(\theta_d))$
  - Two unigram LMs: $\theta_d$ **(the topic of d)**; $\theta_B$ **(background topic)**
  - Mixing weight (topic choice): **$p(\theta_d)+p(\theta_B)=1$**
- **Likelihood** function:

$$p(d \mid \Lambda) = \prod_{i=1}^{|d|} p(x_i \mid \Lambda) = \prod_{i=1}^{|d|} [p(\theta_d)p(x_i \mid \theta_d) + p(\theta_B)p(x_i \mid \theta_B)]$$

$$= \prod_{i=1}^{M} [p(\theta_d)p(w_i \mid \theta_d) + p(\theta_B)p(w_i \mid \theta_B)]^{c(w,d)}$$

- **ML Estimate**: $\Lambda^* = \arg\max_\Lambda p(d \mid \Lambda)$

  **Subject to** $\sum_{i=1}^{M} p(w_i \mid \theta_d) = \sum_{i=1}^{M} p(w_i \mid \theta_B) = 1$    $p(\theta_d) + p(\theta_B) = 1$