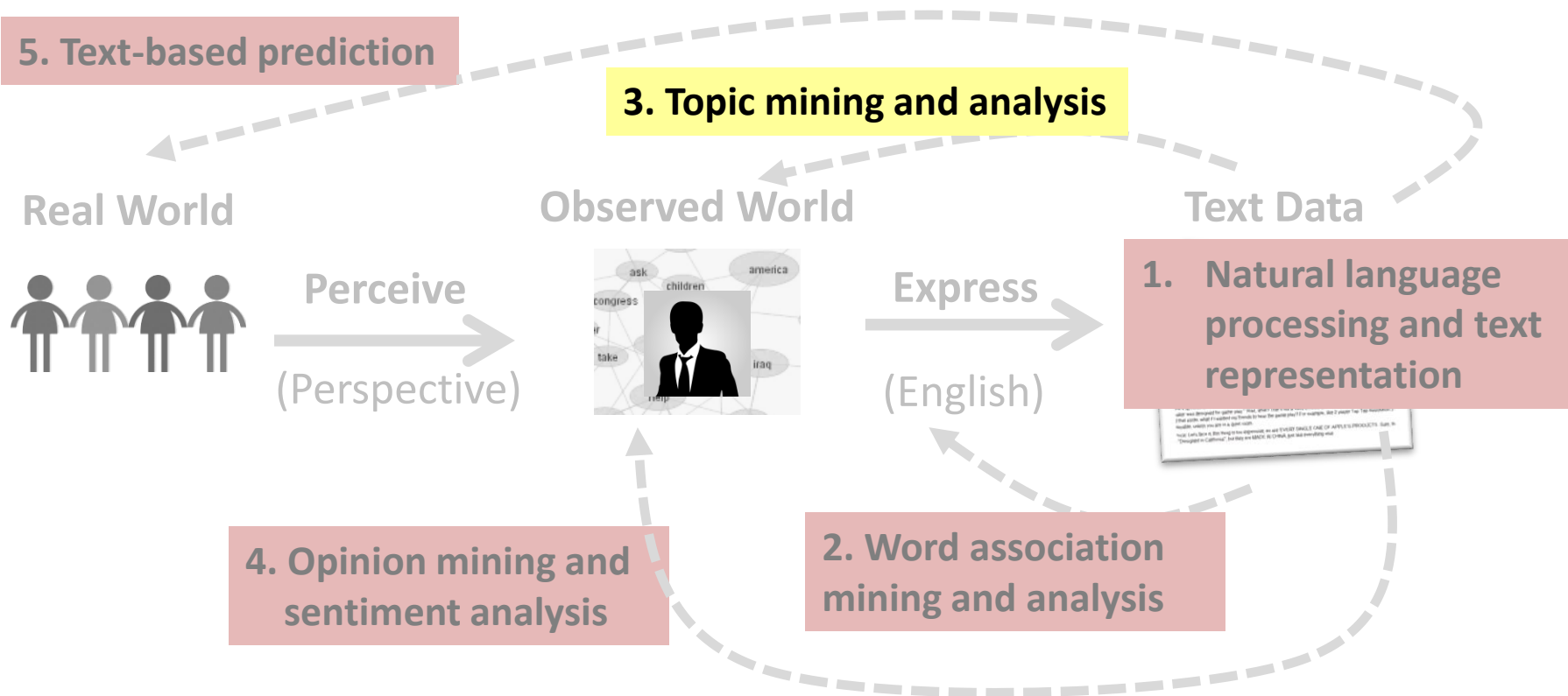# Text Clustering: Generative Probabilistic Models

Part 2

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

# Text Clustering: Generative Probabilistic Models (Part 2)

**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**

# Likelihood Function: p(d)=?

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

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

$d = x_1\ x_2\ \dots\ x_L$

the 0.000001

the 0.03

Topic Choice

**How can we generalize it to include k topics/clusters?**

we 0.01
food 0.003
...
text 0.000006

d

L

# Mixture Model for Document Clustering

- Data: a collection of documents $C=\{d_1, ..., d_N\}$ 在k个分布中选择生成文档的 不晓率

- Model: mixture of k unigram LMs: $\Lambda=(\{\theta_i\}; \{p(\theta_i)\})$, $i \in [1,k]$
  - To generate a document, first **choose a** $\theta_i$ according to $p(\theta_i)$, and then generate all words in the document using $p(w|\theta_i)$ 先选择 $\theta_i$, 再只用 $\theta_i$ 生成

- Likelihood:

$$p(d \mid \Lambda) = \sum_{i=1}^{k} [p(\theta_i) \prod_{j=1}^{|d|} p(x_j \mid \theta_i)]$$

$$= \sum_{i=1}^{k} [p(\theta_i) \prod_{w \in V} p(w \mid \theta_i)^{c(w,d)}]$$

- Maximum Likelihood estimate

用最大ll概率估计

$$\Lambda^* = \arg\max_{\Lambda} p(d \mid \Lambda)$$

# Cluster Allocation After Parameter Estimation

- **Parameters** of the mixture model: $\Lambda=(\{\theta_i\}; \{p(\theta_i)\}), i\in[1,k]$
  - Each $\theta_i$ represents the **content of cluster i** : $p(w|\theta_i)$
  - $p(\theta_i)$ indicates the **size of cluster i** (该类用于生成文档的 所能性、所以 转换成 聚类 的大小)
  - Note that unlike in PLSA, $p(\theta_i)$ doesn't depend on d!
- Which cluster should document d belong to? $c_d$=? $c_d \in [1, k]$
  - **Likelihood only**: Assign d to the cluster corresponding to the topic $\theta_i$ that most likely has been used to generate d
    $$c_d = \arg\max_i p(d|\theta_i)$$
  - **Likelihood + prior $p(\theta_i)$ (Bayesian):** favor large clusters 考虑了聚类的大小..
    $$c_d = \arg\max_i p(d|\theta_i)p(\theta_i)$$
    大聚类 生成 d 的可能性 更大.

5