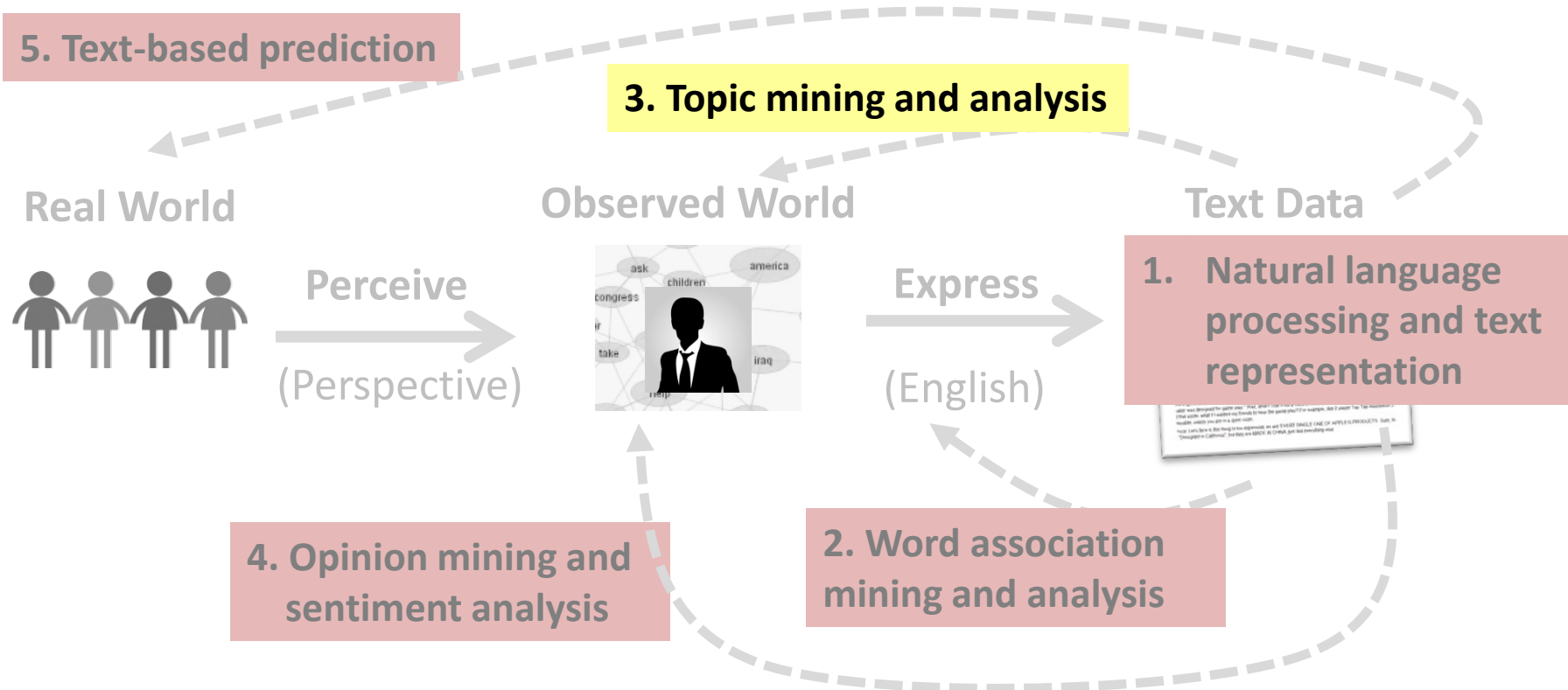# Probabilistic Latent Semantic Analysis (PLSA)

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

# Probabilistic Latent Semantic Analysis (PLSA)



**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**

# Document as a Sample of Mixed Topics

**Topic $\theta_1$**

government 0.3
response 0.2
...

**Topic $\theta_2$**

city 0.2
new 0.1
orleans 0.05
...

. . .

**Topic $\theta_k$**

donate 0.1
relief 0.05
help 0.02
...

**Background $\theta_B$**

the 0.04
a 0.03
...

**Blog article about "Hurricane Katrina"**

[ **Criticism** of **government response** to the **hurricane primarily consisted** of **criticism** of its **response** to the **approach** of the **storm** and its **aftermath, specifically** in the **delayed response** ] to the [ **flooding of New Orleans. … 80%** of the **1.3 million residents** of the **greater New Orleans metropolitan area evacuated** ] …[ **Over seventy countries pledged monetary donations** or other **assistance**]. …

**Many applications are possible if we can "decode" the topics in text…**

# Mining Multiple Topics from Text



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

# Generating Text with Multiple Topics: p(w)=?

$$\sum_{i=1}^{k} \pi_{d,i} = 1$$

$(1-\lambda_B)p(\theta_1)\ p(w|\theta_1)$ —— $p(w|\theta_1)$ —— **Topic $\theta_1$**

government 0.3
response 0.2
...

$p(\theta_1)=\pi_{d,1}$

+

$(1-\lambda_B)p(\theta_2)\ p(w|\theta_2)$ —— $p(w|\theta_2)$ —— **Topic $\theta_2$**

city 0.2
new 0.1
orleans 0.05
...

$p(\theta_2)=\pi_{d,2}$

$\bullet\bullet\bullet$

**W** ...

+

$(1-\lambda_B)p(\theta_k)\ p(w|\theta_k)$ —— $p(w|\theta_k)$ —— **Topic $\theta_k$**

donate 0.1
relief 0.05
help 0.02
...

$p(\theta_k)=\pi_{d,k}$

$1-\lambda_B$

+

$\lambda_B\ p(w|\theta_B)$ —— $p(w|\theta_B)$ —— **Background $\theta_B$**

the 0.04
a 0.03
...

$p(\theta_B)=\lambda_B$

**Topic Choice**

5

# Probabilistic Latent Semantic Analysis (PLSA)

**Percentage of background words (known)** ↘

**Background LM (known)** ↙

↙

**Prob. of word w in topic $\theta_j$** ↙

$$p_d(w) = \lambda_B p(w \mid \theta_B) + (1 - \lambda_B) \sum_{j=1}^{k} \pi_{d,j} p(w \mid \theta_j)$$

$$\log p(d) = \sum_{w \in V} c(w,d) \log[\lambda_B p(w \mid \theta_B) + (1 - \lambda_B) \sum_{j=1}^{k} \pi_{d,j} p(w \mid \theta_j)]$$

$$\log p(C \mid \Lambda) = \sum_{d \in C} \sum_{w \in V} c(w,d) \log[\lambda_B p(w \mid \theta_B) + (1 - \lambda_B) \sum_{j=1}^{k} \pi_{d,j} p(w \mid \theta_j)]$$

**Unknown Parameters: $\Lambda = (\{\pi_{d,j}\}, \{\theta_j\}),\ j=1, \ldots, k$**

**How many unknown parameters are there in total?**

# ML Parameter Estimation

$$p_d(w) = \lambda_B p(w \mid \theta_B) + (1 - \lambda_B) \sum_{j=1}^{k} \pi_{d,j} p(w \mid \theta_j)$$

$$\log p(d) = \sum_{w \in V} c(w, d) \log[\lambda_B p(w \mid \theta_B) + (1 - \lambda_B) \sum_{j=1}^{k} \pi_{d,j} p(w \mid \theta_j)]$$

$$\log p(C \mid \Lambda) = \sum_{d \in C} \sum_{w \in V} c(w, d) \log[\lambda_B p(w \mid \theta_B) + (1 - \lambda_B) \sum_{j=1}^{k} \pi_{d,j} p(w \mid \theta_j)]$$

**Constrained Optimization:** $\quad \Lambda^* = \arg\max_{\Lambda} p(C \mid \Lambda)$ 两个约束体.

$$\forall j \in [1, k], \sum_{i=1}^{M} p(w_i \mid \theta_j) = 1 \qquad \forall d \in C, \sum_{j=1}^{k} \pi_{d,j} = 1$$

# EM Algorithm for PLSA: E-Step

**Hidden Variable (=topic indicator):** $z_{d,w} \in \{B, 1, 2, ..., k\}$

Probability that **w in doc d** is generated from **topic** $\theta_j$

**Use of Bayes Rule**

$$p(z_{d,w} = j) = \frac{\pi_{d,j}^{(n)} p^{(n)}(w \,|\, \theta_j)}{\sum_{j'=1}^{k} \pi_{d,j'}^{(n)} p^{(n)}(w \,|\, \theta_{j'})}$$

$$p(z_{d,w} = B) = \frac{\lambda_B p(w \,|\, \theta_B)}{\lambda_B p(w \,|\, \theta_B) + (1 - \lambda_B) \sum_{j=1}^{k} \pi_{d,j}^{(n)} p^{(n)}(w \,|\, \theta_j)}$$

Probability that **w in doc d** is generated from **background** $\theta_B$

# EM Algorithm for PLSA: M-Step

Hidden Variable (=topic indicator): $z_{d,w} \in \{B, 1, 2, \ldots, k\}$

ML Estimate based on "allocated" word counts to topic $\theta_j$

Re-estimated **probability** of **doc d** covering **topic** $\theta_j$

$$\pi_{d,j}^{(n+1)} = \frac{\sum_{w \in V} c(w,d)(1 - p(z_{d,w} = B))p(z_{d,w} = j)}{\sum_{j'} \sum_{w \in V} c(w,d)(1 - p(z_{d,w} = B))p(z_{d,w} = j')}$$

归一化
所有 topics

$$p^{(n+1)}(w \mid \theta_j) = \frac{\sum_{d \in C} c(w,d)(1 - p(z_{d,w} = B))p(z_{d,w} = j)}{\sum_{w' \in V} \sum_{d \in C} c(w',d)(1 - p(z_{d,w'} = B))p(z_{d,w'} = j)}$$

归一化所有单词.

Re-estimated **probability** of **word w** for **topic** $\theta_j$

# Computation of the EM Algorithm

*手写: $z_{d,w} = j$、$B$ 表示*
*出现中的～单词 是不相同*
*$j$ 或 $B$ (j 是一位题、$B$ 也)*

- Initialize all unknown parameters randomly
- Repeat until likelihood converges
  - E-step $\quad p(z_{d,w} = j) \propto \pi_{d,j}^{(n)} p^{(n)}(w \mid \theta_j)$ $\qquad \sum_{j=1}^{k} p(z_{d,w} = j) = 1$

    $\qquad\qquad p(z_{d,w} = B) \propto \lambda_B p(w \mid \theta_B) \longleftarrow$ What's the normalizer for this one?
  - M-step

$$\pi_{d,j}^{(n+1)} \propto \sum_{w \in V} c(w,d)(1 - p(z_{d,w} = B)) p(z_{d,w} = j) \qquad \forall d \in C, \sum_{j=1}^{k} \pi_{d,j} = 1$$

$$p^{(n+1)}(w \mid \theta_j) \propto \sum_{d \in C} c(w,d)(1 - p(z_{d,w} = B)) p(z_{d,w} = j) \qquad \forall j \in [1,k], \sum_{w \in V} p(w \mid \theta_j) = 1$$

**In general, accumulate counts, and then normalize**

# Summary

- PLSA = mixture model with k unigram LMs (k topics)
- Adding a pre-determined background LM helps discover discriminative topics
- ML estimate "discovers" topical knowledge from text data
  - k word distributions (k topics)
  - proportion of each topic in each document
- The output can enable many applications!
  - Clustering of terms and docs (treat each topic as a cluster)
  - Further associate topics with different contexts (e.g., time periods, locations, authors, sources, etc.)