

Summary of Word Association Mining

- Two basic associations: paradigmatic and syntagmatic
 - Generally applicable to any items in any language (e.g., phrases or entities as units)
- Pure statistical approaches are available for discovering both (can be combined to perform joint analysis).
 - Generally applicable to any text with no human effort
 - Different ways to define “context” and “segment” lead to interesting variations of applications
- Discovered associations can support many other applications.

Recommended Reading

- Chris Manning and Hinrich Schütze, *Foundations of Statistical Natural Language Processing*, MIT Press. Cambridge, MA: May 1999. (Chapter 5 on collocations)
- Chengxiang Zhai, Exploiting context to identify lexical atoms: A statistical view of linguistic context. Proceedings of the International and Interdisciplinary Conference on Modelling and Using Context (CONTEXT-97), Rio de Janeiro, Brazil, Feb. 4-6, 1997. pp. 119-129.
- Shan Jiang and ChengXiang Zhai, Random walks on adjacency graphs for mining lexical relations from big text data. Proceedings of IEEE BigData Conference 2014, pp. 549-554.

Topic Mining and Analysis: Motivation and Task Definition

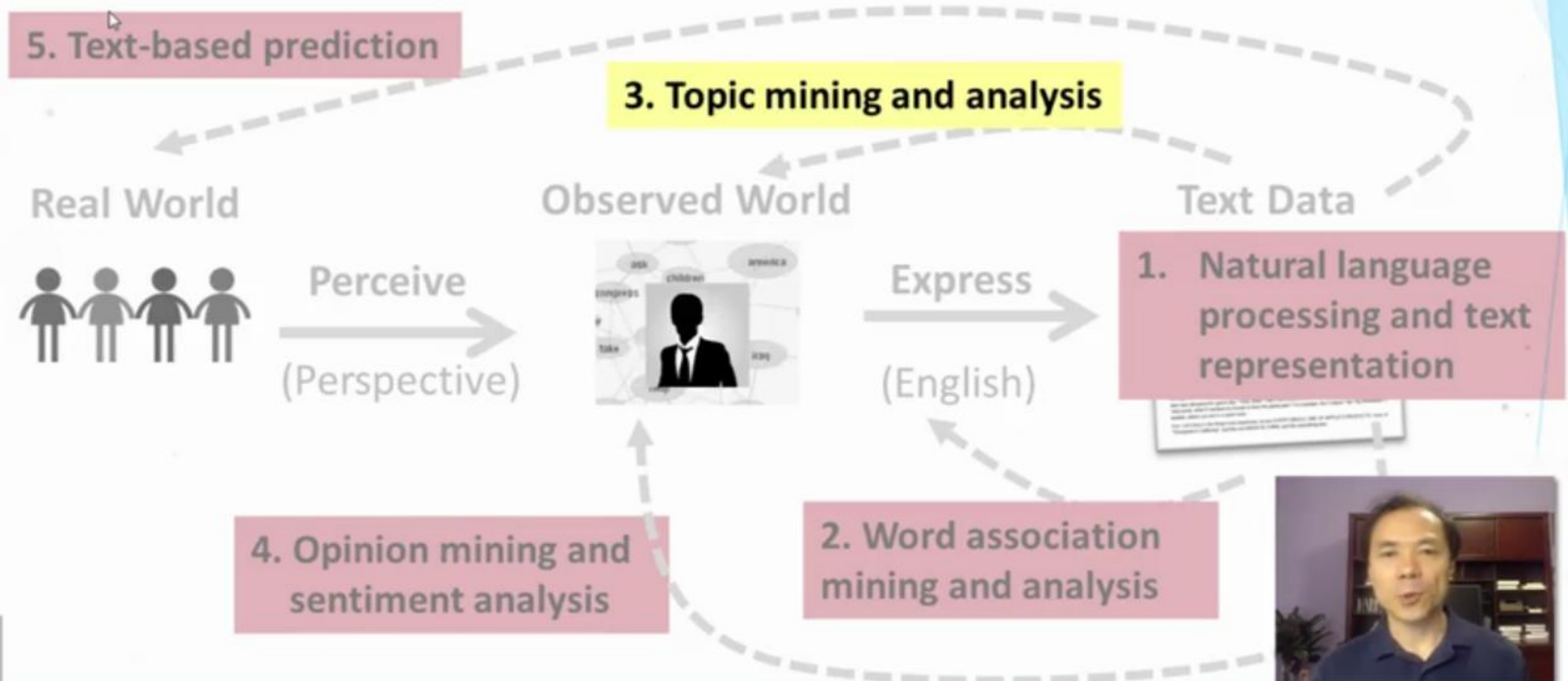
ChengXiang “Cheng” Zhai
Department of Computer Science
University of Illinois at Urbana-Champaign



0:09 / 7:36



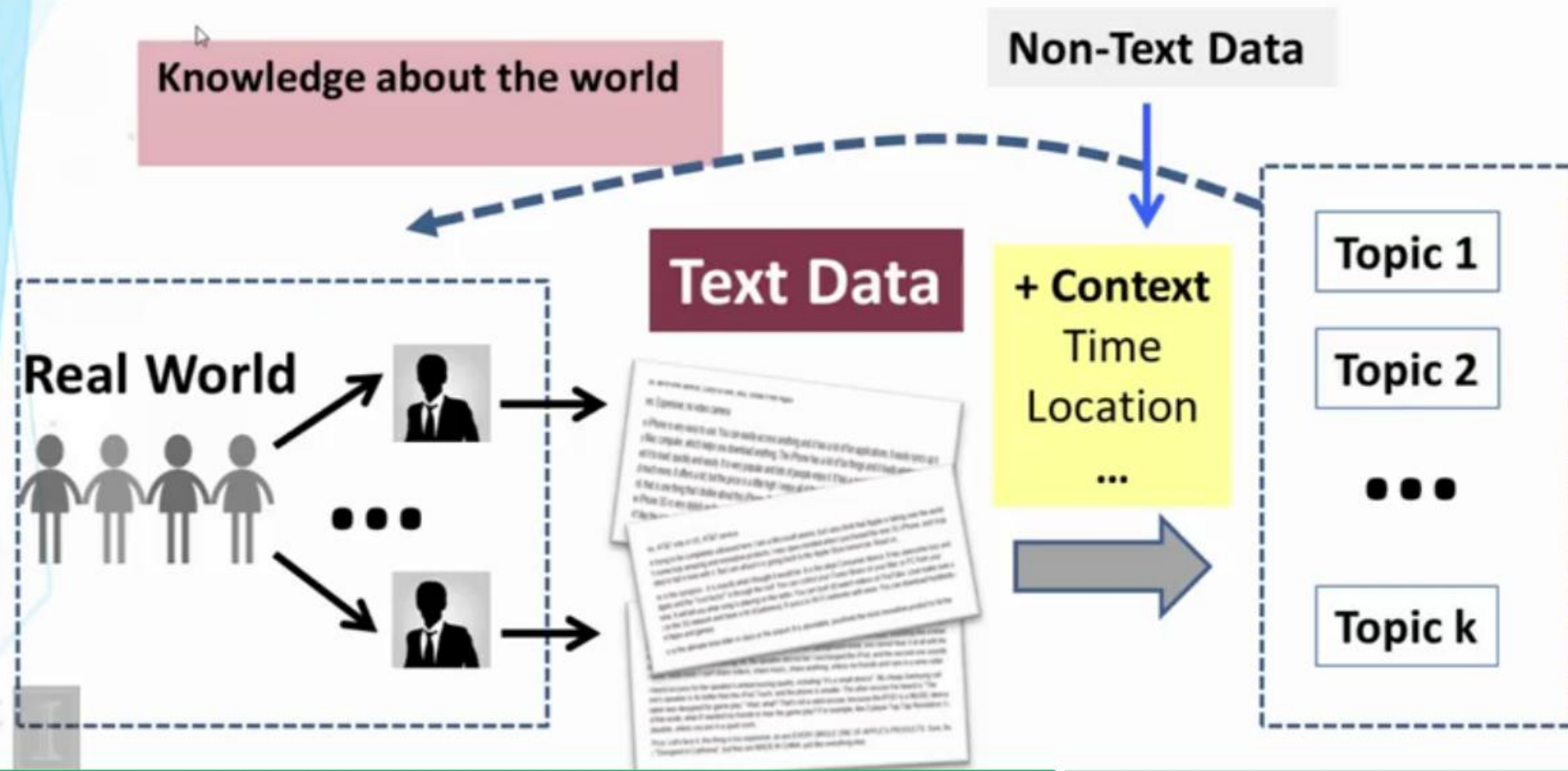
Topic Mining and Analysis: Motivation and Task Definition



Topic Mining and Analysis: Motivation

- Topic \approx main idea discussed in text data
 - Theme/subject of a discussion or conversation
 - Different granularities (e.g., topic of a sentence, an article, etc.)
- Many applications require discovery of topics in text
 - What are Twitter users talking about today?
 - What are the current research topics in data mining? How are they different from those 5 years ago?
 - What do people like about the iPhone 6? What do they dislike?
 - What were the major topics debated in 2012 presidential election?

Topics As Knowledge About the World



Formal Definition of Topic Mining and Analysis

- Input
 - A **collection** of **N** text documents **$C = \{d_1, \dots, d_N\}$**
 - **Number of topics: k**
- Output
 - **k topics: $\{\theta_1, \dots, \theta_k\}$**
 - **Coverage of topics in each d_i : $\{\pi_{i1}, \dots, \pi_{ik}\}$**
 - π_{ij} = prob. of d_i covering topic θ_j

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

How to define θ_i ?

Formal Definition of Topic Mining and Analysis

- Input
 - A **collection** of **N** text documents **$C=\{d_1, \dots, d_N\}$**
 - **Number of topics: k**
- Output
 - **k topics: $\{\theta_1, \dots, \theta_k\}$**
 - **Coverage of topics in each d_i : $\{\pi_{i1}, \dots, \pi_{ik}\}$**
 - π_{ij} = prob. of d_i covering topic θ_j

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

How to define θ_i ?

Topic Mining and Analysis: Term as Topic

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



Formal Definition of Topic Mining and Analysis

- Input
 - A **collection** of **N** text documents **$C=\{d_1, \dots, d_N\}$**
 - **Number of topics: k**
- Output
 - **k topics: $\{\theta_1, \dots, \theta_k\}$**
 - **Coverage of topics in each d_i : $\{\pi_{i1}, \dots, \pi_{ik}\}$**
 - π_{ij} =prob. of d_i covering topic θ_j

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

How to define θ_i ?



00:27 / 11:31



Initial Idea: Topic = Term

Text Data

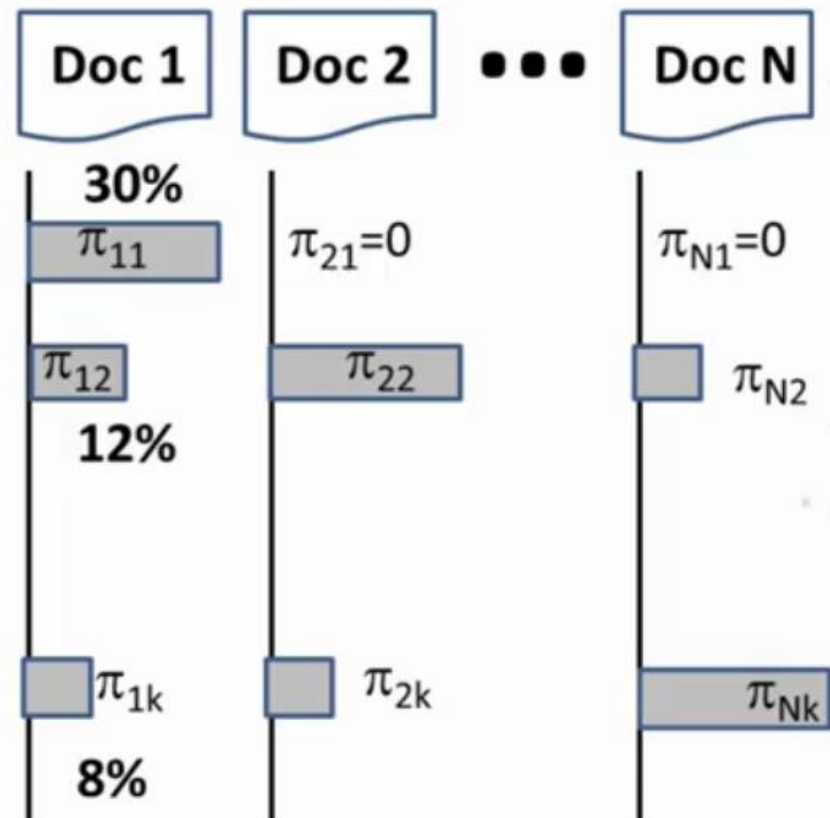


θ_1 "Sports"

θ_2 "Travel"

...

θ_k "Science"



Mining k Topical Terms from Collection C

- Parse text in C to obtain candidate terms (e.g., term = word).



02:27 / 11:31



Mining k Topical Terms from Collection C

- Parse text in C to obtain candidate terms (e.g., term = word).
- Design a scoring function to measure how good each term is as a topic.
 - Favor a representative term (high frequency is favored)
 - Avoid words that are too frequent (e.g., “the”, “a”).
 - TF-IDF weighting from retrieval can be very useful.
 - Domain-specific heuristics are possible (e.g., favor title words, hashtags in tweets).



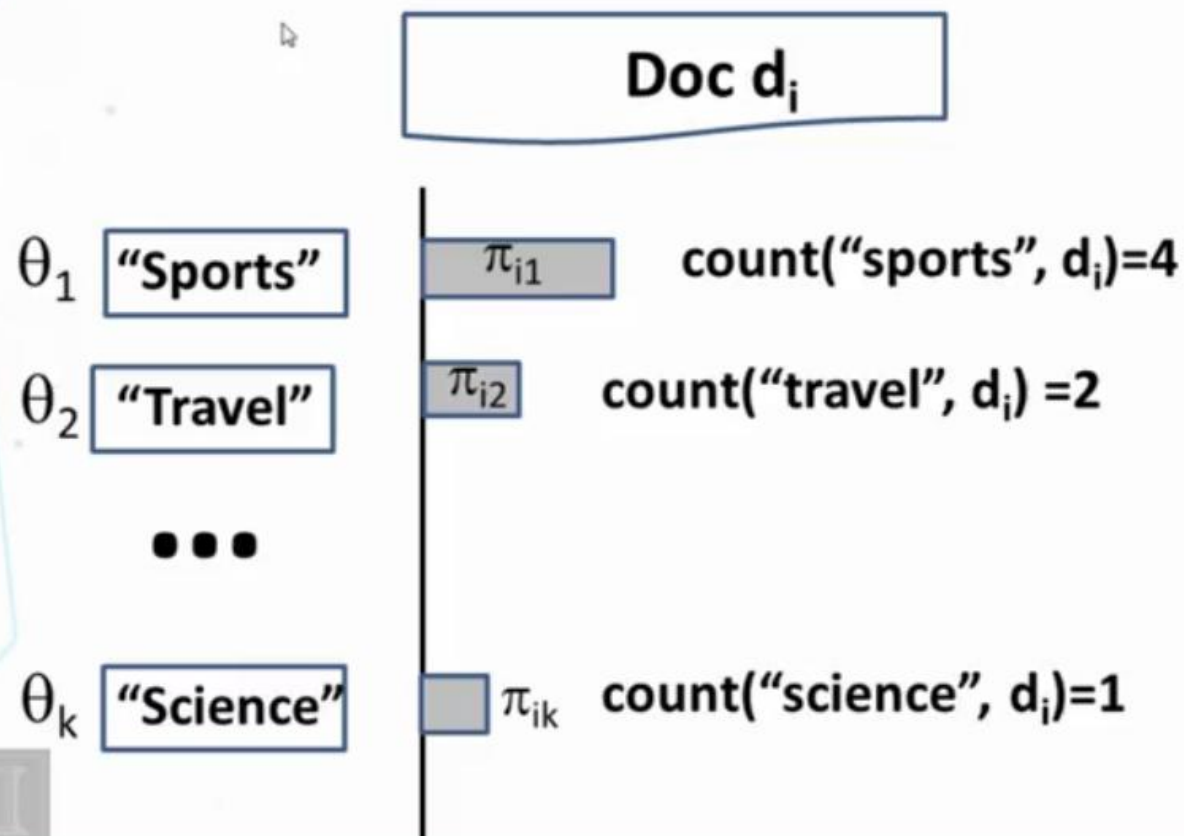
02:53 / 11:31



Mining k Topical Terms from Collection C

- Parse text in C to obtain candidate terms (e.g., term = word).
- Design a scoring function to measure how good each term is as a topic.
 - Favor a representative term (high frequency is favored)
 - Avoid words that are too frequent (e.g., “the”, “a”).
 - TF-IDF weighting from retrieval can be very useful.
 - Domain-specific heuristics are possible (e.g., favor title words, hashtags in tweets).
- Pick k terms with the highest scores but try to minimize redundancy.
 - If multiple terms are very similar or closely related, pick only one of them and ignore others.

Computing Topic Coverage: π_{ij}



$$\pi_{ij} = \frac{\text{count}(\theta_j, d_i)}{\sum_{L=1}^k \text{count}(\theta_L, d_i)}$$

How Well Does This Approach Work?

Doc d_i

Cavaliers vs. Golden State Warriors: NBA playoff finals ...
basketball game ... **travel** to Cleveland ... **star** ...

θ_1 "Sports"

$$\pi_{i1} \propto c(\text{"sports"}, d_i) = 0$$

1. Need to count
related words also!

θ_2 "Travel"

$$\pi_{i2} \propto c(\text{"travel"}, d_i) = 1 > 0$$

...

2. "Star" can be ambiguous (e.g., star in the sky).

θ_k "Science"

$$\pi_{ik} \propto c(\text{"science"}, d_i) = 0$$

3. Mine complicated topics?

Problems with “Term as Topic”

- Lack of expressive power
 - Can only represent simple/general topics
 - Can't represent complicated topics
- Incompleteness in vocabulary coverage
 - Can't capture variations of vocabulary (e.g., related words)
- Word sense ambiguity
 - A topical term or related term can be ambiguous (e.g., basketball star vs. star in the sky)



Topic Mining and Analysis: Probabilistic Topic Models

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



00:07 / 14:17



Topic Mining and Analysis: Probabilistic Topic Models

5. Text-based prediction

3. Topic mining and analysis

Real World



Perceive
(Perspective)

Observed World



Express
(English)

Text Data

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

θ_1 "Sports"

$P(w|\theta_1)$

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

θ_2 "Travel"

$P(w|\theta_2)$

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

...

θ_k "Science"

$P(w|\theta_k)$

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: $V = \{w_1, w_2, \dots\}$

Probabilistic Topic Mining and Analysis

- Input
 - A **collection** of **N** text documents **$C=\{d_1, \dots, d_N\}$**
 - **Vocabulary set**: **$V=\{w_1, \dots, w_M\}$**
 - **Number of topics**: **k**
- Output
 - **k topics**, each a word distribution: **$\{\theta_1, \dots, \theta_k\}$**
 - **Coverage of topics in each d_i** : **$\{\pi_{i1}, \dots, \pi_{ik}\}$**
 - π_{ij} =prob. of d_i covering topic θ_j

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

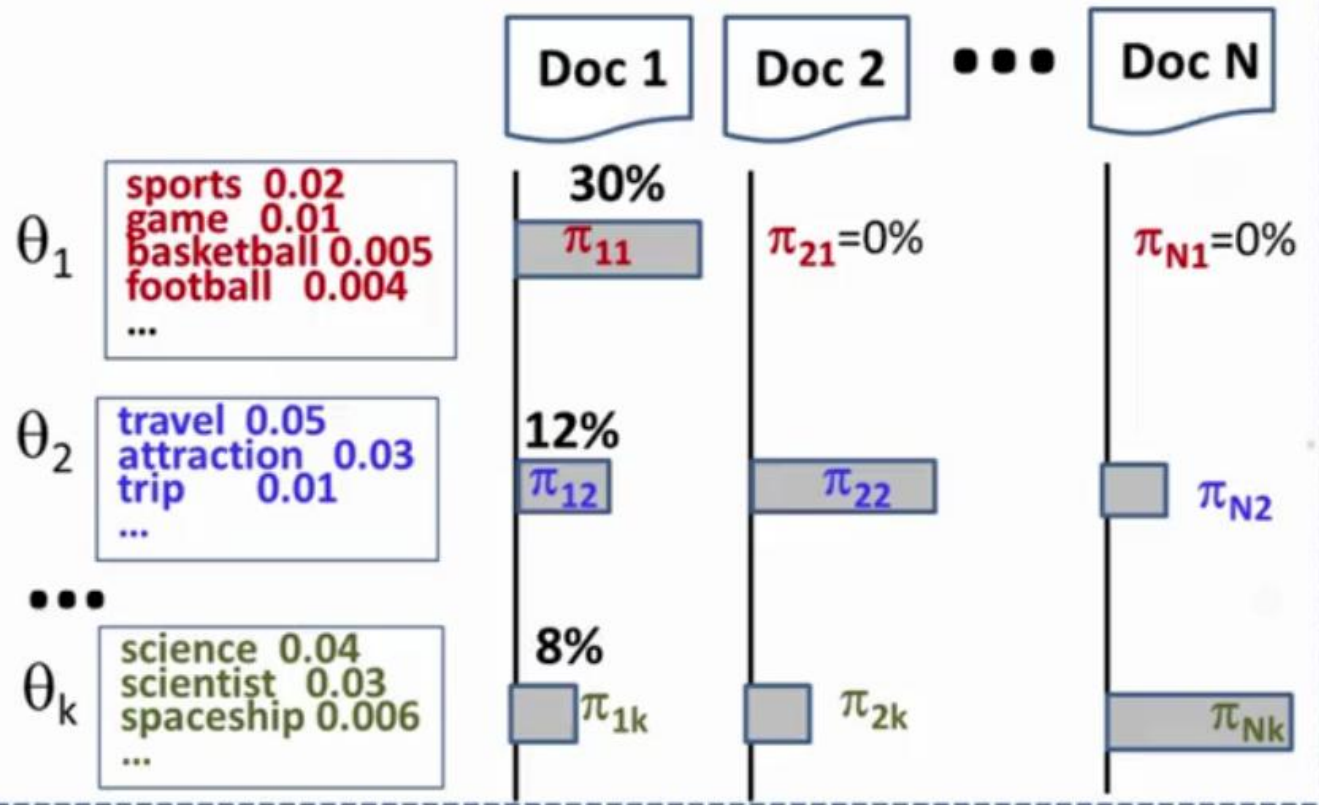
$$\sum_{j=1}^k \pi_{ij} = 1$$

The Computation Task

INPUT: C, k, V

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

Text Data



Generative Model for Text Mining

Modeling of Data Generation: $P(\text{Data} \mid \text{Model}, \Lambda)$

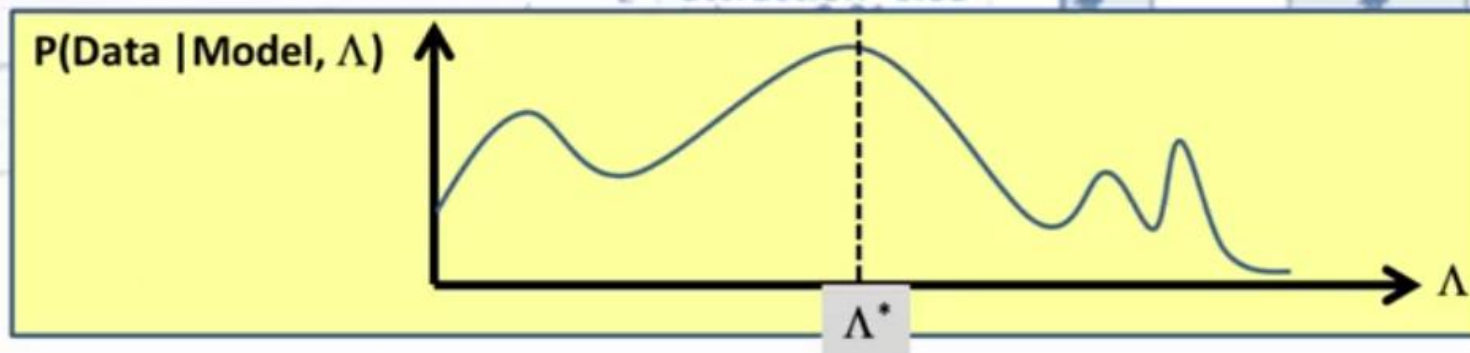
$$\Lambda = (\{\theta_1, \dots, \theta_k\}, \{\pi_{11}, \dots, \pi_{1k}\}, \dots, \{\pi_{N1}, \dots, \pi_{Nk}\})$$

Text Data

How many parameters in total?

Parameter Estimation/ Inferences

$$\Lambda^* = \operatorname{argmax}_{\Lambda} p(\text{Data} \mid \text{Model}, \Lambda)$$



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

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

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

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

Summary (cont.)

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

Topic Mining and Analysis: Overview of Statistical Language Models

Part 1

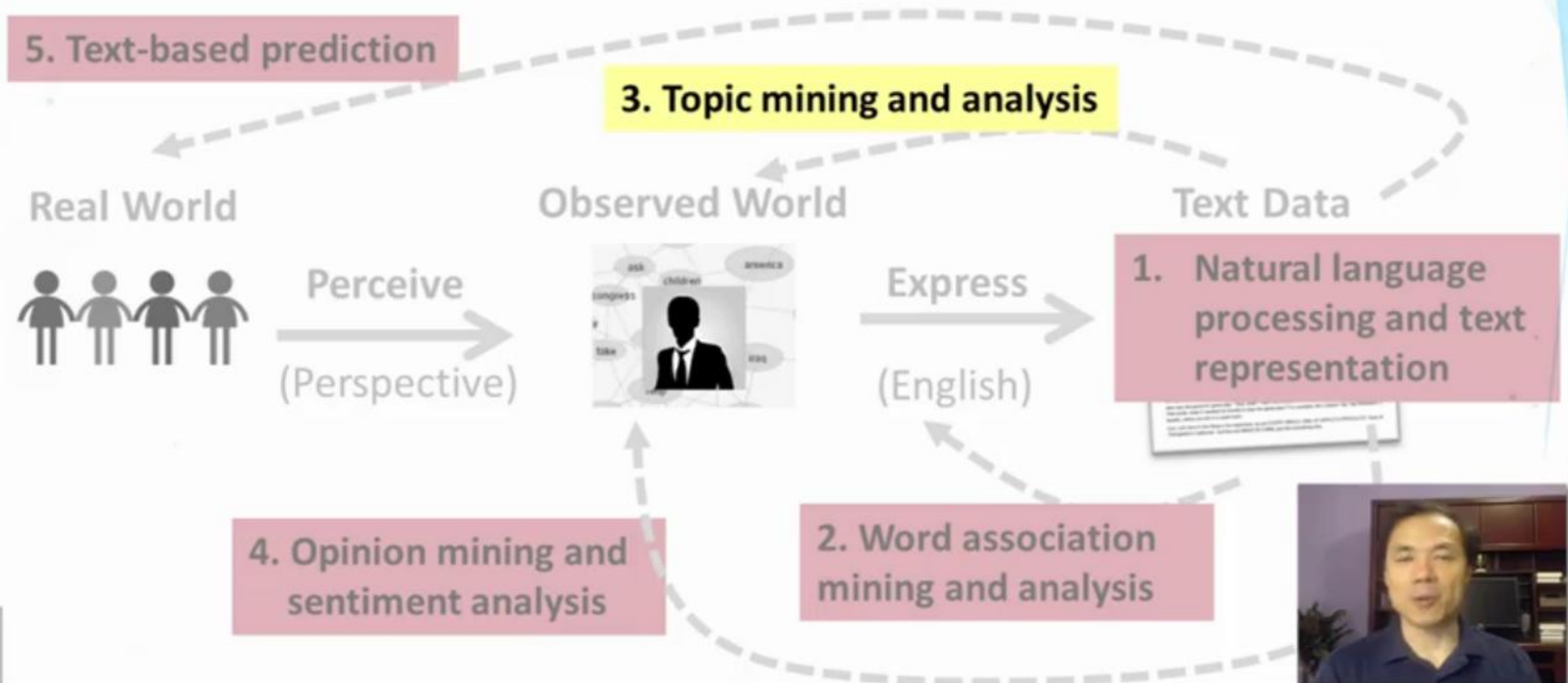
ChengXiang “Cheng” Zhai
Department of Computer Science
University of Illinois at Urbana-Champaign



00:14 / 10:25



Probabilistic Topic Models: Overview of Statistical Language Models



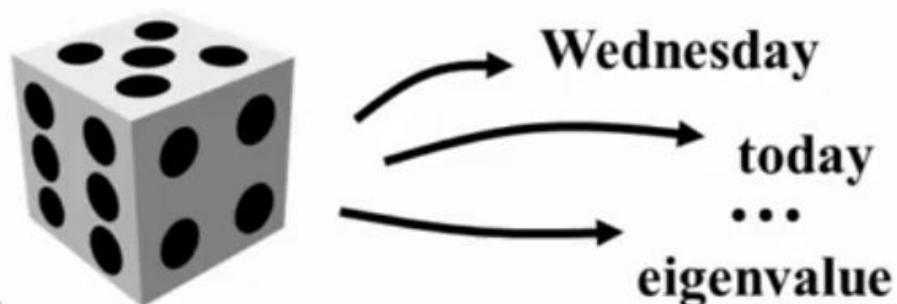
What Is a Statistical Language Model (LM)?

- A probability distribution over word sequences
 - $p(\text{"*Today is Wednesday*"}) \approx 0.001$
 - $p(\text{"*Today Wednesday is*"}) \approx 0.0000000000000001$
 - $p(\text{"*The eigenvalue is positive*"}) \approx 0.00001$
- Context-dependent!
- Can also be regarded as a probabilistic mechanism for "generating" text – thus also called a "generative" model



The Simplest Language Model: Unigram LM

- Generate text by generating each word INDEPENDENTLY
- Thus, $p(w_1 w_2 \dots w_n) = p(w_1)p(w_2)\dots p(w_n)$
- Parameters: $\{p(w_i)\}$ $p(w_1) + \dots + p(w_N) = 1$ (N is voc. size)
- Text = sample drawn according to this **word distribution**



$$\begin{aligned} p(\text{"today is Wed"}) \\ &= p(\text{"today"})p(\text{"is"})p(\text{"Wed"}) \\ &= 0.0002 \times 0.001 \times 0.000015 \end{aligned}$$

04:48 / 10:25



04:48 / 10:25



Text Generation with Unigram LM

Unigram LM $p(w|\theta)$

Sampling

Document d
 $p(d|\theta)=?$

Topic 1:
Text mining

...
text 0.2
mining 0.1
association 0.01
clustering 0.02
...
food 0.00001
...



**Text mining
paper**

Topic 2:
Health

...
food 0.25
nutrition 0.1
healthy 0.05
diet 0.02
...



**Food nutrition
paper**

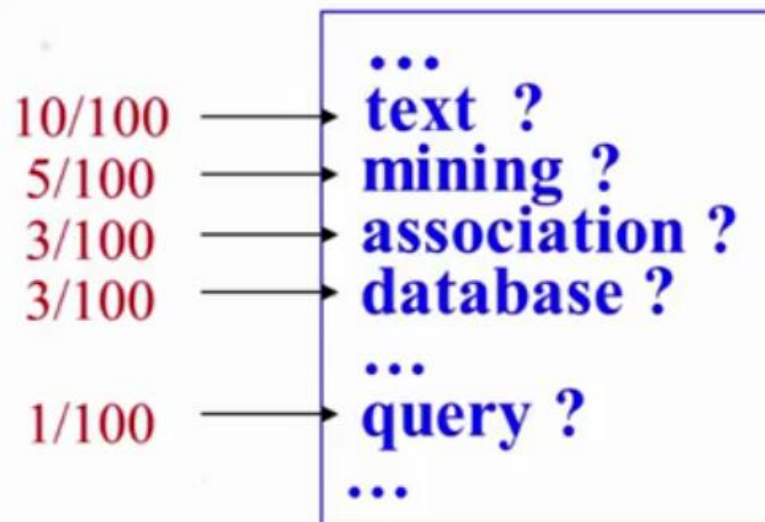
Estimation of Unigram LM

Unigram LM $p(w|\theta)=?$

Estimation

Text Mining Paper d

Total #words=100



Is this our best estimate?
How do we define “best”?

Topic Mining and Analysis: Overview of Statistical Language Models

Part 2

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



Maximum Likelihood vs. Bayesian

- Maximum likelihood estimation

- “Best” means “data likelihood reaches maximum”

$$\hat{\theta} = \arg \max_{\theta} P(X | \theta)$$

- Problem: Small sample

- Bayesian estimation:

Bayes Rule

$$p(X|Y) = \frac{p(Y|X)p(X)}{p(Y)}$$

- “Best” means being consistent with our “prior” knowledge and explaining data well

$$\hat{\theta} = \arg \max_{\theta} P(\theta | X) = \arg \max_{\theta} P(X | \theta) P(\theta)$$

- Problem: How to define prior?



Maximum a Posteriori (MAP) estimate

Illustration of Bayesian Estimation

Bayesian inference: $f(\theta)=?$

$$\hat{f}(\theta) = \sum_{\theta} f(\theta) p(\theta | X)$$

Posterior Mean

$$\hat{\theta} = \sum_{\theta} \theta^* p(\theta | X)$$

Posterior:

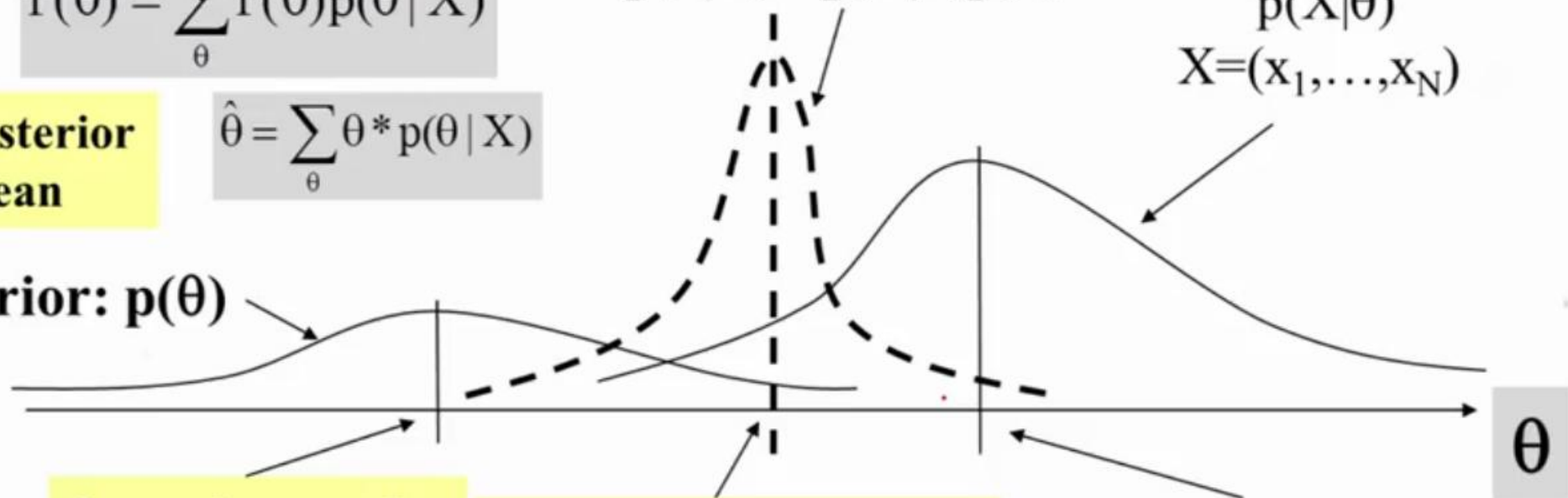
$$p(\theta|X) \propto p(X|\theta)p(\theta)$$

Likelihood:

$$p(X|\theta)$$

$X=(x_1, \dots, x_N)$

Prior: $p(\theta)$



θ_0 : prior mode

θ_1 : posterior mode

θ_{ml} : ML estimate

Summary

- **Language Model** = probability distribution over text = generative model for text data
- **Unigram Language Model** = **word distribution**
- **Likelihood function: $p(X|\theta)$**
 - **Given $\theta \rightarrow$** which X has a higher likelihood?
 - **Given $X \rightarrow$** which θ maximizes $p(X|\theta)$? [**ML estimate**]
- **Bayesian estimation/inference**
 - Must define a **prior: $p(\theta)$**
 - **Posterior distribution: $p(\theta|X) \propto p(X|\theta)p(\theta)$**
 - \rightarrow Allows for inferring any “derived value” from θ !**



Topic Mining and Analysis: Mining One Topic

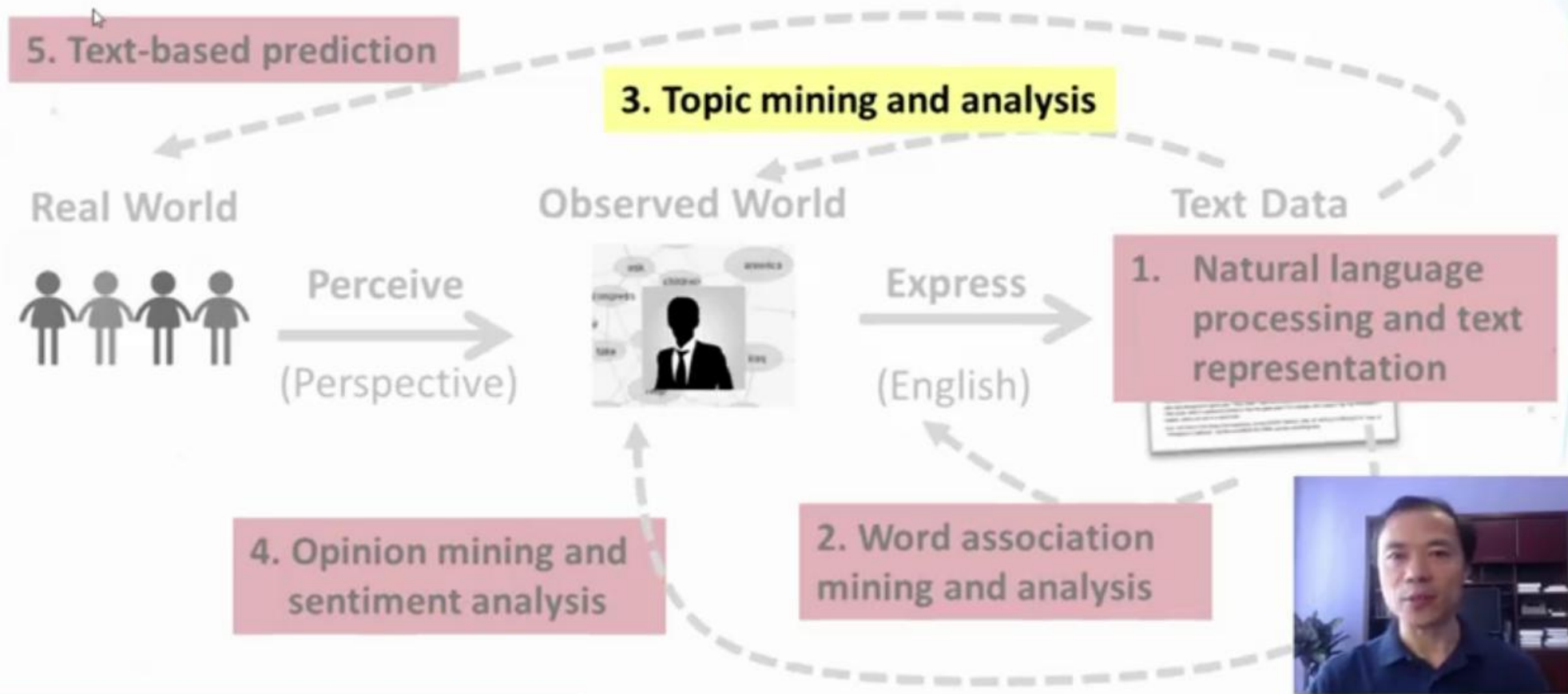
ChengXiang “Cheng” Zhai
Department of Computer Science
University of Illinois at Urbana-Champaign



00:09 / 12:21



Probabilistic Topic Models: Mining One Topic



Simplest Case of Topic Model: Mining One Topic

INPUT: $C=\{d\}, V$

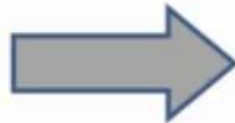
OUTPUT: $\{\theta\}$

Text Data

Speaker quality is ABSOLUTELY CRUCIAL. The speaker is simply not functional, unless you are in perfect recording mode. When you turn it up all the way, because you can't hear it, it becomes fuzzy, sounding like a radio voice. What is the thing, the size of a fist? And there is a RFP background noise, you can't hear it at all with the gain up for some extra recording, so the speaker did not let it exchange the PFD, and the second one sounds some really bad. I can't hear what, then music, then nothing, unless my hands are in a new order.

Headphones to the speaker's ear, including "to a small device" by the way. Setting up your speaker is a lot better than the PFD Touch, and the phone is smaller. The other reason for having "The speaker was designed for the game play." That, what? That's not a valid reason, because the PFD is a RISC device, that's why, what if we had to find a way to hear the game play? For example, the 2 player "Up" Revolution is available, unless you are in a quiet room.

Price is still low if the thing is too expensive, as an EVERYPHASE ONE OF APPLE'S PRODUCTS. See the "Designed in California" button on WAVE IN CHINA, and the recording quality.



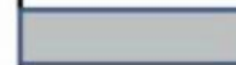
$P(w|\theta)$

θ

text ?
mining ?
association ?
database ?
...
query ?
...

Doc d

100%



Language Model Setup

- **Data:** Document $d = x_1 x_2 \dots x_{|d|}$, $x_i \in V = \{w_1, \dots, w_M\}$ is a word
- **Model:** Unigram LM $\theta (= \text{topic}) : \{\theta_i = p(w_i | \theta)\}$, $i = 1, \dots, M$;
 $\theta_1 + \dots + \theta_M = 1$
- **Likelihood function:** $p(d | \theta) = p(x_1 | \theta) \times \dots \times p(x_{|d|} | \theta)$

$$\begin{aligned} &= p(w_1 | \theta)^{c(w_1, d)} \times \dots \times p(w_M | \theta)^{c(w_M, d)} \\ &= \prod_{i=1}^M p(w_i | \theta)^{c(w_i, d)} = \prod_{i=1}^M \theta_i^{c(w_i, d)} \end{aligned}$$

- **ML estimate:** $(\hat{\theta}_1, \dots, \hat{\theta}_M) = \arg \max_{\theta_1, \dots, \theta_M} p(d | \theta) = \arg \max_{\theta_1, \dots, \theta_M} \prod_{i=1}^M \theta_i^{c(w_i, d)}$

Computation of Maximum Likelihood Estimate

Maximize $p(d | \theta)$ $(\hat{\theta}_1, \dots, \hat{\theta}_M) = \arg \max_{\theta_1, \dots, \theta_M} p(d | \theta) = \arg \max_{\theta_1, \dots, \theta_M} \prod_{i=1}^M \theta_i^{c(w_i, d)}$

Max. Log-Likelihood $(\hat{\theta}_1, \dots, \hat{\theta}_M) = \arg \max_{\theta_1, \dots, \theta_M} \log[p(d | \theta)] = \arg \max_{\theta_1, \dots, \theta_M} \sum_{i=1}^M c(w_i, d) \log \theta_i$

Subject to constraint: $\sum_{i=1}^M \theta_i = 1$

Use Lagrange multiplier approach

Lagrange function: $f(\theta | d) = \sum_{i=1}^M c(w_i, d) \log \theta_i + \lambda (\sum_{i=1}^M \theta_i - 1)$

$$\frac{\partial f(\theta | d)}{\partial \theta_i} = \frac{c(w_i, d)}{\theta_i} + \lambda = 0 \rightarrow \theta_i = -\frac{c(w_i, d)}{\lambda}$$

$$\sum_{i=1}^M -\frac{c(w_i, d)}{\lambda} = 1 \rightarrow \lambda = -\sum_{i=1}^M c(w_i, d) \rightarrow \hat{\theta}_i = p(w_i | \hat{\theta}) = \frac{c(w_i, d)}{\sum_{i=1}^M c(w_i, d)} = \frac{c(w_i, d)}{|d|}$$

Computation of Maximum Likelihood Estimate

Maximize $p(d | \theta)$ $(\hat{\theta}_1, \dots, \hat{\theta}_M) = \arg \max_{\theta_1, \dots, \theta_M} p(d | \theta) = \arg \max_{\theta_1, \dots, \theta_M} \prod_{i=1}^M \theta_i^{c(w_i, d)}$

Max. Log-Likelihood $(\hat{\theta}_1, \dots, \hat{\theta}_M) = \arg \max_{\theta_1, \dots, \theta_M} \log[p(d | \theta)] = \arg \max_{\theta_1, \dots, \theta_M} \sum_{i=1}^M c(w_i, d) \log \theta_i$

Subject to constraint: $\sum_{i=1}^M \theta_i = 1$

Use Lagrange multiplier approach

Lagrange function: $f(\theta | d) = \sum_{i=1}^M c(w_i, d) \log \theta_i + \lambda (\sum_{i=1}^M \theta_i - 1)$

$$\frac{\partial f(\theta | d)}{\partial \theta_i} = \frac{c(w_i, d)}{\theta_i} + \lambda = 0 \rightarrow \theta_i = -\frac{c(w_i, d)}{\lambda}$$

$$\sum_{i=1}^M -\frac{c(w_i, d)}{\lambda} = 1 \rightarrow \lambda = -\sum_{i=1}^M c(w_i, d) \rightarrow \hat{\theta}_i = p(w_i | \hat{\theta}) = \frac{c(w_i, d)}{\sum_{i=1}^M c(w_i, d)} = \frac{c(w_i, d)}{|d|}$$

**Normalized
Counts**



What Does the Topic Look Like?

d

Text mining
paper

$p(w | \theta)$

the 0.031

a 0.018

...

text 0.04

mining 0.035

association 0.03

clustering 0.005

computer 0.0009

...

food 0.000001

...

Can we get rid of
these common words?

