

Syntagmatic Relation Discovery: Entropy

Syntagmatic Relation = Correlated Occurrences

Whenever “**eats**” occurs, what **other words** also tend to occur?

My cat **eats** fish on Saturday
His cat **eats** turkey on Tuesday
My dog **eats** meat on Sunday
His dog **eats** turkey on Tuesday
...

My	_____	eats	_____	on Saturday
His	_____	eats	_____	on Tuesday
My	_____	eats	_____	on Sunday
His	_____	eats	_____	on Tuesday
...	_____		_____	

What words tend to occur
to the **left** of “**eats**”?

What words
are to the
right?

Word Prediction: Intuition and Formal Definition

Prediction questions:

Is word **W** present/ absent in a text segment (sentence, paragraph, document)?

Are some words easier to predict than others:

1) $W = \text{"meat"}$ 2) $W = \text{"the"}$ 3) $W = \text{"unicorn"}$

Binary Random Variable : $X_w \in \{0, 1\}$

$$X_w = \begin{cases} 1 & \text{w is present} \\ 0 & \text{w is absent} \end{cases}$$

$$p(X_w = 1) + p(X_w = 0) = 1$$

The more random X_w , the more difficult the prediction.

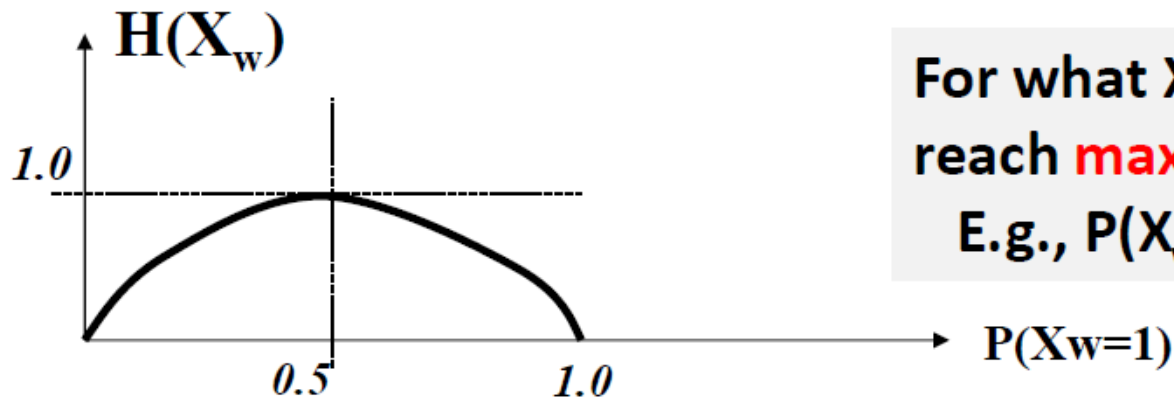
How to quantitatively measure the “randomness” of a random variable like X_w ?

Entropy $H(X)$ Measures Randomness of X

$$H(X_w) = \sum_{v \in \{0,1\}} -p(X_w = v) \log_2 p(X_w = v)$$

$$X_w = \begin{cases} 1 & \text{w is present} \\ 0 & \text{w is absent} \end{cases}$$

$$= -p(X_w = 0) \log_2 p(X_w = 0) - p(X_w = 1) \log_2 p(X_w = 1) \quad \text{Define } 0 \log_2 0 = 0$$



For what X_w , does $H(X_w)$ reach **maximum/minimum**?

E.g., $P(X_w=1)=1$? $P(X_w=1)=0.5$?

or equivalently $P(X_w=0)$ (Why?)

Entropy $H(X)$: Coin Tossing

$$H(X_{\text{coin}}) = -p(X_{\text{coin}} = 0) \log_2 p(X_{\text{coin}} = 0) - p(X_{\text{coin}} = 1) \log_2 p(X_{\text{coin}} = 1)$$

X_{coin} : tossing a coin

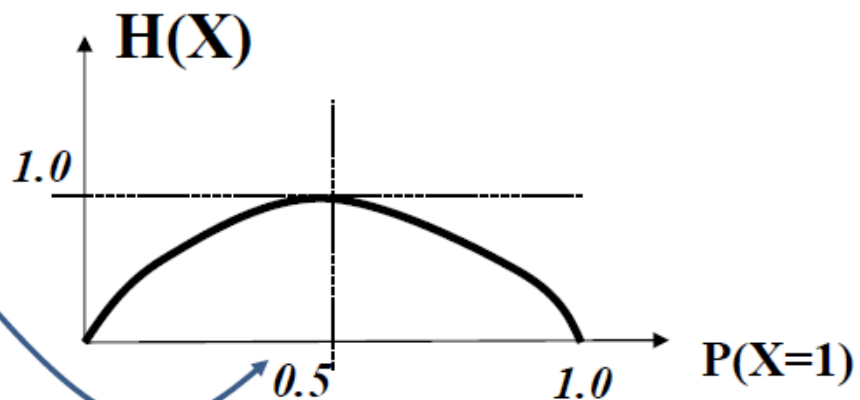
$$X_{\text{coin}} = \begin{cases} 1 & \text{Head} \\ 0 & \text{Tail} \end{cases}$$

Fair coin: $p(X=1)=p(X=0)=1/2$

$$H(X) = -\frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{2} \log_2 \frac{1}{2} = 1$$

Completely biased: $p(X=1)=1$

$$H(X) = -0 * \log_2 0 - 1 * \log_2 1 = 0$$



Entropy for Word Prediction

Is word **W** present (or absent) in this segment?



1) $W = \text{"meat"}$

2) $W = \text{"the"}$

3) $W = \text{"unicorn"}$

Which is **high/low**? $H(X_{\text{meat}})$, $H(X_{\text{the}})$, or $H(X_{\text{unicorn}})$?

$H(X_{\text{the}}) \approx 0 \rightarrow$ no uncertainty since $p(X_{\text{the}}=1) \approx 1$

High entropy words are harder to predict!

Syntagmatic Relation Discovery: Conditional Entropy

What If We Know More About a Text Segment?

Does presence of “**eats**” help predict the presence of “**meat**”?

Does it **reduce** the uncertainty about “meat”, i.e., $H(X_{\text{meat}})$?

What if we know of the absence of “eats”? Does it also help?

Know nothing about the segment

Know “eats” is present ($X_{\text{eats}} = 1$)

$$p(X_{\text{meat}} = 1) \quad \text{-----} \rightarrow \quad p(X_{\text{meat}} = 1 \mid X_{\text{eats}} = 1)$$

$$p(X_{\text{meat}} = 0) \quad \text{-----} \rightarrow \quad p(X_{\text{meat}} = 0 \mid X_{\text{eats}} = 1)$$

$$H(X_{\text{meat}}) = -p(X_{\text{meat}} = 0) \log_2 p(X_{\text{meat}} = 0) - p(X_{\text{meat}} = 1) \log_2 p(X_{\text{meat}} = 1)$$



$$H(X_{\text{meat}} \mid X_{\text{eats}} = 1) = -p(X_{\text{meat}} = 0 \mid X_{\text{eats}} = 1) \log_2 p(X_{\text{meat}} = 0 \mid X_{\text{eats}} = 1) \\ - p(X_{\text{meat}} = 1 \mid X_{\text{eats}} = 1) \log_2 p(X_{\text{meat}} = 1 \mid X_{\text{eats}} = 1)$$

$H(X_{\text{meat}} \mid X_{\text{eats}} = 0)$ can be defined similarly

Conditional Entropy: Complete Definition, Capturing Syntagmatic Relation

$$\begin{aligned} H(X_{meat} | X_{eats}) &= \sum_{u \in \{0,1\}} [p(X_{eats} = u) H(X_{meat} | X_{eats} = u)] \\ &= \sum_{u \in \{0,1\}} [p(X_{eats} = u) \sum_{v \in \{0,1\}} [-p(X_{meat} = v | X_{eats} = u) \log_2 p(X_{meat} = v | X_{eats} = u)]] \end{aligned}$$

In general, for any discrete random variables X and Y , we have $H(X) \geq H(X|Y)$

What's the **minimum** possible value of $H(X|Y)$?

$$H(X_{meat} | X_{eats}) = \sum_{u \in \{0,1\}} [p(X_{eats} = u) H(X_{meat} | X_{eats} = u)]$$

$$H(X_{meat} | X_{meat}) = ?$$

Which is smaller? $H(X_{meat} | X_{the})$ or $H(X_{meat} | X_{eats})$?

For which word w , does $H(X_{meat} | X_w)$ reach its minimum (i.e., 0)?

For which word w , does $H(X_{meat} | X_w)$ reach its maximum, $H(X_{meat})$?

Mining Syntagmatic Relations with Conditional Entropy

- For each word W_1
 - For **every other word** W_2 , compute $H(X_{W_1} | X_{W_2})$
 - Sort all $H(X_{W_1} | X_{W_2})$ in **ascending order**
 - **Top-ranked** candidate words - **potential syntagmatic** relations with W_1
 - Need to use a **threshold** for each W_1
- However, while $H(X_{W_1} | X_{W_2})$ and $H(X_{W_1} | X_{W_3})$ are comparable, $H(X_{W_1} | X_{W_2})$ and $H(X_{W_3} | X_{W_2})$ aren't!

How can we mine the **strongest** K syntagmatic relations from a collection?

Syntagmatic Relation Discovery: Mutual Information $I(X;Y)$

Mutual Information $I(X;Y)$: Entropy Reduction & Sintagm Rel Mining

If we know Y , how much reduction in the entropy of X can we get?

Mutual Information: $I(X; Y) = H(X) - H(X|Y) = H(Y) - H(Y|X)$

Properties:

- **Non-negative**: $I(X;Y) \geq 0$
- **Symmetric**: $I(X;Y) = I(Y;X)$
- $I(X;Y) = 0$ iff X & Y are **independent**

With fixed X and different Y s: same order of $I(X;Y)$ and $H(X|Y)$, but **$I(X;Y)$ allows to compare different (X,Y) pairs.**

Whenever “**eats**” occurs, what **other words** also tend to **occur**?

Which **words** have high mutual information with “**eats**”?

$$I(X_{\text{eats}}; X_{\text{meats}}) = I(X_{\text{meats}}; X_{\text{eats}}) \quad > \quad I(X_{\text{eats}}; X_{\text{the}}) = I(X_{\text{the}}; X_{\text{eats}})$$

$$I(X_{\text{eats}}; X_{\text{eats}}) = H(X_{\text{eats}}) \geq I(X_{\text{eats}}; X_w)$$

Rewriting Mutual Information (MI) Using KL-divergence

The observed joint distribution of X_{w1} and X_{w2}



$$I(X_{w1}; X_{w2}) = \sum_{u \in \{0,1\}} \sum_{v \in \{0,1\}} p(X_{w1} = u, X_{w2} = v) \log_2 \frac{p(X_{w1} = u, X_{w2} = v)}{p(X_{w1} = u)p(X_{w2} = v)}$$



The expected joint distribution of X_{w1} and X_{w2}
if X_{w1} and X_{w2} were independent

MI measures the **divergence** of the **actual joint distribution** from the expected distribution under the independence assumption. The **larger the divergence** is, the **higher the MI** would be.

Mutual Information: Probabilities and Relations Between Them

$$I(X_{w1}; X_{w2}) = \sum_{u \in \{0,1\}} \sum_{v \in \{0,1\}} p(X_{w1} = u, X_{w2} = v) \log_2 \frac{p(X_{w1} = u, X_{w2} = v)}{p(X_{w1} = u)p(X_{w2} = v)}$$

Presence & absence of w1: $p(X_{w1}=1) + p(X_{w1}=0) = 1$

Presence & absence of w2: $p(X_{w2}=1) + p(X_{w2}=0) = 1$

Co-occurrences of w1 and w2:

$$\underline{p(X_{w1}=1, X_{w2}=1)} + \underline{p(X_{w1}=1, X_{w2}=0)} + \underline{p(X_{w1}=0, X_{w2}=1)} + \underline{p(X_{w1}=0, X_{w2}=0)} = 1$$



Both w1 & w2 occur



Only w1 occurs



Only w2 occurs



None of them occurs

Co-occurrences of w1 and w2:

$$p(X_{w1}=1, X_{w2}=1) + p(X_{w1}=1, X_{w2}=0) + p(X_{w1}=0, X_{w2}=1) + p(X_{w1}=0, X_{w2}=0) = 1$$

Constraints:

$$p(X_{w1}=1, X_{w2}=1) + p(X_{w1}=1, X_{w2}=0) = p(X_{w1}=1)$$

$$p(X_{w1}=0, X_{w2}=1) + p(X_{w1}=0, X_{w2}=0) = p(X_{w1}=0)$$

$$p(X_{w1}=1, X_{w2}=1) + p(X_{w1}=0, X_{w2}=1) = p(X_{w2}=1)$$

$$p(X_{w1}=1, X_{w2}=0) + p(X_{w1}=0, X_{w2}=0) = p(X_{w2}=0)$$

Computation of Mutual Information

Presence & absence of w_1 :

$$p(X_{w_1}=1) + p(X_{w_1}=0) = 1$$

Presence & absence of w_2 :

$$p(X_{w_2}=1) + p(X_{w_2}=0) = 1$$

Co-occurrences of w_1 and w_2 :

$$p(X_{w_1}=1, X_{w_2}=1) + p(X_{w_1}=1, X_{w_2}=0) + p(X_{w_1}=0, X_{w_2}=1) + p(X_{w_1}=0, X_{w_2}=0) = 1$$

$$p(X_{w_1}=1, X_{w_2}=1) + p(X_{w_1}=1, X_{w_2}=0) = p(X_{w_1}=1)$$

$$p(X_{w_1}=0, X_{w_2}=1) + p(X_{w_1}=0, X_{w_2}=0) = p(X_{w_1}=0)$$

$$p(X_{w_1}=1, X_{w_2}=1) + p(X_{w_1}=0, X_{w_2}=1) = p(X_{w_2}=1)$$

$$p(X_{w_1}=1, X_{w_2}=0) + p(X_{w_1}=0, X_{w_2}=0) = p(X_{w_2}=0)$$

We only need to know $p(X_{w_1}=1)$, $p(X_{w_2}=1)$, and $p(X_{w_1}=1, X_{w_2}=1)$.

Estimation of Probabilities (Depending on the Data)

$$p(X_{w1} = 1) = \frac{\text{count}(w1)}{N}$$

$$p(X_{w2} = 1) = \frac{\text{count}(w2)}{N}$$

$$p(X_{w1} = 1, X_{w2} = 1) = \frac{\text{count}(w1, w2)}{N}$$

	<u>W1</u>	<u>W2</u>	
Segment_1	1	0	Only W1 occurred
Segment_2	1	1	Both occurred
Segment_3	1	1	Both occurred
Segment_4	0	0	Neither occurred
...			
<u>Segment_N</u>	<u>0</u>	<u>1</u>	Only W2 occurred

Count(w1) = total number segments that contain W1

Count(w2) = total number segments that contain W2

Count(w1, w2) = total number segments that contain both W1 and W2

Smoothing: Accommodating Zero Counts

$$p(X_{w1} = 1) = \frac{\text{count}(w1) + 0.5}{N + 1}$$

$$p(X_{w2} = 1) = \frac{\text{count}(w2) + 0.5}{N + 1}$$

$$p(X_{w1} = 1, X_{w2} = 1) = \frac{\text{count}(w1, w2) + 0.25}{N + 1}$$

	W1	W2
¼ PseudoSeg_1	0	0
¼ PseudoSeg_2	1	0
¼ PseudoSeg_3	0	1
¼ PseudoSeg_4	1	1

Smoothing: Add pseudo data so that
no event has zero counts
(pretend we observed extra data)

Segment_1	1	0
...		
Segment_N	0	1

Actually observed data

Summary of Syntagmatic Relation Discovery

- **Syntagmatic relation** can be discovered by **measuring correlations** between occurrences of two words.
- Three concepts from Information Theory:
 - **Entropy** $H(X)$: measures the uncertainty of a random variable X
 - **Conditional entropy** $H(X|Y)$: entropy of X given we know Y
 - **Mutual information** $I(X;Y)$: entropy reduction of X (or Y) due to knowing Y (or X)
- **Mutual information** provides a principled way for **discovering syntagmatic relations**.

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.

Topic Mining and Analysis: Motivation and Task Definition

Topic Mining and Analysis: Motivation

- Topic \approx **main idea** discussed in text data; **knowledge** about the world
 - 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
 - Summaries from Twitter
 - Current research topics in data mining; their difference from 5 years ago?
 - Likes and dislikes about iPhone 6
 - Major topics debated in 2012 presidential election

Tasks of Topic Mining and Analysis

Task 2: Figure out which documents cover which topics

Text Data



Topic 1

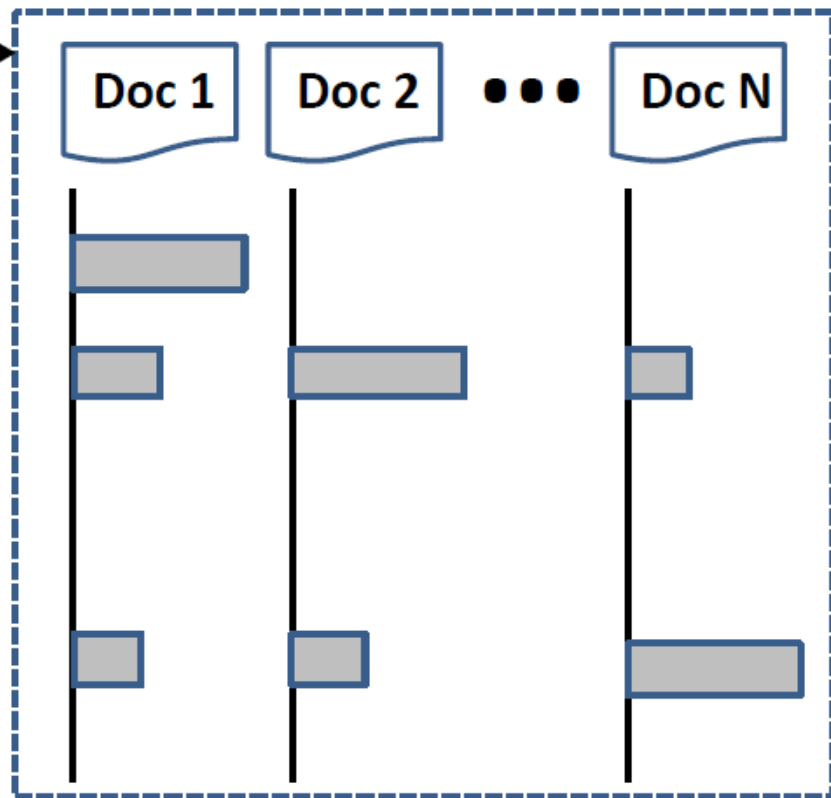
Topic 2

...

Topic k



Task 1: Discover k topics



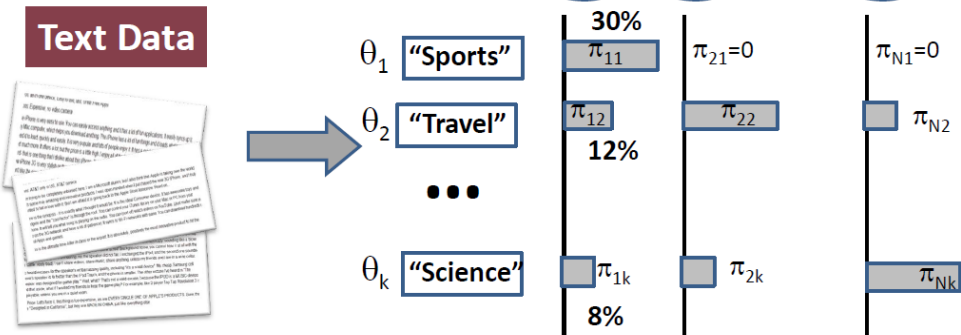
Formal Definition of Topic Mining and Analysis

- 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

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

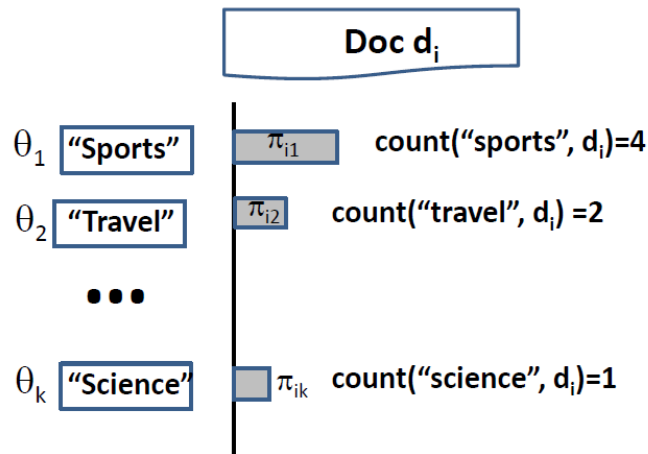
How to define θ_i ?

Topic = Term



Drawbacks:

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



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

Topic Mining and Analysis: Probabilistic Topic Models

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

4

Resolving drawbacks of topic = term:

- Lack of expressive power (no complicated topics) - **Topic = {Multiple Words}**
- Incompleteness in vocabulary coverage (no related words) – introduce **weights on words**
- Word sense ambiguity (star) - **split an ambiguous word** (into different distributions)

A **probabilistic topic model** can do all of these!

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

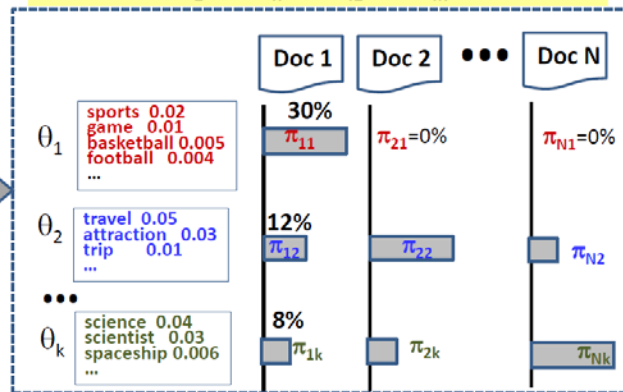
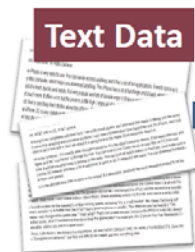
- k topics, each a word distribution: $\{\theta_1, \dots, \theta_k\}$
- Coverage of topics in each d_i : $\{\pi_{i1}, \dots, \pi_{ik}\}$
- π_{ij} =prob. of d_i covering topic θ_j

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

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

INPUT: C, k, V

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



Generative Model for Text Mining

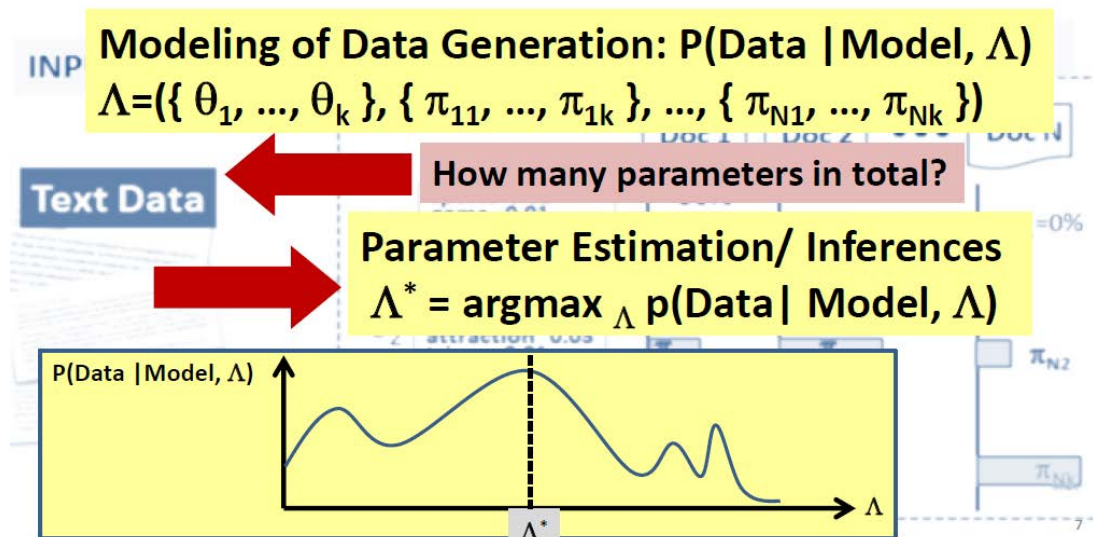
Summary

- **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**

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

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

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



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

Probabilistic Topic Models: Overview of Statistical Language Models

Statistical Language Model (SLM)

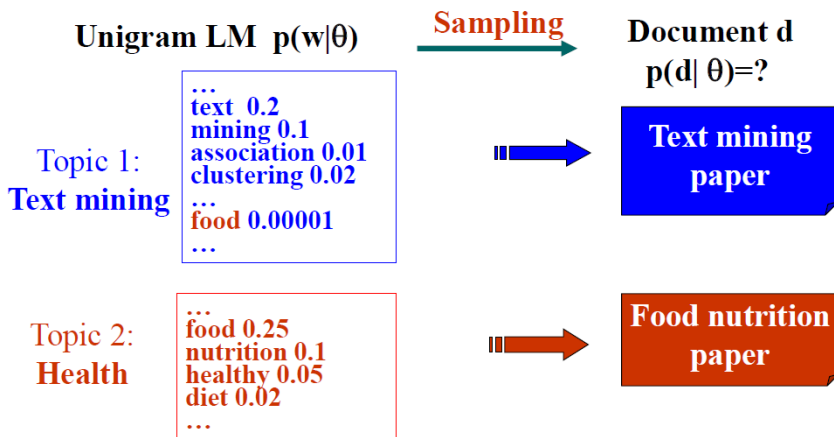
- A probability distribution over word sequences
 - $p(\text{"Today is Wednesday"}) \approx 0.001$
 - $p(\text{"Today Wednesday is"}) \approx 0.00000000000001$
 - $p(\text{"The eigenvalue is positive"}) \approx 0.00001$
- Context-dependent!
- Probabilistic mechanism for "generating" text = **"generative" model**

The Simplest LM: Unigram LM

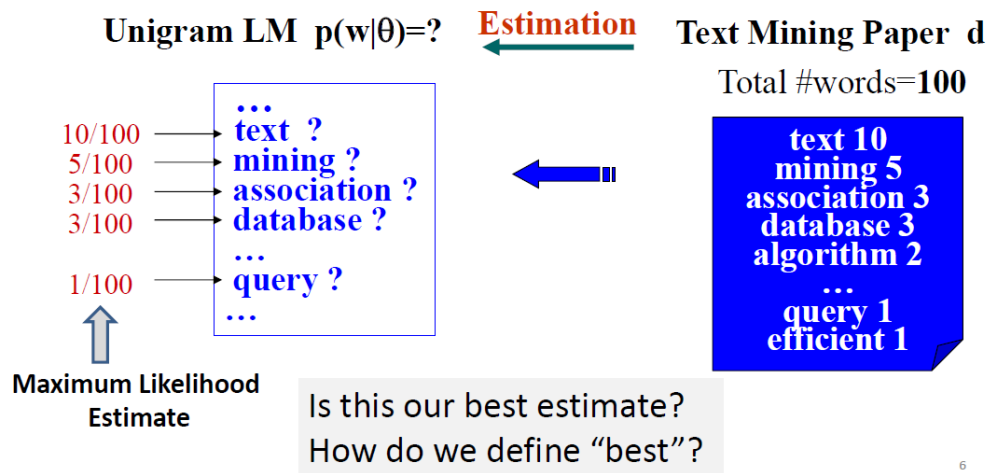
- Each word generated INDEPENDENTLY
- $p(w_1 w_2 \dots w_n) = p(w_1)p(w_2)\dots p(w_n)$
- Parameters: $\{p(w_i)\}$ $p(w_1)+\dots+p(w_N)=1$ (N is voc. size)
- Text = sample drawn according to this word distribution

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

Text Generation with Unigram LM



Estimation of Unigram LM



Maximum Likelihood vs. Bayesian

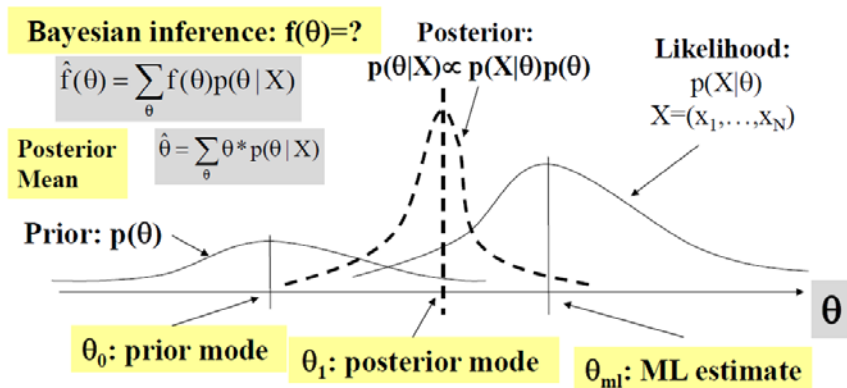
- Maximum likelihood estimation
 - “Best” – maximum data likelihood $\hat{\theta} = \arg \max_{\theta} P(X | \theta)$
 - Problem: Small sample
- Bayesian estimation: **Bayes Rule** $p(X | Y) = \frac{p(Y | X)p(X)}{p(Y)}$
 - “Best” - consistent with “prior” knowledge and explaining data well

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



Maximum a Posteriori (MAP) estimate

Illustration of Bayesian Estimation



Summary

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

Probabilistic Topic Models: Mining One Topic

Simplest Case of Topic Model: Mining One Topic

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

OUTPUT: $\{\theta\}$

Text Data

thing I am still considering taking back. Here's why:

Speaker quality is ABSOLUTELY HORRIFIC! The speaker is simply retuned and, unless you are a perfect coding dork, when you turn it all the way because you can't hear it, it comes across, sounding like a blow dryer. What is this thing, the size of a fly? And if there is ANY background noise, you can't hear it at all with the tone up for just you're wondering, no, the speaker didn't fail. I exchanged the iPod and the second one sounds (same make) but I can't hear videos, show music, than anything unless my hands and arms are in one place.

heard occurs for the speaker's embarrassing quality, including "It's a small device." My cheap Samsung cell isn't a speaker is better than the iPod Touch, and the phone is smaller. The other excuse I've heard is "The speaker was designed for game play." (Wait, what? That's not a valid excuse, because the iPod is a *music* device.) What's the excuse, what if I'm not a gamer? I need to hear the game play? For example, like 2 player Tap Tap Revolution 3 is playable, unless you are in a quiet room.

Price: Let's face it, this thing is *not* expensive, as we EVERY SINGLE ONE OF APPLE'S PRODUCTS. Sure, the i7 "Designed in California," but they are MADE in CHINA, just like everything else.

Navigation 1047E, page 4

$$P(w | \theta)$$
 θ

text ?
mining ?
association ?
database ?

...
query ?
...

Doc d

100%

Language Model Setup

- **Data:** Document $\mathbf{d} = x_1 x_2 \dots x_{|\mathbf{d}|}$, $\mathbf{x}_i \in \mathbf{V} = \{w_1, \dots, w_M\}$ 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:** $p(\mathbf{d} | \theta) = p(x_1 | \theta) \times \dots \times p(x_{|\mathbf{d}|} | \theta)$

$$= p(w_1 | \theta)^{c(w_1, \mathbf{d})} \times \dots \times p(w_M | \theta)^{c(w_M, \mathbf{d})}$$

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

Computation of Maximum Likelihood Estimate

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

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

Subject to constraint:

$$\sum_{i=1}^M \theta_i = 1$$

Use Lagrange multiplier approach

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

$$\frac{\partial f(\square | \mathbf{d})}{\partial \square_i} = \frac{c(w_i, \mathbf{d})}{\square_i} + \square = 0 \rightarrow \square_i = -\frac{c(w_i, \mathbf{d})}{\square}$$

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

Normalized Counts



What Does the Topic Look Like?

