

# Syntagmatic Relation Discovery: Entropy

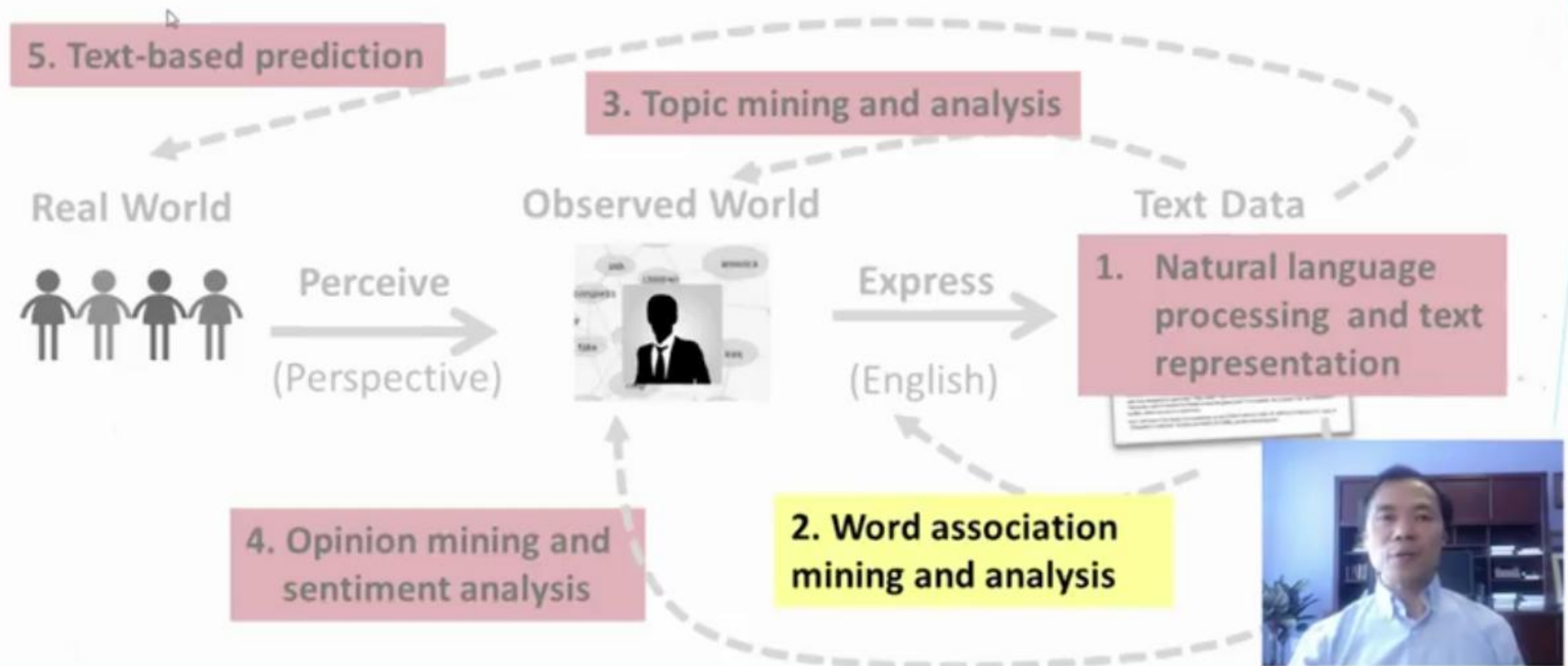
ChengXiang “Cheng” Zhai  
Department of Computer Science  
University of Illinois at Urbana-Champaign



00:09 / 11:00



# 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	_____	<b>eats</b>	_____	on Saturday
His	_____	<b>eats</b>	_____	on Tuesday
My	_____	<b>eats</b>	_____	on Sunday
His	_____	<b>eats</b>	_____	on Tuesday
...	_____		_____	

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

What words  
are to the  
**right**?

# Word Prediction: Intuition

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

Text Segment (any unit, e.g., sentence, paragraph, document)



Are some words easier to predict than others?

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

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

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



01:56 / 11:00



## Word Prediction: Formal Definition

Binary Random Variable :  $X_w = \begin{cases} 1 & \text{w is present} \\ 0 & \text{w is absent} \end{cases}$   
 $X_w \in \{0, 1\}$

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

The more random  $X_w$  is, the more difficult the prediction would be.

How does one quantitatively measure the “randomness” of a random variable like  $X_w$ ?

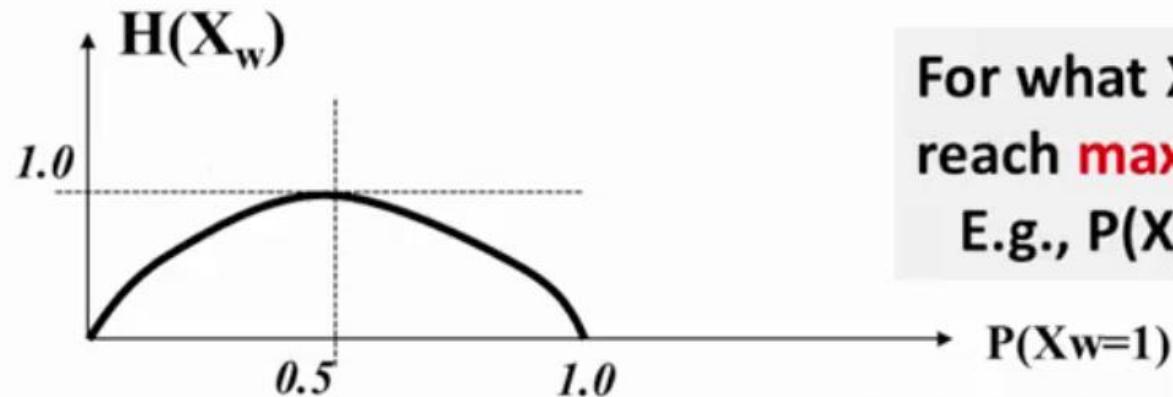


# Entropy $H(X)$ Measures Randomness of $X$

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

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

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



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

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

**$X_{\text{coin}}$ : tossing a coin**

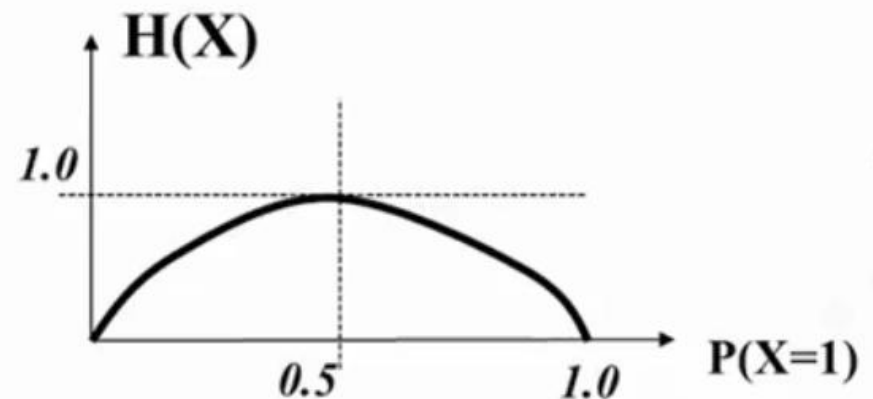
$$X_{\text{coin}} = \begin{cases} 1 & \text{Head} \\ 0 & \text{Tail} \end{cases}$$

**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  $H(X_{\text{unicorn}})$ ?

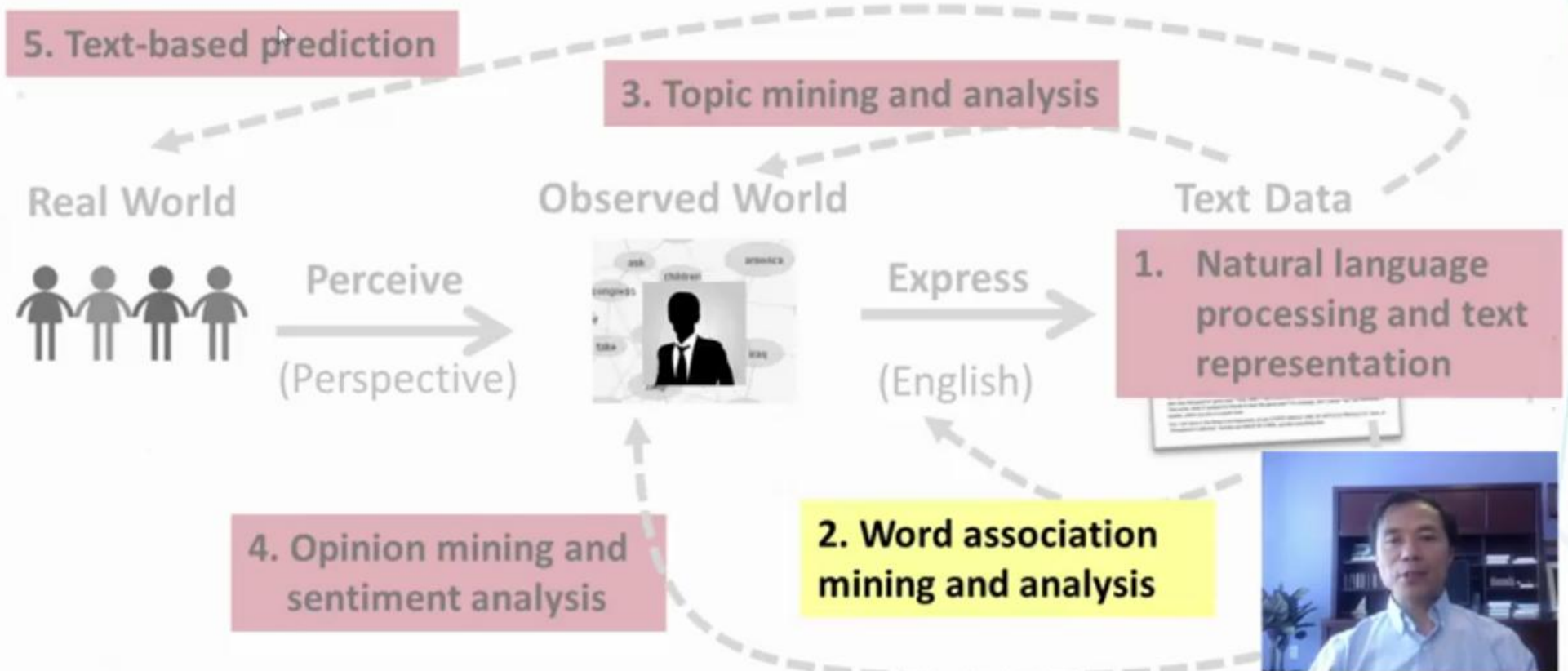
$H(X_{\text{the}}) \approx 0 \rightarrow$  no uncertainty since  $p(X_{\text{the}}=1) \approx 1$

**High entropy words are harder to predict!**





# Syntagmatic Relation Discovery: Conditional Entropy



## What If We Know More About a Text Segment?

Prediction question: Is “**meat**” present (or absent) in this segment?

☐ ☐ ... ☐ eats ☐ ... ☐

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?

# Conditional Entropy

Know nothing about the segment

$$p(X_{meat} = 1)$$



$$p(X_{meat} = 1 | X_{eats} = 1)$$

$$p(X_{meat} = 0)$$



$$p(X_{meat} = 0 | X_{eats} = 1)$$

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



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

$H(X_{meat} | X_{eats} = 0)$  can be defined similarly

# Conditional Entropy: Complete Definition

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

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

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



03:53 / 11:57



5



## Conditional Entropy to Capture Syntagmatic Relation

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

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

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,  $H(X_{meat})$ ?





# Conditional Entropy for Mining Syntagmatic Relations

- For each word  $W1$ 
  - For every other word  $W2$ , compute conditional entropy  $H(X_{W1} | X_{W2})$
  - Sort all the candidate words in ascending order of  $H(X_{W1} | X_{W2})$
  - Take the top-ranked candidate words as words that have potential syntagmatic relations with  $W1$
  - Need to use a threshold for each  $W1$
- However, while  $H(X_{W1} | X_{W2})$  and  $H(X_{W1} | X_{W3})$  are comparable,  $H(X_{W1} | X_{W2})$  and  $H(X_{W3} | X_{W2})$  aren't!

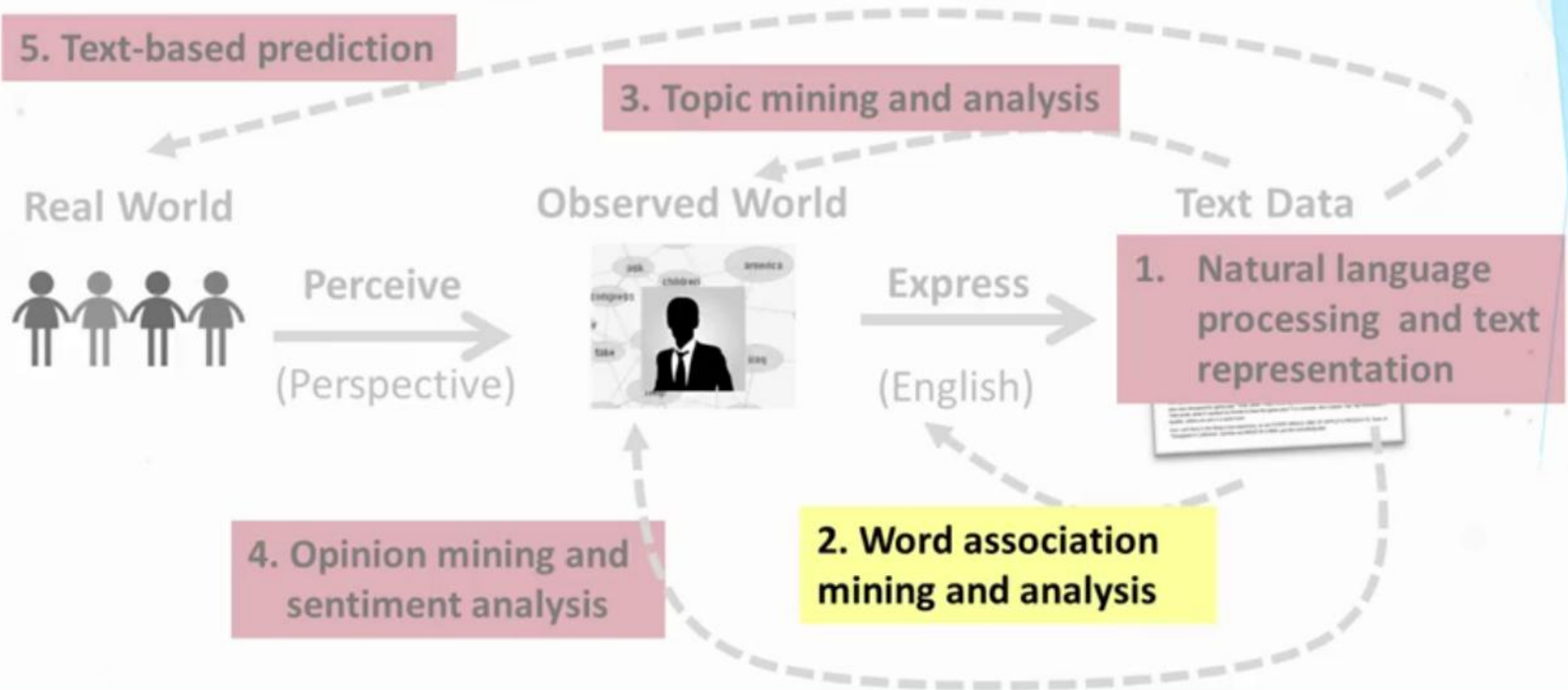
How can we mine the **strongest** K syntagmatic relations from a collection?



09:05 / 11:57



# Syntagmatic Relation Discovery: Mutual Information



# Mutual Information $I(X;Y)$ : Measuring Entropy Reduction

How much reduction in the entropy of  $X$  can we obtain by knowing  $Y$ ?

**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

When we fix  $X$  to rank different  $Y$ s,  $I(X;Y)$  and  $H(X|Y)$  give the same order, but  $I(X;Y)$  allows us to compare different  $(X,Y)$  pairs.

# Mutual Information $I(X;Y)$ for Syntagmatic Relation Mining

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

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

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

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

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



# 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.



# Probabilities Involved in Mutual Information

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

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

# Relations Between Different Probabilities

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

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

**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 and absence of  $w_1$ :

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

Presence and 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  $p(X_{w_1}=1, X_{w_2}=1)$ .

# Syntagmatic Relation Discovery: Mutual Information

Part 2

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





# 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

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

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

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

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

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

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.

