

# Quick Tour of Text Mining

---

Yi-Shin Chen

Institute of Information Systems and Applications

Department of Computer Science

National Tsing Hua University

yishin@gmail.com

## About Speaker

陳宜欣 Yi-Shin Chen

▷ Currently

- 清華大學資訊工程系副教授
- 主持智慧型資料工程與應用實驗室 (IDEA Lab)

▷ Education

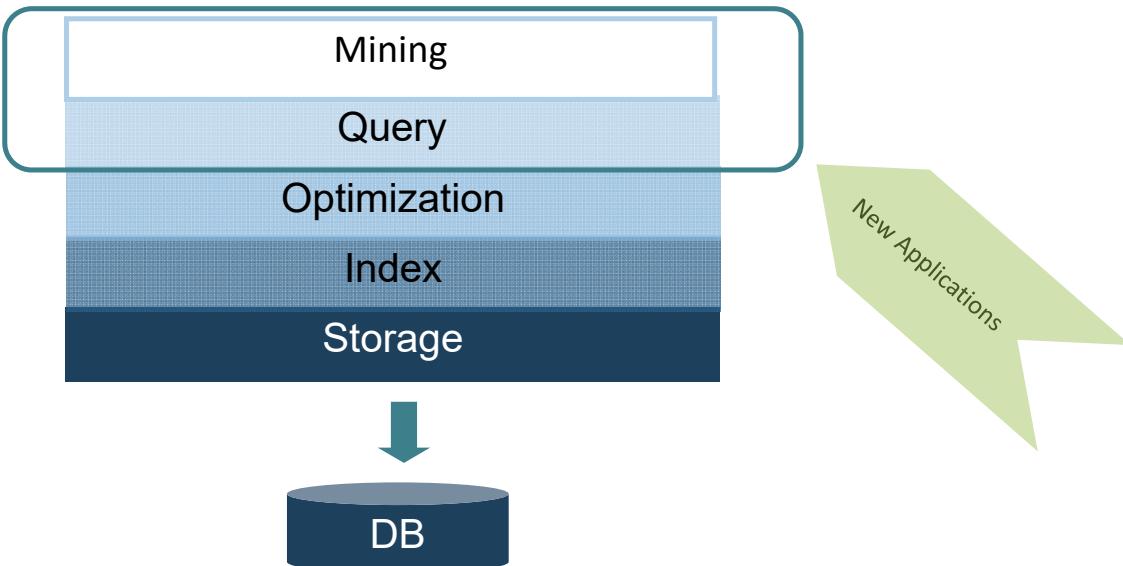
- Ph.D. in Computer Science, USC, USA
- M.B.A. in Information Management, NCU, TW
- B.B.A. in Information Management, NCU, TW

▷ Courses (*all in English*)

- Research and Presentation Skills
- Introduction to Database Systems
- Advanced Database Systems
- Data Mining: Concepts, Techniques, and Applications



# Research Focus from 2000



@ Yi-Shin Chen, Text Mining Overview

3

## Free Resources

### ▷ 免費數據

- 不管公網、私網，能合法下載的資料都是好物



@ Yi-Shin Chen, Text Mining Overview

4

# Past and Current Studies

- ▷ Location identification
- ▷ Interest identification
- ▷ Event identification
- ▷ Extract semantic relationships
- ▷ Unsupervised multilingual sentiment analysis
- ▷ Keyword extraction and summary
- ▷ Emotion analysis
- ▷ Mental illness detection

Text Processing

@ Yi-Shin Chen, Text Mining Overview

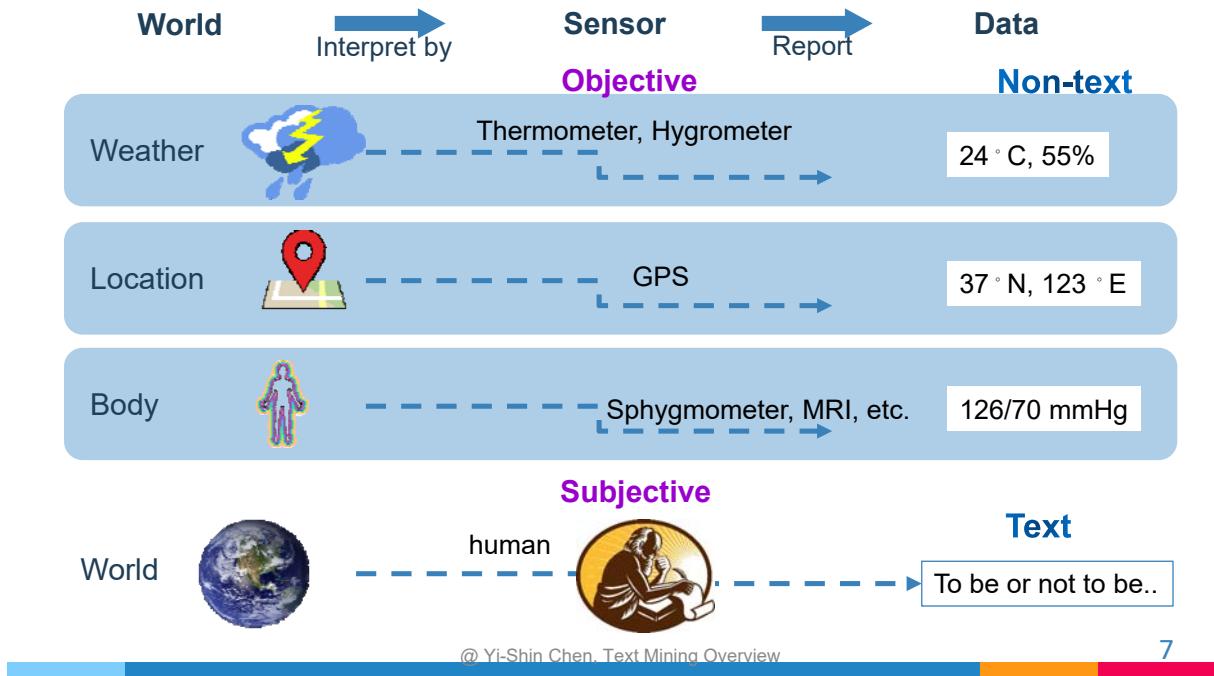
5

# Text Mining Overview

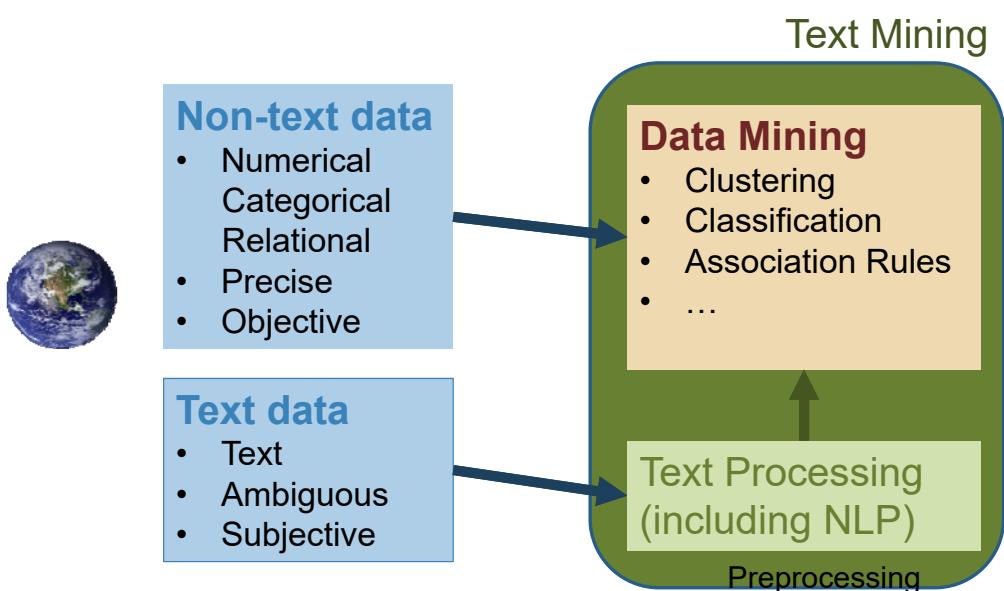
@ Yi-Shin Chen, Text Mining Overview

6

# Data (Text vs. Non-Text)



## Data Mining vs. Text Mining

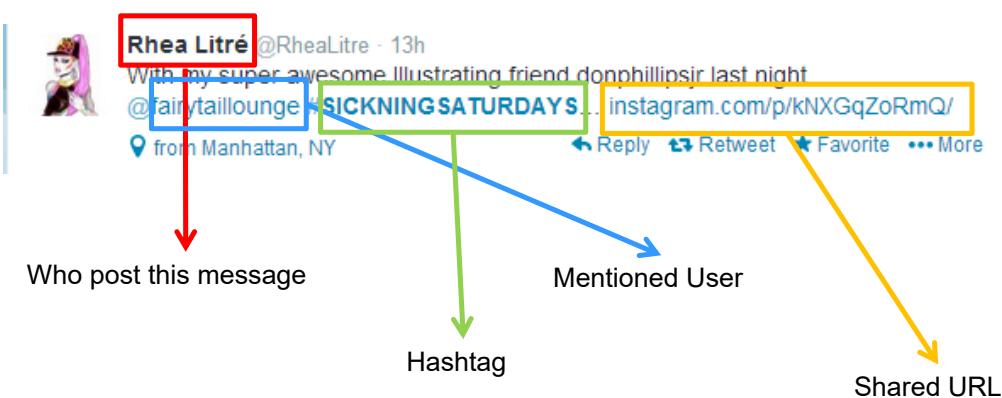


# Preprocessing in Reality

9

## Data Collection

- ▷ Align /Classify the attributes correctly



10

# Language Detection

- ▷ To detect a language (possible languages) in which the specified text is written

你好 現在幾點鐘  
apa kabar sekarang jam berapa ?



繁體中文 (zh-tw)  
印尼文 (id)

- ▷ Difficulties

- Short message
- Different languages in one statement
- Noisy

11

## Wrong Detection Examples

Before / after removing noise

@sayidatynet top song #LailaGhofran  
shokran ya garh new album #listen

en -> id

中華隊的服裝挺特別的，好藍。。。  
#ChineseTaipei #Sochi #2014冬奧

it -> zh-tw

授業前の雪合戦w  
<http://t.co/d9b5peaq7J>

en -> ja

12

# Removing Noise

## ▷ Removing noise before detection

- Html file ->tags
- Twitter -> hashtag, mention, URL

```
<meta name="twitter:description" content="觸犯法國隱私法〔駐歐洲特派記者胡蕙寧、國際新聞中心／綜合報導〕網路搜尋引擎巨擘Google8日在法文版首頁( www.google.fr )張貼悔過書..." />
```

英文  
(en)

觸犯法國隱私法〔駐歐洲特派記者胡蕙寧、國際新聞中心／綜合報導〕網路搜尋引擎巨擘Google8日在法文版首頁( www.google.fr )張貼悔過書...

繁中  
(zh-tw)

13

# Data Cleaning

## ▷ Special character

Unicode emotions	😊, ❤...
Symbol icon	☎, ✉ ...
Currency symbol	€, £, \$...

## ▷ Utilize regular expressions to clean data

Tweet URL

(^|\s\*)http(\S+)?(\s\*|\$)

ube playlist  
mie Riepe

Filter out non-(letters, space,  
punctuation, digit)

(\p{L}+)(\p{Z}+)|  
(\p{Punct}+)(\p{Digit}+)

g ❤ ✉

14

## Japanese Examples

- ▷ Use regular expression remove all special words

\\W

- うふふふふ(\*^-^\*)楽しむ！ありがとうございます ^O^ アイコン、ラブラブ(-\_-)♡
- うふふふふ 楽しむ ありがとうございます アイコン ラブラブ

15

## Part-of-speech (POS) Tagging

- ▷ Processing text and assigning parts of speech to each word
- ▷ Twitter POS tagging
  - Noun (N), Adjective (A), Verb (V), URL (U)...

Happy Easter! I went to work and came home to an empty house now im going for a quick run http://t.co/Ynp0uFp6oZ

Happy\_A Easter\_N !\_, I\_O went\_V to\_P work\_N and\_& came\_V home\_N to\_P an\_D empty\_A house\_N now\_R im\_L going\_V for\_P a\_D quick\_A run\_N http://t.co/Ynp0uFp6oZ\_U

16

# Stemming

love miss be  
RT @kt\_biv : @caycelynne loving and missing you! we are  
still looking for Lucy  
look

- ▷ @DirtyDTran gotta be **caught** up for tomorrow **nights** episode
  - ▷ @ASVP\_Jaykey for some **reasons** I found this very amusing
- ↓
- @DirtyDTran gotta be **catch** up for tomorrow **night** episode
  - @ASVP\_Jaykey for some **reason** I **find** this very amusing

17

# Hashtag Segmentation

- ▷ By using Microsoft Web N-Gram Service (or by using **Viterbi algorithm**)

Wow! explosion at a boston race ... #prayforboston

#pray #for #boston

#citizenscience → #citizen #science

#bostonmarathon → #boston #marathon

#goodthingsarecoming → #good #things #are #coming

#lowbloodpressure → #low #blood #pressure

18

# More Preprocesses for Different Web Data

- ▷ Extract source code without javascript
- ▷ Removing html tags

19

## Extract Source Code Without Javascript

- ▷ Javascript code should be considered as an exception
  - it may contain hidden content

```
<html>
<title>Parakweet - Actionable Signals for Social Media</title>
<script src="//use.typekit.net/ygl8ojq.js" type="text/javascript"></script>
<script>
try
{
    Typekit.load();
}
catch(e){}
</script>
<link href="stylesheets/parakweet.css" media="screen" rel="stylesheet" type="text/css" />
<body><div class='lead'>
    Parakweet is a new platform offering social media analytics, recommendations and metadata to media companies.</div>
    <div class='subtext hide' id='subtext_two'>
        Parakweet's Product Recommendations provide customers with relevant, actionable ideas for brand and product discovery.</div>
</body>
</html>
```

```
<html>
<title>Parakweet - Actionable Signals for Social Media</title>
<link href="stylesheets/parakweet.css" media="screen" rel="stylesheet" type="text/css" />
<body><div class='lead'>
    Parakweet is a new platform offering social media analytics, recommendations and metadata to media companies.</div>
    <div class='subtext hide' id='subtext_two'>
        Parakweet's Product Recommendations provide customers with relevant, actionable ideas for brand and product discovery.</div>
</body>
</html>
```

20

# Remove Html Tags

- ▷ Removing html tags to extract meaningful content

```
<html>
<title>Parakweet - Actionable Signals for Social Media</title>
<link href="stylesheets/parakweet.css" media="screen" rel="stylesheet" type="text/css" />
<body><div class='lead'>
    Parakweet is a new platform offering social media analytics, recommendations and metadata to media companies.</div>
    <div class='subtext hide' id='subtext_two'>
        Parakweet's Product Recommendations provide customers with relevant, actionable ideas for brand and product discovery.</div>
    </body>
</html>
```



```
Parakweet - Actionable Signals for Social Media
Parakweet is a new platform offering social media analytics, recommendations and metadata to media companies.
Parakweet's Product Recommendations provide customers with relevant, actionable ideas for brand and product discovery.
```

21

## More Preprocesses for Different Languages

- ▷ Chinese Simplified/Traditional Conversion
- ▷ Word segmentation

22

# Chinese Simplified/Traditional Conversion

## ▷ Word conversion

- 请乘客从后门落车 → 請乘客從後門下車

## ▷ One-to-many mapping

- @shinrei 出去旅游还是崩坏 → @shinrei 出去旅游還是崩壞  
游 (zh-cn) → 游|遊 (zh-tw)

## ▷ Wrong segmentation

- 人体内存在很多微生物 → 內存: 人體 記憶體 在很多微生物  
→ 存在: 人體內 存在 很多微生物  
  
内存|存在

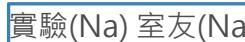
23

# Wrong Chinese Word Segmentation

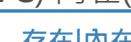
## ▷ Wrong segmentation

- 這(Nep) 地面(Nc) 積(VJ) 還(D) 真(D) 不(D) 小(VH) <http://t.co/QIUbiaz2lz>  
  
地面|面積

## ▷ Wrong word

- @iamzeke 實驗(Na) 室友(Na) 多(Dfa) 危險(VH) 你(Nh) 不(D) 知道(VK) 嘴  
(T) ?  
  
實驗室|室友

## ▷ Wrong order

- 人體(Na) 存(VC) 內在(Na) 很多(Neqa) 微生物(Na)  
  
存在|內在

## ▷ Unknown word

- 半夜(Nd) 逛團(Na) 購(VC) 看到(VE) 太(Dfa) 吸引人(VH) !!  
  
未知詞: 團購

24

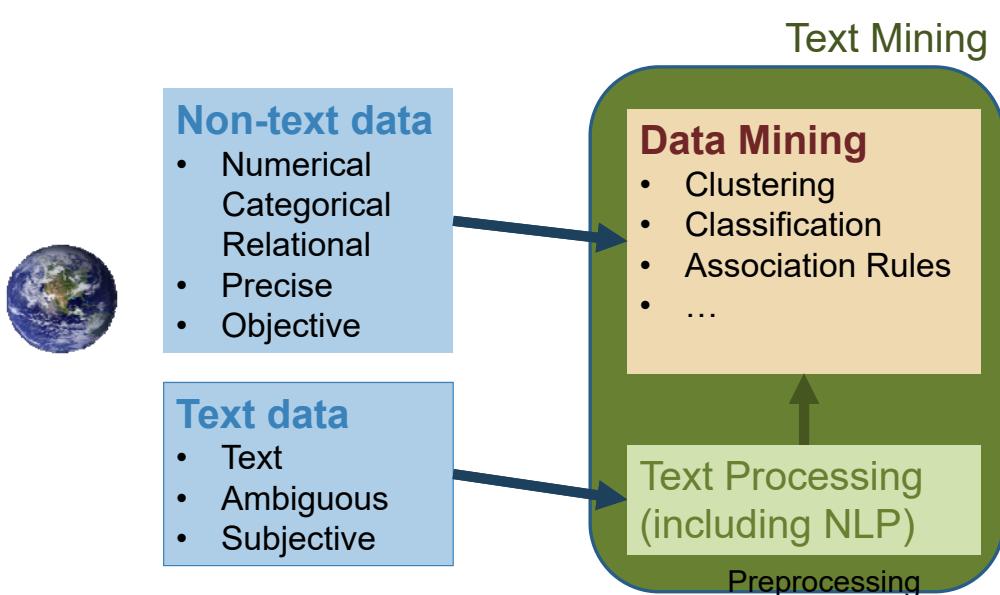
# Back to Text Mining

Let's come back to Text Mining

@ Yi-Shin Chen, Text Mining Overview

25

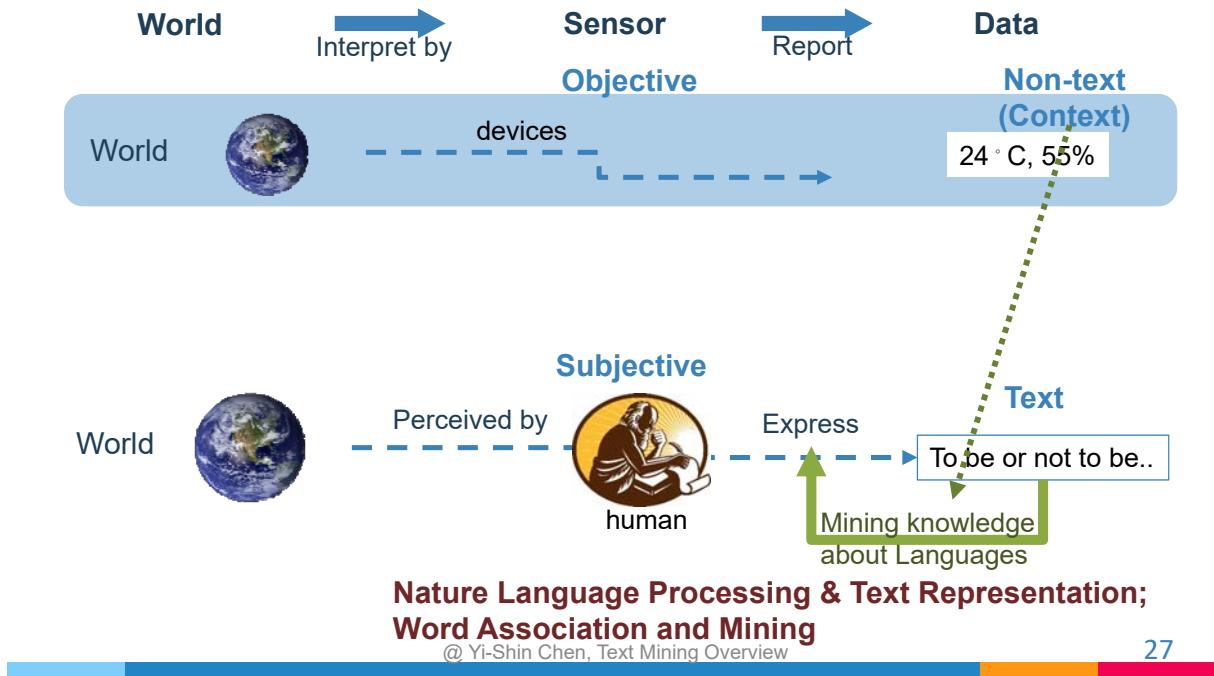
## Data Mining vs. Text Mining



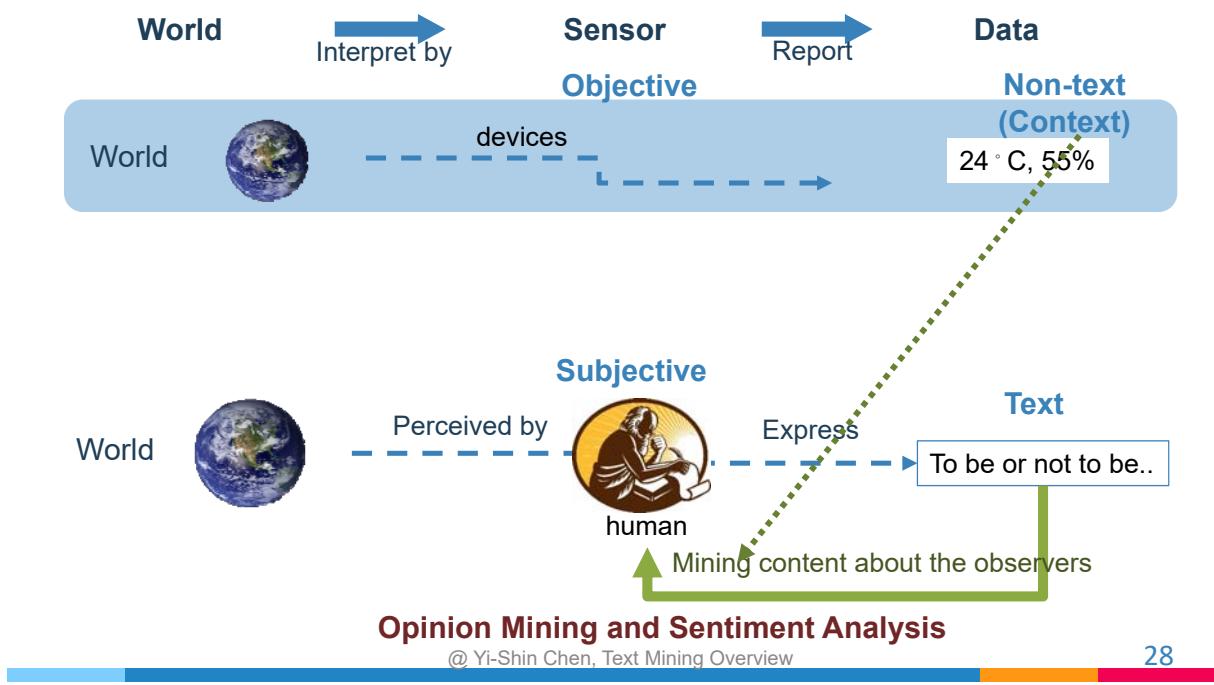
@ Yi-Shin Chen, Text Mining Overview

26

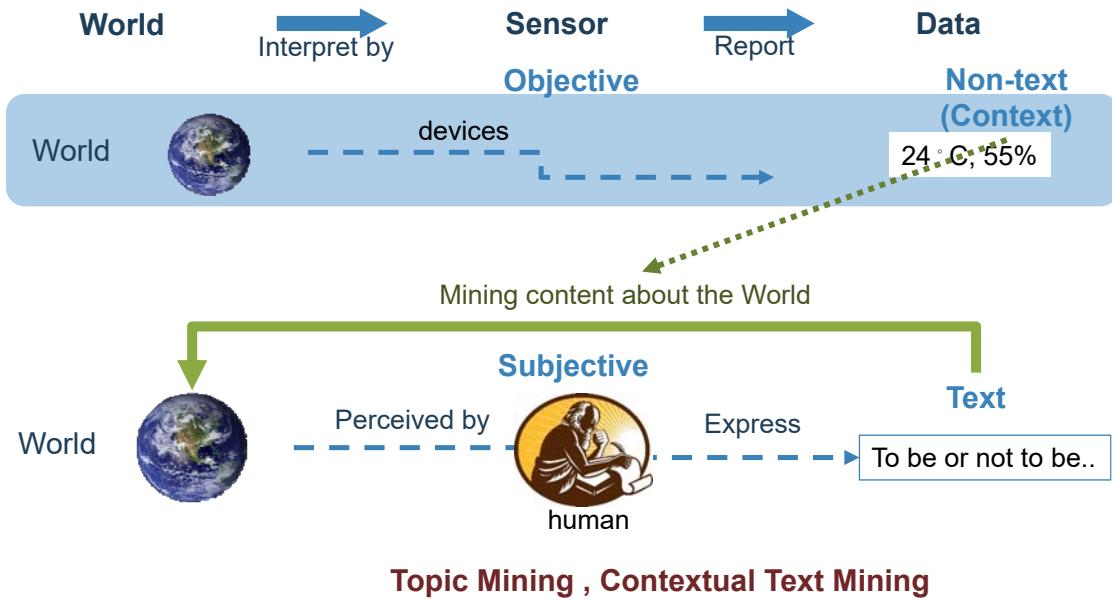
# Landscape of Text Mining



# Landscape of Text Mining



# Landscape of Text Mining



@ Yi-Shin Chen, Text Mining Overview

29

## Basic Concepts in NLP

This is the best thing happened in my life.  
Det. Verb Det. Adj N Verb Pre. PN N

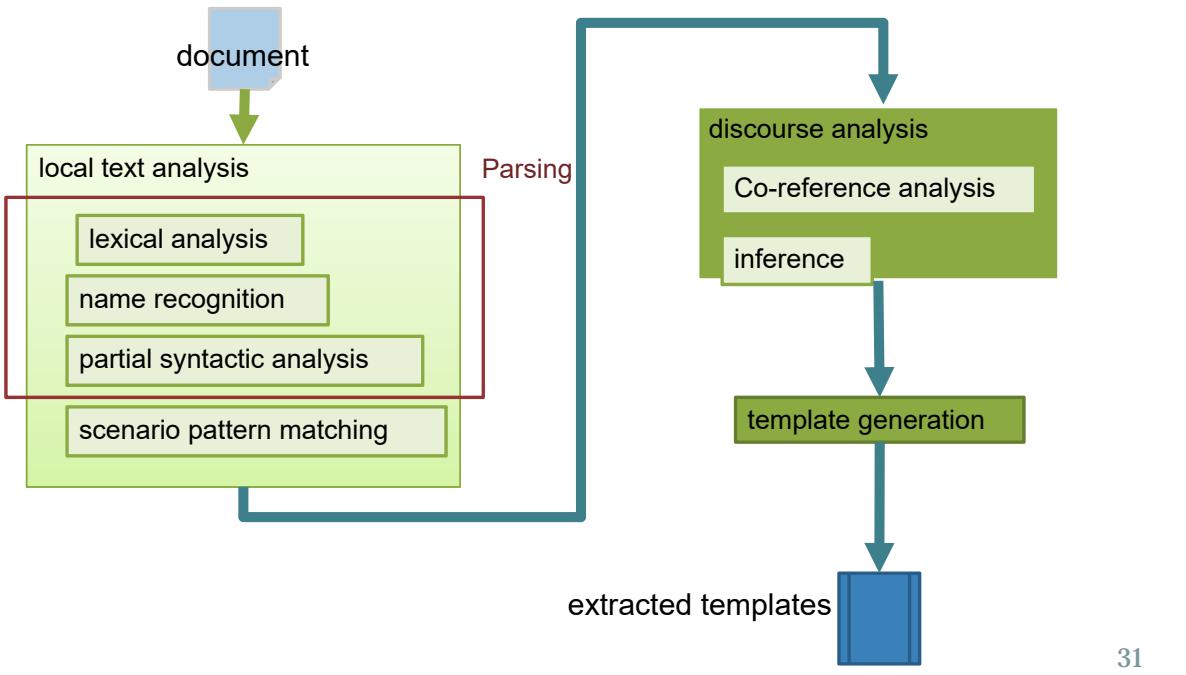
辭彙  
**Lexical analysis**  
**(Part-of Speech Tagging)**

句法  
**Syntactic analysis**  
**(Parsing)**

@ Yi-Shin Chen, Text Mining Overview

30

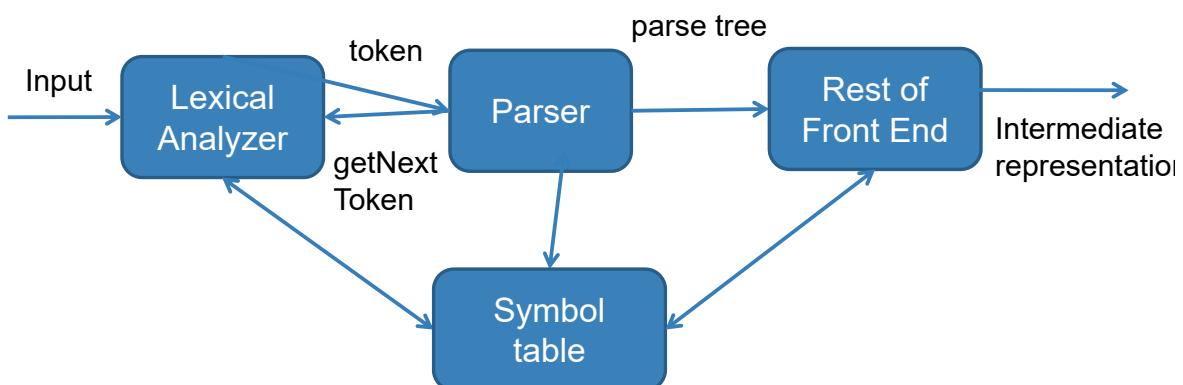
# Structure of Information Extraction System



31

## Parsing

- ▷ Parsing is the process of determining whether a string of tokens can be generated by a grammar



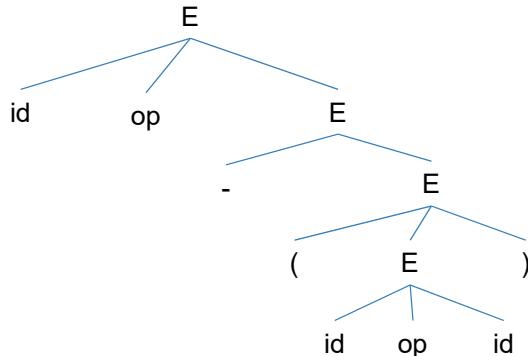
32

# Parsing Example (for Compiler)

▷ Grammar:

- $E ::= E \text{ op } E \mid -E \mid (E) \mid \text{id}$
- $\text{op} ::= + \mid - \mid * \mid /$

a \* - ( b + c )

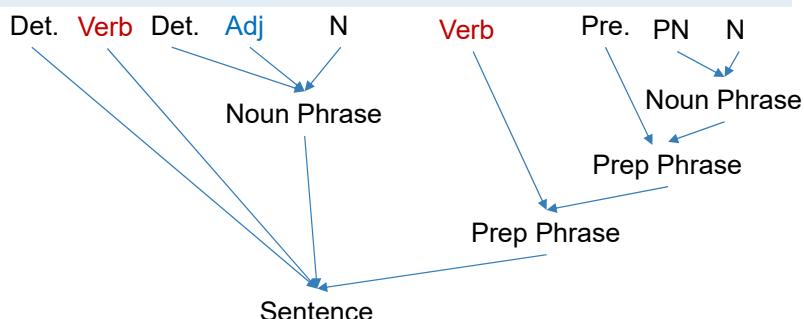


@ Yi-Shin Chen, Text Mining Overview

33

# Basic Concepts in NLP

This is the best thing happened in my life.



辭彙  
**Lexical analysis  
(Part-of Speech  
Tagging)**

句法  
**Syntactic analysis  
(Parsing)**

This? (t1)  
Best thing (t2)  
My (m1)  
Happened (t1, t2, m1)

語意  
**Semantic Analysis**

推理  
**Inference  
(Emotion  
Analysis)**

Happy (x) if Happened (t1, 'Best', m1) → Happy

@ Yi-Shin Chen, Text Mining Overview

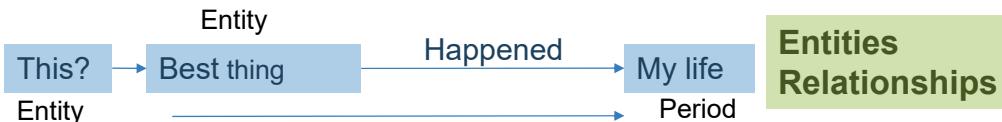
34

# Basic Concepts in NLP

This is the best thing happened in my life. **String of Characters**

This is the best thing happened in my life. **String of Words**

Det. Verb Det. Adj N Verb Pre. PN N **POS Tags**



Happy **Emotion**

The writer loves his new born baby **Understanding (Logic predicates)**

Deeper NLP  
Less accurate  
Closer to knowledge

@ Yi-Shin Chen, Text Mining Overview

35

## NLP vs. Text Mining

### ▷ NLP objectives

- Understanding
- Ability to answer
- Immaculate

### ▷ Text Mining objectives

- Overview
- Know the trends
- Accept noise

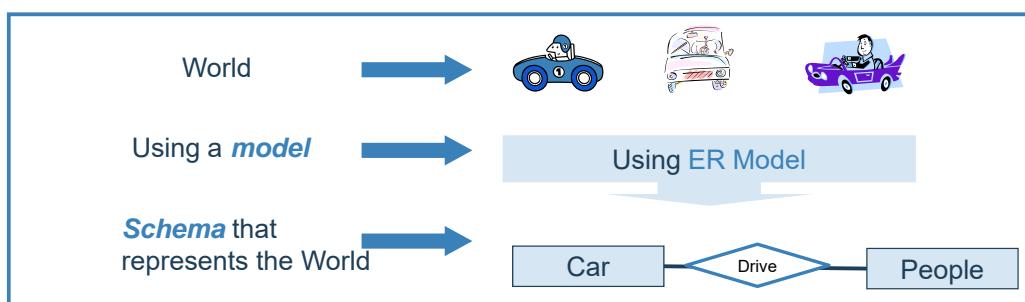
# Basic Data Model Concepts

Let's learn from giants

@ Yi-Shin Chen, Text Mining Overview

37

## Data Models



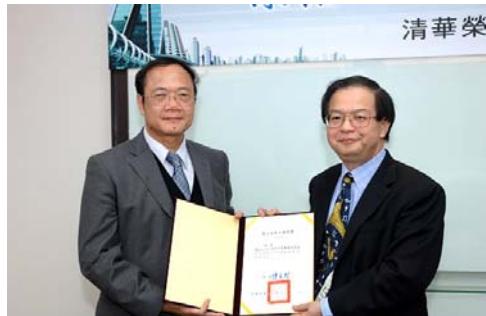
- ▷ Data model/**Language**: a *collection* of concepts for describing data
- ▷ Schema/**Structured observation**: a *description* of a particular collection of data, using the a given data model
- ▷ Data instance/Statements

@ Yi-Shin Chen, Text Mining Overview

38

# E-R Model

- ▷ Introduced by Peter Chen; ACM TODS, March 1976
  - <https://www.linkedin.com/in/peter-chen-43bba0/>



- Additional Readings
  - Peter Chen. "English Sentence Structure and Entity-Relationship Diagram." *Information Sciences*, Vol. 1, No. 1, Elsevier, May 1983, Pages 127-149
  - Peter Chen. "A Preliminary Framework for Entity-Relationship Models." *Entity-Relationship Approach to Information Modeling and Analysis*, North-Holland (Elsevier), 1983, Pages 19 - 28

@ Yi-Shin Chen, Text Mining Overview

39

## E-R Model Basics -Entity

- ▷ Based on a perception of a real world, which consists
  - A set of basic objects ⇒ Entities
  - Relationships among objects
- ▷ Entity: Real-world object distinguishable from other objects
- ▷ Entity Set: A collection of similar entities. E.g., all employees.
  - Presented as:

Animals

Time

People

Things

This is the best thing happened in my life.

Dogs love their owners.

@ Yi-Shin Chen, Text Mining Overview

40

# E-R: Relationship Sets

▷ Relationship: Association among two or more entities

▷ Relationship Set: Collection of similar relationships.

- Relationship set are presented as:



- The relationship cannot exist without having corresponding entities



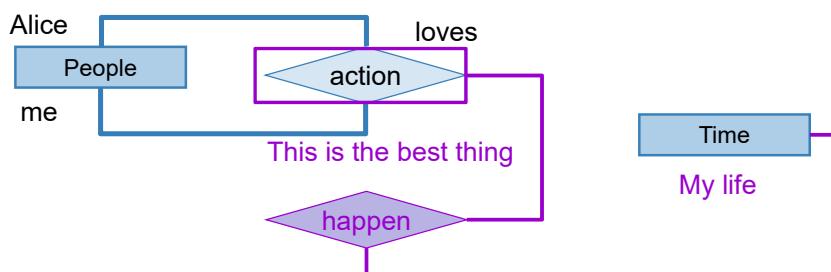
Dogs love their owners.

@ Yi-Shin Chen, Text Mining Overview

41

# High-Level Entity

▷ High-level entity: Abstracted from a group of interconnected low-level entity and relationship types



Alice loves me.

This is the best thing happened in my life.

@ Yi-Shin Chen, Text Mining Overview

42

# Word Relations

[Back to text](#)

@ Yi-Shin Chen, Text Mining Overview

43

## Word Relations

▷ Paradigmatic: can be substituted for each other (similar)

- E.g., Cat & dog, run and walk



▷ Syntagmatic: can be combined with each other (correlated)

- E.g., Cat and fights, dog and barks



→ These two basic and complementary relations can be generalized to describe relations of any times in a language

@ Yi-Shin Chen, Text Mining Overview

44

# Mining Word Associations

## ▷ Paradigmatic

- Represent each word by its context
- Compute context similarity
- Words with high context similarity

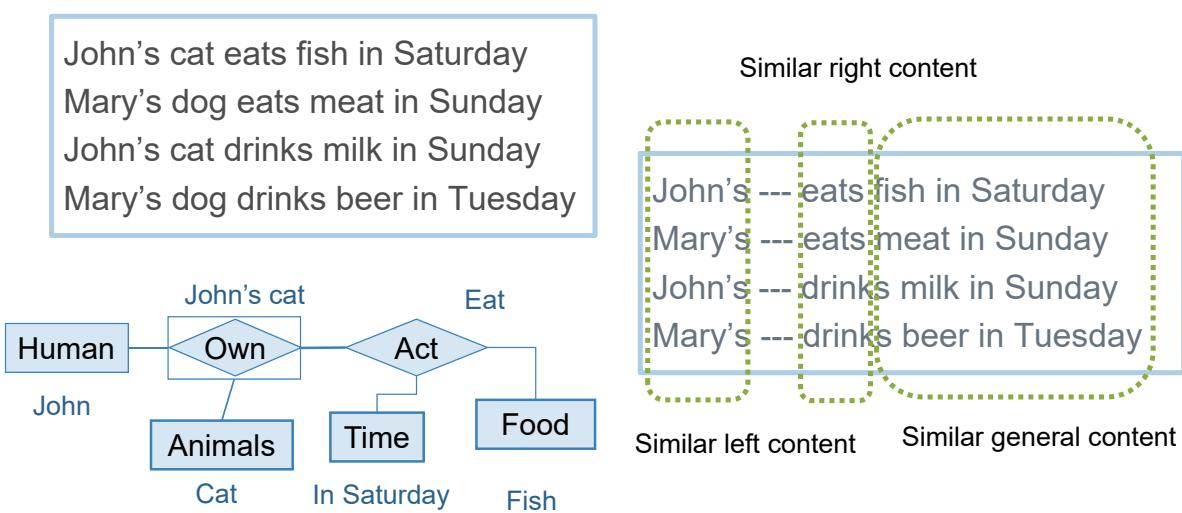
## ▷ Syntagmatic

- Count the number of times two words occur together in a context
- Compare the co-occurrences with the corresponding individual occurrences
- Words with high co-occurrences but relatively low individual occurrence

@ Yi-Shin Chen, Text Mining Overview

45

## Paradigmatic Word Associations



How similar are context ("cat") and context ("dog")?

How similar are context ("cat") and context ("John")?

→ Expected Overlap of Words in Context (EOWC)

Overlap ("cat", "dog")

Overlap ("cat", "John")

46

@ Yi-Shin Chen, Text Mining Overview

## Vector Space Model (Bag of Words)

- ▷ Represent the keywords of objects using a term vector
  - Term: basic concept, e.g., keywords to describe an object
  - Each term represents one dimension in a vector
  - N total terms define an n-element terms
  - Values of each term in a vector corresponds to the importance of that term
- ▷ Measure similarity by the vector distances

	team	catch	y	ball	score	game	u	wi	lost	timeout	season
Document 1	3	0	5	0	2	6	0	2	0	0	2
Document 2	0	7	0	2	1	0	0	0	3	0	0
Document 3	0	1	0	0	1	2	2	0	3	0	0

47

## Common Approach for EOWC: Cosine Similarity

- ▷ If  $d_1$  and  $d_2$  are two document vectors, then
$$\cos(d_1, d_2) = (d_1 \bullet d_2) / \|d_1\| \|d_2\|,$$
where  $\bullet$  indicates vector dot product and  $\|d\|$  is the length of vector  $d.$
- ▷ Example:

$$d_1 = 3 \ 2 \ 0 \ 5 \ 0 \ 0 \ 0 \ 2 \ 0 \ 0$$
$$d_2 = 1 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 1 \ 0 \ 2$$

$$d_1 \bullet d_2 = 3*1 + 2*0 + 0*0 + 5*0 + 0*0 + 0*0 + 0*0 + 2*1 + 0*0 + 0*2 = 5$$
$$\|d_1\| = (3^2 + 2^2 + 0^2 + 5^2 + 0^2 + 0^2 + 0^2 + 2^2 + 0^2 + 0^2)^{0.5} = (42)^{0.5} = 6.481$$
$$\|d_2\| = (1^2 + 0^2 + 0^2 + 0^2 + 0^2 + 0^2 + 0^2 + 1^2 + 0^2 + 2^2)^{0.5} = (6)^{0.5} = 2.245$$

$$\cos(d_1, d_2) = .3150$$

→ Overlap ("John", "Cat") = .3150

## Quality of EOWC?

- ▷ The more overlap the two context documents have, the higher the similarity would be
- ▷ However:
  - It favors matching one frequent term very well over matching more distinct terms
  - It treats every word equally (overlap on “the” should not be as meaningful as overlap on “eats”)

## Term Frequency and Inverse Document Frequency (TFIDF)

- ▷ Since not all objects in the vector space are equally important, we can weight each term using its occurrence probability in the object description
  - Term frequency:  $TF(d, t)$ 
    - number of times  $t$  occurs in the object description  $d$
  - Inverse document frequency:  $IDF(t)$ 
    - to scale down the terms that occur in many descriptions

# Normalizing Term Frequency

▷  $n_{ij}$  represents the number of times a term  $t_i$  occurs in a description  $d_j$ .  $tf_{ij}$  can be normalized using the total number of terms in the document

- $tf_{ij} = \frac{n_{ij}}{NormalizedValue}$

▷ Normalized value could be:

- Sum of all frequencies of terms
- Max frequency value
- Any other values can make  $tf_{ij}$  between 0 to 1
- BM25\*:  $tf_{ij} = \frac{n_{ij} \times (k+1)}{n_{ij} + k}$

# Inverse Document Frequency

▷ IDF seeks to scale down the coordinates of terms that occur in many object descriptions

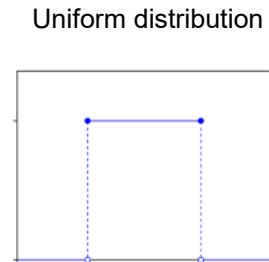
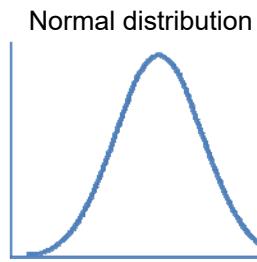
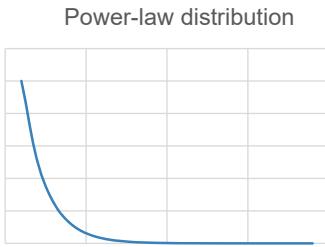
- For example, some stop words(the, a, of, to, and...) may occur many times in a description. However, they should be considered as non-important in many cases

- $idf_i = \log \left( \frac{N}{df_i} + 1 \right)$

→ where  $df_i$  (document frequency of term  $t_i$ ) is the number of descriptions in which  $t_i$  occurs

▷ IDF can be replaced with ICF (inverse class frequency) and many other concepts based on applications

# Reasons of Log



▷ Each distribution can indicate the hidden force

- Time
- Independent
- Control

# Mining Word Associations

## ▷ Paradigmatic

- Represent each word by its context
- Compute context similarity
- Words with high context similarity

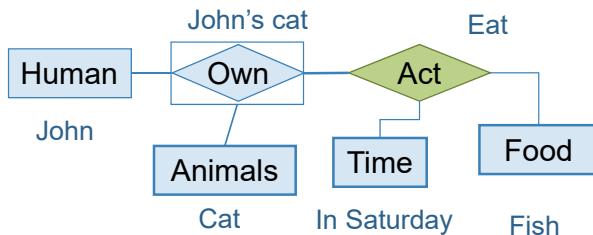
## ▷ Syntagmatic

- Count the number of times two words occur together in a context
- Compare the co-occurrences with the corresponding individual occurrences
- Words with high co-occurrences but relatively low individual occurrence

# Syntagmatic Word Associations

Correlated occurrences

John's cat eats fish in Saturday  
Mary's dog eats meat in Sunday  
John's cat drinks milk in Sunday  
Mary's dog drinks beer in Tuesday



What words to the right?

John's \*\*\* eats \*\*\* in Saturday  
Mary's \*\*\* eats \*\*\* in Sunday  
John's \*\*\* drinks \*\*\* in Sunday  
Mary's \*\*\* drinks \*\*\* in Tuesday

What words tend to occur to the left of "eats"

Whenever "eats" occurs, what other words also tend to occur?

$$P(\text{dog} \mid \text{eats}) = ? ; P(\text{cats} \mid \text{eats}) = ?$$

@ Yi-Shin Chen, Text Mining Overview

55

## Word Prediction

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

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



Apple designs and creates the iPhone, iPad, Mac 1. [Apple Inc.](#) - Official site, with details of products and services.  
iOS 8, OS X, iPod and iTunes, and the new Apple' -- <http://www.apple.com/> Computers: Apple Inc. (1)

Predict the occurrence of word

W1 = 'meat'

W2 = 'a'

W3 = 'unicorn'

@ Yi-Shin Chen, Text Mining Overview

56

# Word Prediction: Formal Definition

▷ Binary random variable {0,1}

- $x_w = \begin{cases} 1 & w \text{ is present} \\ 0 & w \text{ is absent} \end{cases}$
- $P(x_w = 1) + P(x_w = 0) = 1$

▷ The more random  $x_w$  is, the more difficult the prediction is

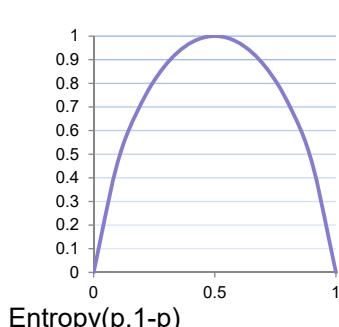
▷ How do we quantitatively measure the **randomness**?

## Entropy

▷ Entropy measures the amount of randomness or surprise or uncertainty

▷ Entropy is defined as:

$$H(p_1, \dots, p_n) = \sum_{i=1}^n \left( p_i \times \log \frac{1}{p_i} \right) = - \sum_{i=1}^n (p_i \times \log p_i)$$



$$\text{where } \sum_{i=1}^n (p_i) = 1$$

- entropy = 0  
easy
- entropy = 1  
difficult

# Conditional Entropy

$$H(Y|X) = \sum_{x \in X} (p(x)H(Y|X=x))$$

Know nothing about the segment

$$\begin{aligned} p(x_{meat} = 1) \\ p(x_{meat} = 0) \end{aligned}$$

Know "eats" is present ( $X_{eat}=1$ )

$$\begin{aligned} p(x_{meat} = 1|x_{eats} = 1) \\ p(x_{meat} = 0|x_{eats} = 1) \end{aligned}$$

$$H(x_{meat})$$

$$= -p(x_{meat} = 0) \times \log_2(p(x_{meat} = 0)) - p(x_{meat} = 1) \times \log_2(p(x_{meat} = 1))$$

$$H(x_{meat}|x_{eats} = 1)$$

$$\begin{aligned} &= -p(x_{meat} = 0|x_{eats} = 1) \times \log_2(p(x_{meat} = 0|x_{eats} = 1)) - p(x_{meat} = 1|x_{eats} = 1) \\ &\quad \times \log_2(p(x_{meat} = 1|x_{eats} = 1)) \end{aligned}$$

## Mining Syntagmatic Relations

▷ For each word W1

- For every 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 with some given threshold

▷ However

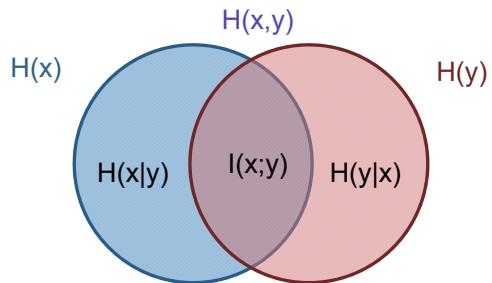
- $H(x_{w1}|x_{w2})$  and  $H(x_{w1}|x_{w3})$  are comparable
- $H(x_{w1}|x_{w2})$  and  $H(x_{w3}|x_{w2})$  are not
  - Because the upper bounds are different

▷ Conditional entropy is not symmetric

- $H(x_{w1}|x_{w2}) \neq H(x_{w2}|x_{w1})$

# Mutual Information

▷  $I(x; y) = H(x) - H(x|y) = H(y) - H(y|x)$



▷ Properties:

- Symmetric
- Non-negative
- $I(x;y)=0$  iff  $x$  and  $y$  are independent

▷ Allow us to compare different  $(x,y)$  pairs

@ Yi-Shin Chen, Text Mining Overview

61

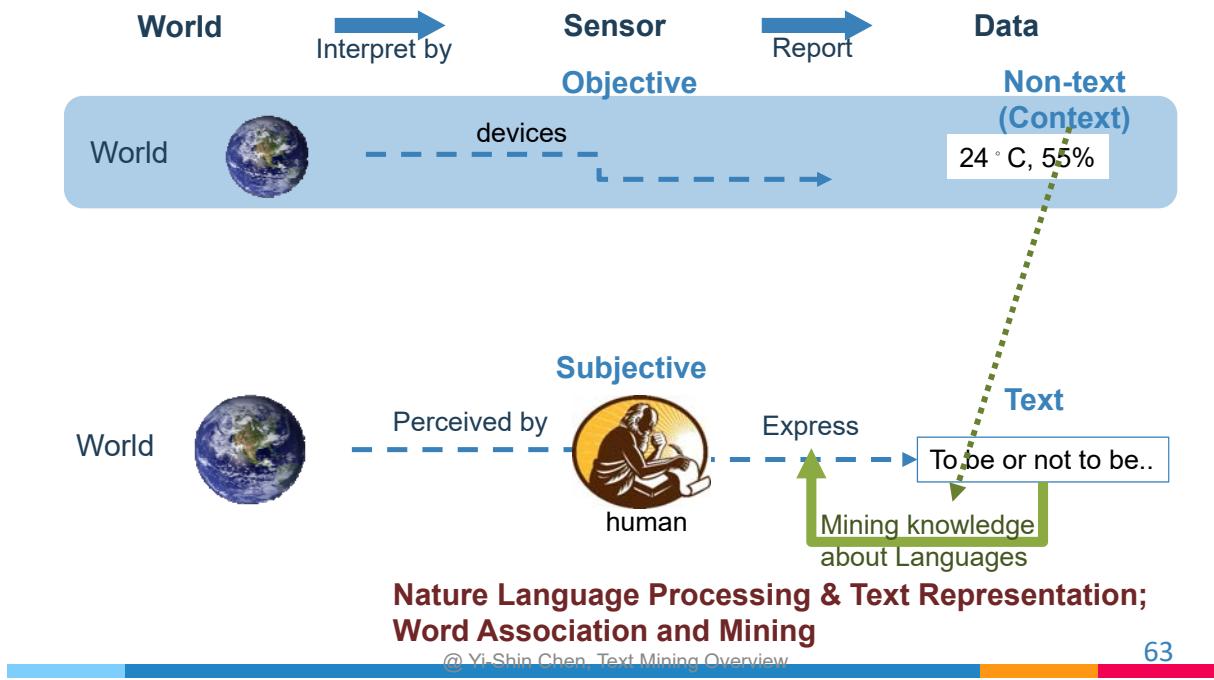
# Topic Mining

Assume we already know the word relationships

@ Yi-Shin Chen, Text Mining Overview

62

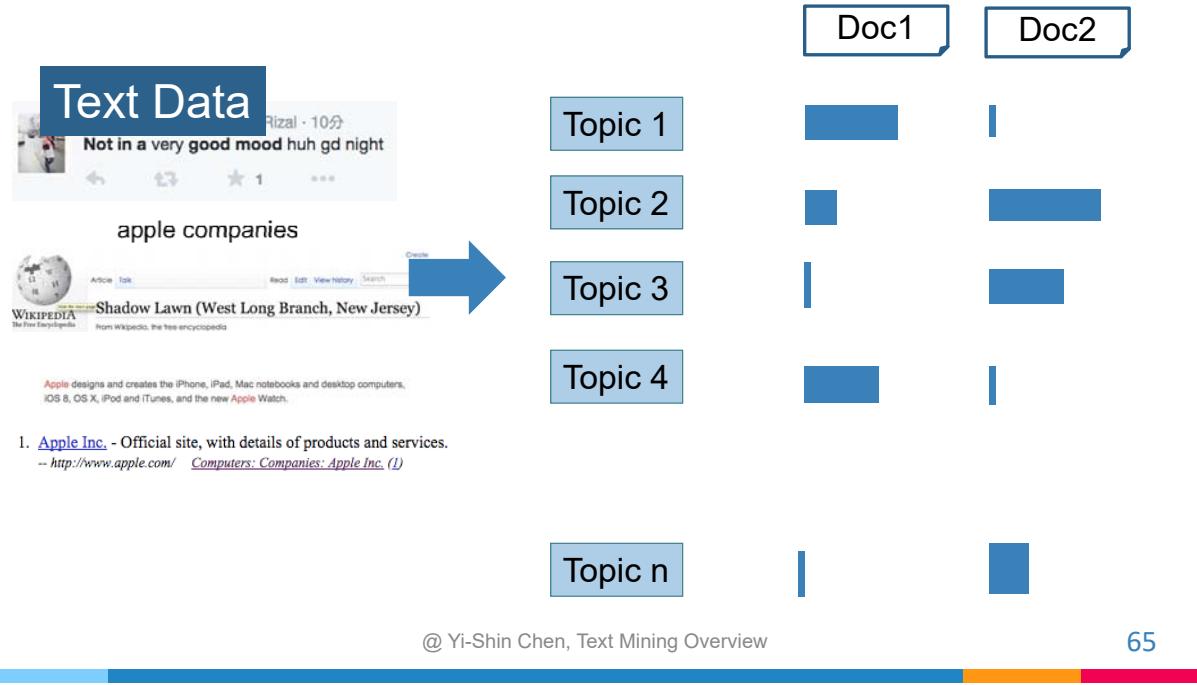
# Landscape of Text Mining



## Topic Mining: Motivation

- ▷ Topic: key idea in text data
  - Theme/subject
  - Different granularities (e.g., sentence, article)
- ▷ Motivated applications, e.g.:
  - Hot topics during the debates in 2016 presidential election
  - What do people like about Windows 10
  - What are Facebook users talking about today?
  - What are the most watched news?

# Tasks of Topic Mining



## Formal Definition of Topic Mining

### ▷ Input

- A collection of **N** text documents  $S = \{d_1, d_2, d_3, \dots d_n\}$
- Number of topics: **k**

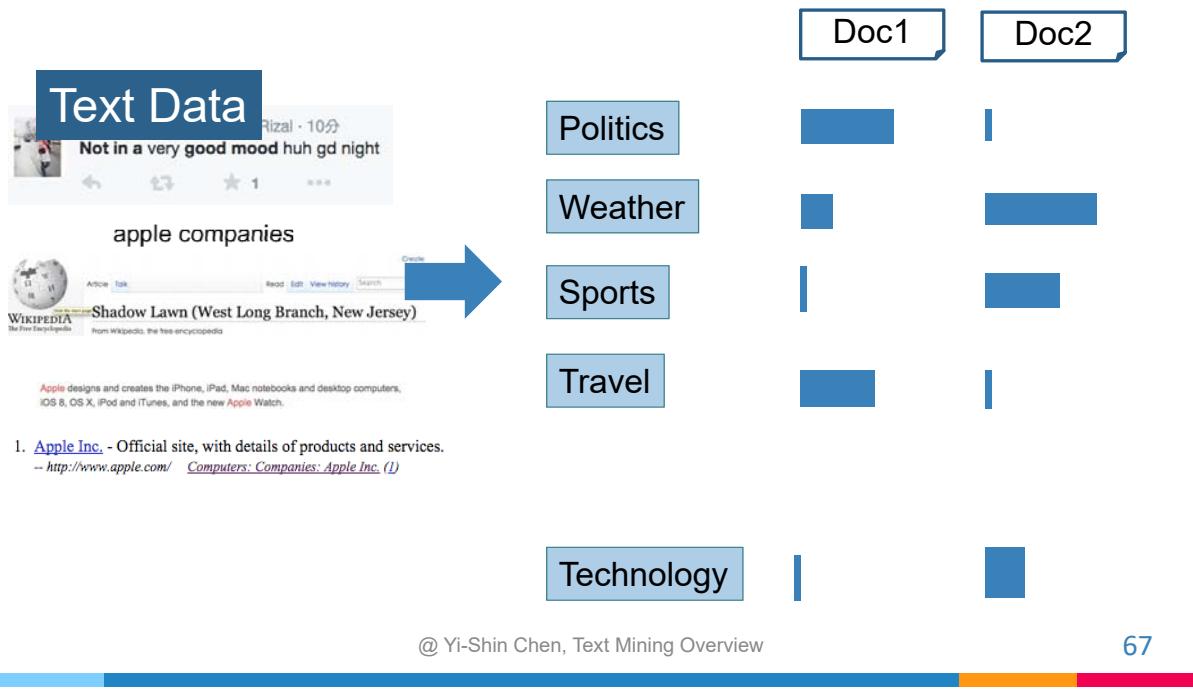
### ▷ Output

- k topics:  $\{\theta_1, \theta_2, \theta_3, \dots \theta_n\}$
- Coverage of topics in each  $d_i$ :  $\{\mu_{i1}, \mu_{i2}, \mu_{i3}, \dots \mu_{in}\}$

### ▷ How to define topic $\theta_i$ ?

- Topic=term (word)?
- Topic= classes?

# Tasks of Topic Mining (Terms as Topics)



## Problems with “Terms as Topics”

- ▷ Not generic
  - Can only represent simple/general topic
  - Cannot represent complicated topics  
→ E.g., “uber issue”: political or transportation related?
- ▷ Incompleteness in coverage
  - Cannot capture variation of vocabulary
- ▷ Word sense ambiguity
  - E.g., Hollywood star vs. stars in the sky; apple watch vs. apple recipes

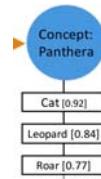
# Improved Ideas

▷ Idea1 (Probabilistic topic models): topic = word distribution

- E.g.: Sports = {(Sports, 0.2), (Game 0.01), (basketball 0.005),  
(play, 0.003), (NBA,0.01)...}
- ✓: generic, easy to implement

▷ Idea 2 (Concept topic models): topic = concept

- Maintain concepts (manually or automatically)  
→ E.g., ConceptNet



## Possible Approaches for Probabilistic Topic Models

▷ Bag-of-words approach:

- Mixture of unigram language model
- Expectation-maximization algorithm
- Probabilistic latent semantic analysis
- Latent Dirichlet allocation (LDA) model

▷ Graph-based approach :

- TextRank (Mihalcea and Tarau, 2004)
- Reinforcement Approach (Xiaojun et al., 2007)
- CollabRank (Xiaojun er al., 2008)

# Bag-of-words Assumption

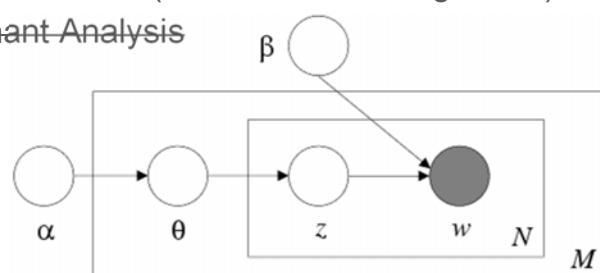
- ▷ Word order is *ignored*
- ▷ “bag-of-words” – exchangeability
- ▷ **Theorem (De Finetti, 1935)** – if  $(x_1, x_2, \dots, x_n)$  are infinitely exchangeable, then the joint probability  $p(x_1, x_2, \dots, x_n)$  has a representation as a mixture:
- ▷  $p(x_1, x_2, \dots, x_n) = \int d\theta p(\theta) \prod_{i=1}^N p(x_i|\theta)$   
for some random variable  $\theta$

@ Yi-Shin Chen, Text Mining Overview

71

# Latent Dirichlet Allocation

- ▷ Latent Dirichlet Allocation (D. M. Blei, A. Y. Ng, 2003)  
~~Linear Discriminant Analysis~~



$$(2) p(\theta, z, w | \alpha, \beta) = p(\theta | \alpha) \prod_{n=1}^N p(z_n | \theta) p(w_n | z_n, \beta)$$

$$(3) p(w | \alpha, \beta) = \int p(\theta | \alpha) \left( \prod_{n=1}^N \sum_{z_n} p(z_n | \theta) p(w_n | z_n, \beta) \right) d^k \theta$$

$$p(D | \alpha, \beta) = \prod_{d=1}^M \int p(\theta_d | \alpha) \left( \prod_{n=1}^{N_d} \sum_{z_{dn}} p(z_{dn} | \theta_d) p(w_{dn} | z_{dn}, \beta) \right) d^k \theta_d$$

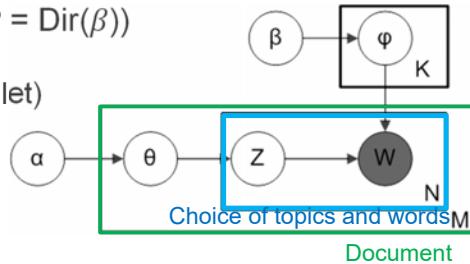
@ Yi-Shin Chen, Text Mining Overview

72

# LDA Assumption

▷ Assume:

- When writing a document, you
  1. Decide how many words
  2. Decide distribution( $P = \text{Dir}(\alpha)$ )  $P = \text{Dir}(\beta)$ )
  3. Choose topics (Dirichlet)
  4. Choose words for topics (Dirichlet)
  5. Repeat 3
- Example
  1. 5 words in document
  2. 50% food & 50% cute animals
  3. 1st word - food topic, gives you the word "bread".
  4. 2nd word - cute animals topic, "adorable".
  5. 3rd word - cute animals topic, "dog".
  6. 4th word - food topic, "eating".
  7. 5th word - food topic, "banana".    "bread adorable dog eating banana"



@ Yi-Shin Chen, Text Mining Overview

73

# LDA Learning (Gibbs)

- ▷ How many topics you think there are ?
- ▷ Randomly assign words to topics
- ▷ Check and update topic assignments (**Iterative**)
  - $p(\text{topic } t \mid \text{document } d)$
  - $p(\text{word } w \mid \text{topic } t)$
  - Reassign  $w$  a new topic,  $p(\text{topic } t \mid \text{document } d) * p(\text{word } w \mid \text{topic } t)$

I eat fish and vegetables.  
Dog and fish are pets.  
My kitten eats fish.

#Topic: 2

$p(\text{red}|1)=0.67; p(\text{purple}|1)=0.33; p(\text{red}|2)=0.50; p(\text{purple}|2)=0.50; p(\text{red}|3)=0.67; p(\text{purple}|3)=0.33;$

$p(\text{eat}|\text{red})=0.20; p(\text{eat}|\text{purple})=0.33; p(\text{fish}|\text{red})=0.20; p(\text{fish}|\text{purple})=0.33;$   
 $p(\text{vegetable}|\text{red})=0.20; p(\text{dog}|\text{purple})=0.33; p(\text{pet}|\text{red})=0.20; p(\text{kitten}|\text{red})=0.20;$

$p(\text{purple}|2)*p(\text{fish}|\text{purple})=0.5*0.33=0.165; p(\text{red}|2)*p(\text{fish}|\text{red})=0.5*0.2=0.1;$

@ Yi-Shin Chen, Text Mining Overview

74

## Related Work – Topic Model (LDA)

- ▷ I eat fish and vegetables.
- ▷ Dog and fish are pets.
- ▷ My kitten eats fish.

LDA

Topic 1	
0.268	fish
0.210	pet
0.210	dog
0.147	kitten

- Sentence 1:** 14.67% Topic 1, 85.33% Topic 2  
**Sentence 2:** 85.44% Topic 1, 14.56% Topic 2  
**Sentence 3:** 19.95% Topic 1, 80.05% Topic 2

Topic 2	
0.296	eat
0.265	fish
0.189	vegetable
0.121	kitten

@ Yi-Shin Chen, Text Mining Overview

75

## Possible Approaches for Probabilistic Topic Models

- ▷ Bag-of-words approach:
  - Mixture of unigram language model
  - Expectation-maximization algorithm
  - Probabilistic latent semantic analysis
  - Latent Dirichlet allocation (LDA) model

- ▷ Graph-based approach :
  - TextRank (Mihalcea and Tarau, 2004)
  - Reinforcement Approach (Xiaojun et al., 2007)
  - CollabRank (Xiaojun er al., 2008)

@ Yi-Shin Chen, Text Mining Overview

76

# Construct Graph

- ▷ Directed graph
- ▷ Elements in the graph
  - Terms
  - Phrases
  - Sentences

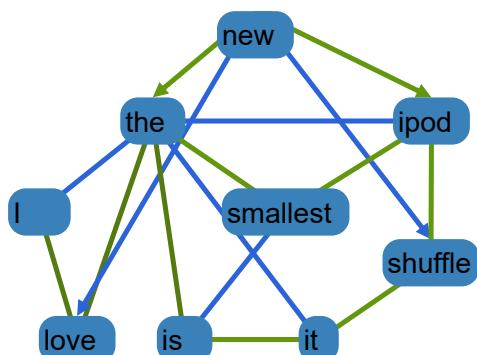
@ Yi-Shin Chen, Text Mining Overview

77

## Connect Term Nodes

- ▷ Connect terms based on its slop.

I love the new ipod shuffle.  
It is the smallest ipod.



@ Yi-Shin Chen, Text Mining Overview

78

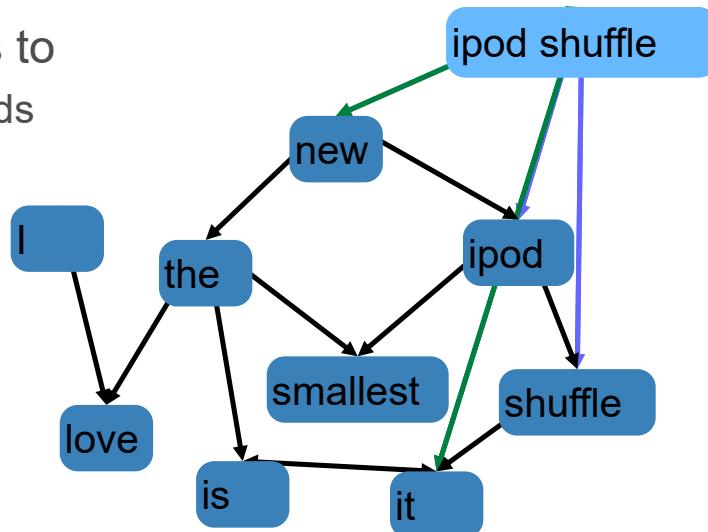
## Connect Phrase Nodes

### ▷ Connect phrases to

- Compound words
- Neighbor words

I love the new  
ipod shuffle.

It is the smallest  
ipod.



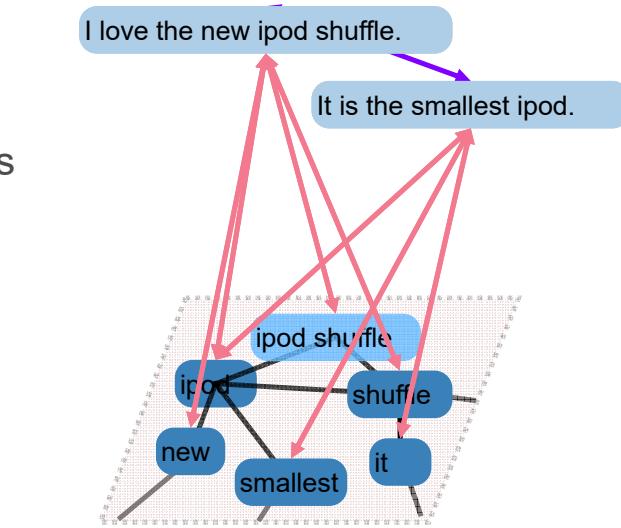
@ Yi-Shin Chen, Text Mining Overview

79

## Connect Sentence Nodes

### ▷ Connect to

- Neighbor sentences
- Compound terms
- Compound phrase

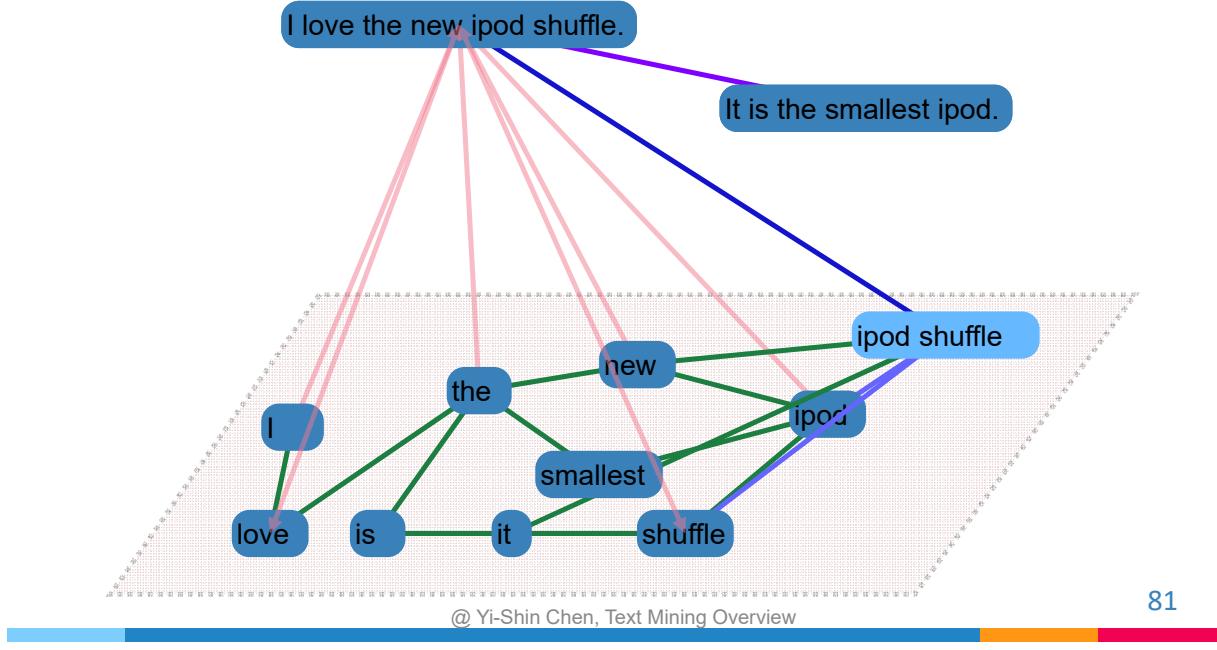


I love the new ipod shuffle.  
It is the smallest ipod.

@ Yi-Shin Chen, Text Mining Overview

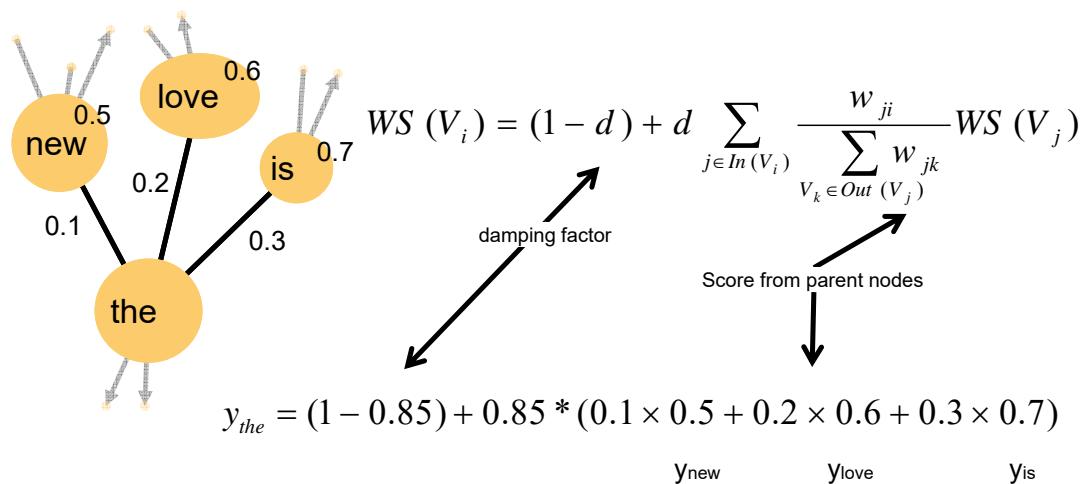
80

# Edge Weight Types



# Graph-base Ranking

▷ Scores for each node (TextRank 2004)



# Result

the	1.45
ipod	1.21
new	1.02
is	1.00
shuffle	0.99
it	0.98
smallest	0.77
love	0.57

@ Yi-Shin Chen, Text Mining Overview

83



# Graph-based Extraction

- Pros
  - Structure and syntax information
  - Mutual influence
- Cons
  - Common words get higher scores

@ Yi-Shin Chen, Text Mining Overview

84



# Summary: Probabilistic Topic Models

▷ Probabilistic topic models): topic = word distribution

- E.g.: Sports = {(Sports, 0.2), (Game 0.01), (basketball 0.005),  
(play, 0.003), (NBA,0.01)...}
- ✓: generic, easy to implement
- ?: Not easy to understand/communicate
- ?: Not easy to construct semantic relationship between topics



Topic= {(Crooked, 0.02), (dishonest, 0.001), (News, 0.0008), (totally, 0.0009), (total, 0.000009), (failed, 0.0006), (bad, 0.0015), (failing, 0.00001), (presidential, 0.000008), (States, 0.000004), (terrible, 0.000085),(failed, 0.000021), (lightweight,0.00001),(weak, 0.000075), .....}

@ Yi-Shin Chen, Text Mining Overview

85

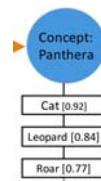
## Improved Ideas

▷ Idea1 (Probabilistic topic models): topic = word distribution

- E.g.: Sports = {(Sports, 0.2), (Game 0.01), (basketball 0.005),  
(play, 0.003), (NBA,0.01)...}
- ✓: generic, easy to implement

▷ Idea 2 (Concept topic models): topic = concept

- Maintain concepts (manually or automatically)  
→ E.g., ConceptNet



@ Yi-Shin Chen, Text Mining Overview

86

## NLP Related Approach: Named Entity Recognition

- ▷ Find and classify all the named entities in a text.
- ▷ What's a named entity?
  - A mention of an entity using its name.
    - *Kansas Jayhawks*
  - This is a subset of the possible mentions...
    - *Kansas, Jayhawks, the team, it, they*
- ▷ Find means identify the exact span of the mention
- ▷ Classify means determine the category of the entity being referred to

87

## Named Entity Recognition Approaches

- ▷ As with partial parsing and chunking there are two basic approaches (and hybrids)
  - Rule-based (regular expressions)
    - Lists of names
    - Patterns to match things that look like names
    - Patterns to match the environments that classes of names tend to occur in.
  - Machine Learning-based approaches
    - Get annotated training data
    - Extract features
    - Train systems to replicate the annotation

88

# Rule-Based Approaches

- ▷ Employ regular expressions to extract data
- ▷ Examples:
  - Telephone number:  $(\d\{3\}[-.\s]\{1,2\})\d\{4\}$ .
    - 800-865-1125
    - 800.865.1125
    - (800)865-CARE
  - Software name extraction:  $([A-Z][a-z]^{*}\s^{*})^{+}$ 
    - Installation Designer v1.1

A horizontal progress bar consisting of four colored segments: light blue, dark blue, orange, and red. The red segment is the shortest.

89

# Relations

- ▷ Once you have captured the entities in a text you might want to ascertain how they relate to one another.
  - Here we're just talking about explicitly stated relations

A horizontal progress bar consisting of four colored segments: light blue, dark blue, orange, and red. The red segment is the shortest.

90

# Relation Types

- ▷ As with named entities, the list of relations is application specific. For generic news texts...

Relations	Examples	Types
Affiliations	Personal	<i>married to, mother of</i>
	Organizational	<i>spokesman for, president of</i>
	Artifactual	<i>owns, invented, produces</i>
Geospatial	Proximity	<i>near, on outskirts</i>
	Directional	<i>southeast of</i>
Part-Of	Organizational	<i>a unit of, parent of</i>
	Political	<i>annexed, acquired</i>

91

# Bootstrapping Approaches

- ▷ What if you don't have enough annotated text to train on.
- But you might have some seed tuples
  - Or you might have some patterns that work pretty well
- ▷ Can you use those seeds to do something useful?
- Co-training and active learning use the seeds to train classifiers to tag more data to train better classifiers...
  - Bootstrapping tries to learn directly (populate a relation) through direct use of the seeds

92

## Bootstrapping Example: Seed Tuple

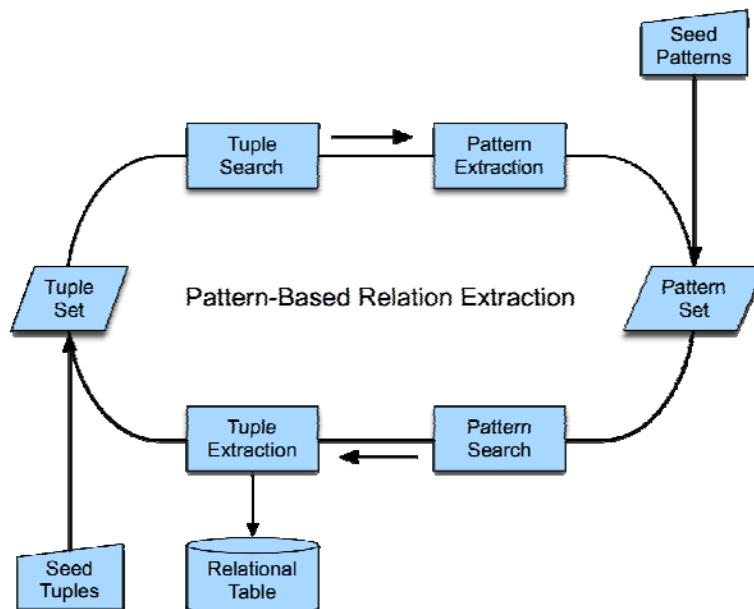
▷ <Mark Twain, Elmira> Seed tuple

- Grep (google)
- “Mark Twain is buried in Elmira, NY.”  
→ X is buried in Y
- “The grave of Mark Twain is in Elmira”  
→ The grave of X is in Y
- “Elmira is Mark Twain’s final resting place”  
→ Y is X’s final resting place.

▷ Use those patterns to grep for new tuples that you don’t already know

93

## Bootstrapping Relations



94

# Wikipedia Infobox

- ▷ Infoboxes are kept in a namespace separate from articles
  - Namespce example: Special:SpecialPages; Wikipedia>List of infoboxes
  - Example:
  -

```
{ {Infobox person
|name = Casanova
|image = Casanova_self_portrait.jpg
|caption = A self portrait of Casanova
...
|website = }}
```

95

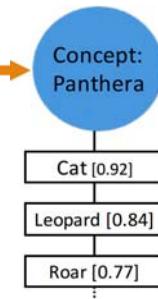
## Concept-based Model

- ▷ ESA (Egozi, Markovitch, 2011)  
Every Wikipedia article represents a concept

From Wikipedia, the free encyclopedia

*Panthera* is a genus of the family Felidae (the **cats**, which contains four well-known living **species**: the lion, tiger, jaguar, and **leopard**). The genus comprises about half of the big **cats**. One meaning of the word **panther** is to designate **cats** of this family. Only these four **cat** species have the anatomical changes enabling them to **roar**. The primary reason for this was assumed to be the incomplete **ossification** of the hyoid bone. However, new studies show that the ability to **roar** is due to other morphological features, especially of the larynx. The snow leopard (*Uncia uncia*), which is sometimes included within *Panthera*, does not **roar**. Although it has an incomplete ossification of the hyoid bone, it lacks the special morphology of the larynx, which is typical for lions, tigers, jaguars and leopards.<sup>[1]</sup>

Species and subspecies



Manually-maintained knowledge base

# Yago

- ▷ YAGO: A Core of Semantic Knowledge Unifying WordNet and Wikipedia, WWW 2007
- ▷ Unification of Wikipedia & WordNet
- ▷ Make use of rich structures and information
  - Infoboxes, Category Pages, etc.



@ Yi-Shin Chen, Text Mining Overview

97

## Mining Concepts from User-Generated Web Data

- ▷ Concept: in a word sequence, its meaning is known by a group of people and everyone within the group refers to the same meaning



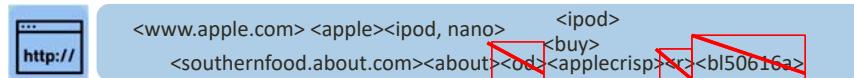
Pei-Ling Hsu, Hsiao-Shan Hsieh, Jheng-He Liang, and Yi-Shin Chen\*, Mining Various Semantic Relationships from Unstructured User-Generated Web Data, Journal of Web Semantics, 2015

@ Yi-Shin Chen, Text Mining Overview

98

# Concepts in Word Sequences

- Word sequences are likely meaningless and noise



→ A word sequence as a candidate to be a concept

→ About noun

- A noun e.g., Basketball
- A sequence of noun e.g., Christmas Eve
- A noun and a number e.g., 311 earthquake
- An adjective and a noun e.g., Desperate housewives

→ Special format

- A word sequence contains "of" e.g., Cover of harry potter

# Concept Modeling

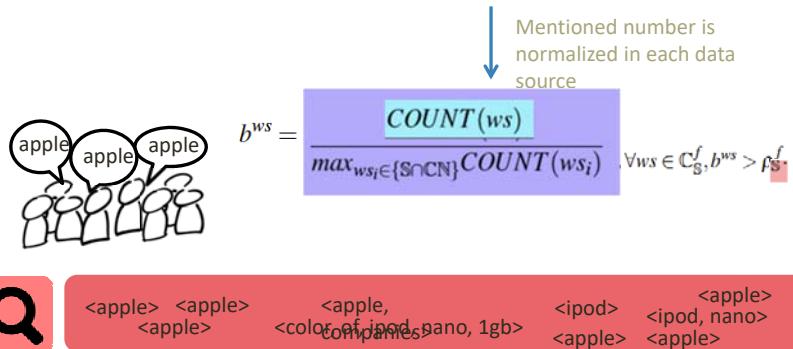
- Concept: in a word sequence, **its meaning is known by a group of people** and everyone within the group refers to the same meaning.

→ Try to detect a word sequence that is known by a group of people.

# Concept Modeling

- Frequent Concept:

- A word sequence mentioned frequently from a single source.



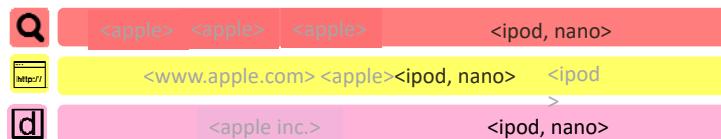
@ Yi-Shin Chen, Text Mining Overview

101

# Concept Modeling

- ▷ Some data sources are public, sharing data with other users.
- ▷ These data sources with frequently seen word sequences.
  - These data sources provide concepts with higher confidence.
- ▷ Confident Value: every data source has a confident value
- ▷ Confident concept: a word sequence is from data sources with higher confident values.

$$ac_{ws} = \sum_{S_k \in SS, ws \in (S_k \cap CN)} (cv_k), \forall ws \in \mathbb{C}^c, ac_{ws} > \rho^c.$$



@ Yi-Shin Chen, Text Mining Overview

102

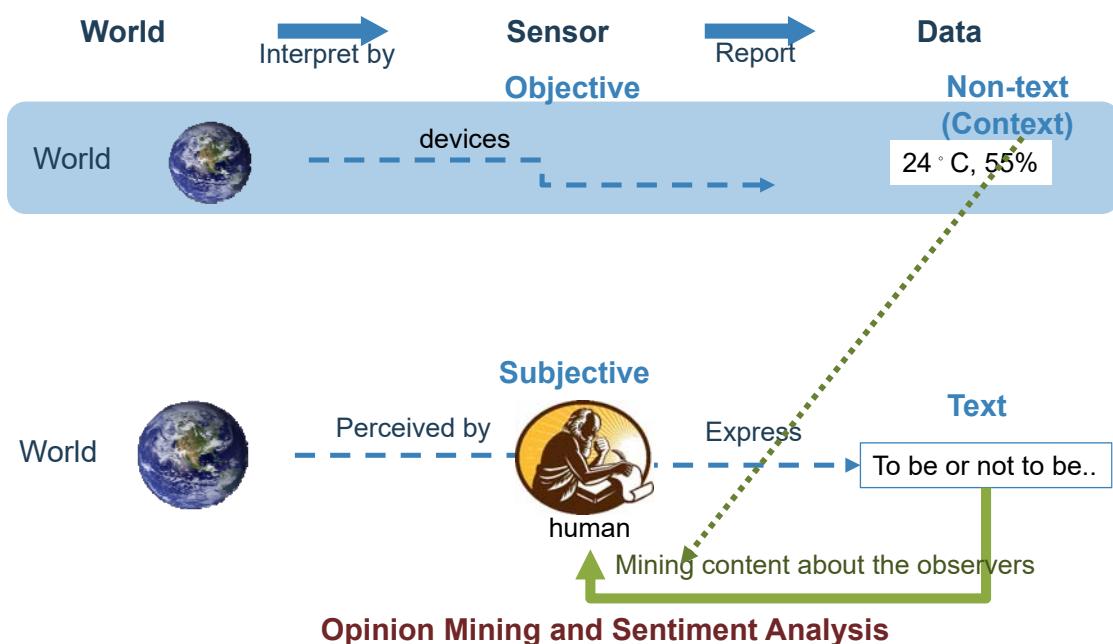
# Opinion Mining

How people feel?

@ Yi-Shin Chen, Text Mining Overview

103

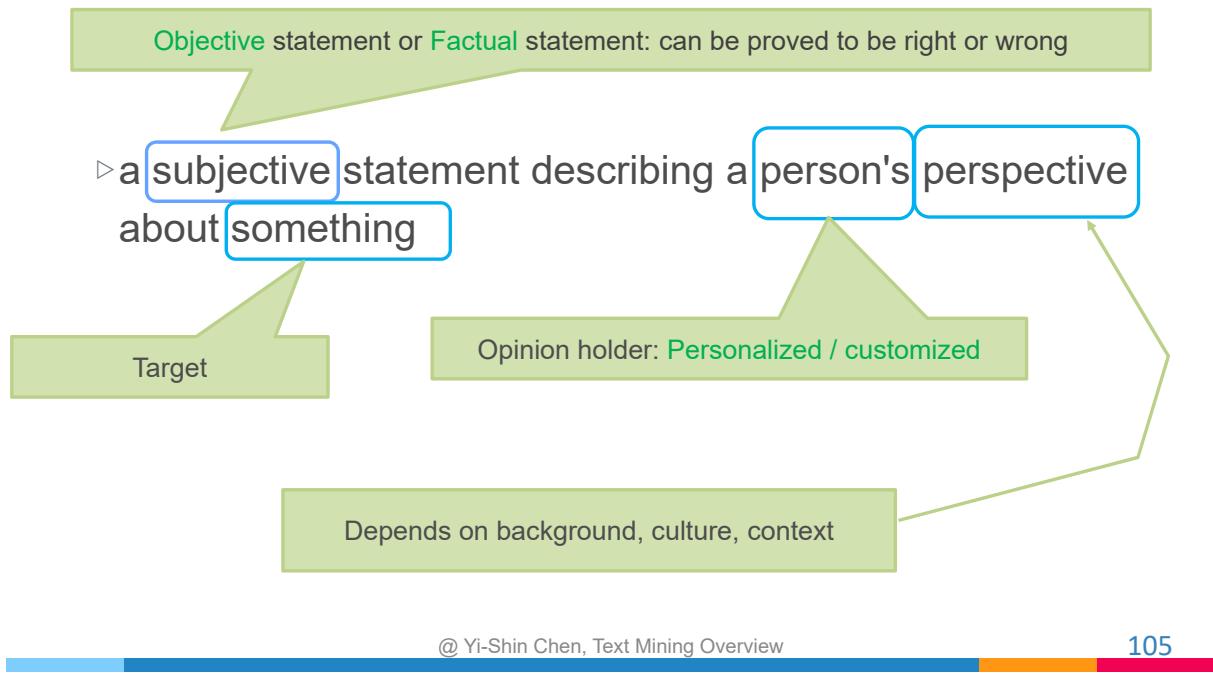
## Landscape of Text Mining



@ Yi-Shin Chen, Text Mining Overview

104

# Opinion



@ Yi-Shin Chen, Text Mining Overview

105

# Opinion Representation

- ▷ Opinion holder: user
- ▷ Opinion target: object
- ▷ Opinion content: keywords?
- ▷ Opinion context: time, location, others?
- ▷ Opinion sentiment (emotion): positive/negative, happy or sad

@ Yi-Shin Chen, Text Mining Overview

106

# Sentiment Analysis

- ▷ Input: An opinionated text object
- ▷ Output: Sentiment tag/Emotion label
  - Polarity analysis: {positive, negative, neutral}
  - Emotion analysis: happy, sad, anger
- ▷ Naive approach:
  - Apply classification, clustering for extracted text features

# Text Features

- ▷ Character  $n$ -grams
  - Usually for spelling/recognition proof
  - Less meaningful
- ▷ Word  $n$ -grams
  - $n$  should be bigger than 1 for sentiment analysis
- ▷ POS tag  $n$ -grams
  - Can mixed with words and POS tags
    - E.g., “adj noun”, “sad noun”

# More Text Features

## ▷ Word classes

- Thesaurus: LIWC
- Ontology: WordNet, Yago, DBpedia
- Recognized entities: DBpedia, Yago

## ▷ Frequent patterns in text

- Could utilize pattern discovery algorithms
- Optimizing the tradeoff between *coverage* and *specificity* is essential

@ Yi-Shin Chen, Text Mining Overview

109

# LIWC

## ▷ Linguistic Inquiry and word count

- LIWC2015

## ▷ Home page: <http://liwc.wpengine.com/>

## ▷ >70 classes

▷ Developed by researchers with interests in social, clinical, health, and cognitive psychology

## ▷ Cost: US\$89.95



110

# Emotion Analysis: Pattern Approach

- ▷ Carlos Argueta, Fernando Calderon, and Yi-Shin Chen,  
Multilingual Emotion Classifier using Unsupervised Pattern  
Extraction from Microblog Data, Intelligent Data Analysis - An  
International Journal, 2016

@ Yi-Shin Chen, Text Mining Overview

111

## Collect Emotion Data

@ Yi-Shin Chen, Text Mining Overview

112

Mike Duggan (WMAS)	Stephen Parry	Violet Blue	Donna Fry
@wmasmikeduggan · 14小時 Another day another assault on staff! Female paramedic assaulted by a male now has a possible fractured wrist. #disgusted & #angry	@StephenParry80 · 15小時 Thanks for the last few weeks, Labour. Unprecedented situation of no opposition party for last 2-3 weeks. #angry	@violetblue · 7月11日 Discovered today that enough people are so far behind on paying me that I almost bounced my rent check this month. It's a first. #angry	@Ds fry67 · 7月10日 If you're going to camp & light fires, learn how to extinguish before leaving the fire unattended or STAY HOME! #angry #ColdSpringsFire

# Collect Emotion Data

Twitter 主頁 關於 搜尋

PUGGY @PuggyBand - 7月11日  
On our way to @FestivalEteQc canada #happy

jussi69official @jussi69official - 7月11日  
Thank you so much for all the birthday wishes!!!! I love you all!!! 🎉❤️ #happy  
instagram.com/p/BHtylhhDyV/

Sheena @m\_sheena\_ - 7月10日  
To each their own but I'm glad I'm not someone who needs to party every weekend. I love my quiet life. #Happy

HRH Mordi Joel (M.J) @MordiOfficial - 7月10日  
To everyone going through any struggle or pain right now, keep your head up, and keep on smiling, you'll get through it soon. #Happy

@ Yi-Shin Chen, Text Mining Overview

113

# Not-Emotion Data

Twitter 主頁 請助 @latimes 搜尋 已經擁有帳戶？登入 +

新加入 Twitter？  
立即註冊，取得你的個人化時間軸！

註冊

熱門  
#3YearsWithoutCory  
7,128 推文  
#ChevTheCola  
前一小時開始的流行趨勢  
#ArunachalPradesh  
前一小時開始的流行趨勢  
#WednesdayWisdom  
前一小時開始的流行趨勢  
#AllStarGame  
7,236 推文  
Owen Smith  
6,760 推文  
Jennifer Aniston  
3,676 推文  
Miss Spice - Voltage Video  
793 推文  
Steinhoff International  
剛開始的流行趨勢  
NBCC  
剛開始的流行趨勢

帳戶

Los Angeles Times @latimes  
News from Los Angeles and the world. Staffed by latimes.com editors.  
Los Angeles, CA • latimes.com

已喜歡 88 次  
Los Angeles Times @latimes - 12小時  
Obama and Bush join crowd at Dallas Memorial in singing "The Battle Hymn of the Republic" lat.ms/29vb820

@ Yi-Shin Chen, Text Mining Overview

114

# Preprocessing Steps

## ▷ Hints: Remove troublesome ones

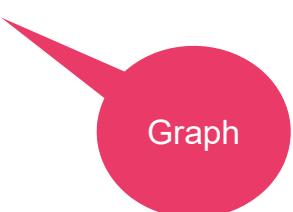
- Too short
  - Too short to get important features
- Contain too many hashtags
  - Too much information to process
- Are retweets
  - Increase the complexity
- Have URLs
  - Too trouble to collect the page data
- Convert user mentions to <usermention> and hashtags to <hashtag>
  - Remove the identification. We should not peek answers!

@ Yi-Shin Chen, Text Mining Overview

115

# Basic Guidelines

- ▷ Identify the common and differences between the experimental and control groups
  - Analyze the frequency of words
    - TF•IDF (Term frequency, inverse document frequency)
  - Analyze the co-occurrence between words/patterns
    - Co-occurrence
  - Analyze the importance between words
    - Centrality



Graph

@ Yi-Shin Chen, Text Mining Overview

116

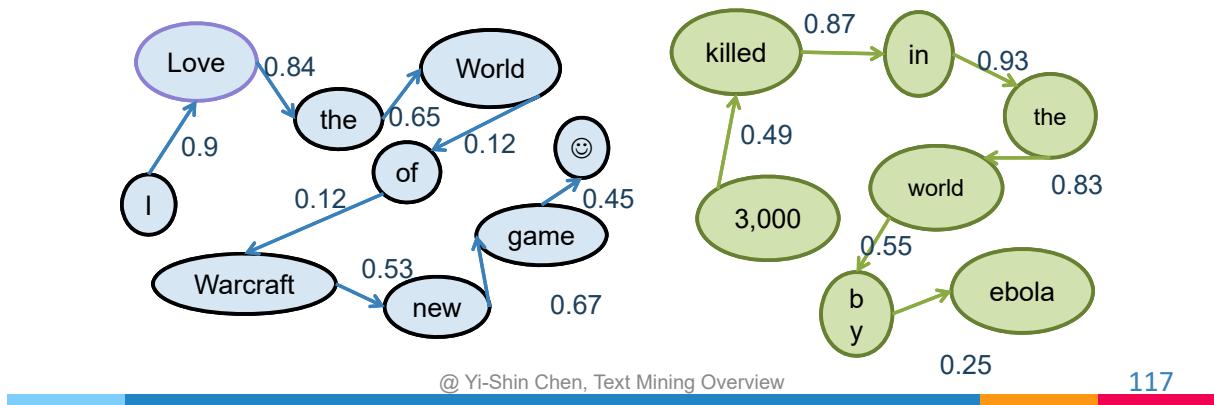
# Graph Construction

## ▷ Construct two graphs

- E.g.

→ Emotion one: I love the World of Warcraft new game 😊

→ Not-emotion one: 3,000 killed in the world by ebola



# Graph Processes

## ▷ Remove the common ones between two graphs

- Leave the significant ones only appear in the emotion graph

## ▷ Analyze the centrality of words

- Betweenness, Closeness, Eigenvector, Degree, Katz
  - Can use the free/open software, e.g., Gaphi, GraphDB

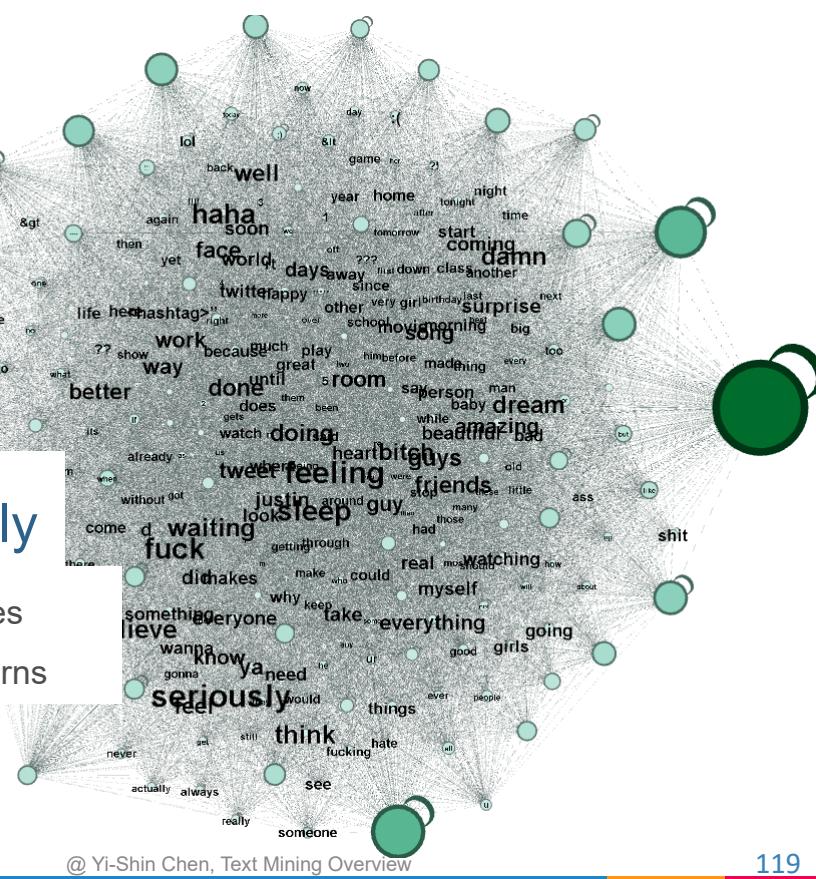
## ▷ Analyze the cluster degrees

- Clustering Coefficient

Key patterns

## Essence Only

Only key phrases  
→ emotion patterns



119

## Emotion Patterns Extraction

### ○ The goal:

- Language independent extraction – not based on grammar or manual templates
- More representative set of features - balance between generality and specificity
- High recall/coverage – adapt to unseen words
- Requiring only a relatively small number – high reliability
- Efficient—fast extraction and utilization
- Meaningful - even if there are no recognizable emotion words in it

# Patterns Definition

- Constructed from two types of elements:
  - Surface tokens: hello, 😊, lol, house, ...
  - Wildcards: \* (matches every word)
- Contains at least 2 elements
- Contains at least one of each type of element

Examples:

Pattern	Matches
* this *	“Hate this weather”, “love this drawing”
* * 😊	“so happy 😊”, “to me 😊”
luv my *	“luv my gift”, “luv my price”
* that	“want that”, “love that”, “hate that”

@ Yi-Shin Chen, Text Mining Overview

121

# Patterns Construction

- Constructed from instances
- An instance is a sequence of 2 or more words from CW and SW
- Contains at least one CW and one SW

Examples

Connector Words	SubjectWords	Instances
this	love	“hate this weather”
luv	hate	“so happy 😊”
my	gift	“luv my gift”
😊	weather	“love this drawing”
...	...	“luv my price”

@ Yi-Shin Chen, Text Mining Overview

122

Connector Words
this
luv
my
☺
...

## Patterns Construction (2)

- Find all instances in a corpus with their frequency
- Aggregate counts by grouping them based on length and position of matching CW

Instances	Count
"hate this weather"	5
"so happy ☺"	4
"luv my gift"	7
"love this drawing"	2
"luv my price"	1
"to me ☺"	3
"kill this idiot"	1
"finish this task"	4

Groups	Count
"Hate this weather", "love this drawing", "kill this idiot", "finish this task"	12
"so happy ☺", "to me ☺"	7
"luv my gift", "luv my price"	8
...	...

@ Yi-Shin Chen, Text Mining Overview

123

Connector Words
this
got
my
pain
...

## Patterns Construction (3)

- Replace all the SWs by a wildcard \* and keep the CWs to convert all instances into the representing pattern

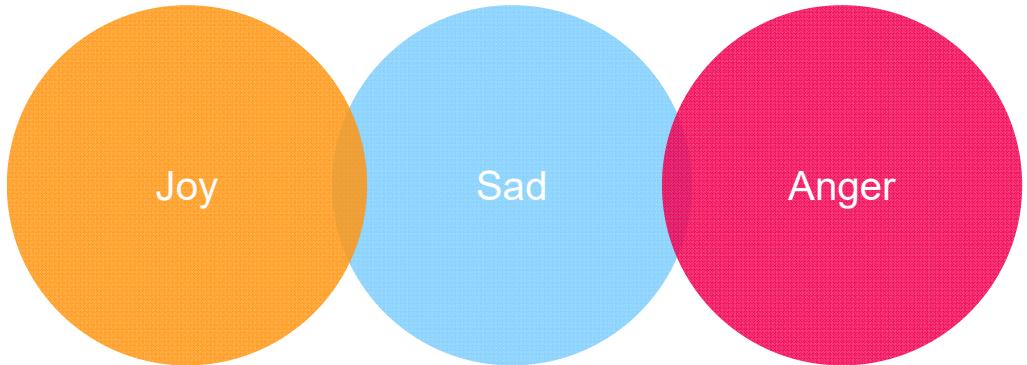
Pattern	Groups	Count
* this *	"Hate this weather", "love this drawing", "kill this idiot", "finish this task"	12
* * ☺	"so happy ☺", "to me ☺"	7
luv my *	"luv my gift", "luv my price"	8
...	...	...

- The wildcard matches any word and is used for term generalization
- Infrequent patterns are filtered out

@ Yi-Shin Chen, Text Mining Overview

124

## Ranking Emotion Patterns



- ▷ Ranking the emotion patterns for each emotion
  - Frequency, exclusiveness, diversity
  - One ranked list for each emotion

@ Yi-Shin Chen, Text Mining Overview

125

## Contextual Text Mining

Basic Concepts

# Context

- ▷ Text usually has rich context information
  - Direct context (meta-data): time, location, author
  - Indirect context: social networks of authors, other text related to the same source
  - Any other related text
- ▷ Context could be used for:
  - Partition the data
  - Provide extra features

@ Yi-Shin Chen, Text Mining Overview

127

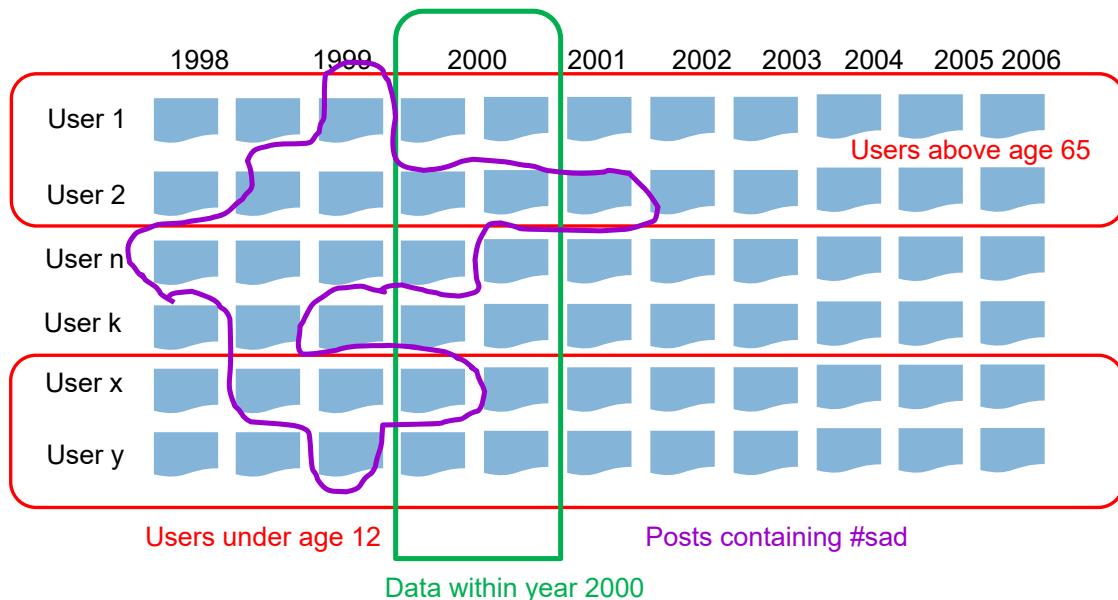
# Contextual Text Mining

- ▷ Query log + **User** = Personalized search
- ▷ Tweet + **Time** = Event identification
- ▷ Tweet + **Location-related patterns** = Location identification
- ▷ Tweet + **Sentiment** = Opinion mining
- ▷ Text Mining +Context → Contextual Text Mining

@ Yi-Shin Chen, Text Mining Overview

128

# Partition Text

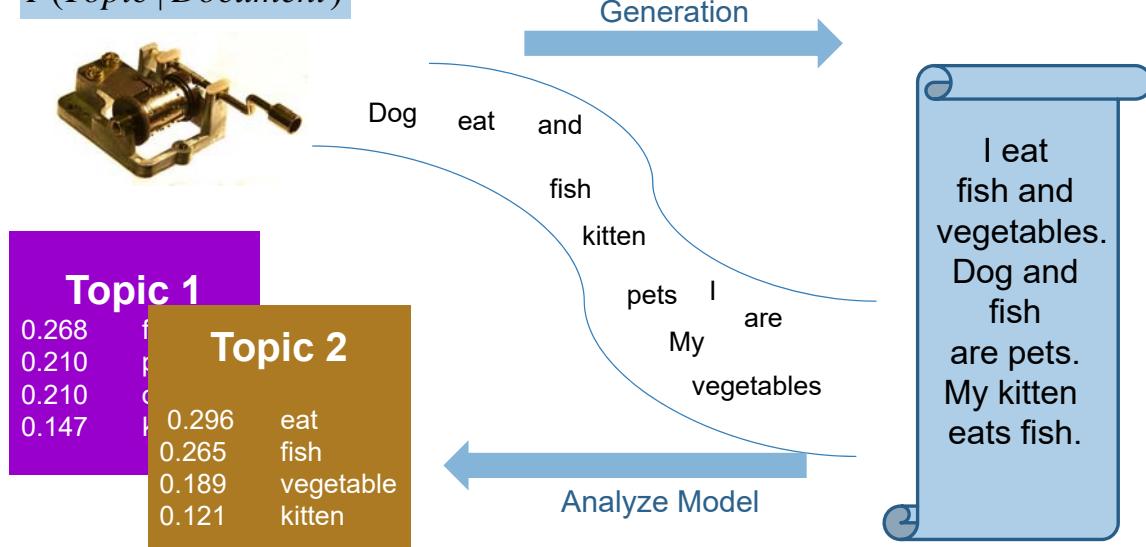


@ Yi-Shin Chen, Text Mining Overview

129

# Generative Model of Text

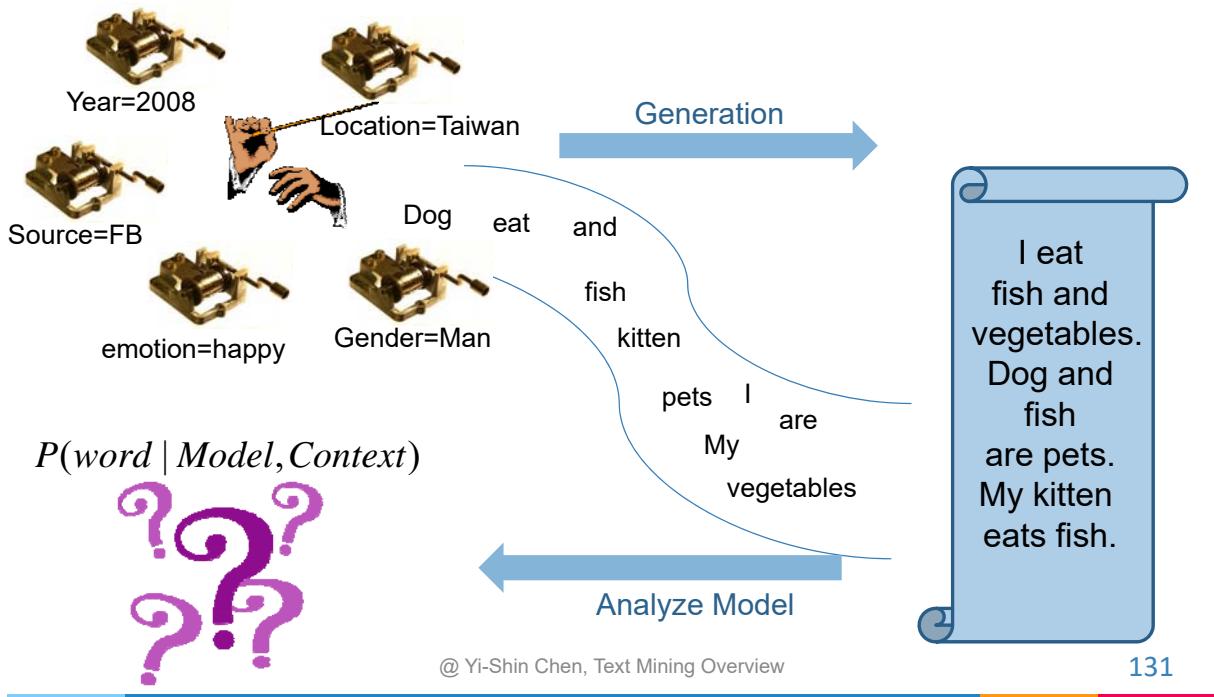
$$P(\text{word} \mid \text{Topic})$$
$$P(\text{Topic} \mid \text{Document})$$



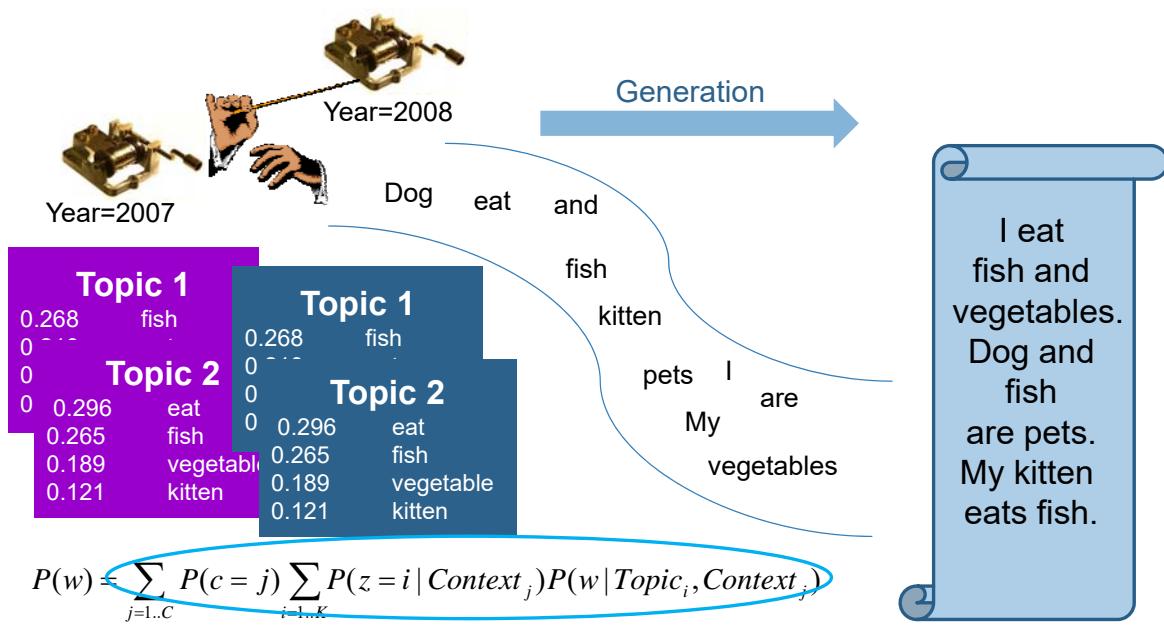
@ Yi-Shin Chen, Text Mining Overview

130

# Contextualized Models of Text



# Naïve Contextual Topic Model



How do we estimate it? → Different approaches for different contextual data and problems 132

# Contextual Probabilistic Latent Semantic Analysis (CPLAS) (Mei, Zhai, KDD2006)

▷ An extension of PLSA model ([Hofmann 99]) by

- Introducing context variables
- Modeling views of topics
- Modeling coverage variations of topics

▷ Process of contextual text mining

- Instantiation of CPLSA (*context, views, coverage*)
- Fit the model to text data (EM algorithm)
- Compare a topic from different views
- Compute strength dynamics of topics from coverages
- Compute other probabilistic topic patterns

@ Yi-Shin Chen, Text Mining Overview

133

## The Probabilistic Model

- A probabilistic model explaining the generation of a document D and its context features C: if an author wants to write such a document, he will
  - Choose a view  $v_i$  according to the view distribution  $p(v_i|D, C)$
  - Choose a coverage  $\kappa_j$  according to the coverage distribution  $p(k_j|D, C)$
  - Choose a theme  $\theta_{il}$  according to the coverage  $\kappa_j$
  - Generate a word using  $\theta_{il}$
  - The likelihood of the document collection is:

$$\log p(D) = \sum_{(D,C) \in D} \sum_{w \in V} c(w, D) \log \left( \sum_{i=1}^n p(v_i | D, C) \sum_{j=1}^m p(\kappa_j | D, C) \sum_{l=1}^k p(l | \kappa_j) p(w | \theta_{il}) \right)$$

@ Yi-Shin Chen, Text Mining Overview

134

# Contextual Text Mining Example 1

## Event Identification

@ Yi-Shin Chen, Text Mining Overview

135

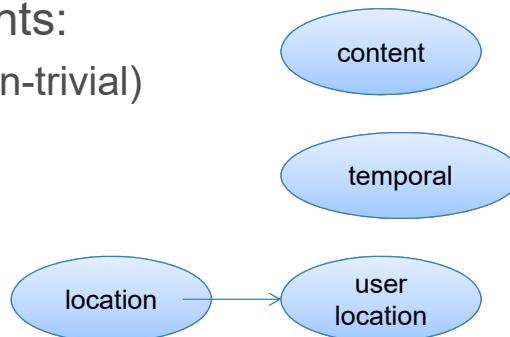
## Introduction

### ▷ Event definition:

- Something (non-trivial) happening in a certain place at a certain time (Yang et al. 1998)

### ▷ Features of events:

- Something (non-trivial)
- Certain time
- Certain place

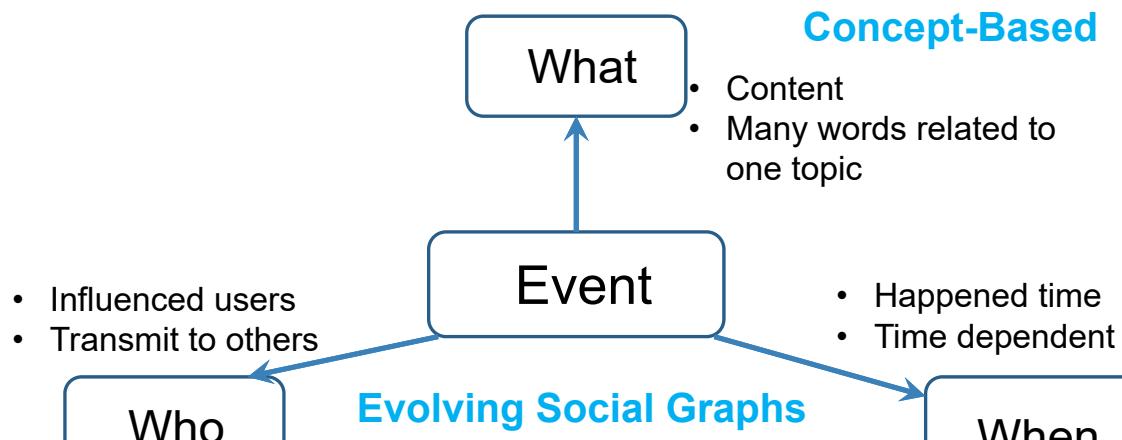


@ Yi-Shin Chen, Text Mining Overview

136

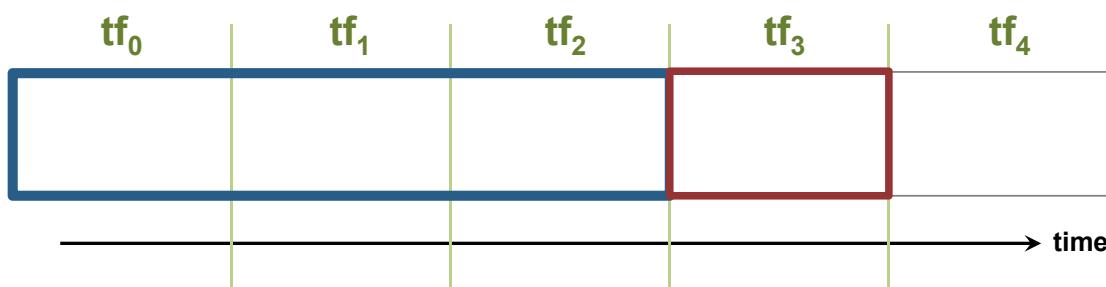
# Goal

- ▷ Identify events from social streams
- ▷ Events contain these characteristics



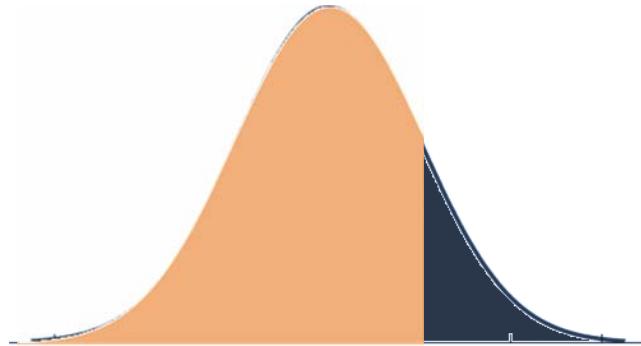
# Keyword Selection

- ▷ Well-noticed criterion
  - Compared to the past, if a word suddenly be mentioned by many users, it is well-noticed
  - Time Frame – a unit of time period
  - Sliding Window – a certain number of past time frames



# Keyword Selection

- ▷ Compare the probability proportion of this word count with the past sliding window



$$Z_{tf_j,s}(w_i) = \frac{\kappa_{tf_j}(w_i) - AVG_{tf_j,s}(w_i)}{SD_{tf_j,s}(w_i)}$$

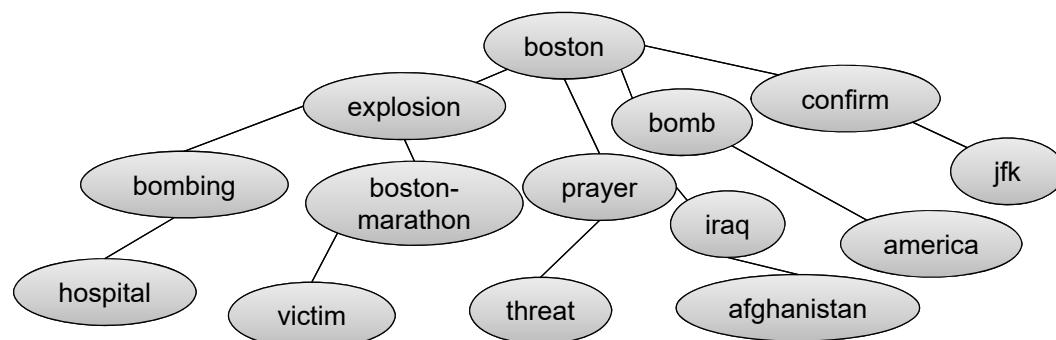
$$Prob(Z \leq Z_{tf_j,s}(w_i)) > \delta_{category}$$

139



# Event Candidate Recognition

- ▷ Concept-based event: all keywords about the same event
  - Use the co-occurrence of words in tweets to connect keywords

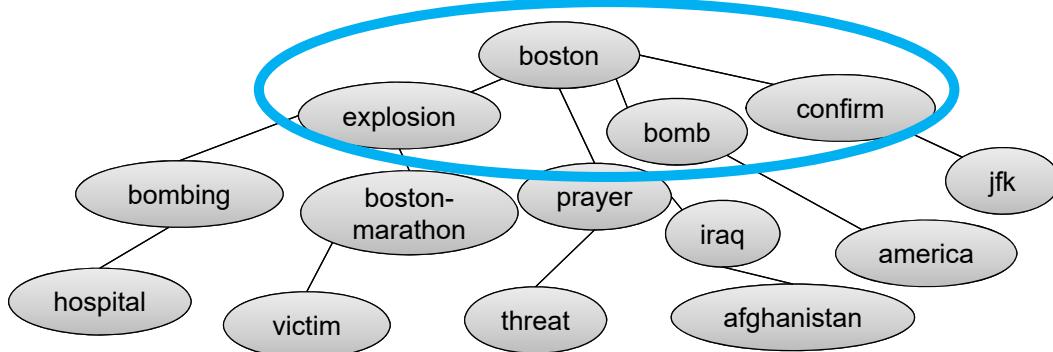


Huge amount of keywords connected as one event



## Event Candidate Recognition

- ▷ Idea: group one keyword with its most relevant keywords into one event candidate



- ▷ How do we decide which keyword to start grouping?

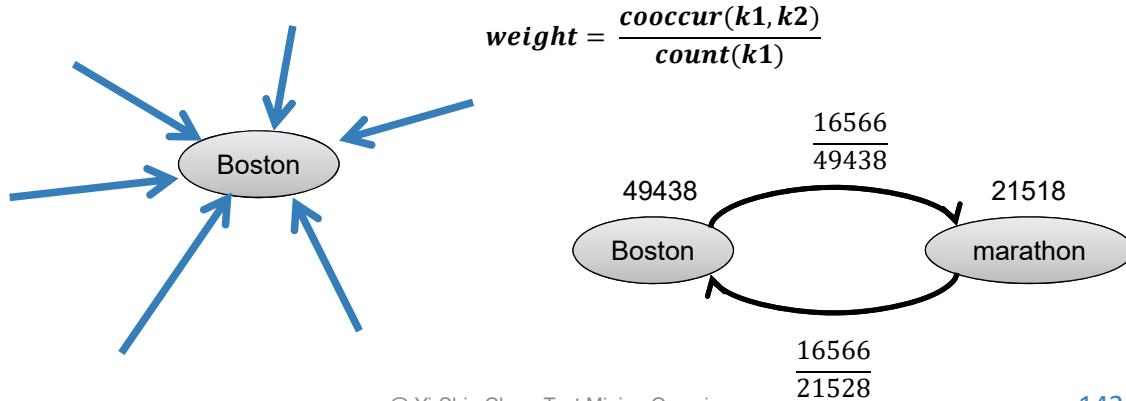
@ Yi-Shin Chen, Text Mining Overview

141

## Event Candidate Recognition

- ▷ TextRank: rank importance of each keyword

- Number of edges
- Edge weights (word count and word co-occurrence)



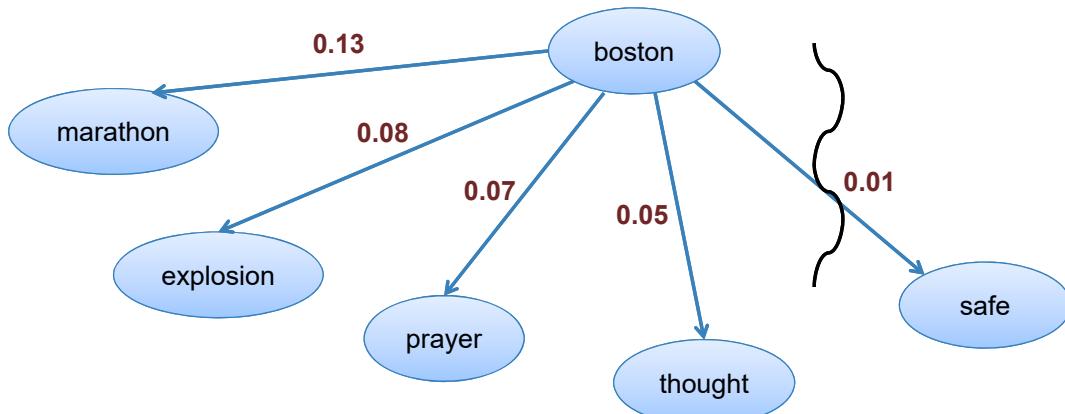
@ Yi-Shin Chen, Text Mining Overview

142

# Event Candidate Recognition

1. Pick the most relevant neighbor keywords

**Normalized weight**

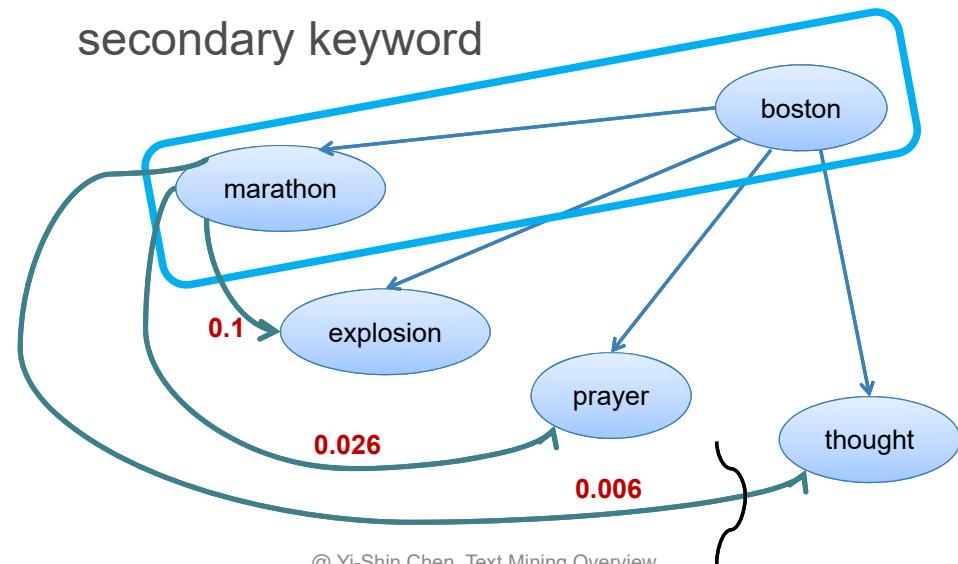


@ Yi-Shin Chen, Text Mining Overview

143

# Event Candidate Recognition

2. Check neighbor keywords with the secondary keyword

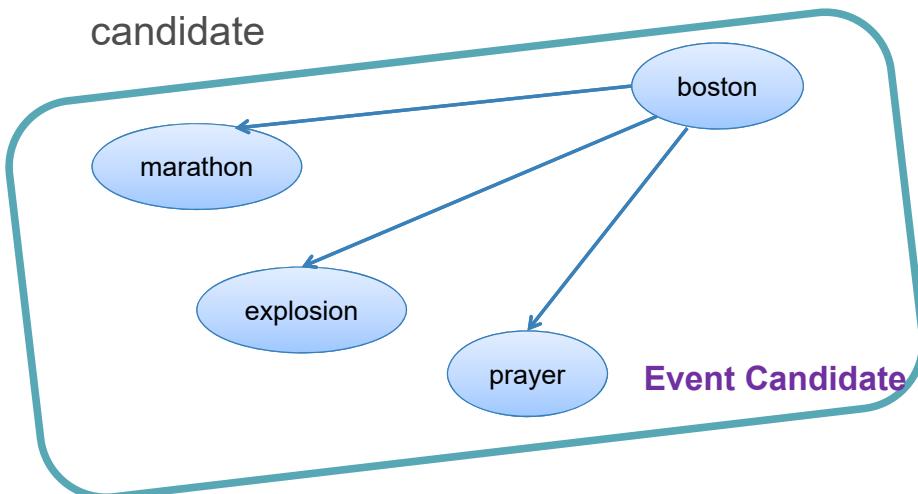


@ Yi-Shin Chen, Text Mining Overview

144

# Event Candidate Recognition

3. Group the remaining keywords as one event candidate

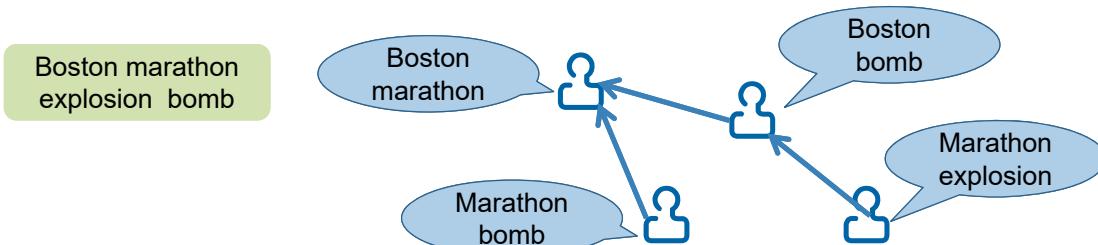


@ Yi-Shin Chen, Text Mining Overview

145

# Evolving Social Graph Analysis

- ▷ Bring in social relationships to estimate information propagation
- ▷ Social Relation Graph
  - Vertex: users that mentioned one or more keywords in the candidate event
  - Edge: the following relationship between users

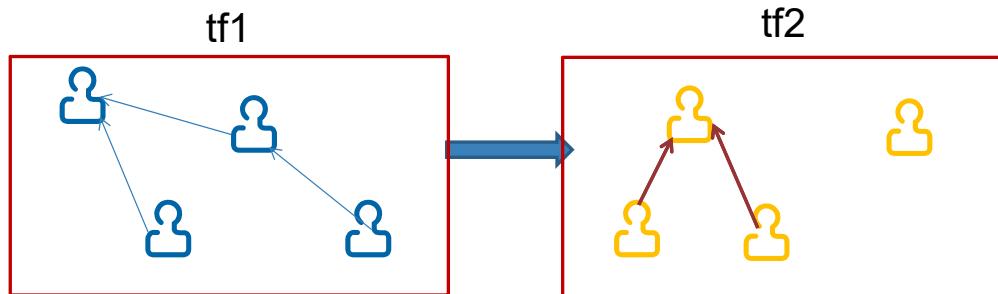


@ Yi-Shin Chen, Text Mining Overview

146

# Evolving Social Graph Analysis

- ▷ Add in evolving idea (time frame increment): graph sequences

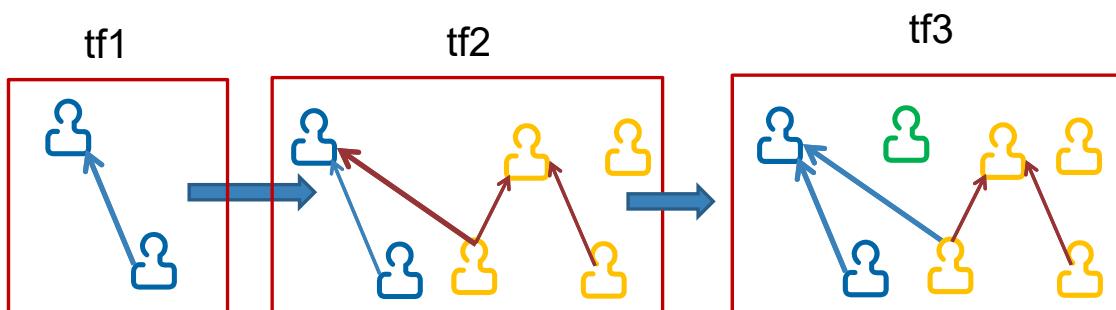


@ Yi-Shin Chen, Text Mining Overview

147

# Evolving Social Graph Analysis

- ▷ Information decay:
  - Vertex weight, edge weight
  - Decay mechanism



- ▷ Concept-Based Evolving Graph Sequences (cEGS): a sequence of directed graphs that demonstrate information propagation

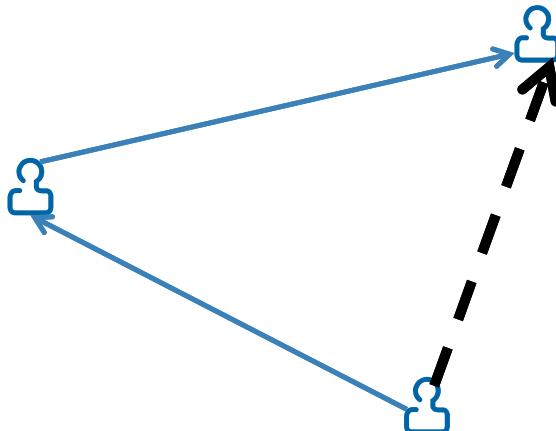
@ Yi-Shin Chen, Text Mining Overview

148

## Methodology – Evolving Social Graph Analysis

### ▷ Hidden link – construct hidden relationship

- To model better interaction between users
- Sample data



@ Yi-Shin Chen, Text Mining Overview

149

## Evolving Social Graph Analysis

### ▷ Analysis of cEGS

- Number of vertices( $nV$ )
  - The number of users mentioned this event candidate
- Number of edges( $nE$ ):
  - The number of following relationship in this cEGS
- Number of connected components( $nC$ )
  - The number of communities in this cEGS
- Reciprocity( $R$ )
  - The degree of mutual connections is in this cEGS

$$\text{Reciprocity} = \frac{\text{number of edges}}{\text{number of possible edges}}$$

@ Yi-Shin Chen, Text Mining Overview

150

# Event Identification

## ▷ Type1 Event: One-shot event

- An event that receives attention in a short period of time
  - Number of users, number of followings, number of connected components suddenly increase

## ▷ Type2 Event: Long-run event

- An event that attracts many discussion for a period of time
  - Reciprocity

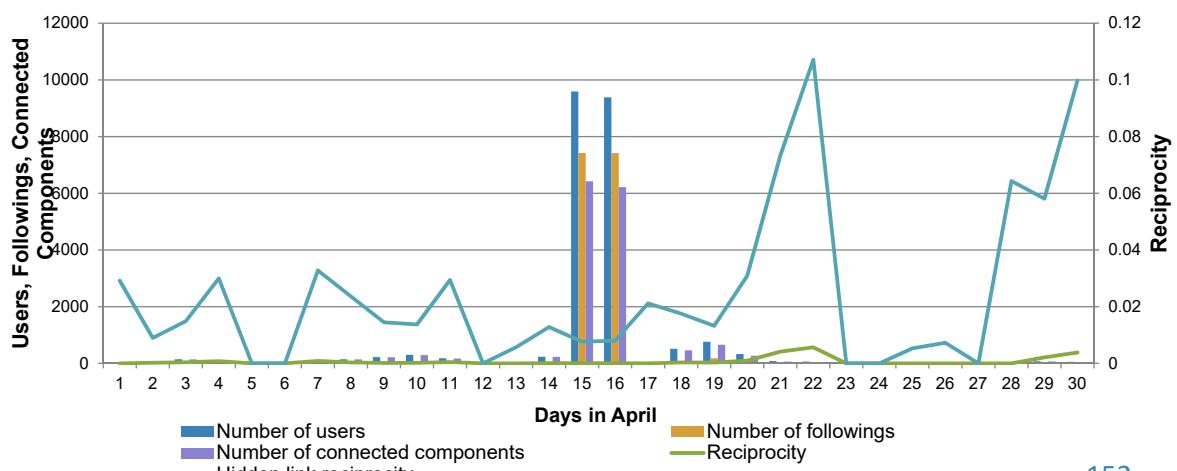
## ▷ Non-event

@ Yi-Shin Chen, Text Mining Overview

151

# Experimental Results

▷ April 15, “library jfk blast bostonmarathon prayforboston pd possible area boston police social explosion bomb bombing marathon confirm incident”



152

## Contextual Text Mining Example 2

### Location Identification

153

## Twitter and Geo-Tagging

- ▷ Geospatial features in Twitter have been available
  - User's Profile location
  - Per tweet geo-tagging
- ▷ Users have been slow to adopt these features
  - On a 1-million-record sample, 0.30% of tweets were geo-tagged
- ▷ Locations are not specific

# Our Goal

- ▷ Identify the location of a particular Twitter user at a given time
  - Using exclusively the content of his/her tweets



@ Yi-Shin Chen, Text Mining Overview

155

# Data Sources

- ▷ Twitter Dataset
  - 13 million tweets, 1.53 million profiles
  - From Nov. '11 to Apr. '12

## ▷ Local Cache

- Geospatial Resources
- GeoNames Spatial DB
- Wikiplaces DB
- Lexical Resources
- WordNet

## ▷ Web Resources

@ Yi-Shin Chen, Text Mining Overview

156

# Baseline Classification

▷ We are interested only in tweets that might suggest a location.

## ▷ Tweet Classification

- Direct Subject
- Has a first person personal pronoun – “I love New York”
  - (I, me, myself,...)
- Anonymous Subject
  - Starts with Verbs – “Rocking the bar tonight!”
  - Composed of only nouns and adjectives – “Perfect day at the beach”
- Others

# Rule Generation

▷ By identifying certain association rules, we can identify certain locations

▷ Combinations of certain verbs and nouns (bigrams) imply some locations

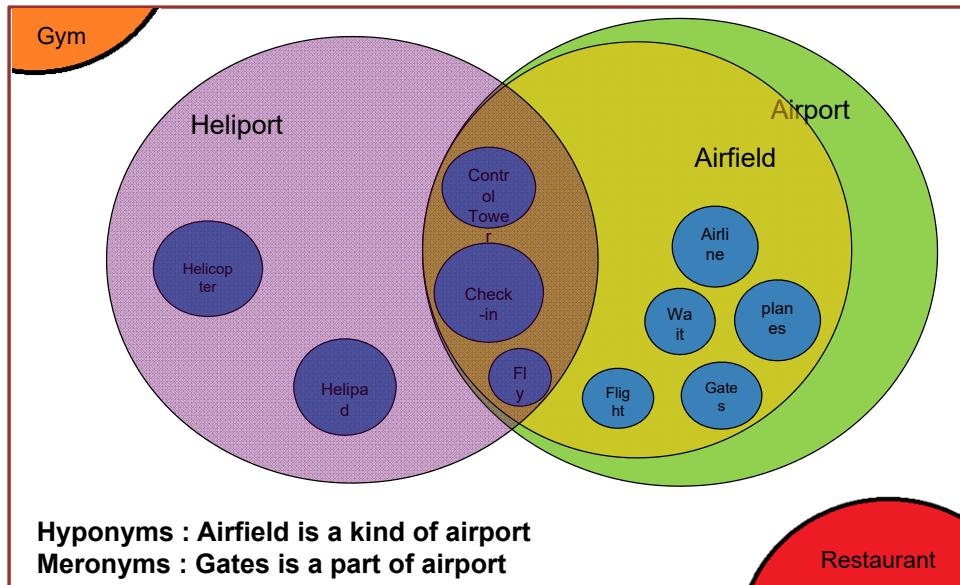
- “Wait” + “Flight” → Airport
- “Go” + “Funeral” → Funeral Home

▷ Create combinations of all synonyms, meronyms and hyponyms related to a location

▷ Number of occurrences must be greater than K

▷ Each rule does not overlap with other categories

# Examples



@ Yi-Shin Chen, Text Mining Overview

159

## N-gram Combinations

- ▷ Identify likely location candidates
- ▷ Contiguous sequences of n-items
- ▷ Only nouns and adjectives
- ▷ Extracted from Direct Subject and Anonymous Subject categories.

Field	Original Tweet	After Processing	1-gram Seq	2-gram Seq	3-gram Seq
Vernacular name			Unigram	Bigram	Trigram
Order of resulting Markov model			0	1	2
Example	... I love the beautiful city of Hsinchu ...	...beautiful city Hsinchu...	...,beautiful, city, Hsinchu,...	...,beautiful city city Hsinchu,...	...,beautiful city Hsinchu,...

@ Yi-Shin Chen, Text Mining Overview

160

# Tweet Types

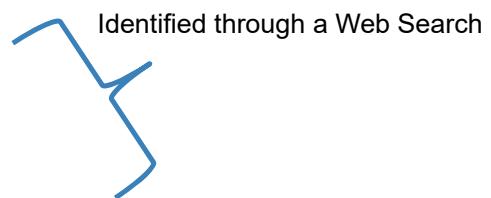
## ▷ Coordinates

- Tweet has geographical coordinates

## ▷ Explicit Specific

- “I Love Hsinchu”

→ Toponyms



Identified through a Web Search

## ▷ Explicit General

- “Going to the Gym”

→ Through Hyponyms

## ▷ Implicit

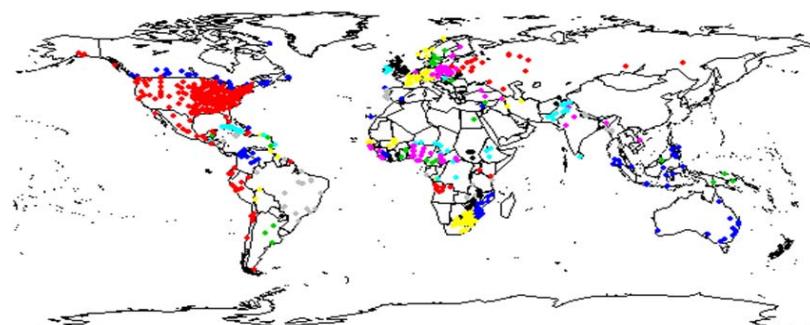
- “Buying a new scarf” -> Department Store
- Emphasize on actions

@ Yi-Shin Chen, Text Mining Overview

161

# Country Discovery

- ▷ Identify the user's country
- ▷ Identified with OPTICS algorithm
- ▷ Cluster all previously marked n-grams
- ▷ Most significant cluster is retrieved
  - User's country



162

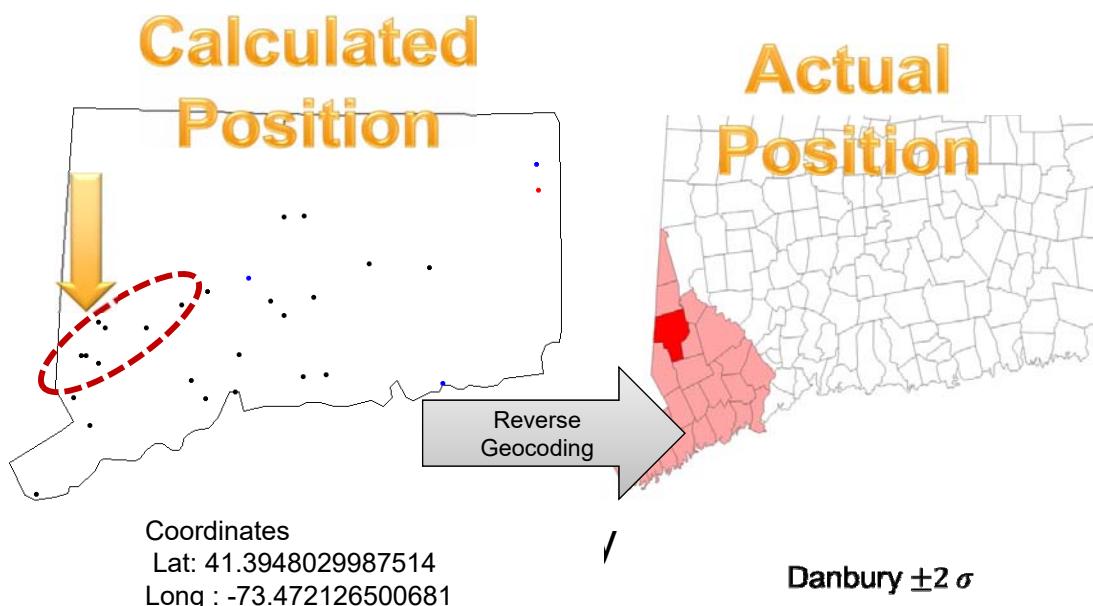
## Inner Region Discovery

- ▷ Identify user's Hometown
- ▷ Remove toponyms
- ▷ Clustered with OPTICS
- ▷ Only locations within  $\pm 2 \sigma$

@ Yi-Shin Chen, Text Mining Overview

163

## Inner Region Discovery



@ Yi-Shin Chen, Text Mining Overview

164

# Web Search

▷ Use web services

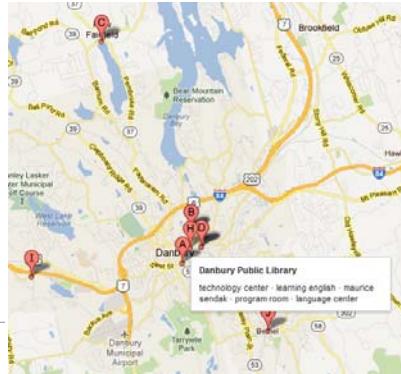
▷ Example :

- “What's the weather like outside? I haven't left the library in three hours”
- Search: Library near Danbury, Connecticut , US

Google places

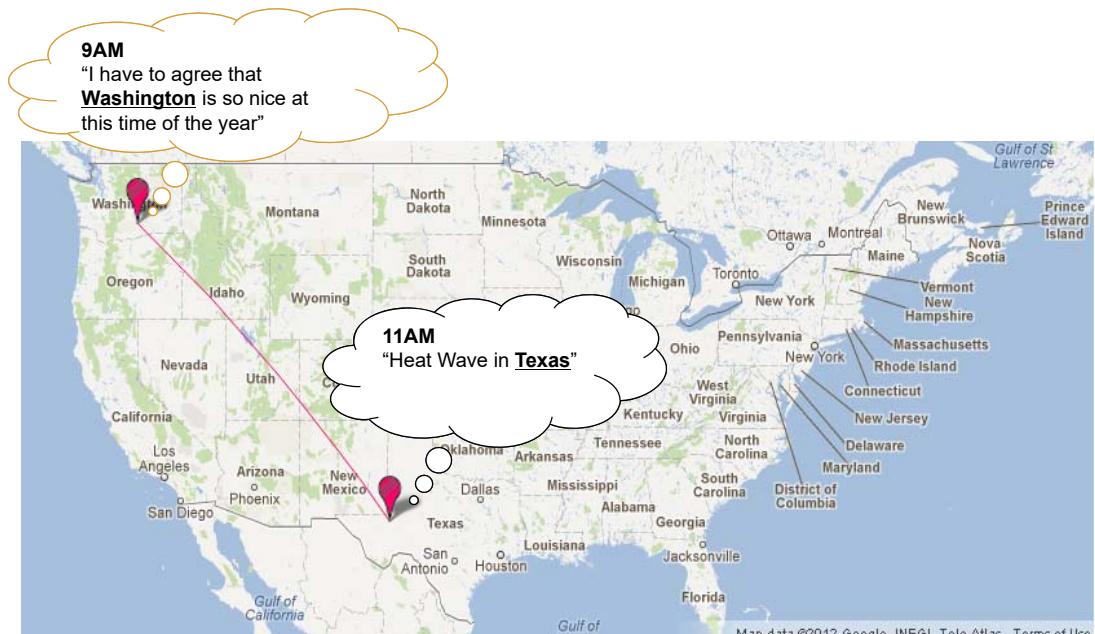


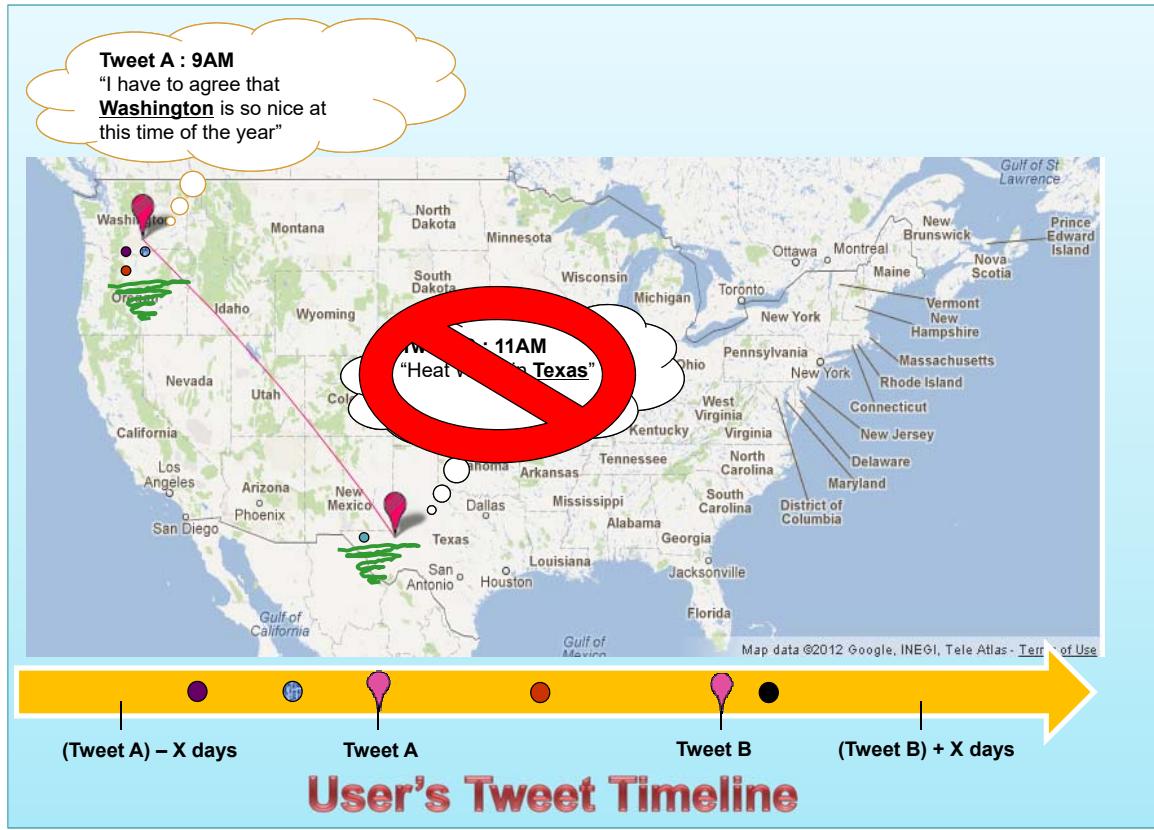
Danbury Public Library  
Danbury Health Sciences Library  
New Fairfield Free Public Library  
Ruth A Haas Library  
Danbury Hospital Library  
Robert S Young Business Library  
Long Ridge Library  
Superior Court Law Library



165

## Timeline Sorting





## Location Inferred

▷ Users are classified according to the particular findings

:

- No Country (No Information)
- Just Country
- Timeline
  - Current and past locations
- Timeline with Hometown
  - Current and past locations
  - User's Hometown
  - General Locations

# General Statistics

Country	Profiles	Factual	Empty Fictional	MAE Factual	MAE Empty Fictional
US	26	16	10	6.26	21.31
GB	22	19	3	6.18	6.45
CA	29	25	4	11.18	12.63
AU	8	8	0	6.47	-
IN	15	13	2	7.76	18.63
Others	5	4	1	2.06	14.13
<b>Total</b>	<b>105</b>	<b>85</b>	<b>20</b>	<b>6.65</b>	<b>14.62</b>

Country	ACC100	WMAE	Tw MAE
US	96%	12.05	9.51
GB	100%	6.22	2.02
CA	100%	11.38	3.41
AU	100%	6.47	2.42
IN	93%	9.21	1.89
Others	100%	4.47	27.54

@ Yi-Shin Chen, Text Mining Overview

169

## Case Studies

@ Yi-Shin Chen, Text Mining Overview

170

# Reddit Data

@ Yi-Shin Chen, Text Mining Overview

171

## Reddit Data

<https://www.kaggle.com/c/reddit-comments-may-2015/data>

The screenshot shows a list of top comments from a search query. The comments are listed in descending order of upvotes.

- 1 How many of you have watched the TV Show "Misfits"? (self.AskReddit)  
submitted 12 分前 by Spyder to /r/AskReddit  
5 留言 分享
- 2 What are the "Beats headphones" of your hobby? What makes you cringe to see others flexing? (self.AskReddit)  
submitted 9小時前 by MadeANewAccountUgh to /r/AskReddit  
21708 留言 分享
- 3 Sony, Steam and Amazon issuing full refunds for No Man's Sky. [Misleading Title] (inquisitr.com)  
submitted 10小時前 by IxWoodstockX1 to /r/gaming  
2837 留言 分享
- 4 My 90 years old french grandmother gave me the WW2 medal of my grandfather after he died (imgur.com)  
submitted 4小时前 by Fitzcalm to /r/pics  
108 留言 分享
- 5 The reason I like Reddit so much is because no one from my family uses it (self.Showerthoughts)  
submitted 10小時前 by Chouplina to /r/Showerthoughts  
670 留言 分享
- 6 Chromatophores in a squid react to light. (i.imgur.com)  
submitted 10小時前 by natsdorf to /r/gifs  
287 留言 分享

172

# Reddit: The Front Page of the Internet

MY SUBREDDITS ▾

- ANNOUNCEMENTS
- ART
- ASKREDDIT
- ASKSCIENCE
- AWW
- BLOG
- BOOKS
- CREEPY
- DATAISBEAUTIFUL
- DIY
- DOCUMENTARIES
- EARTHPORN
- EXPLAINLIKEIMFIVE

50k+ on  
this set

↑ [-] redundant6939 1285 points 5 hours ago  
↓ The way she moved her head and winked is disturbing af.  
permalink save report give gold reply

↑ [-] KayPike [S] 871 points 5 hours ago  
↓ Thanks I was going for a cute kind of disturbing and weird ^  
permalink save parent report give gold reply

↑ [-] Suckonmyfatvagina 314 points 5 hours ago  
↓ That's you? Also are you the artist as well?  
permalink save parent report give gold reply

↑ [-] KayPike [S] 692 points 4 hours ago  
↓ Yppers, this is my first try at facepainting myself :)



## Subreddit Categories

- ▷ Reddit's structure may already provide a baseline similarity

gadgets

related\_subreddits [view](#) [history](#) [talk](#)

RELATED SUBREDDITS

Hacked Gadgets	Products
Tech	TechNews
Technology	TechProducts
Technawlogy	

music

hot new rising controversial top gilded [wiki](#) promoted

Menu Listen Social Music Subreddits Info Filter

Are you a musician? [Read our guide to promoting yo](#)

musicsubreddits [view](#) [history](#) [talk](#)

Welcome to the official list of [/r/Music](#) recommended subreddits!  
In bold are the subreddits with over 10,000 readers.  
Please add a very descriptive, objective and useful comment after each subreddit name !  
Enjoy!

Subreddits by Genre

- Classical Music
- Electronic Music
- Rock / Metal
- Hip Hop
- Some decades.
- By country/region/culture
- Other



# Provided Data

- created\_utc
- ups
- subreddit\_id
- link\_id
- name
- score\_hidden
- author\_flair\_css\_class
- author\_flair\_text
- subreddit
- id
- removal\_reason
- gilded
- downs
- archived
- author
- score
- retrieved\_on
- body
- distinguished
- edited
- controversiality
- parent\_id



# Recover Structure

DjSoundBlast	2015-05-01 10:10:30	2 points
Will the stems be available, or will we have to work with the acceppalla on Zedd's Sound cloud?		

JustTheNorm	2015-05-01 10:14:50	6 points
Stems provided.		

DjSoundBlast	2015-05-01 10:23:13	6 points
Do you know if zedd is looking for people to completely transform the song like he did for his first remix contests, or a remix that still stays true to the original?		

rycar	2015-05-01 14:09:35	2 points
Watch the video on the contest page.		

Nawabi	2015-05-01 20:30:02	1 points
the video only explains how the contest works, not what hes looking for. Id say transform.		



# Facebook Fanpage

@ Yi-Shin Chen, Text Mining Overview

177

## Fan Page?



@ Yi-Shin Chen, Text Mining Overview

178

## Fan Page Post



@ Yi-Shin Chen, Text Mining Overview

179

## Fan Page Feedback



@ Yi-Shin Chen, Text Mining Overview

180

# Twitter

@ Yi-Shin Chen, Text Mining Overview

181

## Tweet

Music Ed Now @musicednow · Feb 23  
What an amazing #piano #recital last night at #uofm featuring @KotaroFukuma!!!  
So much #energy #expression #confidence and #excellence! 🌟🌟🌟



Music Ed Now @musicednow · Feb 23  
Thank you to all our new #followers!!! #music #musiced #homeed #homeschool  
#unschool #parents #teens #kids #children #musical #musical

@ Yi-Shin Chen, Text Mining Overview

182

# Reddit DM Examples

@ Yi-Shin Chen, Text Mining Overview

183

Data Mining Final Project  
*Role-Playing Games Sales  
Prediction*

Group 6

## Dataset

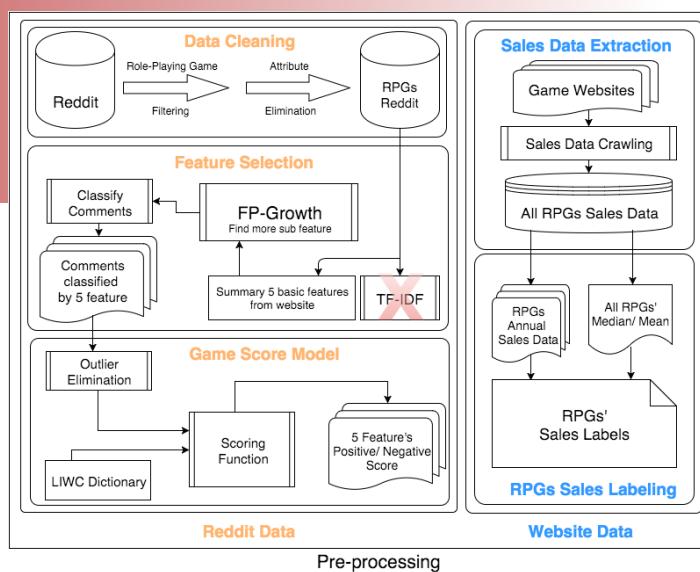
▷ Use Reddit comment dataset

▷ Data Selection

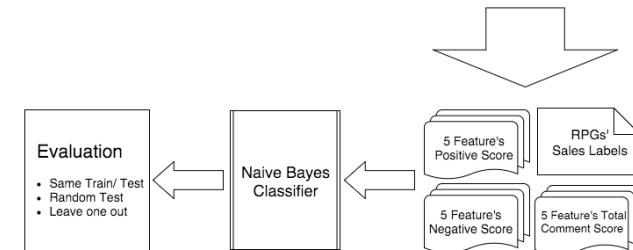
- Choose 30 Role-playing games from the subreddits.

mountandblade	witcher	Diablo	assassinscreed	metalgearsolid	HalfLife
blackops3	battlefield_4	bloodborne	Borderlands2	FF14	Warhammer
GrandTheftAutoV	thelastofusfactions	halo	DarkSouls2	KingdomHearts	StreetFighter
CodAW	BF_Hardline	dragonage	masseffect	darksouls	Shadowrun
StarWarsBattlefront	skyrim	batman	Fallout	FinalFantasy	arma

- Choose useful attributes form their comment
  - Body
  - Score
  - Subreddit



Pre-processing



Role-Playing Games (RPGs) Sales Prediction Model



## Preprocessing

### ▷ Common Game Features

- Drop trivial words
  - We manually build the trivial words dictionary by ourselves.
- Ex : "is" "the" "I" "you" "we" "a".
- TF-IDF
  - We use Tf-IDF to calculate the importance of each words.



## Preprocessing

- TF-IDF result
  - But the TF-IDF value are too small to find the significant importance.

Words	TFIDF	Words	TFIDF	Words	TFIDF
souls	0.005845	ancient	0.003847	city	0.00209
bloodborn	0.00536	solo	0.003843	guns	0.001887
boss	0.00531	group	0.003536	rifle	0.001883
blood	0.005161	blizzard	0.003419	nuclear	0.001785
damage	0.004916	weapon	0.003242	sounds	0.00158
weapon	0.004386	speed	0.003196	shot	0.001375
beast	0.004062	character	0.003131	building	0.001351

- Define Comment Feature
  - We manually define the common gaming words in five categories by referring several game website.



## Preprocessing

### ▷ Filtering useful comments

- Common keywords of 5 categories:

GamePlay	music	story	graphic	criticize
Balance	orchestral	main story	graphic	Downgrade
Combat	BGM	character story	graphical	engine
Atmosphere	soundtrack	story	texture	performance
Gameplay	voice act	storyline	Framerate	load%time

- Filtering useful comments

→ Using the frequent keywords to find out the useful comments in each category.



## Preprocessing

- Using FP-Growth to find other feature  
→ To Find other frequent keywords

{**time**, play, people, gt, **limited**}

{**time**, up, good, now, damage, gear, way, need, build, better, d3, something, right, being, gt, **limited**}

- Then, manually choose some frequent keywords

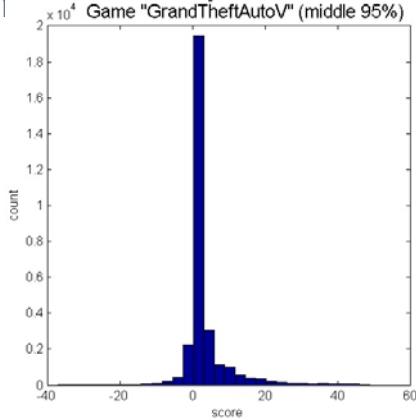


## Preprocessing

### ▷ Comment Emotion

- Filtering outliers

→ First filter out those comment whose “score” separate in top and bottom 2.5% which indicate that they m<sup>~</sup> Game "GrandTheftAutoV" (middle 95%)



## Preprocessing

- Emotion detection

→ We use LIWC dictionary to calculate each comment's positive and negative emotion percentage.

Ex: I like this game. → **positive**

EX: I hate this character. → **negative**

→ Then find each category's positive and negative emotion score.

- $Score_{each\_category\_P} = \sum_{j=0}^n score_{each\_comment}_j * (positive\ words)$
- $Score_{each\_category\_N} = \sum_{j=0}^n score_{each\_comment}_j * (negative\ words)$
- $TotalScore_{each\_category} = \sum_{i=0}^m score_{each\_comment}_m$
- $FinalScore_{each\_category\_P} = Score_{each\_category\_P} / TotalScore_{each\_category}$
- $FinalScore_{each\_category\_N} = Score_{each\_category\_N} / TotalScore_{each\_category}$



## Preprocessing

- Emotion detection

→ We use LIWC dictionary to calculate each comment's positive and negative emotion percentage.

Ex: I like this game. → **positive**

EX: I hate this character. → **negative**

→ Then find each category's positive and negative emotion score.

The meaning is :

Finding each category's emotion score for each game

Calculating TotalScore of each category's comments

Having the average emotion score FinalScore of each category

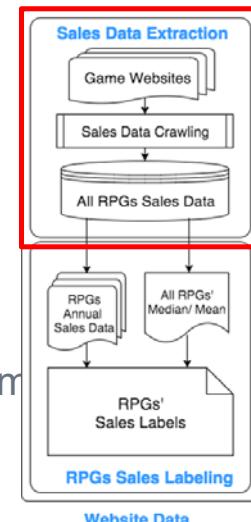


## Preprocessing

### ▷ Sales Data Extraction

- Crawling website's data
- Find the games' sales on each platform

1		Call of Duty: Black Ops 3 (PS4) Activision, Shooter	6	8,303,244	8,303,244
2		FIFA 16 (PS4) Electronic Arts, Sports	12	6,451,242	6,451,242
3		Star Wars: Battlefront (2015) (PS4) Electronic Arts, Shooter	4	4,832,058	4,832,058
4		Call of Duty: Black Ops 3 (XOne) Activision, Shooter	6	4,423,401	4,423,401
5		Fallout 4 (PS4) Bethesda Softworks, Role-Playing	5	4,418,849	4,418,849
6		Grand Theft Auto V (PS4) Take-Two Interactive, Action	56	3,684,126	8,026,560
7		Splatoon (WiiU) Nintendo, Shooter	29	3,335,728	3,335,728
8		Batman: Arkham Knight (PS4) Warner Bros. Interactive Entertainment, Action	25	3,044,664	3,044,664
9		The Witcher 3: Wild Hunt (PS4) Warner Bros. Interactive Entertainment, Role-Playing	30	2,980,543	2,980,543



Website Data



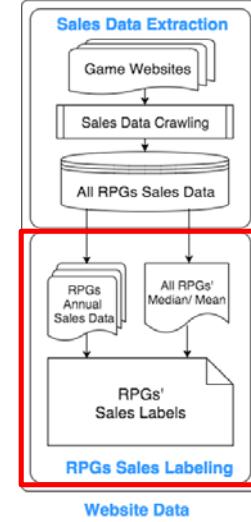
## Preprocessing

- Annually sales for each game
- Median of each platform
- Mean of all platform

Median : **0.17** -> 30 games (26H 4L)

Mean : **0.533** -> 30 games (18H 12L)

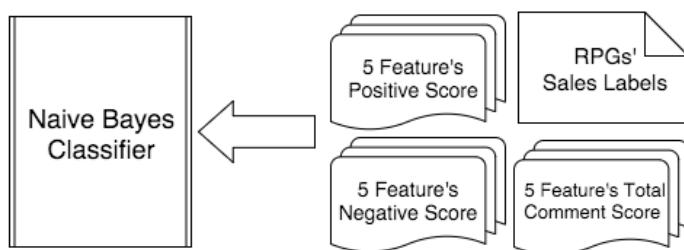
- Try both of Median/ Mean to do the prediction.



## Sales Prediction

### ▷ Model Construction

- We use the 4 outputs from pre-processing step to be the input of the Naïve Bayes classifier



197

## Evaluation

We use naïve Bayes to evaluate our result

1. training 70% & test30%
2. Leave-one-out

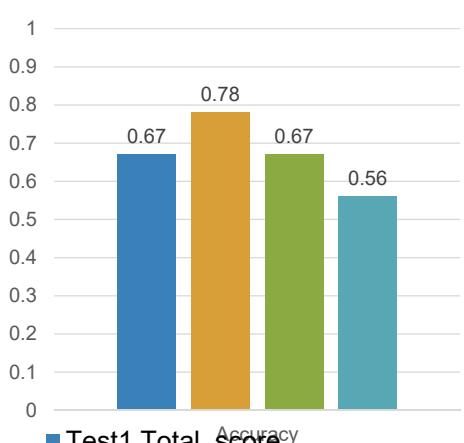


198

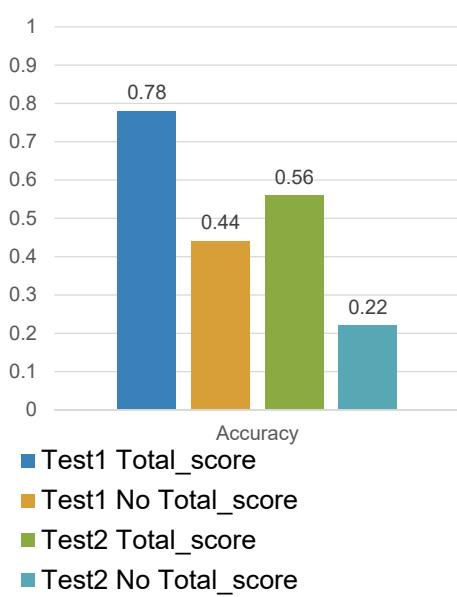
## Evaluation (train70% & test30%)

Median 0.17

Mean 0.533



- Test1 Total\_score
- Test1 No Total\_score
- Test2 Total\_score
- Test2 No Total\_score



- Test1 Total\_score
- Test1 No Total\_score
- Test2 Total\_score
- Test2 No Total\_score



199

## Evaluation(Leave-one-out)

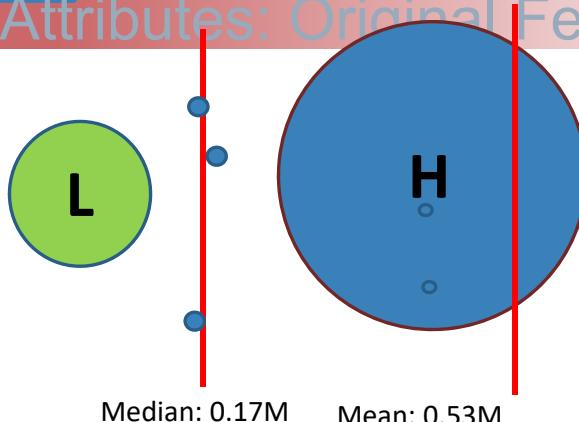
Choose 29 games as training data and the other one as test data, and totally run 30times to get the accuracy.

	Total 30 times	Accuracy
Median Total_score	7 times wrong	77%
Median NO Total_score	5 times wrong	83%
Mean Total_score	6 times wrong	80%
Mean No Total_score	29 times wrong	3%



200

## Attributes: Original Features Scores



Evaluation: Leave-one-Out  
Total Sample size: 30

Error Times	Accuracy
Median 5	83%
Mean 29	3%



Attribute distribution to target(sales) domain



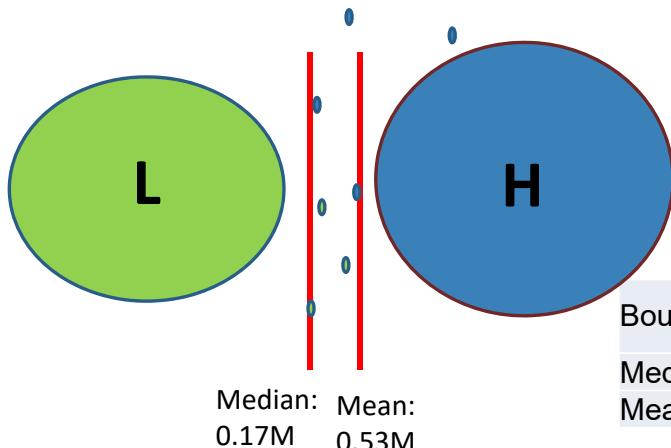
Sales boundary for H and L class

Sales class



201

## Attributes: Transformed Features Scores Better



Evaluation: Leave-one-Out  
Total Sample size: 30

Boundary	Error Times	Accuracy
Median	7	77%
Mean	6	80%

Attribute distribution project to target(sales) domain  
 Sales boundary for H and L class  
 Sales class

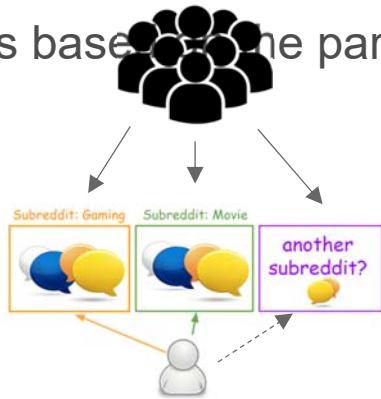


Finding related subreddits based on the gamers' participation

Data Mining Final Project  
Group#4

## Goal & Hypothesis

- To suggest new subreddits to the Gaming subreddit users based on the participation of other users



203

## Data Exploration

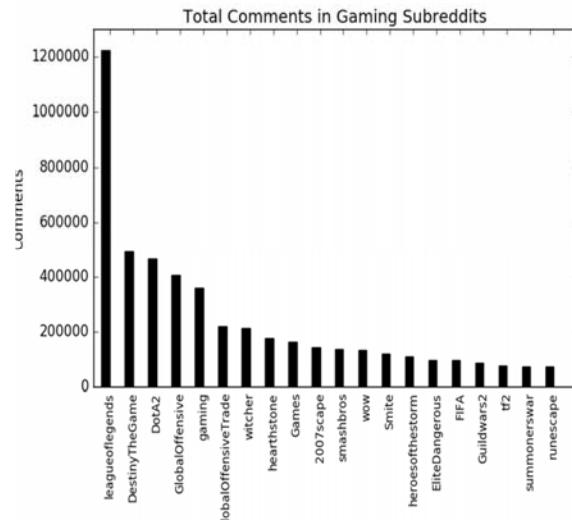
- ▷ We extracted the Top 20 gaming related subreddits to be focused on.

leagueoflegends	smashbros
DestinyTheGame	wow
DotA2	Smite
GlobalOffensive	heroesofthestorm
Gaming	EliteDangerous
GlobalOffensiveTrade	FIFA
witcher	Guildwars2
hearthstone	tf2
Games	summonerswar
2007scape	runescape

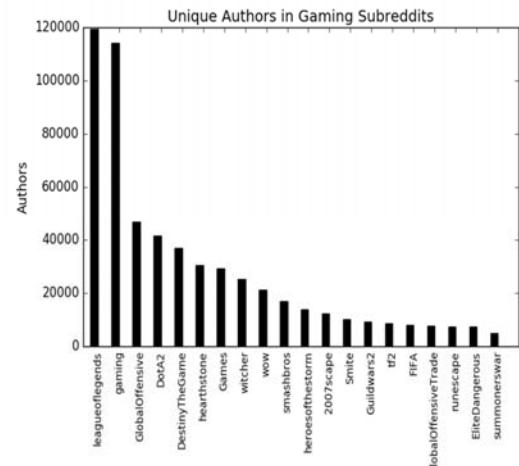
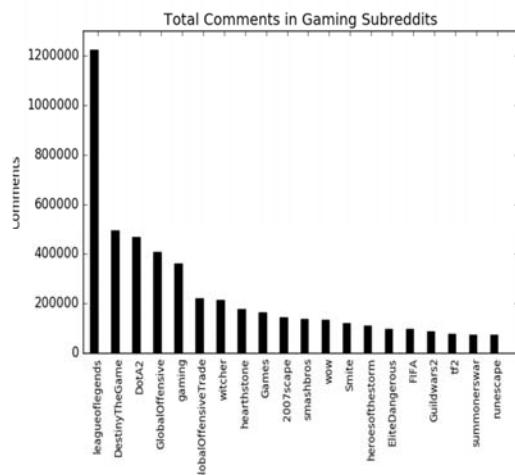
204

## Data Exploration (Cont.)

- ▷ We sorted the subreddits by users' comment activities.
- ▷ We found out that the most active subreddit is Leagueoflegends.



## Data Exploration (Cont.)



# Data Pre-processing

- Removed comments from moderators and bots (>2,000 comments)
- Removed comments from default subreddits
- Extracted all the users that commented on at

```
 subreddit
3ch, 4chan, AdviceAnimals, AndroidQuestions, AutoModerator, Boxing, Ca
ndidFashionPolice, CoonTown, CuteFemaleCorpses, DotA2, DumpsterDivin
g, Ellenpaiaaction, GasTheKikes, JusticePorn, KotakuInAction, Public
Freakout, RedditFacelift, SRSSucks, SexWithDogs, SubredditDrama, Tumb
lrInAction, WhiteSmite, WhitesWinFights, beta, circlejerk, cssnews, da
taiisbeautiful, modhelp, pussypassdenied, redditdev, siliconvalley, sn
ew, subredditcancer, techsupport, userstyles
4chan, Agario, Drugs, GlobalOffensive, TheWeeknd, fatpeoplehate, killi
ngfloor, leagueoflegends, millionairemakers
Animesuggest, BlackPeopleTwitter, CringeAnarchy, PussyPass, Robocraf
t, TumblrInAction, WTF, awnnverts, comics, ftlgame, homeopathy, mantids
,pokemon, tf2
```

"ket Basket"  
DestinyTheGame, Fireteams, theflash, FIU

207

# Processing

## ▷ Sorted the frequent items (sort.js)

▷ Eg. DotA2, Neverwinter, atheism → DotA2, atheism, Neverwinter

### Original Transactions

```
107452 Dirtybomb,Games,Warframe
107453 4chan,BlackPeopleTwitter,FlashTV,Games,Guildwars2,Hig
107454 KotakuInAction,NotTimAndEric,PublicFreakout,thatHappened
107455 AdviceAnimals,AskCulinary,Boise,CampingandHiking,DotA
107456 DunderMifflin,GlobalOffensive,HotlineMiami,paydaytheheist
107457 Cynicalbrit,DotA2,DrinkingGames,Pathfinder,Pathfinder
107458 DotA2,Neverwinter,atheism
107459 Brazil,Fallout,GlobalOffensive,Minecraft,Steam,brasil
107460 Games,bloodborne,firstworldproblems,pcgaming,witcher
107461 Berserk,Eyebleach,Games,OnePiece,OutOfTheLoop,TokyoGh
107462 Berserk,DnD,EDH,Guildwars2,GundamExVs,Gunpla,fatpeopl
107463 AdviceAnimals,battlestations,gameofthrones,hardwaresw
107464 Starwarscommander,TumblrInAction,fireemblemcasual,fri
107465 Games,Mechanicalkeyboards,churning,gamecollecting,wit
107466 EliteDangerous,TalesFromRetail,WTF,brum,patientgamers
107467 DestinyTheGame,Fireteams,Frugal,nfl
```

### Sorted Transactions

```
107452 Games,Warframe,Dirtybomb
107453 leagueoflegends,Games,WTF,thebutton,TumblrInAction,PS
107454 witcher,KotakuInAction,thatHappened,PublicFreakout,No
107455 WTF,AdviceAnimals,DotA2,gameofthrones,buildapc,pcgami
107456 pcmasterrace,GlobalOffensive,paydaytheheist,DunderMif
107457 DotA2,Cynicalbrit,Pathfinder_RPG,Pathfinder,DrinkingG
107458 DotA2,atheism,Neverwinter
107459 pcmasterrace,GlobalOffensive,science,woahdude,Steam,M
107460 Games,witcher,pcgaming,bloodborne,firstworldproblems
107461 leagueoflegends,Games,reactiongifs,OutOfTheLoop,OnePi
107462 thebutton,fatpeoplehate,Guildwars2,magicTCG,DnD,EDH,v
107463 pcmasterrace,AdviceAnimals,wow,gameofthrones,pcgaming
107464 hearthstone,TumblrInAction,fireemblemcasual,friendsa
107465 Games,witcher,Mechanicalkeyboards,gamecollecting,chur
107466 WTF,EliteDangerous,patientgamers,TalesFromRetail,brum
107467 DestinyTheGame,nfl,Fireteams,Frugal
```

208

# Processing

- Choosing the minimum support from 159,475 Transactions
- Max possible support 25.71% (at least in 41,004 Transactions)
- Min possible support  
Min. Support as  
 $0.05\% = 79.73$   
transactions

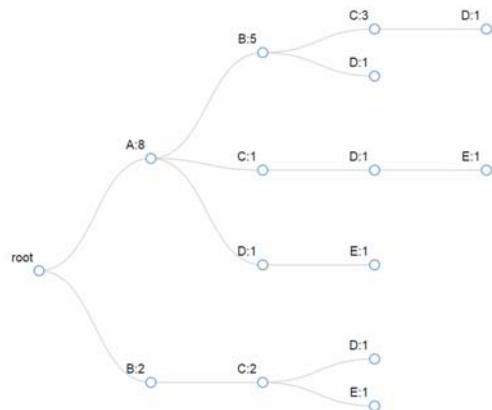
```
Testing if: <DILLY> -> <canine>
19 items remaining
Testing if: <DILLY> -> <politics>
18 items remaining
Testing if: <gameofthrones,technology,UTF,pennaterrace> -> <politics>
18 items remaining
Testing if: <Charyn> -> <thebutton>
16 items remaining
Testing if: <NexusNewbies> -> <gameofthrones>
14 items remaining
Testing if: <Metroid> -> <leagueoflegends>
13 items remaining
Association Rules Generated in: 01:01:24m:26s:66ms
Saved in: data\nosort\basket_trans_sort-tree
Saved Association Rules in: 0ms
Saved in: data\nosort\basket_trans-sort-tree.json
Saved Tree File in: 47ms
Done!
C:\xampp\htdocs\DM>
```

n 1

209

# Processing

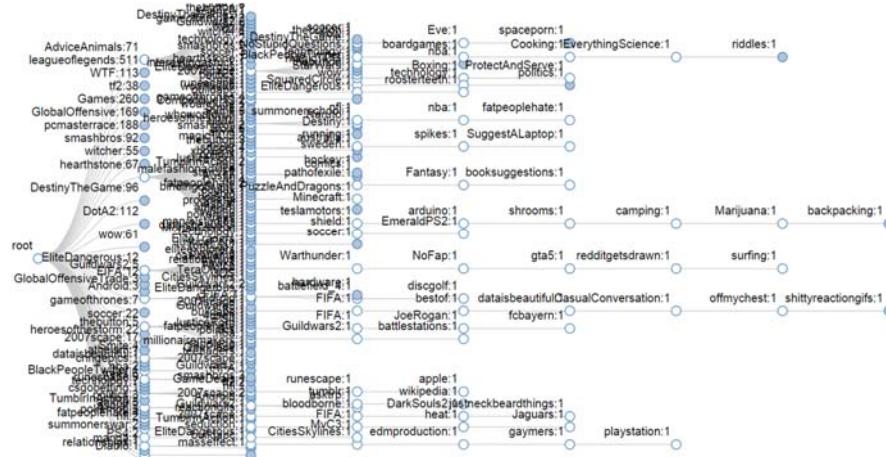
- How a FP-Growth tree looks like



210

# Processing

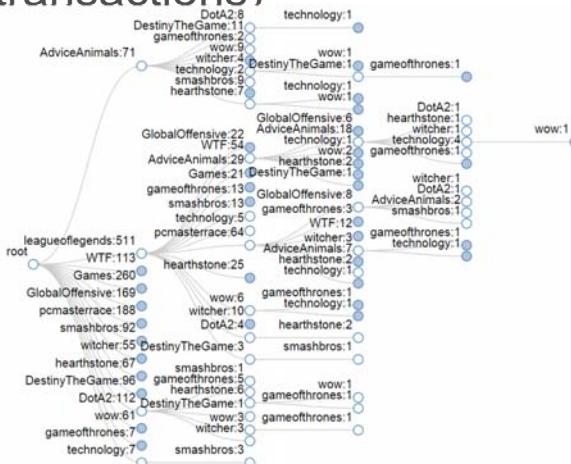
- Our FP-Growth (minimum support = 1% at least)



211

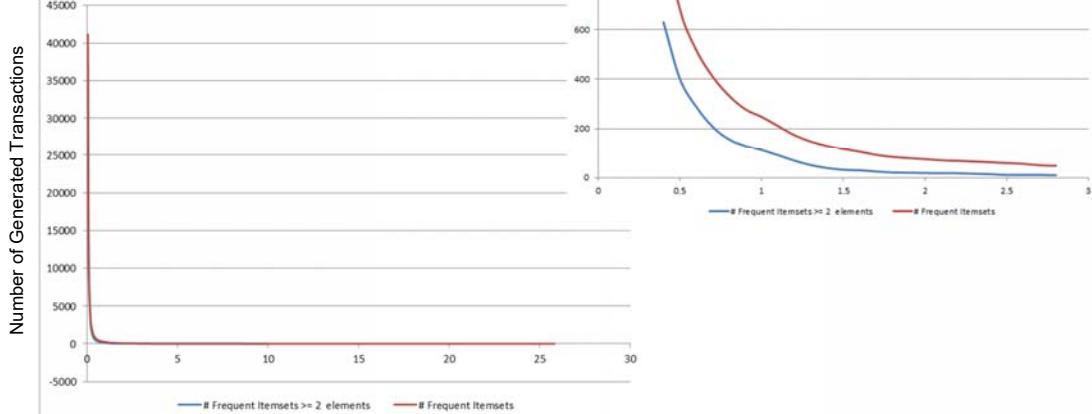
# Processing

- Our FP-Growth (minimum support = 5% at least  
in 7,974 transactions)



212

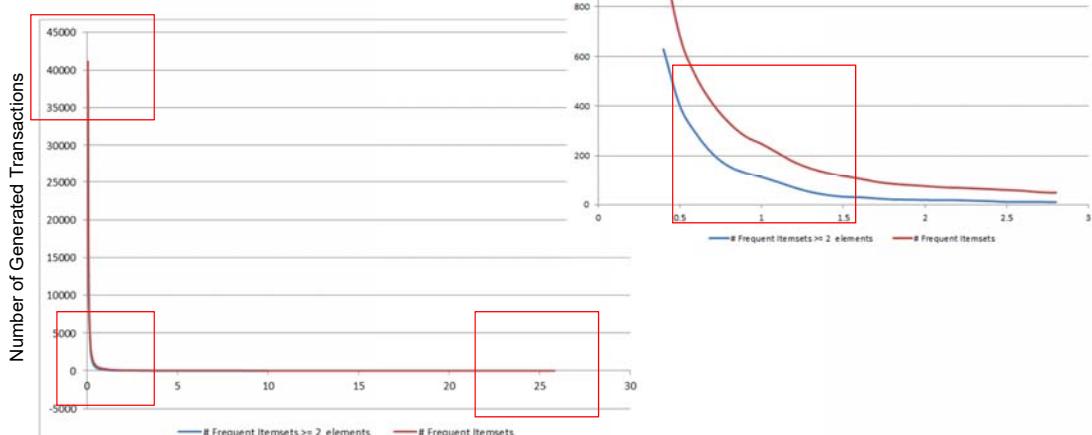
# Processing



213



# Processing



214



# Processing

- Our minimum support = 0.8% (at least in 1,276 trans)

Association Rules , min-support = 79.7375 items (0.05%)

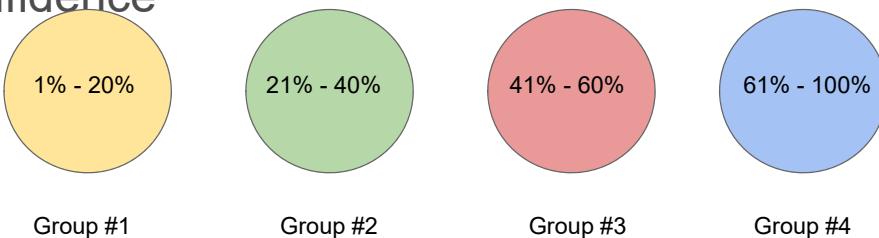
Limit to:	100	Search By Item:	item name...	Min Support	0.8	Max Support	100	Min
Confidence	0	Max Confidence	100	Min Items	1	Max Items	10	Showing 100 of 41169 itemsets
1	[ 1.51% , 47.70% , 5034 ]	{ pcmasterrace , AdviceAnimals }	-> WTF					
2	[ 1.51% , 46.68% , 5144 ]	{ pcmasterrace , WTF }	-> AdviceAnimals					
3	[ 1.51% , 30.70% , 7821 ]	{ WTF , AdviceAnimals }	-> pcmasterrace					
4	[ 1.46% , 38.85% , 5984 ]	{ leagueoflegends , WTF }	-> AdviceAnimals					
5	[ 1.46% , 39.17% , 5935 ]	{ leagueoflegends , AdviceAnimals }	-> WTF					
6	[ 1.46% , 29.73% , 7821 ]	{ WTF , AdviceAnimals }	-> leagueoflegends					
7	[ 1.28% , 26.02% , 7821 ]	{ WTF , AdviceAnimals }	-> technology					

215

- Games -> politics Support : 1.30% (2,073 users)

# Post-processing

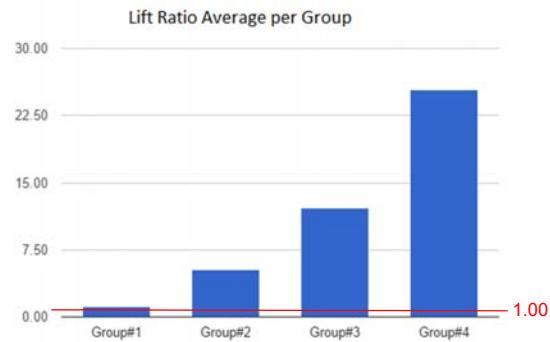
- We classified the rules in 4 groups based on their confidence



216

## Post-processing

- {Antecedent} → {Consequent}
- ▷ eg. {FIFA, Soccer} -



217

- **Lift ratio** = confidence / Benchmark confidence
- **Benchmark confidence** = # consequent element

## Post-processing

We created **8 surveys** with **64 rules** (2 rules per group) and **2 questions** per rules.

A screenshot of a "Reddit Survey" questionnaire. The title is "Subreddit Survey" and the subtitle is "Questionnaire for identifying the interestingness between subreddits". It says "Required".  
Question (1/8)  
Subreddit No.1 :: /r/DotA2 : an action RTS game developed by Valve Corporation  
DOTA 2  
Subreddit No.2 :: /r/AdviceAnimals: Reddit's Gold Mine  
Advice Animals  
1.1 Do you think /r/AdviceAnimals and /r/DotA2 are related? \*  
 Yes  
 No  
1.2 If you are subscriber of /r/AdviceAnimals, will you also be interested in /r/DotA2? \*  
1 2 3 4 5  
Definitely No                Definitely Yes  
NEXT

218

## Post-processing

▷ Q1: Do you think **subreddit A** and **subreddit B** are related? → Dependency

▷ [ yes | no ]

▷ Q2: If you are subscriber of **subreddit A**, will you also be interested in **subreddit B**?

▷ Definitely No [ 1 , 2 , 3 , 4 , 5 ] Definitely Yes  
Expected Confidence

219



## Post-processing

● We got response from 52 persons

● Data Confidence vs Expected Confidence → Lift

Subreddit No.1	Form#	# of response	Related Yes	Related No	0	25	50	75	100	% Not related	Expected Con (1=0)	Unrelatedness	Lift (1=0)
AdviceAnimals	1	4	0	4	1		3			100.00%	37.50%	68.75%	0
pcmasterrace	1	4	2	2	1	2		1		50.00%	31.25%	40.63%	0
AdviceAnimals	2	7	0	7			5	2		100.00%	32.14%	66.07%	0
WTF	3	7	1	6	5	2				85.71%	7.14%	46.43%	2.36%
AdviceAnimals	3	7	2	5	2	1	3		1	71.43%	39.29%	55.36%	0
Games	4	6	4	2		1	3	1	1	33.33%	58.33%	45.83%	0
WTF	4	6	3	3			4	2		50.00%	58.33%	54.17%	0
leagueoflegends	5	6	4	2		1	2	3		33.33%	58.33%	45.83%	0
Games	5	6	1	5		3	2	1		83.33%	41.67%	62.50%	0
leagueoflegends	6	7	1	6		4	3			85.71%	35.71%	60.71%	0
Games	4	7	6.5	0.5	0	0	1.5	4.5	1	7.14%	73.21%	40.18%	0
wow	7	4	1	3	2		2			75.00%	25.00%	50.00%	0
Games	7	4	3	1			2	2		25.00%	62.50%	43.75%	0
pcmasterrace	8	11	0	11	1	2	5	3		100.00%	47.73%	73.86%	0
WTF	8	11	0	11	2	6	3			100.00%	27.27%	63.64%	0
politics	1	4	0	4	4					100.00%	0.00%	50.00%	28.70%
pcgaming	1	4	4	0	1		3			0.00%	37.50%	18.75%	0
technology	2	7	7	0			1	6		0.00%	71.43%	35.71%	0



## Post-processing

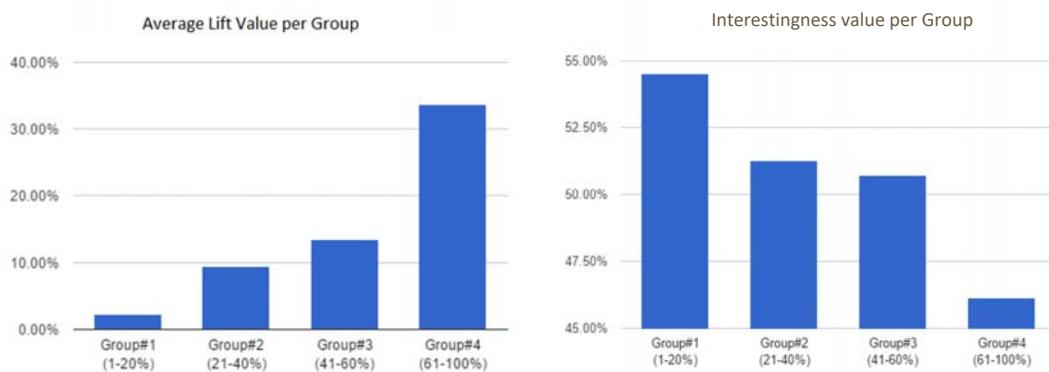
- % Non-related was proportion of the “No” answer to entire response in first question.
- Interestingness = Average of Expected confidence

# of response	Related Yes	Related No	% Not related	Expected Con (1=0)	Interestingness
4	0	4	100.00%	0.375	0.6875
4	2	2	50.00%	0.3125	0.40625
7	0	7	100.00%	0.3214285714	0.6607142857
7	1	6	85.71%	0.07142857143	0.4642857143
7	2	5	71.43%	0.3928571429	0.5535714286
6	4	2	33.33%	0.5833333333	0.4583333333
6	3	3	50.00%	0.5833333333	0.5416666667
6	4	2	33.33%	0.5833333333	0.4583333333
6	1	5	83.33%	0.4166666667	0.625
7	3	6	85.71%	0.3571428571	0.6071428571
7	6.5	0.5	7.14%	0.7321428571	0.4017857143
4	1	3	75.00%	0.25	0.5
4	3	1	25.00%	0.625	0.4375
11	0	11	100.00%	0.4772727273	0.7386363636
11	0	11	100.00%	0.2772727272	0.6363636364
4	0	4	100.00%	0	0.5
4	4	0	0.00%	0.375	0.1875
7	7	0	0.00%	0.7142857143	0.3571428571
7	4	3	42.86%	0.5357142857	0.4821428571
7	6	1	14.29%	0.75	0.4464285714

221



## Post-processing



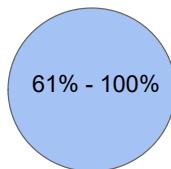
222



## Results

- How we can suggest new subreddits then?

High Confidence  
Less Interestingness



Group #4

Less Confidence  
High Interestingness



Group #1

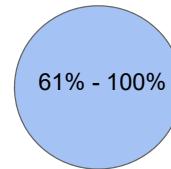
223



## Results

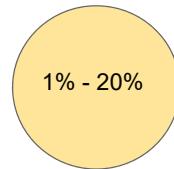
- Suggest them As:

Other People also  
talks about....



Group #4

Maybe you could be  
Interested in...



Group #1

224



# Results

## ● Results from Group# 4 ( 9 Rules )

```
1 [ 0.94% , 63.05% , 2387 ] { AdviceAnimals , BlackPeopleTwitter } -> WTF
2 [ 0.94% , 61.71% , 2439 ] { WTF , BlackPeopleTwitter } -> AdviceAnimals
3 [ 2.80% , 96.36% , 4640 ] { Fireteams } -> DestinyTheGame
4 [ 2.77% , 96.43% , 4589 ] { summonerschool } -> leagueoflegends
5 [ 1.64% , 84.54% , 3086 ] { csgobetting } -> GlobalOffensive
6 [ 1.03% , 73.12% , 2251 ] { GlobalOffensiveTrade } -> GlobalOffensive
7 [ 1.00% , 89.48% , 1787 ] { CompetitiveHS } -> hearthstone
8 [ 0.91% , 63.76% , 2285 ] { truegaming } -> Games
9 [ 0.81% , 96.19% , 1338 ] { CLG } -> leagueoflegends
```

225



# Results

## ● Results from Group# 1 ( 173 Rules )

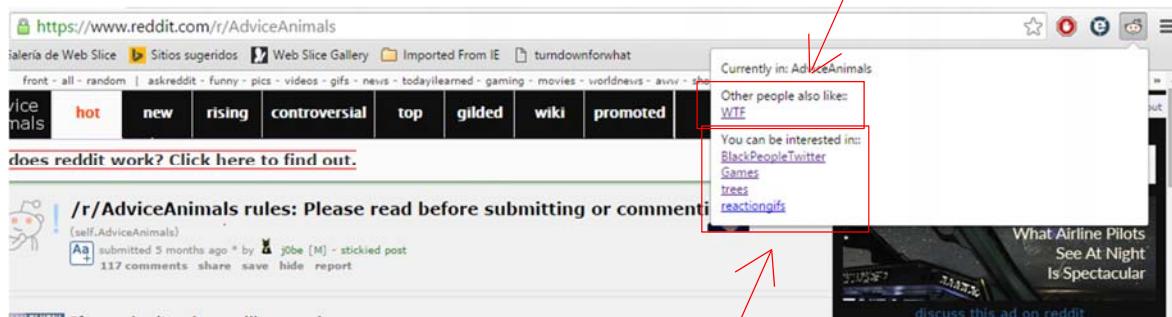
```
13 [ 2.31% , 14.68% , 25063 ] { Games } -> AdviceAnimals
14 [ 2.31% , 17.85% , 20605 ] { AdviceAnimals } -> Games
15 [ 2.31% , 17.85% , 20605 ] { AdviceAnimals } -> technology
16 [ 2.29% , 14.56% , 25063 ] { Games } -> technology
17 [ 2.28% , 16.52% , 22022 ] { GlobalOffensive } -> leagueoflegends
18 [ 2.28% , 8.87% , 41004 ] { leagueoflegends } -> GlobalOffensive
19 [ 2.01% , 14.26% , 22501 ] { WTF } -> GlobalOffensive
20 [ 2.01% , 14.57% , 22022 ] { GlobalOffensive } -> WTF
21 [ 1.96% , 12.49% , 25063 ] { Games } -> witcher
22 [ 1.90% , 12.32% , 24560 ] { pcmasterrace } -> pcgaming
23 [ 1.79% , 11.62% , 24560 ] { pcmasterrace } -> DotA2
24 [ 1.79% , 17.19% , 16598 ] { DotA2 } -> pcmasterrace
25 [ 1.78% , 13.81% , 20605 ] { AdviceAnimals } -> politics
26 [ 1.77% , 12.51% , 22501 ] { WTF } -> politics
```

13 GlobalOffensive,->,leagueoflegends,2.28%,16.52%,22022  
14 leagueoflegends,->,GlobalOffensive,2.28%,8.87%,41004  
15 WTF,->,GlobalOffensive,2.01%,14.26%,22501  
16 GlobalOffensive,->,WTF,2.01%,14.57%,22022  
17 Games,->,witcher,1.96%,12.49%,25063  
18 pcmasterrace,->,pcgaming,1.98%,12.32%,24560  
19 pcmasterrace,->,DotA2,1.79%,11.62%,24560  
20 DotA2,->,pcmasterrace,1.79%,17.19%,16598  
21 AdviceAnimals,->,politics,1.78%,13.81%,20605  
22 WTF,->,politics,1.77%,12.51%,22501  
23 AdviceAnimals,->,GlobalOffensive,1.73%,13.42%,20605  
24 GlobalOffensive,->,AdviceAnimals,1.73%,12.56%,22022  
25 pcmasterrace,->,technology,1.70%,11.02%,24560  
26 Games,->,pcgaming,1.64%,10.41%,25063

226



# Results



227

# Challenges

- Reducing scope of data for decreasing computational time
- Defining and calculating the interestingness value
- How to suggest the rules that we got to reddit users

228

## Conclusion

- We cannot be sure that our recommendation system will be 100% useful to user since the interestingness can vary depending on the purpose of the experiment
- To getting more accurate result, we need to ask all generated association rules, more than 300 rules in survey

229



## References

- Liu, B. (2000). Analyzing the subjective interestingness of association rules. 15(5), 47-55.  
doi:10.1109/5254.889106
- Calculating Lift, How We Make Smart Online Product Recommendations  
<https://www.youtube.com/watch?v=DeZFe1LerA>

Q)

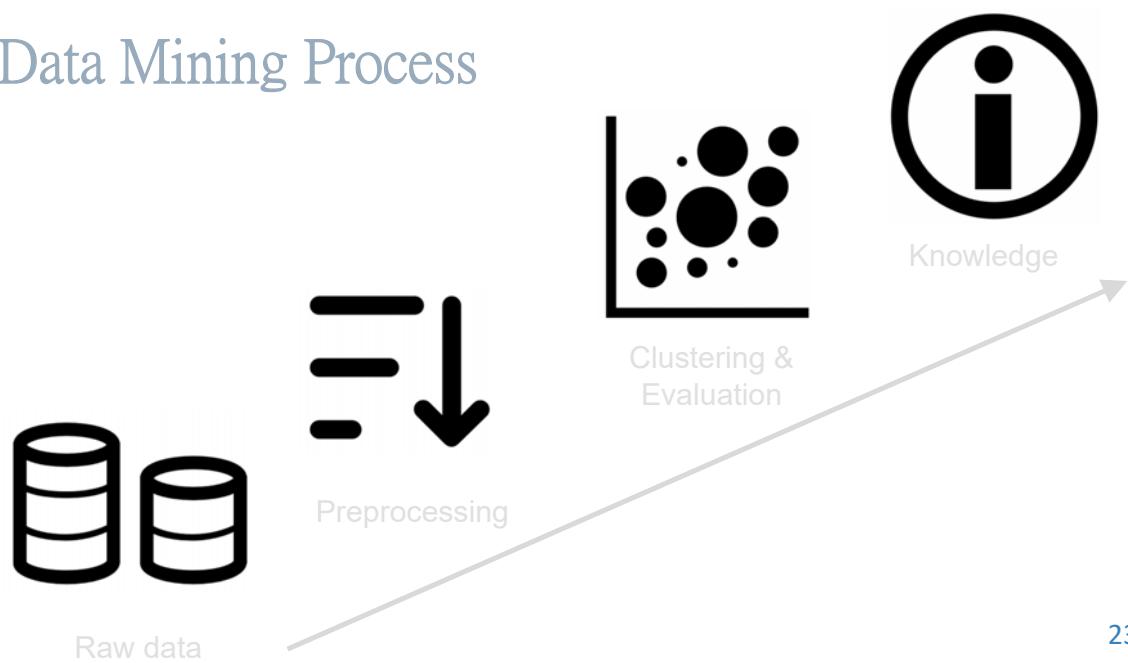


230

# Application of Data Mining Techniques on Active Users in Reddit

231

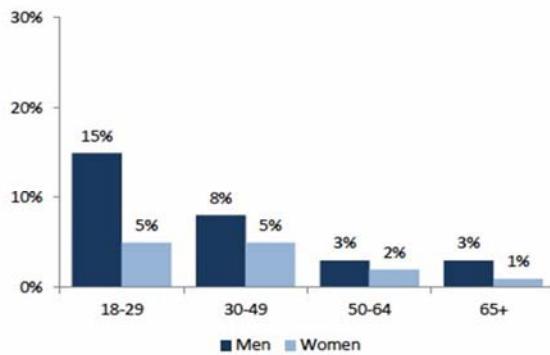
## Data Mining Process



232

## Facts - Why Active Users?

Young males are especially likely to use reddit  
% of internet users in each age/gender grouping who use reddit



Source: Pew Research Center's Internet & American Life Project Spring Tracking Survey, April 17 – May 19, 2013. N=2,252 adults ages 18+. Interviews were conducted in English and Spanish and on landline and cell phones. The margin of error for results based on all internet users is +/- 2.5 percentage points.

233

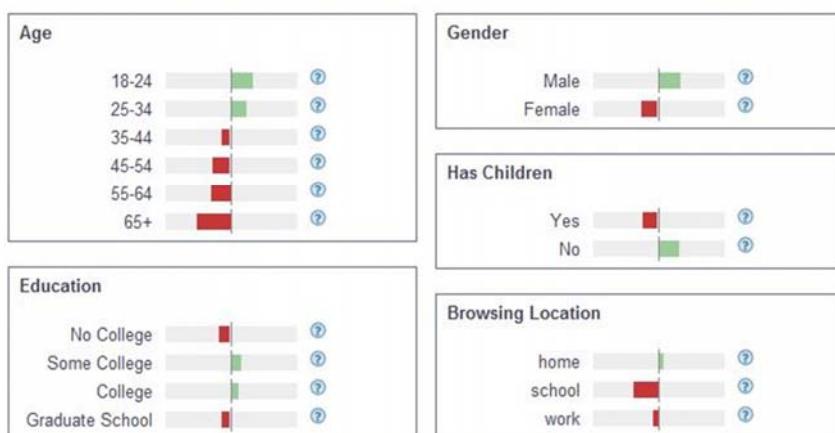
\*\* <http://www.pewinternet.org/2013/07/03/6-of-online-adults-are-reddit-users/>



## Facts - Why Active Users?

### Audience Demographics for Reddit.com

Relative to the general internet population how popular is reddit.com with each audience below?



\*\* <http://www.snoosecret.com/statistics-about-reddit.html>

234



## Active Users Definitions

- ▷ We define active users as people who :
  1. who had **posted or commented in at least 5 subreddits**
  2. who has **at least 5 Posts or Comments** in each of the **subreddits**
  3. # of Users' comments **above Q3**
  4. **Average Score of the User**, who satisfies three criteria **above > Q3**

235



## Preprocessing

- ▷ Total **posts** in **May2015**: 54,504,410
- ▷ Total **distinct authors** in **May2015**: 2,611,449
- ▷ After deleting BOTs, [deleted], non-English, only-URL posts, and length < 3 posts, we got 46,936,791 rows and 2,514,642 distinct authors.

- ▷ Finally, we extracted **25,298 “active users”** (0.97%)
- ▷ and **5,007,845 posts (9.18%)** by our active users

*definitions*

236



# Clustering: K-means, K- prototype

- ▷ # of clusters (K) = 10
- ,  $k = \sqrt{n/2}$ ,  $k = \sqrt{\sqrt{n/2}}$
- ▷ using python sklearn module - *KMeans()*, *open source* -

Attributes	Type	Description
subreddit	nominal	The title of the subreddit.
time	nominal	The median time that the author posted comments. It ranges from 0 to 23.
activity	ratio	(number of user's comment in the subreddit they commented most) / (number of user's comment in total) It indicates how much the user devoted their time in the subreddit they commented most.
assurance	ratio	(user's median score within the subreddit) / (median score of the subreddit) It indicates their similarity with others in the subreddit they commented most.

237

## Attributes

author	subreddit	frequency	median score
1	A	11	5
1	B	6	2
1	C	27	3
1	D	17	0
1	E	3	1
2	C	19	1
2	D	5	1



author = 1

subreddit = C (frequency: 27>others)

activity = 27/(11+6+27+17+3) = 0.42

assurance = 3/1 = 3

subreddit	A	B	C	D	E	F
median score	2	3	1	1	0	4



238

## Clustering: K-means, K-prototype

Attributes	Type	Description
subreddit	nominal	The title of the subreddit.
time	ratio	The median time that the author posted comments. It ranges from 0 to 23.
activity	ratio	(number of user's comment in the subreddit they commented most) / (number of user's comment in total) It indicates how much the user devoted their time in the subreddit they commented most.
assurance	ratio	(user's median score within the subreddit) / (median score of the subreddit) It indicates their similarity with others in the subreddit they commented most.

239

## Clustering: K-means, K-prototype

Attributes	Type	Description
subreddit	nominal	The title of the subreddit.
time	nominal	The median time that the author posted comments. It ranges from 0 to 23.
activity	ratio	(number of user's comment in the subreddit they commented most) / (number of user's comment in total) It indicates how much the user devoted their time in the subreddit they commented most.
assurance	ratio	(user's median score within the subreddit) / (median score of the subreddit) It indicates their similarity with others in the subreddit they commented most.

240

## Evaluation

data	A	B	C
data	A	B	C

- ▷ Separate data into 3 parts,
- ▷ K-means && K-prototype for AB, AC
- ▷ Compare labels of A which was from clustering result of AB, AC.

### ▷ Measurement:

- **Adjusted Rand index:** measure similarity between two list. (Ex. AB-A & AC-A)
- **Homogeneity:** each cluster contains only members of a single class
- **Completeness:** all members of a given class are assigned to the same cluster. 241
- **V-measure:** harmonic mean of homogeneity & completeness



## Evaluation: K-means v.s. K-prototype

	K-means	K-prototype
Adjusted Rand index	0.510904	0.193399
Homogeneity	0.803911	0.326116
Completeness	0.681669	0.298049
V-measure	0.737761	0.311451

242



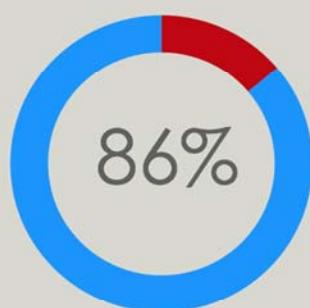
# VISUALIZATION

243



## Reddit Active Users Analysis for May 2015

Percentage of Reddit May Users



Summary

Total Posts During May 2015

54504410

Posts from the Normal users

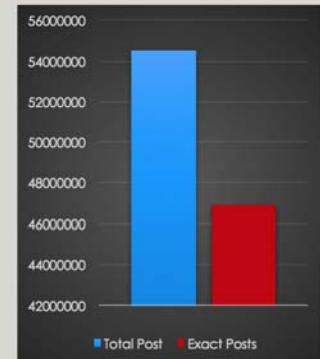
46936791

BOT / Spam Users

401

Total Registered Users

2611449



Monthly Users

Total Users ▾ Amount ▾

Total Registered Users 2,611,449

Spam / Bot Users 410

Normal Users 2,611,048

244



# The Active Users in Reddit

Active Users



Summary

TOTAL USERS

2,611,449

ACTIVE USERS

25,298

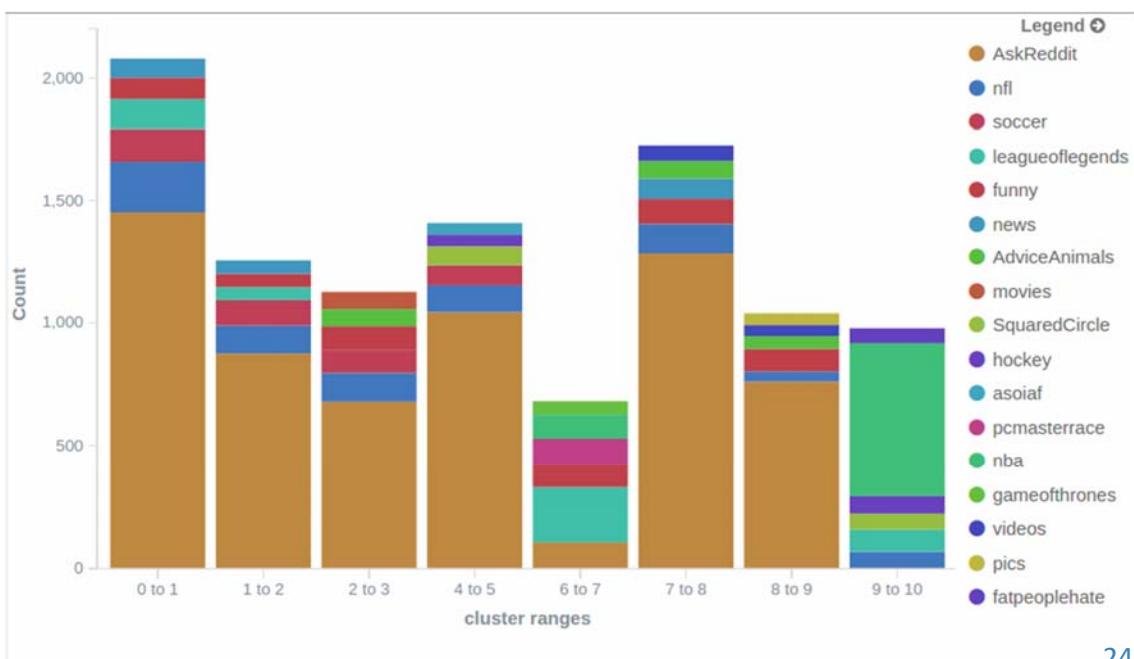
ACTIVE USERS RATIO

1%

USERS WHO SHOWN UP DURING MAY 2015

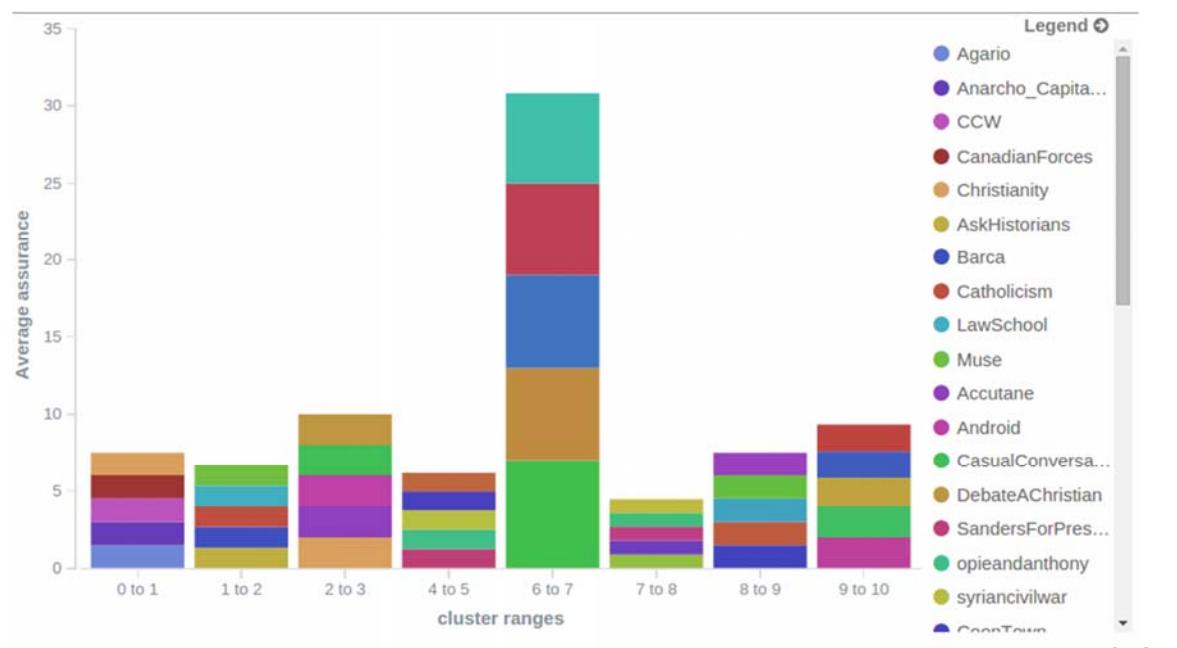
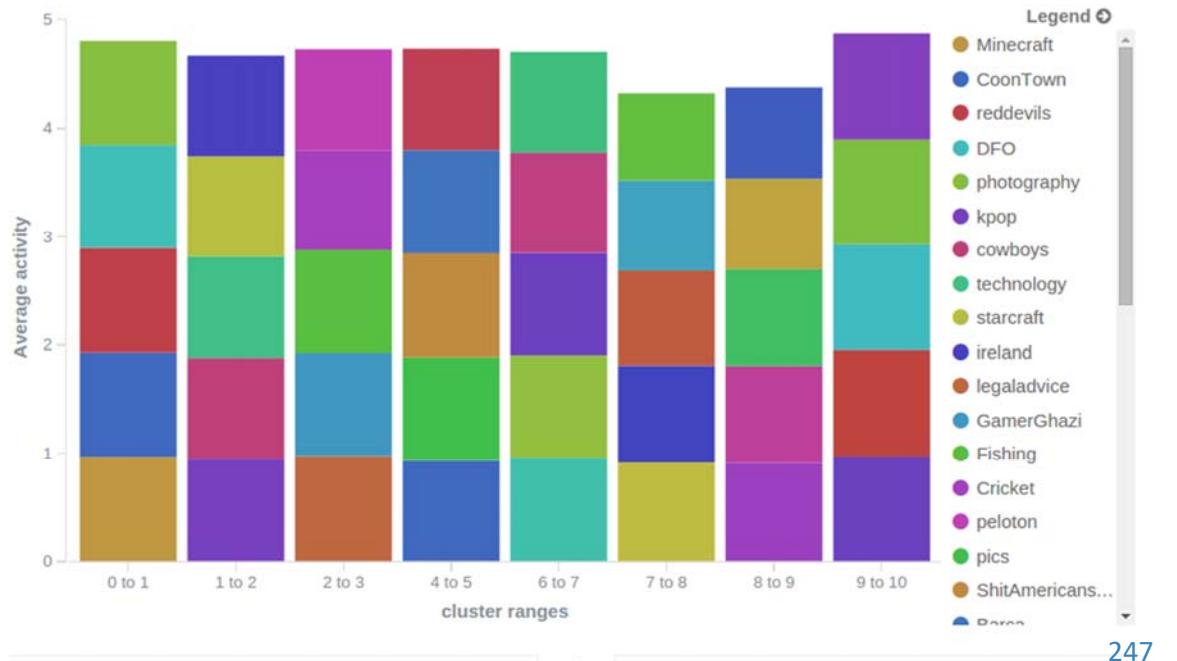
99%

245



246





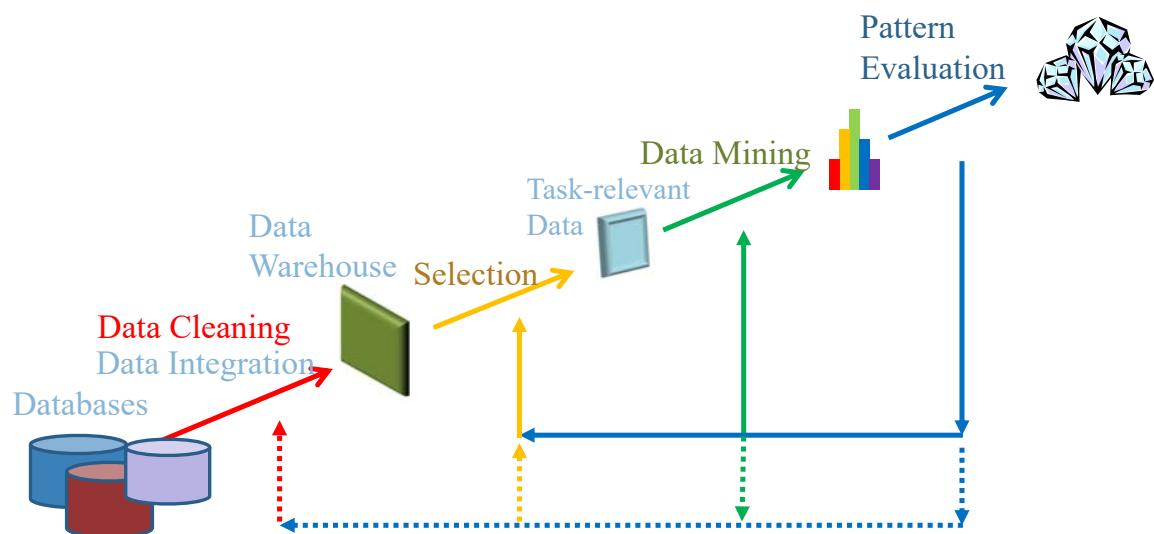
# Some Remarks

Some but not all

@ Yi-Shin Chen, Text Mining Overview

249

## Knowledge Discovery (KDD) Process



250

# Always Remember

▷ Have a good and solid objective

- No goal no gold
- Know the relationships between them