



# TwisSent: A Multi-Stage System for Analyzing Sentiment in Twitter

Subhabrata Mukherjee, Akshat Malu, Balamurali A.R.  
and Pushpak Bhattacharyya

Dept. of Computer Science and Engineering, IIT Bombay

21st ACM Conference on Information and Knowledge Management

**CIKM 2012,**

Hawai, Oct 29 - Nov 2, 2012

# Social Media Analysis

2

*Had Hella fun today with the team. Y'all are hilarious! & Yes, i do need more black homies.....*

# Social Media Analysis

3

- Social media sites, like Twitter, generate around 250 million tweets daily

*Had Hella fun today with the team. Y'all are hilarious! &Yes, i do need more black homies.....*

# Social Media Analysis

4

- Social media sites, like Twitter, generate around 250 million tweets daily
- This information content could be leveraged to create applications that have a social as well as an economic value

*Had Hella fun today with the team. Y'all are hilarious! &Yes, i do need more black homies.....*

# Social Media Analysis

5

- Social media sites, like Twitter, generate around 250 million tweets daily
- This information content could be leveraged to create applications that have a social as well as an economic value
- Text limit of 140 characters per tweet makes Twitter a noisy medium
  - ▣ Tweets have a poor syntactic and semantic structure
  - ▣ Problems like *slangs, ellipses, nonstandard vocabulary etc.*

*Had Hella fun today with the team. Y'all are hilarious! &Yes, i do need more black homies.....*

# Social Media Analysis

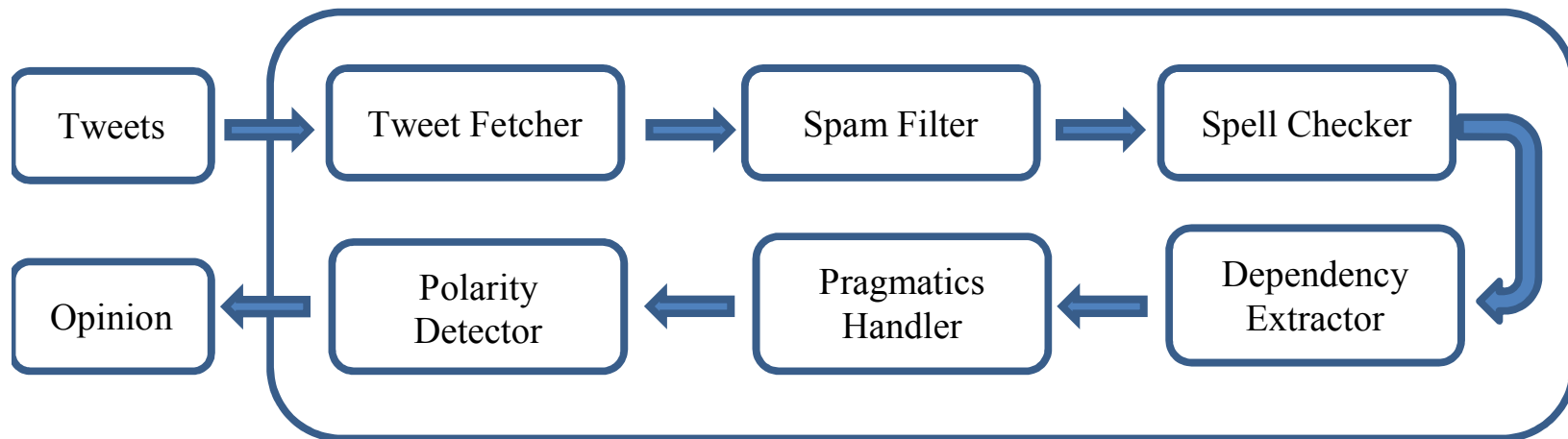
6

- Social media sites, like Twitter, generate around 250 million tweets daily
- This information content could be leveraged to create applications that have a social as well as an economic value
- Text limit of 140 characters per tweet makes Twitter a noisy medium
  - ▣ Tweets have a poor syntactic and semantic structure
  - ▣ Problems like *slangs, ellipses, nonstandard vocabulary etc.*
- Problem is compounded by increasing number of *spams* in Twitter
  - ▣ *Promotional tweets, bot-generated tweets, random links to websites etc.*  
In fact Twitter contains around 40% tweets as pointless babble

*Had Hella fun today with the team. Y'all are hilarious! &Yes, i do need more black homies.....*

# TwisSent: Multi-Stage System Architecture

7



# Spam Categorization and Features

8

- *Re-tweets*
- *Promotional tweets for some entity*
- *Tweets containing links to some other websites*
- *Tweets in languages other than English*
- *Tweets with incomplete text*
- *Automatically generated tweets by bots*
- *Tweets built primarily for search engines or tweets with excessive off-topic keywords*
- *Multiple tweets offering substantially the same content*

<ol style="list-style-type: none"><li>1. <i>Number of Words per Tweet</i></li><li>2. <i>Average Word Length</i></li><li>3. <i>Frequency of “?” and “!”</i></li><li>4. <i>Frequency of Numeral Characters</i></li><li>5. <i>Frequency of hashtags</i></li><li>6. <i>Frequency of @users</i></li><li>7. <i>Extent of Capitalization</i></li><li>8. <i>Frequency of the First POS Tag</i></li></ol>	<ol style="list-style-type: none"><li>9. <i>Frequency of Foreign Words</i></li><li>10. <i>Validity of First Word</i></li><li>11. <i>Presence / Absence of links</i></li><li>12. <i>Frequency of POS Tags</i></li><li>13. <i>Strength of Character Elongation</i></li><li>14. <i>Frequency of Slang Words</i></li><li>15. <i>Average Positive and Negative Sentiment of Tweets</i></li></ol>
--	---



# Algorithm for Spam Filter

9

Input: Build an initial naive bayes classifier NB-  $C$ , using the tweet sets  $M$  (mixed unlabeled set containing spams and non-spams) and  $P$  (labeled non-spam set)

- 1: Loop while classifier parameters change
- 2: for each tweet  $t_i \in M$  do
- 3: Compute  $\Pr[c_1 | t_i]$ ,  $\Pr[c_2 | t_i]$  using the current NB //  $c_1$  - non-spam class ,  $c_2$  - spam class
- 4:  $\Pr[c_2 | t_i] = 1 - \Pr[c_1 | t_i]$
- 5: Update  $\Pr[f_{i,k} | c_1]$  and  $\Pr[c_1]$  given the probabilistically assigned class for all  $t_i$  ( $\Pr[c_1 | t_i]$ ).  
(a new NB- $C$  is being built in the process)
- 6: end for
- 7: end loop

$$\Pr[c_j | t_i] = \frac{\Pr[c_j] \prod_k \Pr[f_{i,k} | c_j]}{\sum_r \Pr[c_r] \prod_k \Pr[f_{i,k} | c_r]}$$

# Categorization of Noisy Text

10

- Dropping of Vowels - *btfl* (*beautiful*), *lvng* (*loving*)
- Vowel Exchange - *good* vs. *gud* (*o, u*)
- Mis-spelt words - *redicule* (*ridicule*), *magnificant* (*magnificent*)
- Text Compression - *shok* (*shock*), *terorism* (*terrorism*)
- Phonetic Transformation - *be8r* (*better*), *gud* (*good*), *fy9* (*fine*), *gr8* (*great*)
- Normalization and Pragmatics - *hapyyyyyy* (*happy*), *guuuuud* (*good*)
- Segmentation with Punctuation - *beautiful*, (*beautiful*)
- Segmentation with Compound Words - *breathtaking* (*breath-taking*), *eyecatching* (*eye-catching*), *good-looking* (*good looking*)
- Hashtags and Segmentation - *#notevenkidding*, *#worthawatch*
- Combination of all - *#awsummm* (*awesome*), *gr88888* (*great*), *amzng, btfl* (*amazing, beautiful*).

# Spell-Checker Algorithm

# Spell-Checker Algorithm

12

- Heuristically driven to resolve the *identified errors* with a *minimum edit distance based spell checker*

# Spell-Checker Algorithm

13

- Heuristically driven to resolve the *identified errors* with a *minimum edit distance based spell checker*
- A *normalize* function takes care of *Pragmatics* and *Number Homophones*
  - ▣ Replaces *happppyyy* with *hapy*, '2' with 'to', '8' with 'eat', '9' with 'ine'

# Spell-Checker Algorithm

14

- Heuristically driven to resolve the *identified errors* with a *minimum edit distance based spell checker*
- A *normalize* function takes care of *Pragmatics* and *Number Homophones*
  - ▣ Replaces *happppyyy* with *hapy*, '2' with 'to', '8' with 'eat', '9' with 'ine'
- A *vowel\_dropped* function takes care of the *vowel dropping* phenomenon

# Spell-Checker Algorithm

15

- Heuristically driven to resolve the *identified errors* with a *minimum edit distance based spell checker*
- A *normalize* function takes care of *Pragmatics* and *Number Homophones*
  - ▣ Replaces *happppyyyy* with *hapy*, '2' with 'to', '8' with 'eat', '9' with 'ine'
- A *vowel\_dropped* function takes care of the *vowel dropping* phenomenon
- The parameters *offset* and *adv* are determined empirically

# Spell-Checker Algorithm

16

- Heuristically driven to resolve the *identified errors* with a *minimum edit distance based spell checker*
- A *normalize* function takes care of *Pragmatics* and *Number Homophones*
  - ▣ Replaces *happppyyy* with *hapy*, '2' with 'to', '8' with 'eat', '9' with 'ine'
- A *vowel\_dropped* function takes care of the *vowel dropping* phenomenon
- The parameters *offset* and *adv* are determined empirically
- Words are marked during normalization, to preserve their pragmatics  
*happppyyyyy*, normalized to *hapy* and thereafter spell-corrected to *happy*, is marked so as to not lose its pragmatic content



# Spell-Checker Algorithm

17

- **Input:** For string  $s$ , let  $S$  be the set of words in the lexicon starting with the initial letter of  $s$ .
- */\* Module Spell Checker \*/*
- **for** each word  $w \in S$  **do**
- $w' = \text{vowel\_dropped}(w)$
- $s' = \text{normalize}(s)$
- */\*diff(s,w) gives difference of length between s and w\*/*
- **if**  $\text{diff}(s', w') < \text{offset}$  **then**
- $\text{score}[w] = \min(\text{edit\_distance}(s, w), \text{edit\_distance}(s, w'), \text{edit\_distance}(s', w))$
- **else**
- $\text{score}[w] = \text{max\_centinel}$
- **end if**
- **end for**

# Spell-Checker Algorithm

## Contd..

18

- Sort score of each  $w$  in the Lexicon and retain the top  $m$  entries in suggestions( $s$ ) for the original string  $s$
- **for** each  $t$  in suggestions( $s$ ) **do**
- $\text{edit}_1 = \text{edit\_distance}(s', s)$
- */\*t.replace(char1,char2) replaces all occurrences of char1 in the string t with char2\*/*
- $\text{edit}_2 = \text{edit\_distance}(t.\text{replace}(a, e), s')$
- $\text{edit}_3 = \text{edit\_distance}(t.\text{replace}(e, a), s')$
- $\text{edit}_4 = \text{edit\_distance}(t.\text{replace}(o, u), s')$
- $\text{edit}_5 = \text{edit\_distance}(t.\text{replace}(u, o), s')$
- $\text{edit}_6 = \text{edit\_distance}(t.\text{replace}(i, e), s')$
- $\text{edit}_7 = \text{edit\_distance}(t.\text{replace}(e, i), s')$
- $\text{count} = \text{overlapping\_characters}(t, s')$
- $\text{min\_edit} =$
- $\min(\text{edit}_1, \text{edit}_2, \text{edit}_3, \text{edit}_4, \text{edit}_5, \text{edit}_6, \text{edit}_7)$
- **if** ( $\text{min\_edit} == 0$  or  $\text{score}[s] == 0$ ) **then**
- $\text{adv} = -2$  */\* for exact match assign advantage score \*/*
- **else**
- $\text{adv} = 0$
- **end if**
- $\text{final\_score}[t] = \text{min\_edit} + \text{adv} + \text{score}[w] - \text{count};$
- **end for**
- return  $t$  with minimum final\_score;

# Feature Specific Tweet Analysis

- *I have an **ipod** and it is a great buy but I'm probably the only person that dislikes the **iTunes software**.*
- ▣ *Here the sentiment w.r.t ipod is positive whereas that respect to software is negative*

# Opinion Extraction Hypothesis

---

- *“More closely related words come together to express an opinion about a feature”*

# Hypothesis Example

- *“I want to use Samsung which is a great product but am not so sure about using Nokia”.*
  - ▣ *Here “great” and “product” are related by an adjective modifier relation, “product” and “Samsung” are related by a relative clause modifier relation. Thus “great” and “Samsung” are transitively related.*
  - ▣ ***Here “great” and “product” are more related to Samsung than they are to Nokia***
  - ▣ *Hence “great” and “product” come together to express an opinion about the entity “Samsung” than about the entity “Nokia”*

# Hypothesis Example

- *“I want to use Samsung which is a great product but am not so sure about using Nokia”.*
- ▣ *Here “great” and “product” are related by an adjective modifier relation, “product” and “Samsung” are related by a relative clause modifier relation. Thus “great” and “Samsung” are transitively related.*
- ▣ ***Here “great” and “product” are more related to Samsung than they are to Nokia***
- ▣ *Hence “great” and “product” come together to express an opinion about the entity “Samsung” than about the entity “Nokia”*

# Hypothesis Example

□ *“I want to use Samsung which is a great product but am not so sure about using Nokia”.*

- ▣ *Here “great” and “product” are related by an adjective modifier relation, “product” and “Samsung” are related by a relative clause modifier relation. Thus “great” and “Samsung” are transitively related.*
- ▣ ***Here “great” and “product” are more related to Samsung than they are to Nokia***
- ▣ *Hence “great” and “product” come together to express an opinion about the entity “Samsung” than about the entity “Nokia”*

# Hypothesis Example

□ “I want to use Samsung which is a great product but am not so sure about using Nokia”.

- Here “great” and “product” are related by an adjective modifier relation, “product” and “Samsung” are related by a relative clause modifier relation. Thus “great” and “Samsung” are transitively related.
- **Here “great” and “product” are more related to Samsung than they are to Nokia**
- Hence “great” and “product” come together to express an opinion about the entity “Samsung” than about the entity “Nokia”



# Hypothesis Example

- *"I want to use Samsung which is a great product but am not so sure about using Nokia".*

Adjective Modifier

- ▣ *Here "great" and "product" are related by an adjective modifier relation, "product" and "Samsung" are related by a relative clause modifier relation. Thus "great" and "Samsung" are transitively related.*
- ▣ ***Here "great" and "product" are more related to Samsung than they are to Nokia***
- ▣ *Hence "great" and "product" come together to express an opinion about the entity "Samsung" than about the entity "Nokia"*

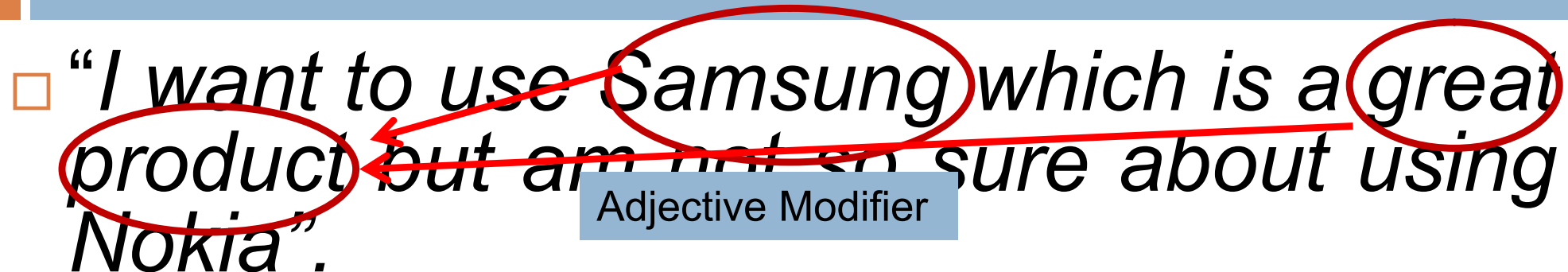
# Hypothesis Example

- “I want to use **Samsung** which is a **great** **product** but am not so sure about using Nokia”.

Adjective Modifier

- ▣ Here “great” and “product” are related by an adjective modifier relation, “product” and “Samsung” are related by a relative clause modifier relation. Thus “great” and “Samsung” are transitively related.
- ▣ **Here “great” and “product” are more related to Samsung than they are to Nokia**
- ▣ Hence “great” and “product” come together to express an opinion about the entity “Samsung” than about the entity “Nokia”

# Hypothesis Example

- “I want to use **Samsung** which is a **great** **product** but am not so sure about using Nokia”.
- Adjective Modifier
- 

- Here “great” and “product” are related by an adjective modifier relation, “product” and “Samsung” are related by a relative clause modifier relation. Thus “great” and “Samsung” are transitively related.
- **Here “great” and “product” are more related to Samsung than they are to Nokia**
- Hence “great” and “product” come together to express an opinion about the entity “Samsung” than about the entity “Nokia”

# Hypothesis Example

□ “I want **use Samsung** which is a **great** product but am not so sure about using Nokia”.

Relative Clause Modifier

Adjective Modifier

- Here “great” and “product” are related by an adjective modifier relation, “product” and “Samsung” are related by a relative clause modifier relation. Thus “great” and “Samsung” are transitively related.
- **Here “great” and “product” are more related to Samsung than they are to Nokia**
- Hence “great” and “product” come together to express an opinion about the entity “Samsung” than about the entity “Nokia”

# Example of a Review

---

- *I have an **ipod** and it is a great buy but I'm probably the only person that dislikes the **iTunes software**.*

# Example of a Review

---

- *I have an **ipod** and it is a great buy but I'm probably the only person that dislikes the iTunes **software**.*

# Example of a Review

---

- *I have an **ipod** and it is a great buy but I'm probably the only person that dislikes the **iTunes software**.*

# Example of a Review

---

- *I have an **ipod** and it is a **great** buy but I'm probably the only person that dislikes the **iTunes software**.*



# Example of a Review

---

- *I have an **ipod** and it is a great buy but I'm probably the only person that dislikes the **iTunes software**.*

# Feature Extraction : Domain Info Not Available



# Feature Extraction : Domain Info Not Available

- Initially, all the Nouns are treated as features and added to the *feature list* ***F***.

# Feature Extraction : Domain Info Not Available

- Initially, all the Nouns are treated as features and added to the *feature list* ***F***.
- ***F*** = { *ipod, buy, person, software* }

# Feature Extraction : Domain Info Not Available

- Initially, all the Nouns are treated as features and added to the *feature list*  $F$ .
- $F = \{ ipod, buy, person, software \}$
- Pruning the feature set
  - ▣ Merge 2 features if they are **strongly related**

# Feature Extraction : Domain Info Not Available

- Initially, all the Nouns are treated as features and added to the *feature list*  $F$ .
- $F = \{ ipod, buy, person, software \}$
- Pruning the feature set
  - ▣ Merge 2 features if they are **strongly related**
- “*buy*” merged with “*ipod*”, when target feature = “*ipod*”,
  - ▣ “*person, software*” will be ignored.

# Feature Extraction : Domain Info Not Available

- Initially, all the Nouns are treated as features and added to the *feature list*  $F$ .
- $F = \{ ipod, buy, person, software \}$
- Pruning the feature set
  - ▣ Merge 2 features if they are **strongly related**
- “*buy*” merged with “*ipod*”, when target feature = “*ipod*”,
  - ▣ “*person, software*” will be ignored.
- “*person*” merged with “*software*”, when target feature = “*software*”
  - ▣ “*ipod, buy*” will be ignored.

# Relations

- Direct Neighbor Relation
  - Capture **short range dependencies**
  - Any 2 consecutive words (such that none of them is a StopWord) are directly related
  - Consider a sentence  $S$  and 2 consecutive words .
  - If  $w_i, w_{i+1} \in \text{Stopwords}$  , then they are directly related.  $w_i, w_{i+1} \in S$
- Dependency Relation
  - Capture **long range dependencies**
  - Let *Dependency\_Relation* be the list of **significant relations**.
- Any 2 words  $w_i$  and  $w_j$  in  $S$  are directly related, if  $\exists D_i \text{ s.t. } D_i(w_i, w_j) \in \text{Dependency\_Relation}$



# Graph representation

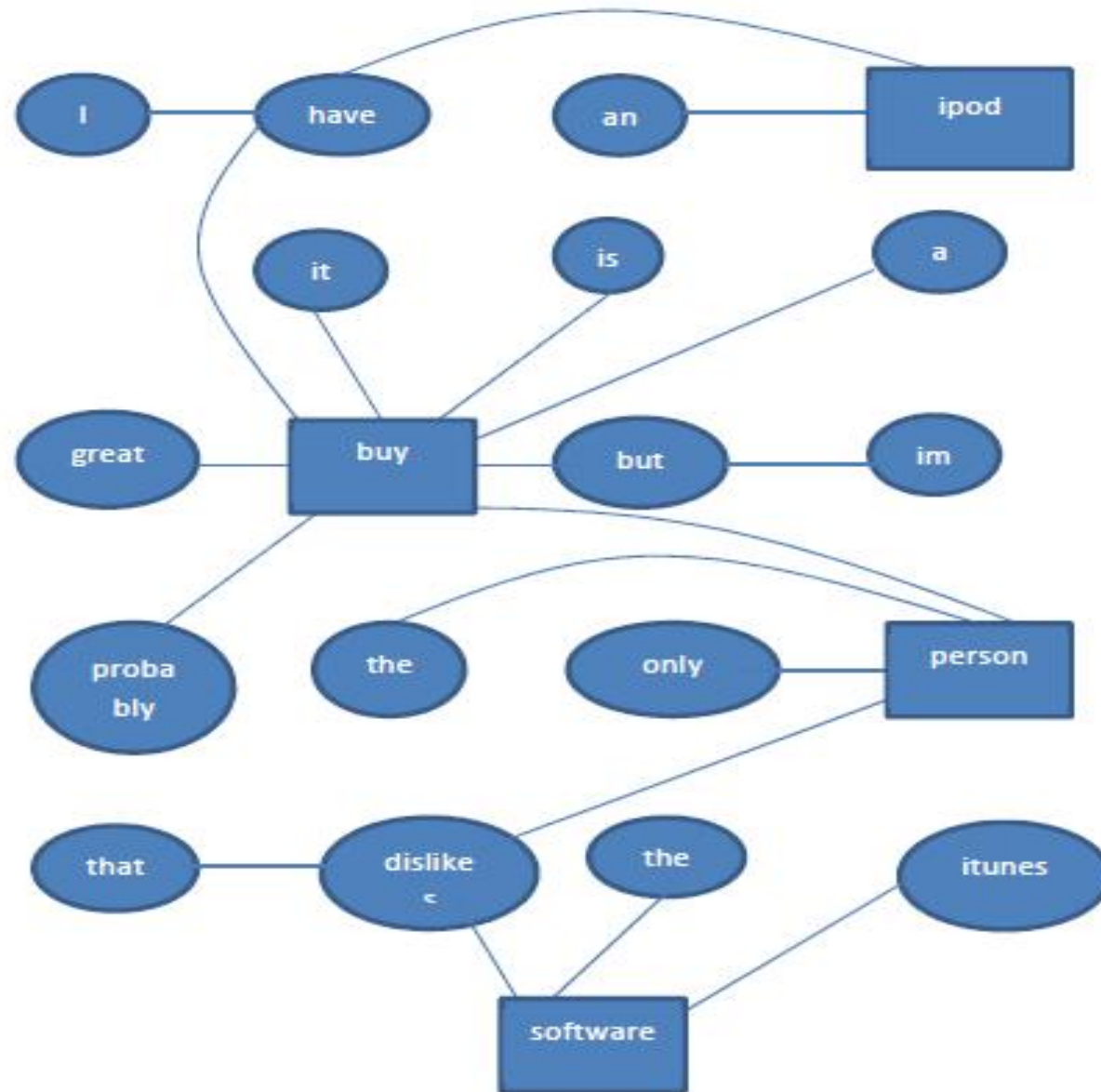
Given a sentence  $S$ , let  $W$  be the set of all words in the sentence  $S$ .

A Graph  $G(W, E)$  is constructed such that any

$w_i, w_j \in W$  are directly connected by  $e_k \in E$  ,if

$\exists R_l$  s.t.  $R_l(w_i, w_j) \in R$  .

# Graph



# Algorithm

- i. Initialize  $n$  clusters  $C_i \quad \forall i = 1..n$
- ii. Make each  $f_i \in F$  the clusterhead of  $C_i$ . The target feature  $f_t$  is the clusterhead of  $C_t$ . Initially, each cluster consists only of the clusterhead.

# Algorithm

# Contd...

iii. Assign each word  $w_j \in S$  to cluster  $C_k$

$$\text{s.t. } k = \arg \min_{i \in n} \text{dist}(w_j, f_i),$$

*Where  $\text{dist}(w_j, f_i)$  gives the number of edges, in the shortest path, connecting  $w_j$  and  $f_i$  in  $G$ .*

# Algorithm

## Contd...

- iv. Merge any cluster  $C_i$  with  $C_t$  if,  
 $dist(f_i, f_t) < \theta$ ,

Where  $\theta$  is some threshold distance.

- v. Finally the set of words  $w_i \in C_t$  gives  
the opinion expression regarding the  
target feature  $f_t$ .

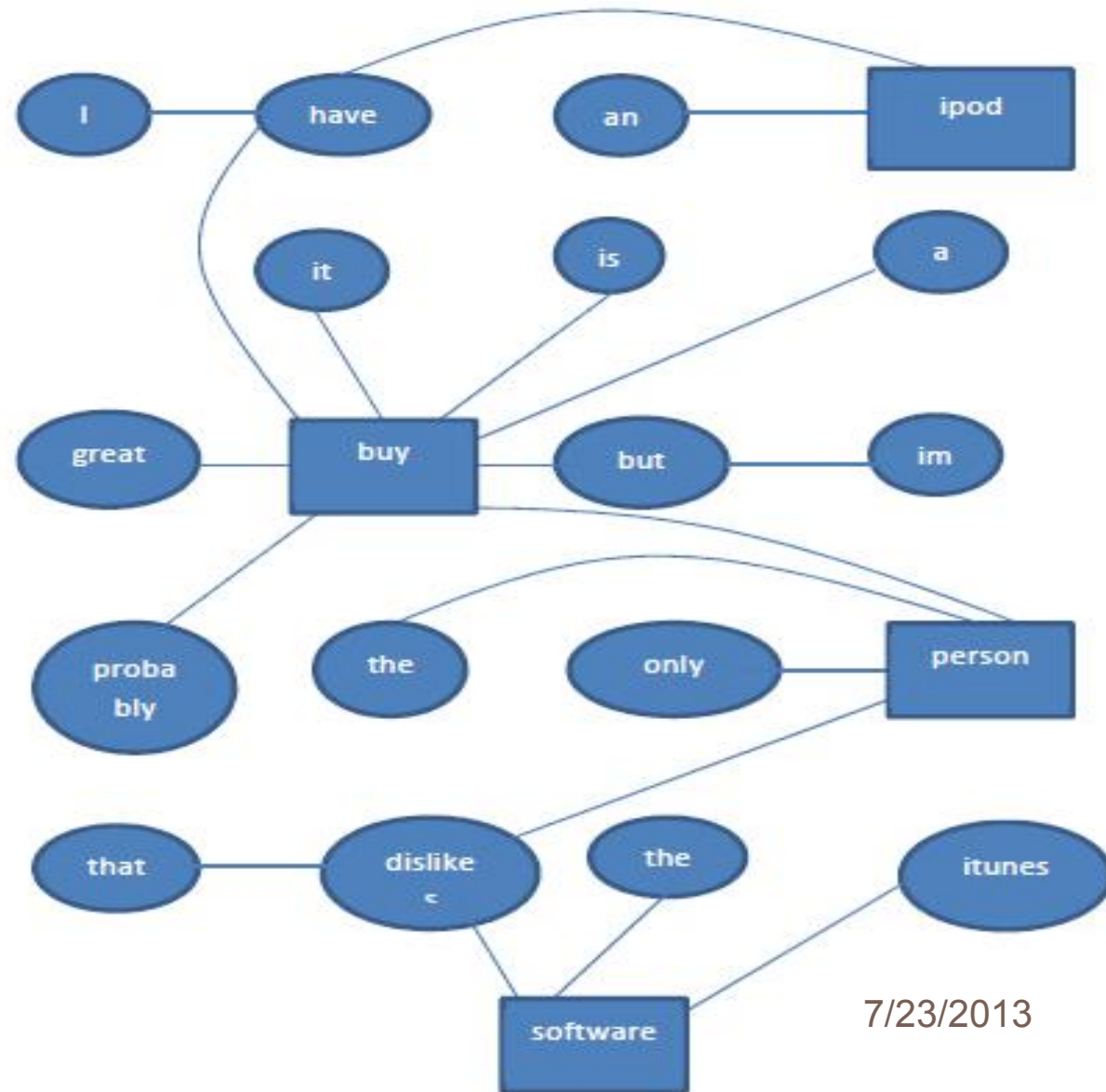
# Clustering

46

7/23/2013

# Clustering

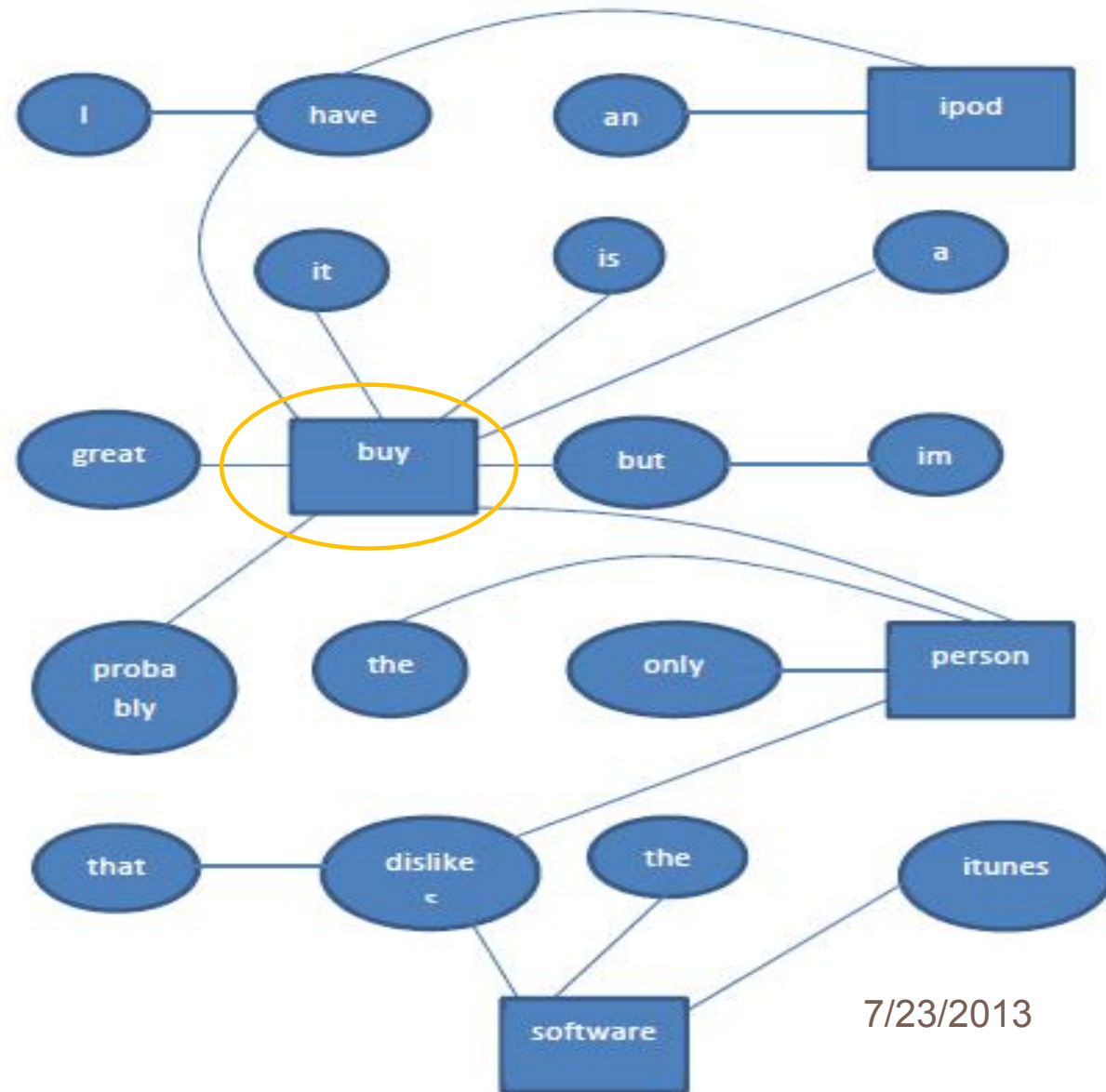
47



7/23/2013

# Clustering

48

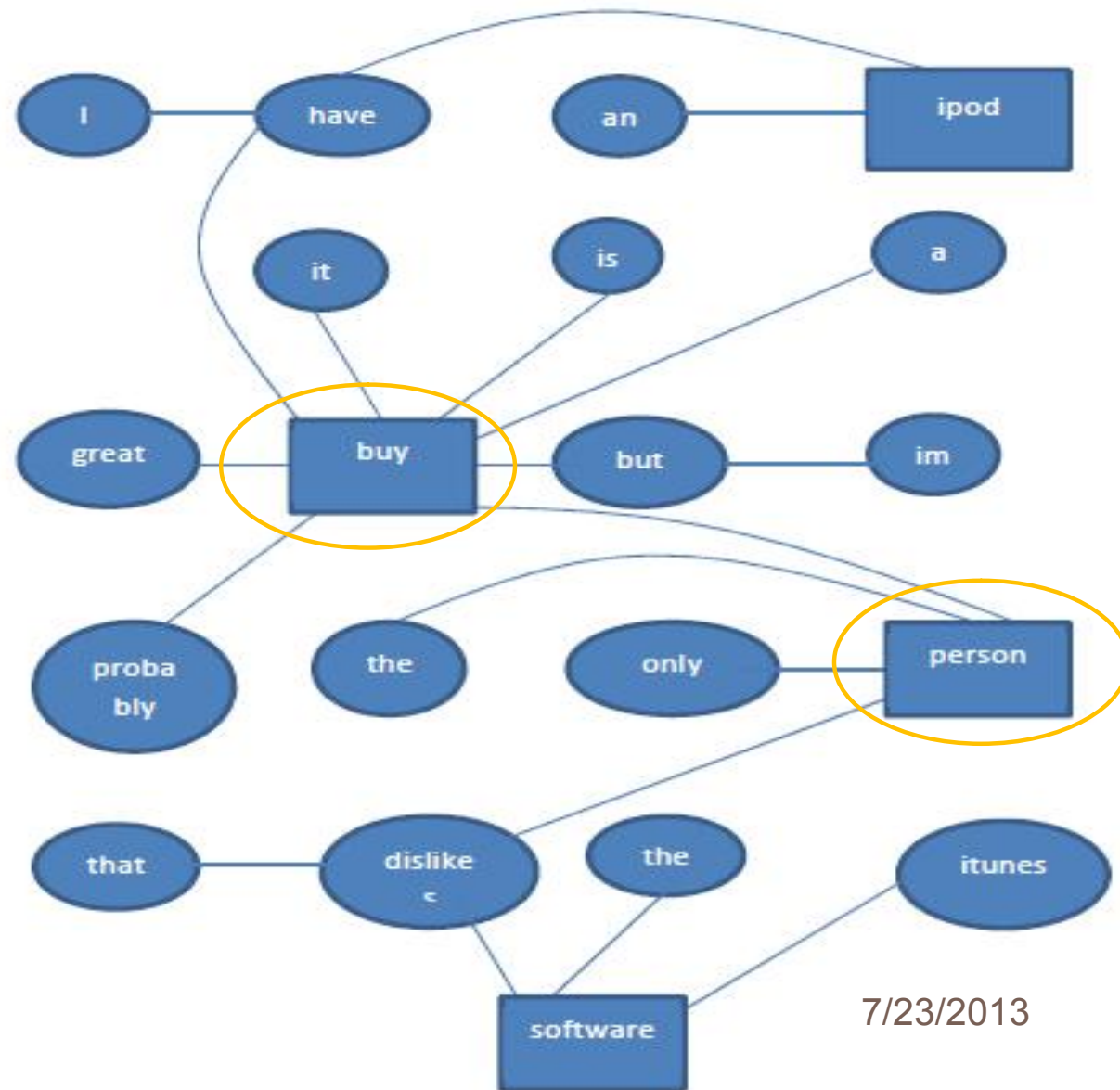


7/23/2013



# Clustering

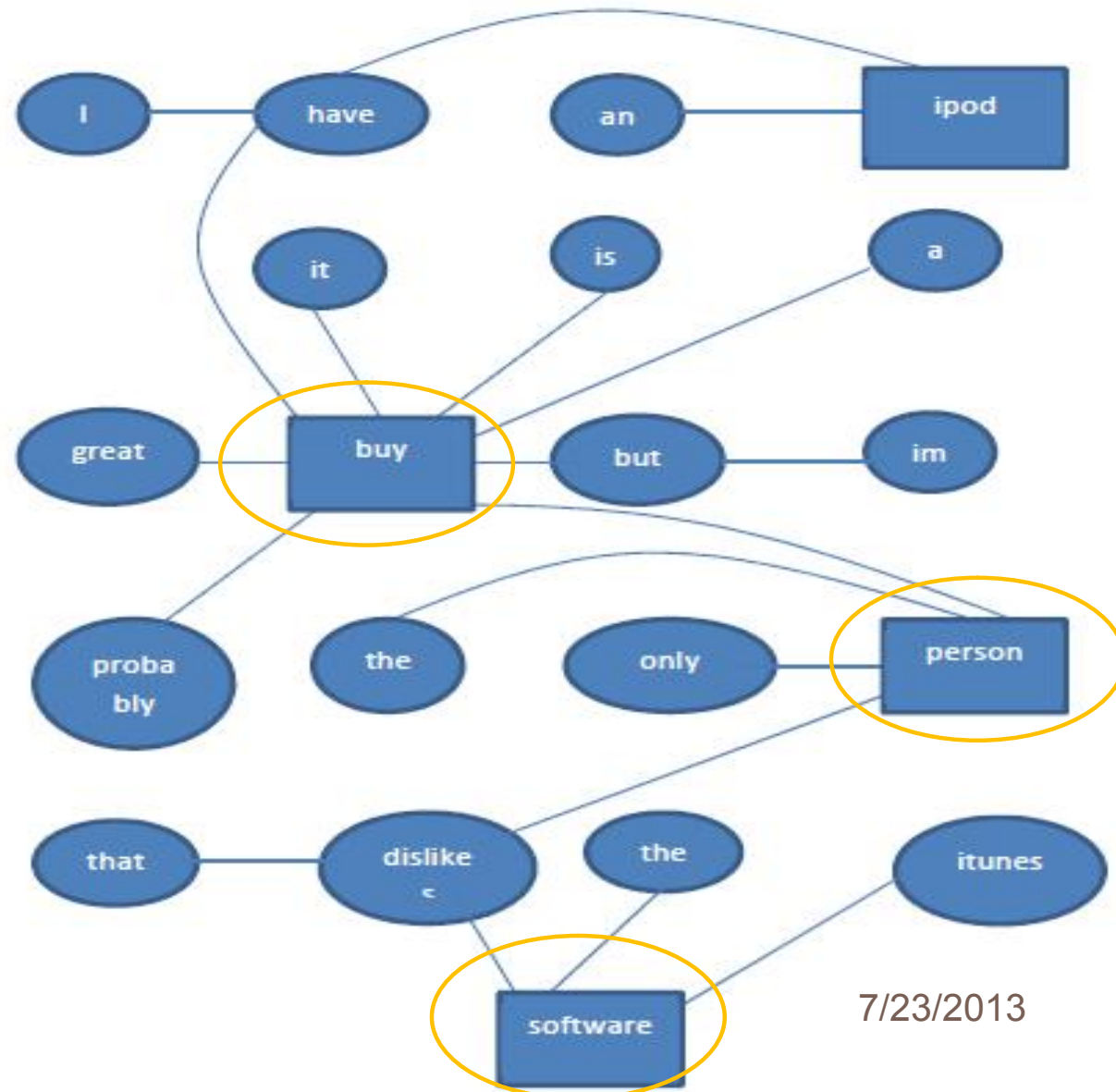
49



7/23/2013

# Clustering

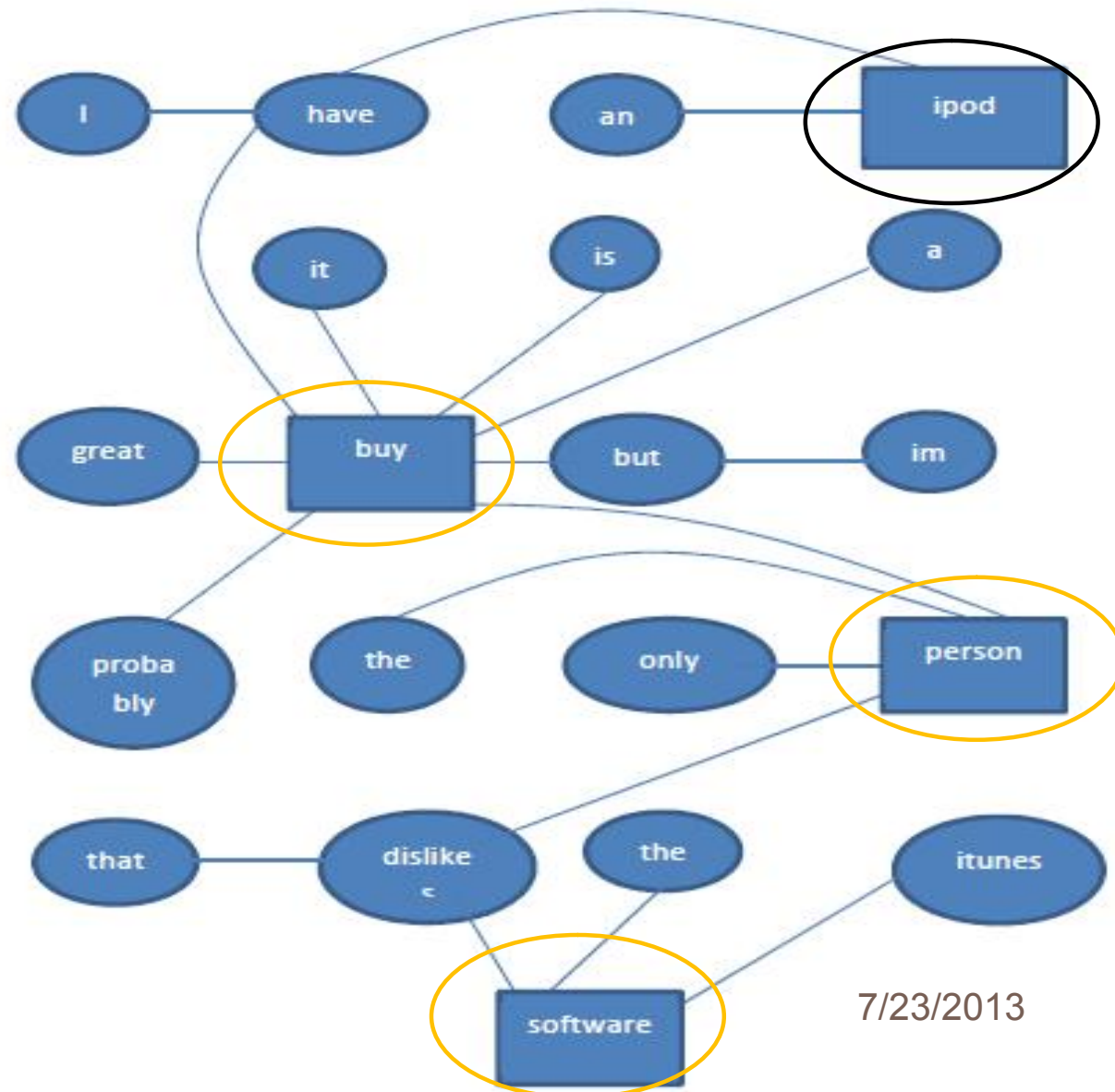
50



7/23/2013

# Clustering

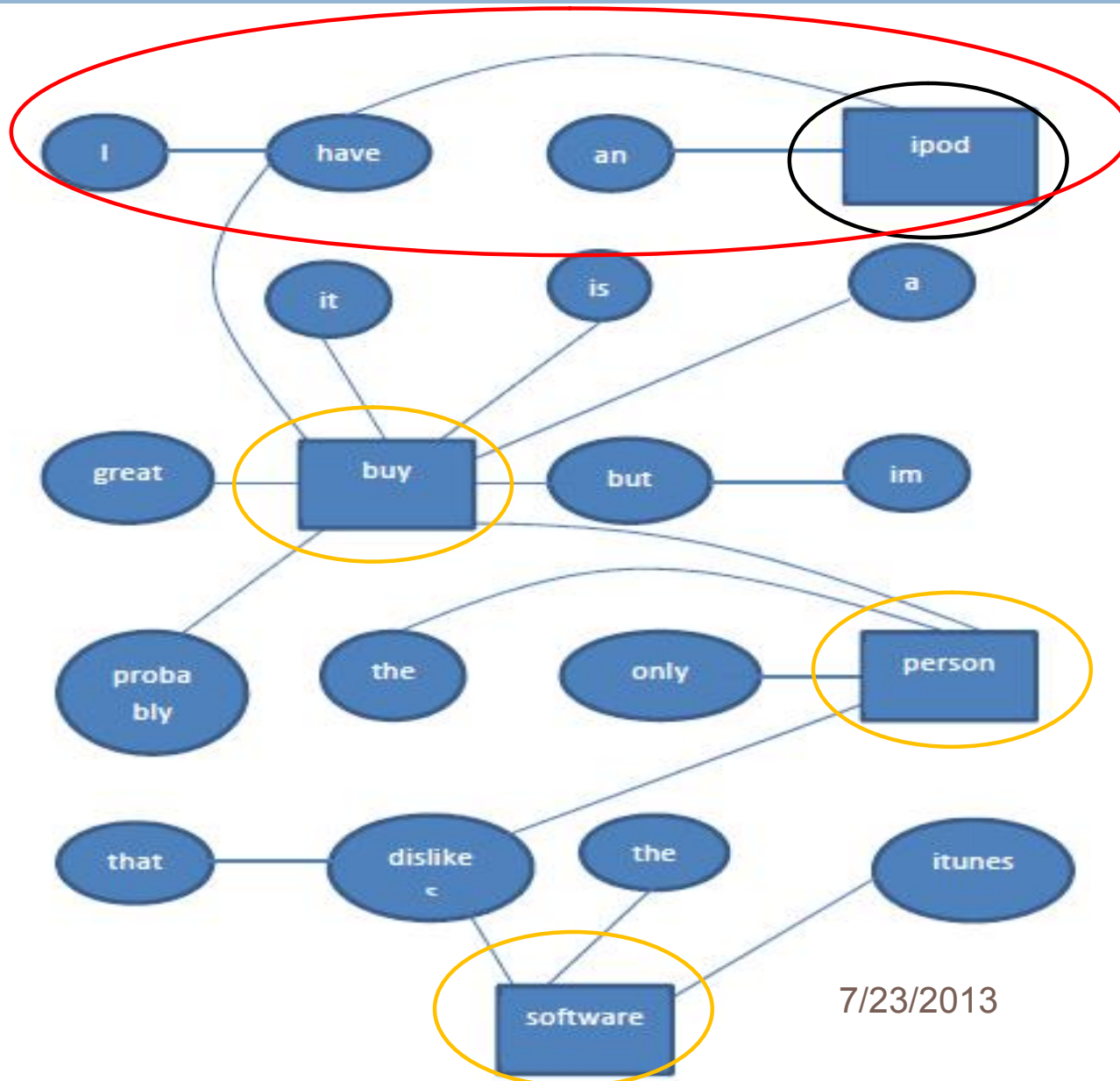
51



7/23/2013

# Clustering

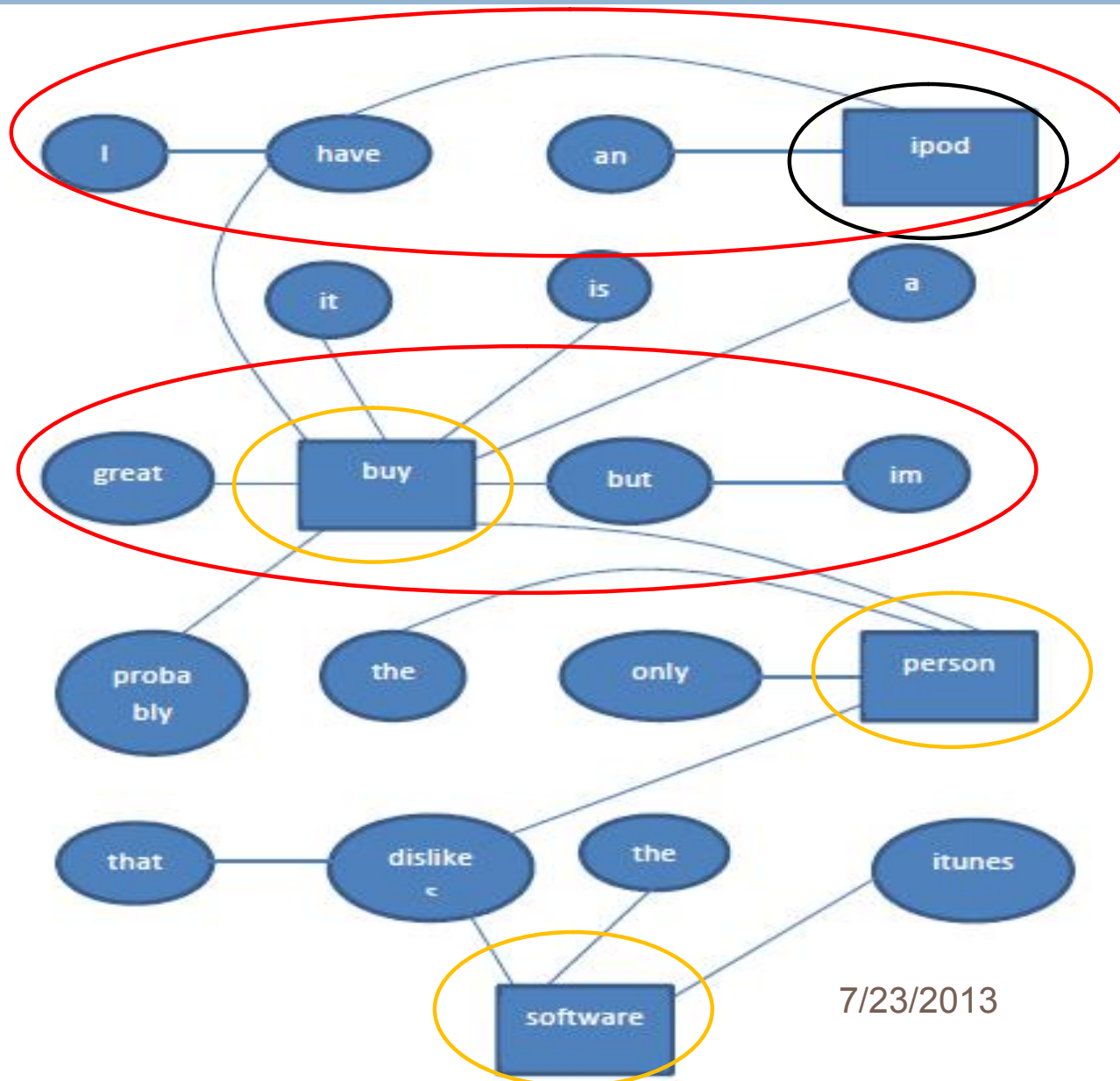
52



7/23/2013

# Clustering

53

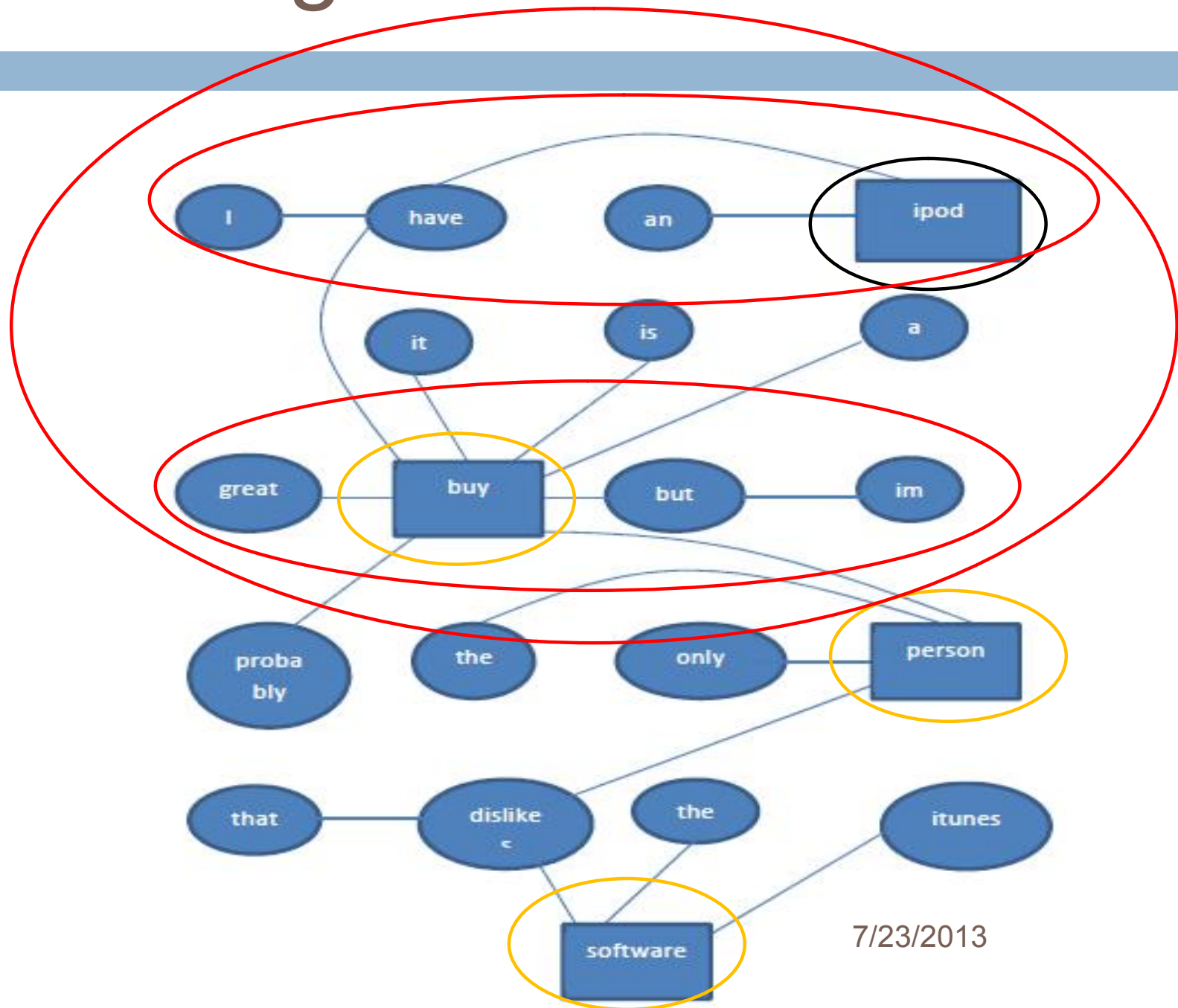


7/23/2013



# Clustering

54



7/23/2013

# Pragmatics

# Pragmatics

56

- *Elongation of a word, repeating alphabets multiple times - Example: **happpppyyyyyy, goooooood**. More weightage is given by repeating them twice*



# Pragmatics

57

- *Elongation of a word, repeating alphabets multiple times* - Example: *happpppyyyyyyy, goooooood*. More weightage is given by repeating them twice
- *Use of Hashtags* - *#overrated, #worthawatch*. More weightage is given by repeating them thrice

# Pragmatics

58

- *Elongation of a word, repeating alphabets multiple times* - Example: *happppppyyyyyy, goooooood*. More weightage is given by repeating them twice
- *Use of Hashtags* - *#overrated, #worthawatch*. More weightage is given by repeating them thrice
- *Use of Emoticons* - 😊 (*happy*), ☹️ (*sad*)

# Pragmatics

59

- *Elongation of a word, repeating alphabets multiple times* - Example: *happpppyyyyyy, goooooood*. More weightage is given by repeating them twice
- *Use of Hashtags* - *#overrated, #worthawatch*. More weightage is given by repeating them thrice
- *Use of Emoticons* - 😊 (*happy*), ☹️ (*sad*)
- *Use of Capitalization* - where words are written in capital letters to express intensity of user sentiments
  - ▣ *Full Caps* - Example: *I HATED that movie*. More weightage is given by repeating them thrice
  - ▣ *Partial Caps*- Example: *She is a Loving mom*. More weightage is given by repeating them twice

# Spam Filter Evaluation

60

## 2-Class Classification

Tweets	Total Tweets	Correctly Classified	Misclassified	Precision (%)	Recall (%)
All	7007	3815	3192	54.45	55.24
<b>Only spam</b>	<b>1993</b>	<b>1838</b>	<b>155</b>	<b>92.22</b>	<b>92.22</b>
Only non-spam	5014	2259	2755	45.05	-

## 4-Class Classification

Tweets	Total Tweets	Correctly Classified	Misclassified	Precision (%)	Recall (%)
All	7007	5010	1997	<b>71.50</b>	<b>54.29</b>
<b>Only spam</b>	1993	1604	389	80.48	80.48
Only non-spam	5014	4227	787	84.30	-

# TwisSent Evaluation

# TwisSent Evaluation

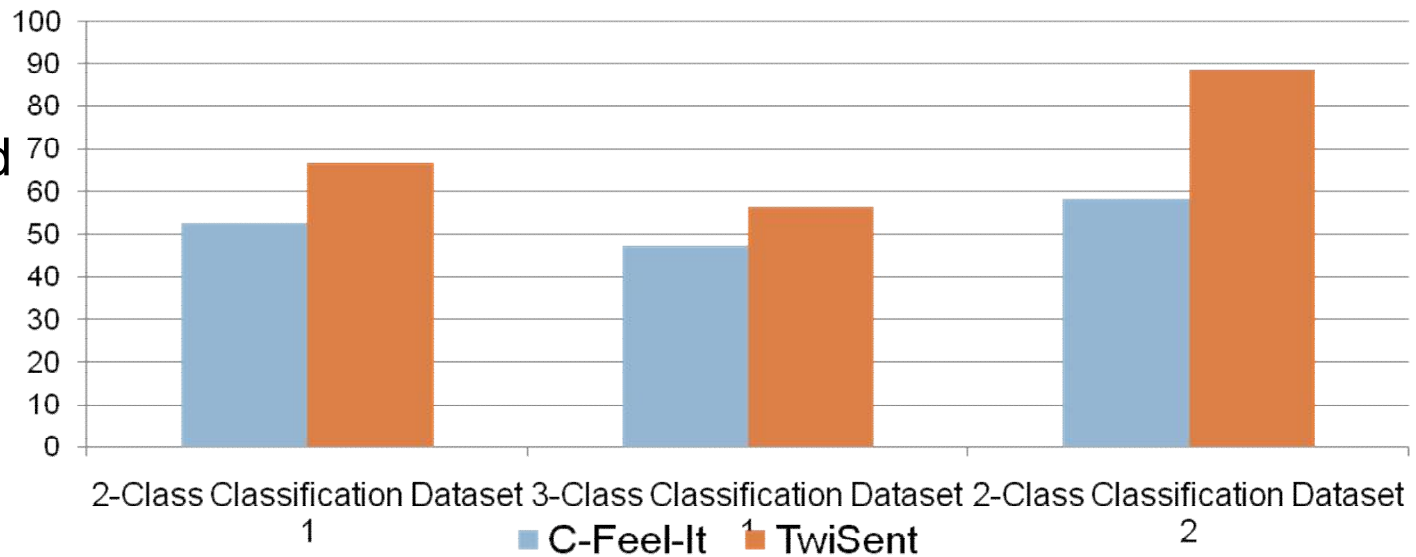
62

Lexicon-based  
Classification

# TwisSent Evaluation

63

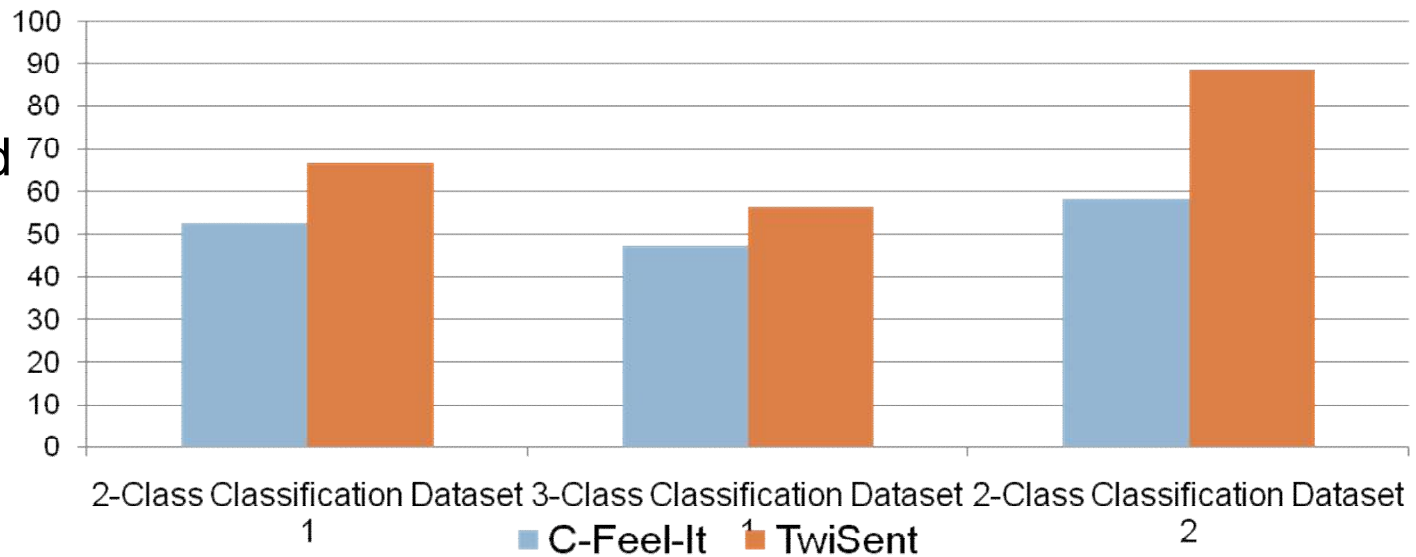
Lexicon-based  
Classification



# TwisSent Evaluation

64

Lexicon-based  
Classification



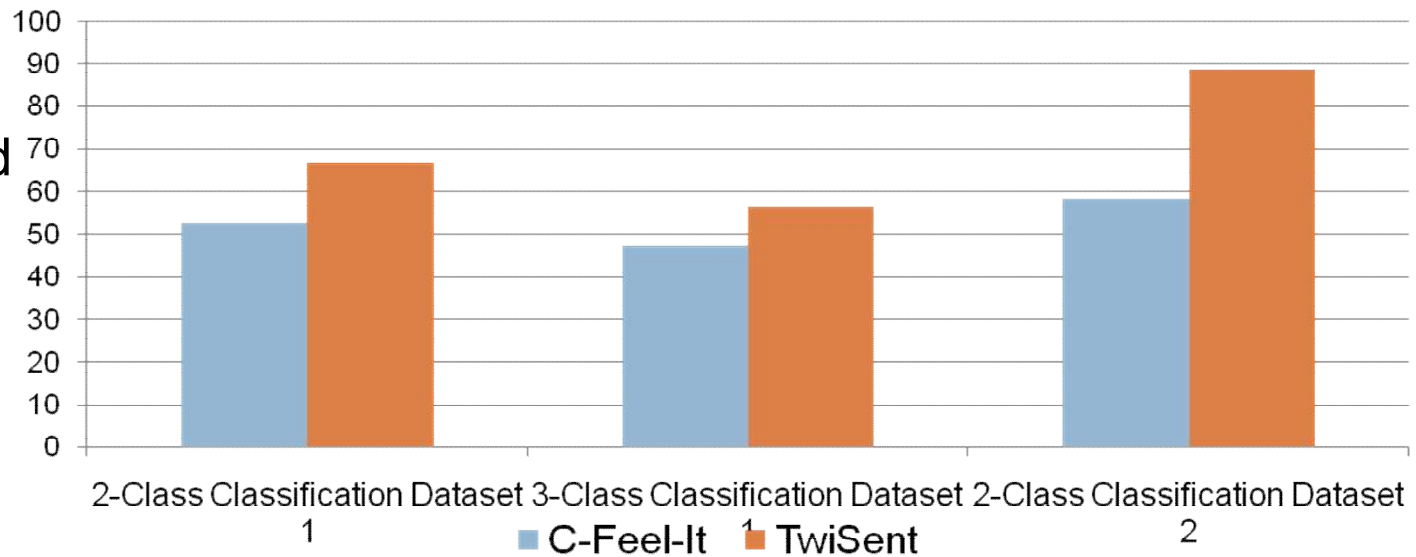
Supervised  
Classification



# TwisSent Evaluation

65

Lexicon-based  
Classification



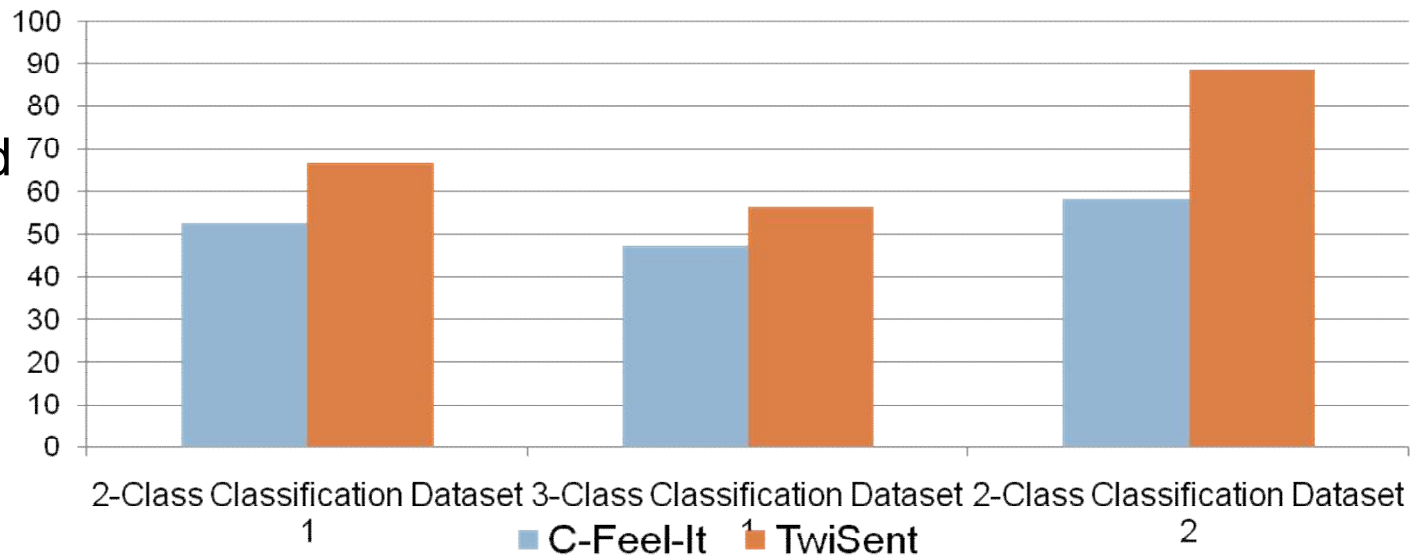
Supervised  
Classification

System	2-class Accuracy	Precision/Recall
C-Feel-It	50.8	53.16/72.96
<b>TwisSent</b>	<b>68.19</b>	<b>64.92/69.37</b>

# TwisSent Evaluation

66

Lexicon-based  
Classification



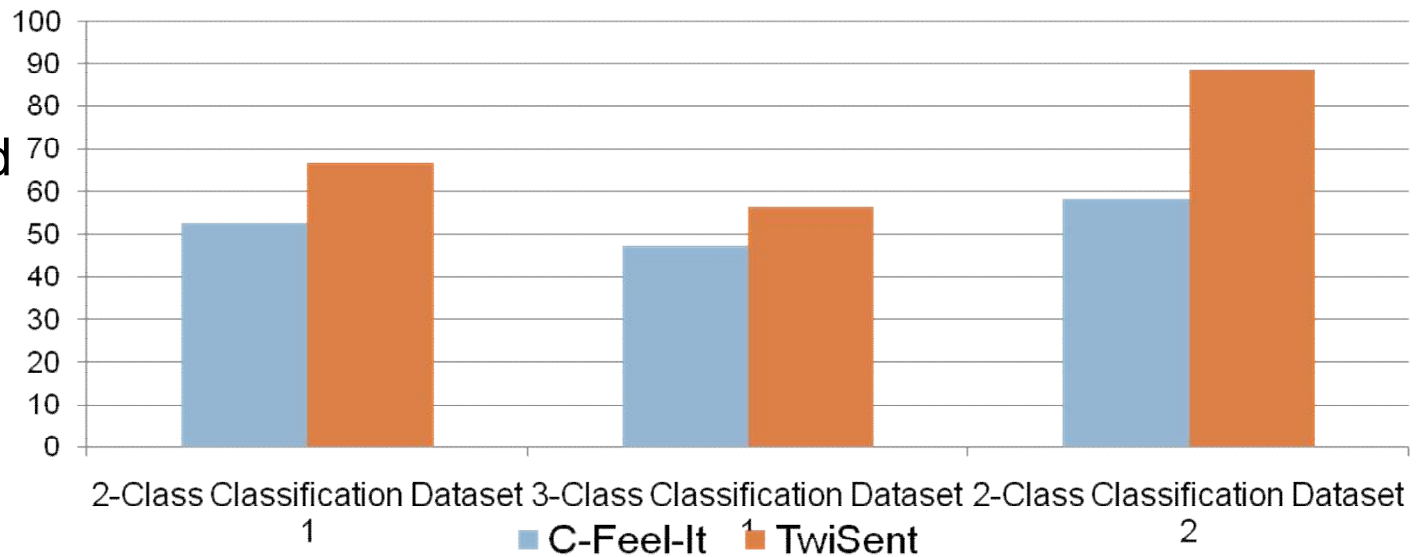
Supervised  
Classification

System	2-class Accuracy	Precision/Recall
C-Feel-It	50.8	53.16/72.96
<b>TwisSent</b>	<b>68.19</b>	<b>64.92/69.37</b>

# TwisSent Evaluation

67

Lexicon-based  
Classification



Ablation  
Test

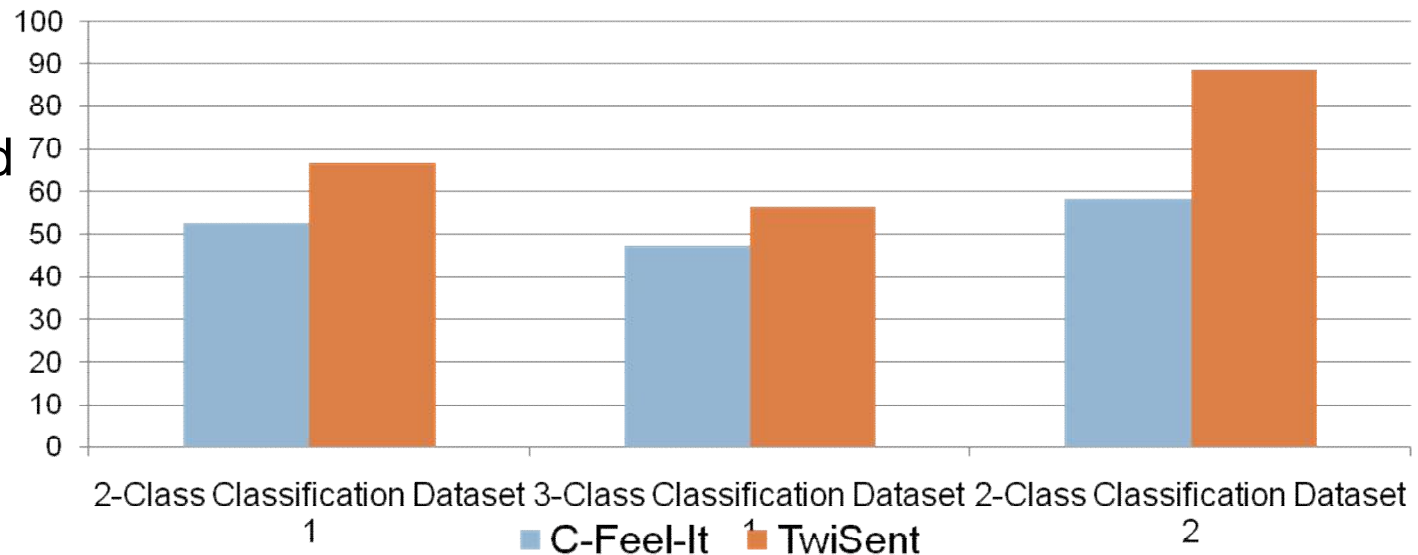
1

System	2-class Accuracy	Precision/Recall
C-Feel-It	50.8	53.16/72.96
<b>TwiSent</b>	<b>68.19</b>	<b>64.92/69.37</b>

# Twisent Evaluation

68

Lexicon-based  
Classification



Ablation  
Test

Module Removed	Accuracy	Statistical Significance Confidence (%)
Entity-Specificity	65.14	95
Spell-Checker	64.2	99
Pragmatics Handler	63.51	99
<b>Complete System</b>	<b>66.69</b>	-