

Adaptation of Sentiment Analysis to New Linguistic Features, Informal Language Form and World Knowledge

Subhabrata Mukherjee
Master's Thesis

Guide: Dr. Pushpak Bhattacharyya
Department of Computer Science and Engineering,
IIT Bombay

Roadmap

2

- Motivation
- Role of Feature Specificity in Sentiment Analysis
- Role of Discourse Specificity in Sentiment Analysis
- Role of World Knowledge in Sentiment Analysis
- Role of Social Media Content and Informal Language Form in Sentiment Analysis
- Applications
 - Role of Social Media and World Knowledge in IR
 - Role of Sentiment in Semantic Similarity Measure

The movie
was fabulous!

The movie
stars Mr. X

The movie
was horrible!

The movie
was fabulous!

[Sentimental]

The movie
stars Mr. X

[Factual]

The movie
was horrible!

[Sentimental]



The movie
was fabulous!

[Sentimental]



The movie
stars Mr. X

[Factual]



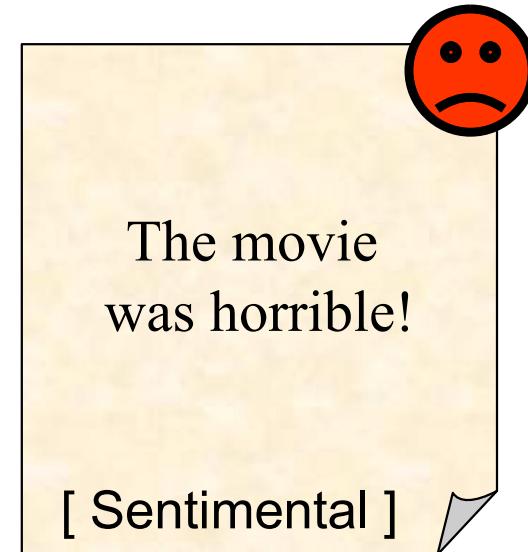
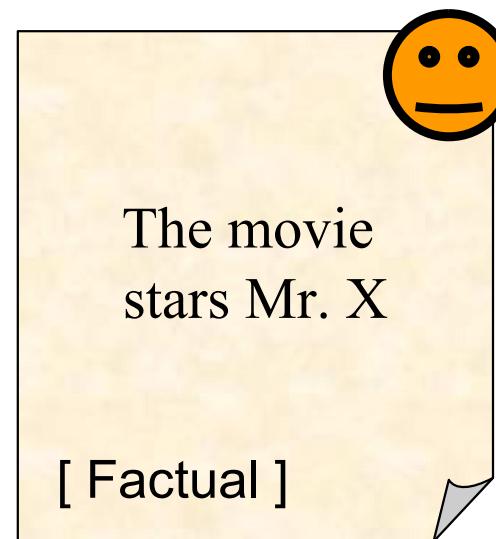
The movie
was horrible!

[Sentimental]

What is Sentiment Analysis

6

- Identify the orientation of opinion in a piece of text



Can be generalized to a wider set of emotions

Department of Computer Science and Engineering, IIT Bombay

Slide Courtesy: Aditya Joshi

7/23/2013

Motivation

Contd...

7

Motivation

Contd...

8

□ Discourse Coherency

Motivation

Contd...

9

□ Discourse Coherency

*He won the election **despite** a lot of smearing and slanderizing by the opposition.*

Motivation

Contd...

10

- Discourse Coherency
- Feature Specificity

*He won the election **despite** a lot of smearing and slanderizing by the opposition.*

Motivation

Contd...

11

- Discourse Coherency
- Feature Specificity

I like Samsung's multimedia features but that of Nokia is not that good.

Motivation

Contd...

12

- Discourse Coherency
- Feature Specificity
- Informal Language Form and Social Media Content

I like Samsung's multimedia features but that of Nokia is not that good.

Motivation

Contd...

13

- Discourse Coherency
- Feature Specificity
- Informal Language Form and Social Media Content

It really was. :) RT @AlissonRicardo: Chernobyl lukd dumb but
Men in Black 3 was prety guuud ☺

Motivation

Contd...

14

- Discourse Coherency
- Feature Specificity
- Informal Language Form and Social Media Content

World Knowledge

It really was. :) RT @AlissonRicardo: Chernobyl lukd dumb but
Men in Black 3 was prety guuud ☺

Motivation

Contd...

15

- Discourse Coherency
- Feature Specificity
- Informal Language Form and Social Media Content

World Knowledge

He is behaving like a Frankenstein.

*“I am looking forward to spend a nice evening with my parents”
– part of a movie review*

□ Feature Specific Sentiment Analysis

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 to 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”

Relations

25

- 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 $w_i, w_{i+1} \notin \text{Stopwords}$, then they are directly related. $w_i, w_{i+1} \in S$
- Dependency Relation
 - Capture **long range dependencies**
 - Let $\text{Dependency_Relation}$ be the list of **significant relations**.
- Any 2 words w_i and w_j in S are directly related, if
 $\exists D_i \text{ s.t. } D_i(w_i, w_j) \in \text{Dependency_Relation}$

Algorithm

26

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. $R_l(w_i, w_j) \in R$.

- i. Initialize n clusters C_i $\forall i = 1\dots n$
- ii. Make each $f_t \in F$ the clusterhead of C_i . The target feature f_t is the clusterhead of C_t . Initially, each cluster consists only of the clusterhead.
- iii. Assign each word $w_j \in S$ to cluster C_k s.t., $k = \operatorname{argmin}_{i \in n} \operatorname{dist}(w_j, f_i)$, Where $\operatorname{dist}(w_j, f_i)$ gives the number of edges, in the shortest path, connecting w_j and f_i in G .
- iv. Merge any cluster C_i with C_t if $\operatorname{dist}(w_j, f_i) < \theta$, Where θ 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 .

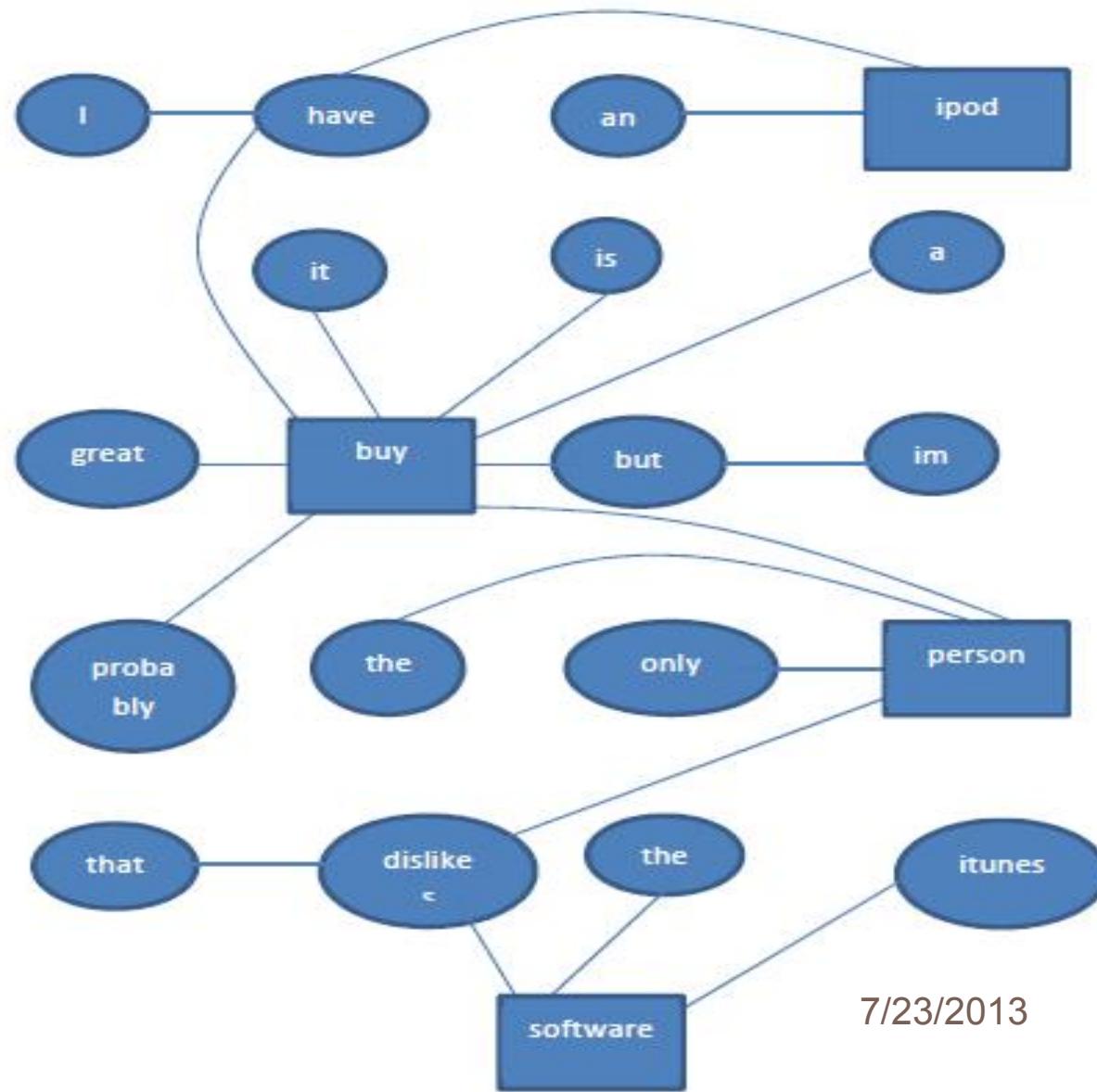
Graph

27

7/23/2013

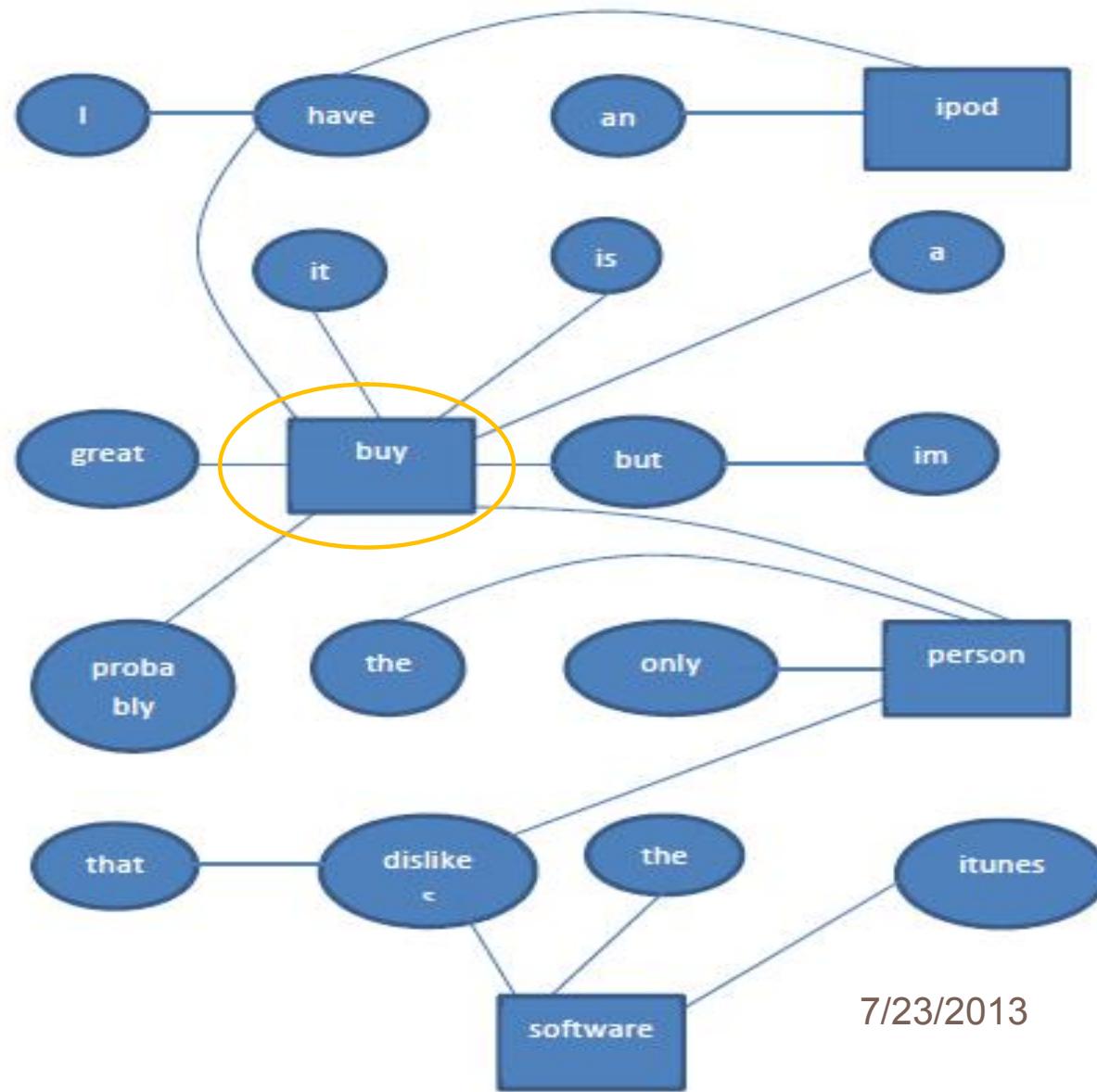
Graph

28



Graph

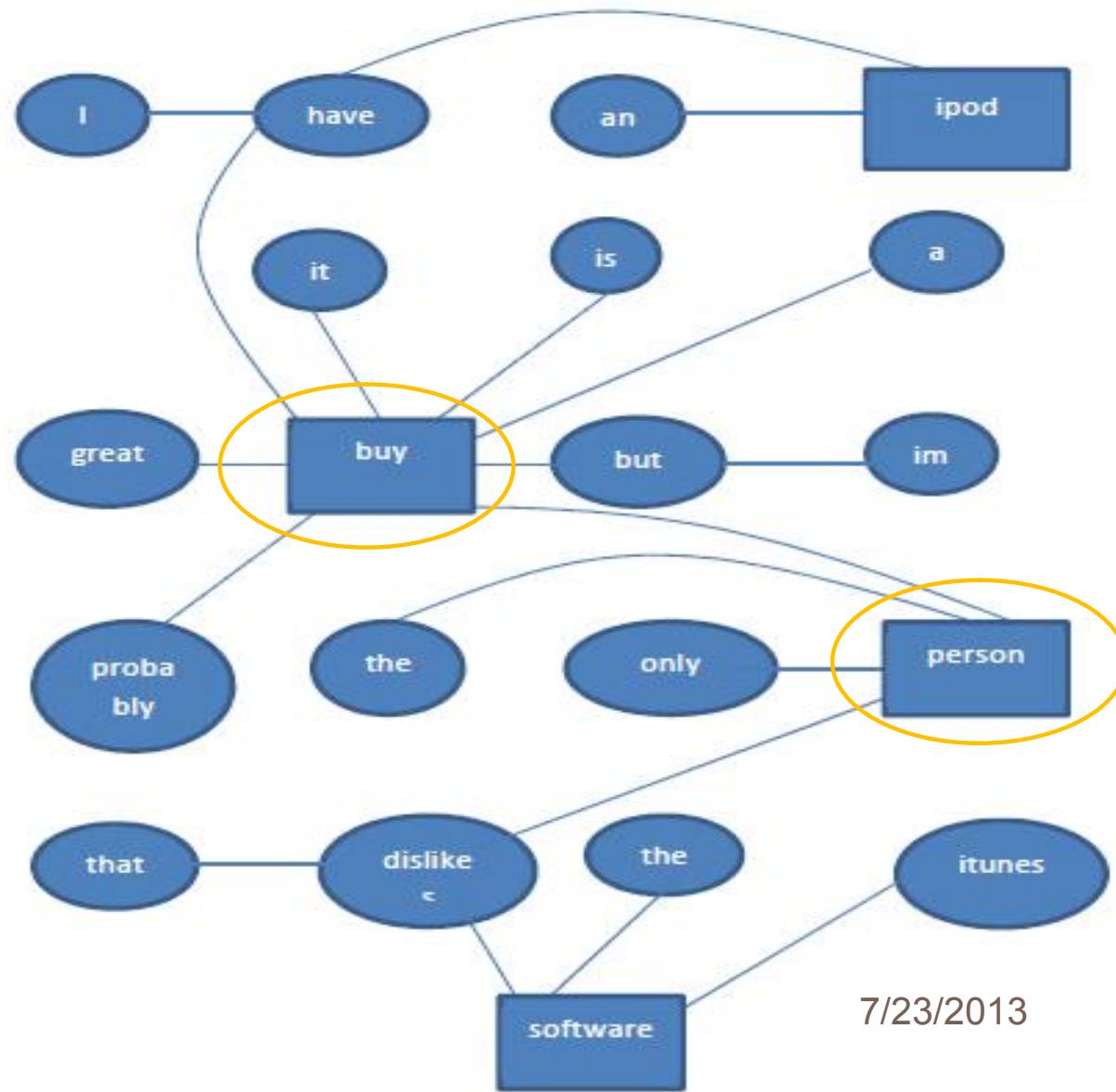
29



7/23/2013

Graph

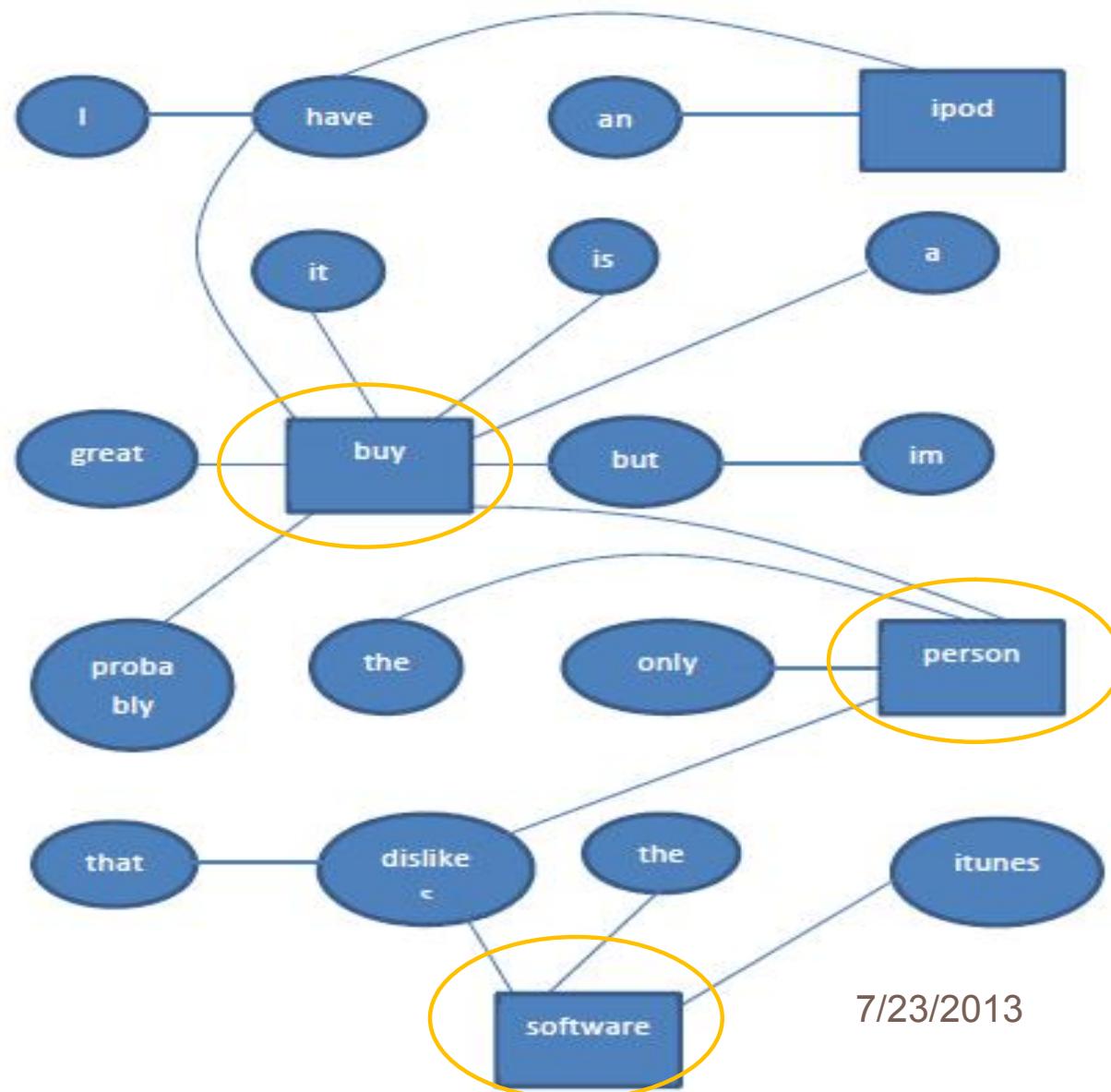
30



7/23/2013

Graph

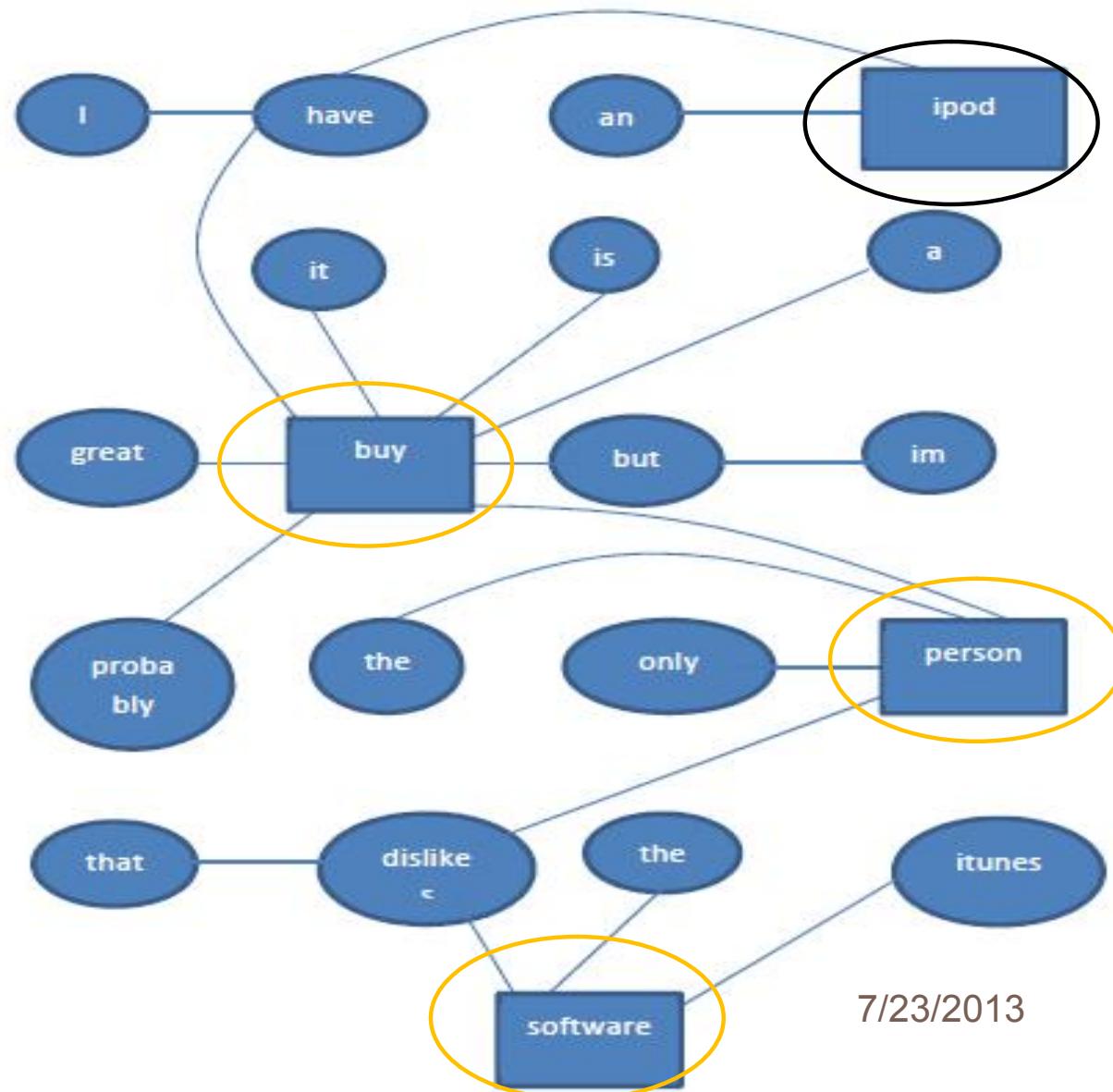
31



7/23/2013

Graph

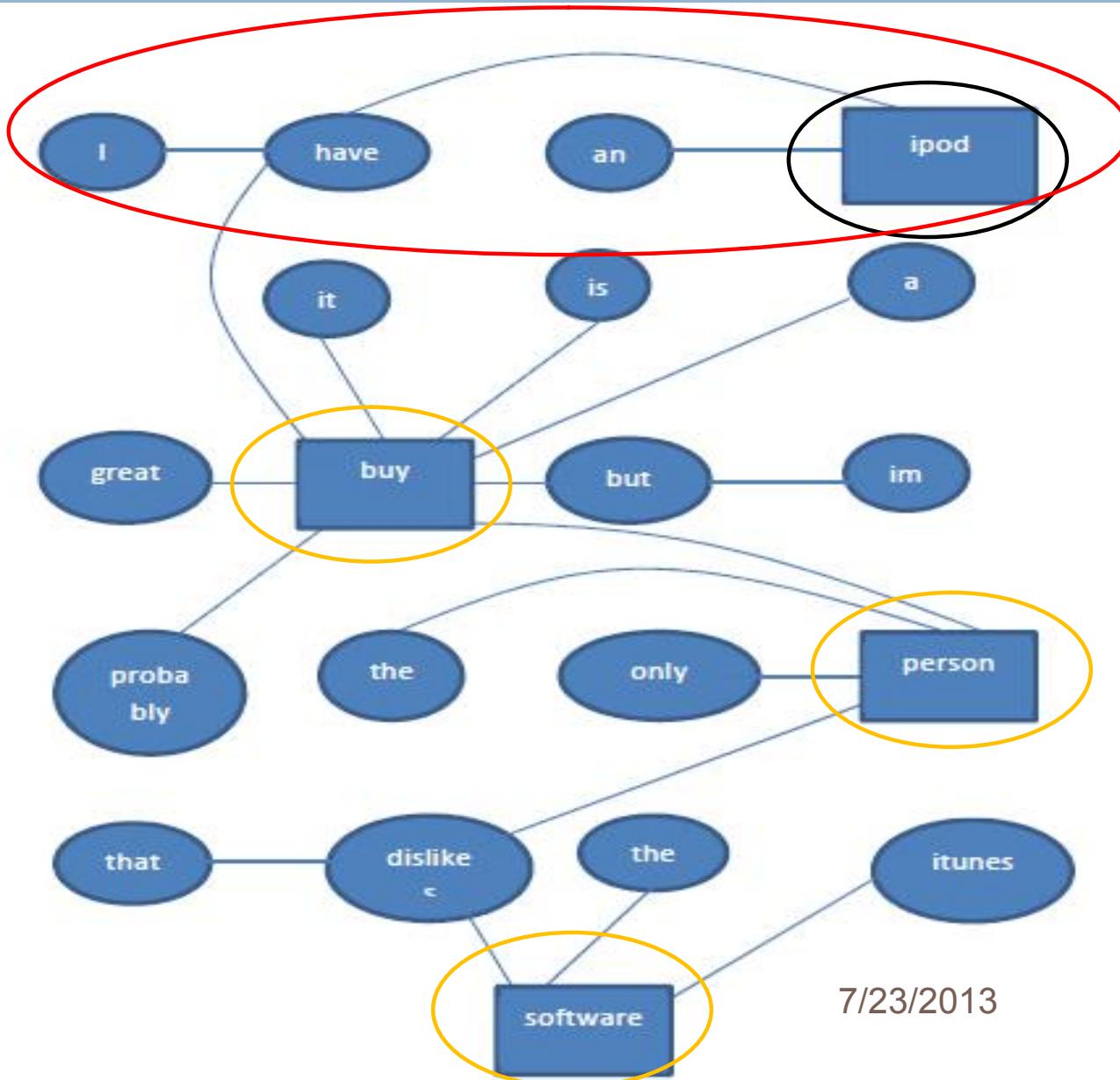
32



7/23/2013

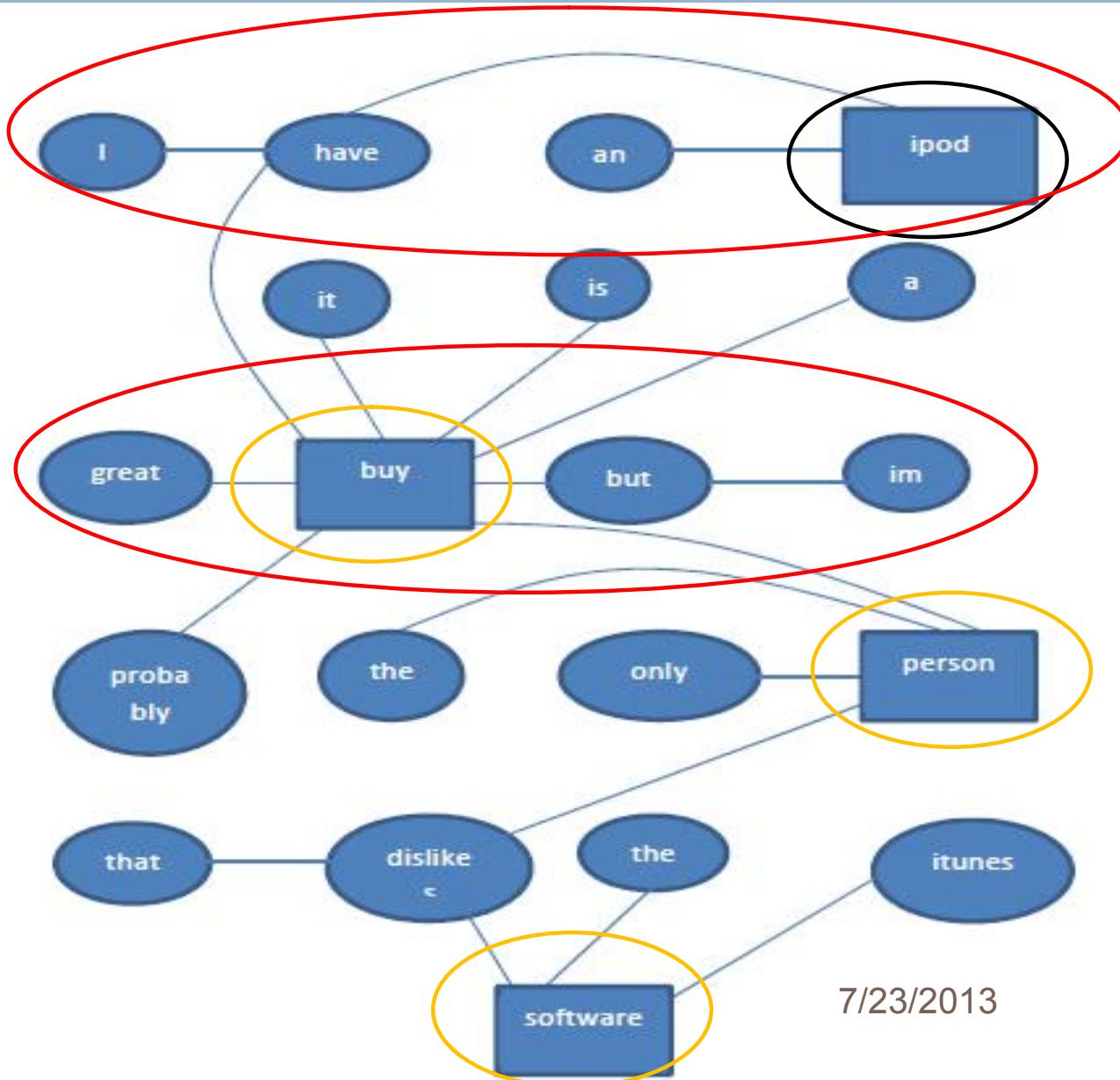
Graph

33



Graph

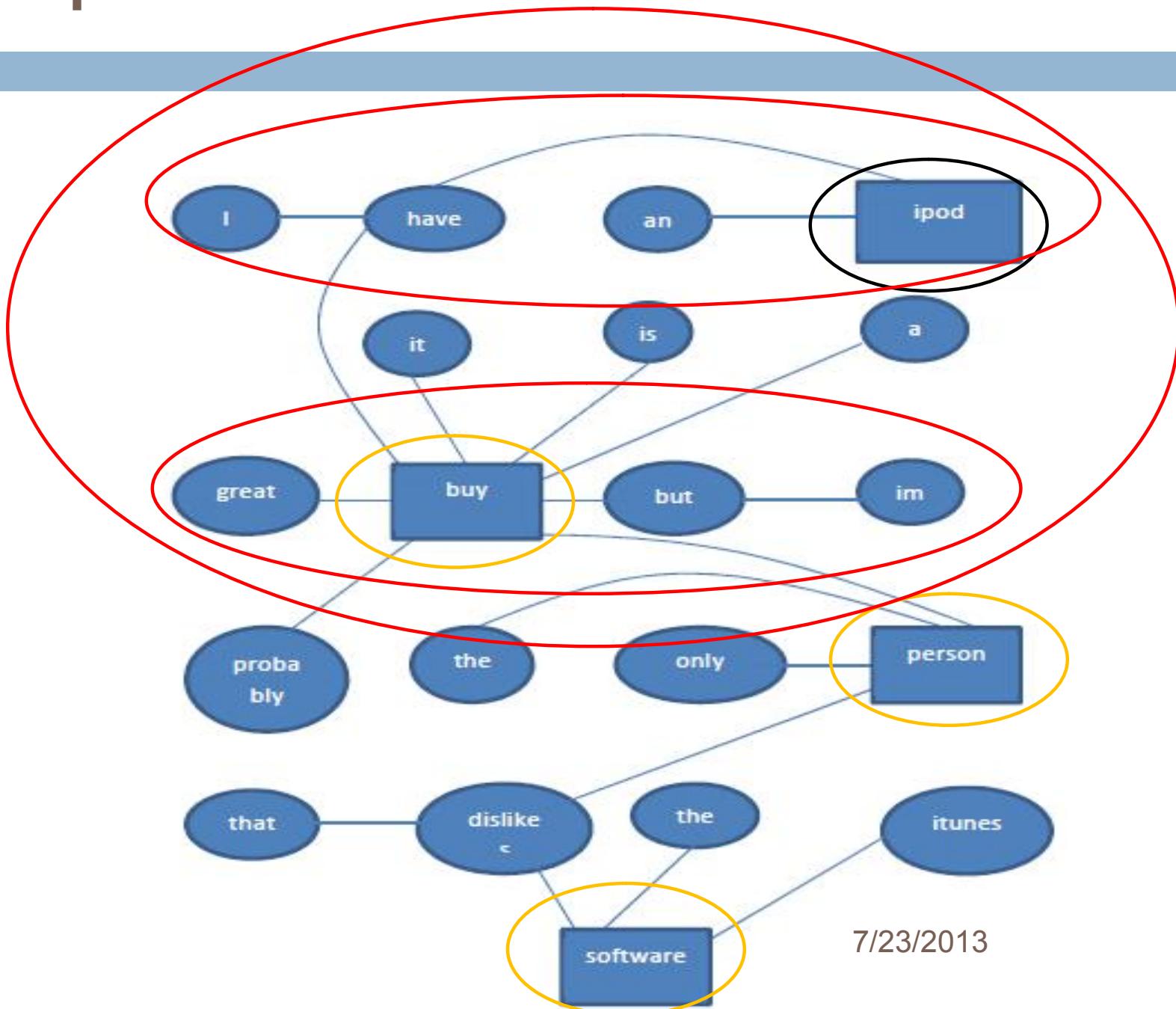
34



7/23/2013

Graph

35



Classification and Datasets

36

Classification and Datasets

37

- Lexicon Classification

Classification and Datasets

38

- Lexicon Classification

Bing Liu Sentiment Lexicon for finding the polarity of the words in the target cluster (majority voting)

Classification and Datasets

39

- Lexicon Classification
- Supervised Classification

Bing Liu Sentiment Lexicon for finding the polarity of the words in the target cluster (majority voting)

Classification and Datasets

40

- Lexicon Classification
- Supervised Classification

*Words in the target cluster as features (unigram bag-of-words)
SVM's as classifier*

Classification and Datasets

41

- Lexicon Classification
- Supervised Classification
- Dataset 1 - Lakkaraju *et al.* (SDM 2011)

*Words in the target cluster as features (unigram bag-of-words)
SVM's as classifier*

Classification and Datasets

42

- Lexicon Classification
- Supervised Classification
- Dataset 1 - Lakkaraju *et al.* (SDM 2011)

Domains - camera, laptop, mobile, printer

Classification and Datasets

43

- Lexicon Classification
- Supervised Classification
- Dataset 1 - Lakkaraju *et al.* (SDM 2011)
- Dataset 2 - Hu and Liu *et al.* (SIGKDD 2004)

Domains - camera, laptop, mobile, printer

Classification and Datasets

44

- Lexicon Classification
- Supervised Classification
- Dataset 1 - Lakkaraju *et al.* (SDM 2011)
- Dataset 2 - Hu and Liu *et al.* (SIGKDD 2004)

Domains - antivirus, camera, dvd, ipod, music player, router, mobile

Classification and Datasets

45

- Lexicon Classification
- Supervised Classification
- Dataset 1 - Lakkaraju *et al.* (SDM 2011)
- Dataset 2 - Hu and Liu *et al.* (SIGKDD 2004)
- Baseline 1

Domains - antivirus, camera, dvd, ipod, music player, router, mobile

Classification and Datasets

46

- Lexicon Classification
- Supervised Classification
- Dataset 1 - Lakkaraju *et al.* (SDM 2011)
- Dataset 2 - Hu and Liu *et al.* (SIGKDD 2004)
- Baseline 1

Majority Voting on the number of positive and negative opinion words

Classification and Datasets

47

- Lexicon Classification
- Supervised Classification
- Dataset 1 - Lakkaraju *et al.* (SDM 2011)
- Dataset 2 - Hu and Liu *et al.* (SIGKDD 2004)
- Baseline 1
- Baseline 2

Majority Voting on the number of positive and negative opinion words

Classification and Datasets

48

- Lexicon Classification
- Supervised Classification
- Dataset 1 - Lakkaraju *et al.* (SDM 2011)
- Dataset 2 - Hu and Liu *et al.* (SIGKDD 2004)
- Baseline 1
- Baseline 2

Opinion word nearest to the target entity decides the polarity

Parameter Learning

49

Parameter Learning

50

- Dependency Parsing uses approx. 40 relations.

Parameter Learning

51

- Dependency Parsing uses approx. 40 relations.
 - Relation Space – $(2^{40} - 1)$

Parameter Learning

52

- Dependency Parsing uses approx. 40 relations.
 - Relation Space – $(2^{40} - 1)$
- *Fix relations certain to be significant*
 - *nsubj, nsubjpass, dobj, amod, advmod, nn, neg*

Parameter Learning

53

- Dependency Parsing uses approx. 40 relations.
 - Relation Space – $(2^{40} - 1)$
- *Fix relations certain to be significant*
 - *nsubj, nsubjpass, dobj, amod, advmod, nn, neg*
- *Reject relations certain to be non-significant*
 - *copula, det, predet etc.*
- Compute Leave-One-Relation out accuracy over the remaining 21 relations over a development set
 - Find the relations for which there is ***significant accuracy change***
- Find inter-clusters distance threshold using development set

Ablation test

54

Ablation test

55

Relations	Accuracy (%)
All	63.5
Dep	67.3
Rmod	65.4
xcomp, conj_and ccomp, iobj	61.5
advcl , appos, csubj, abbrev, infmod, npavmod, rel, acomp, agent, csubjpass, partmod, pobj, purpcl, xsubj	63.5

Ablation test

56

Relations	Accuracy (%)
All	63.5
Dep	67.3
Rcmod	65.4
xcomp, conj_and ccomp, iobj	61.5
advcl , appos, csubj, abbrev, infmod, npavmod, rel, acomp, agent, csubjpass, partmod, pobj, purpcl, xsubj	63.5

Relation Set	Accuracy
With Dep+Rcmod	66
Without Dep	69
Without Rcmod	67
Without Dep+Rcmod	68

Ablation test

57

Relations	Accuracy (%)
All	63.5
Dep	67.3
Rcmod	65.4
xcomp, conj_and ccomp, iobj	61.5
advcl , appos, csubj, abbrev, infmod, npavmod, rel, acomp, agent, csubjpass, partmod, pobj, purpcl, xsubj	63.5

Relation Set	Accuracy
With Dep+Rcmod	66
Without Dep	69
Without Rcmod	67
Without Dep+Rcmod	68

θ	Accuracy (%)
2	67.85
3	69.28
4	68.21
5	67.40

Lexicon based classification (Dataset 2)

58

Domain	Baseline 1 (%)	Baseline 2 (%)	Proposed System (%)
Antivirus	50.00	56.82	63.63
Camera 1	50.00	61.67	78.33
Camera 2	50.00	61.76	70.58
Camera 3	51.67	53.33	60.00
Camera 4	52.38	57.14	78.57
Diaper	50.00	63.63	57.57
DVD	52.21	63.23	66.18
IPOD	50.00	57.69	67.30
Mobile 1	51.16	61.63	66.28
Mobile 2	50.81	65.32	70.96
Music Player 1	50.30	57.62	64.37
Music Player 2	50.00	60.60	67.02
Router 1	50.00	58.33	61.67
Router 2	50.00	59.72	70.83

Overall Classification Accuracy (Dataset 2)

59

Lexicon-based
Classification

Supervised
Classification

Overall Classification Accuracy (Dataset 2)

60

Lexicon-based
Classification

System	Accuracy (%)
Baseline ₁	50.35
Baseline ₂	58.93
Proposed System	70.00

Supervised
Classification

Overall Classification Accuracy (Dataset 2)

61

Lexicon-based
Classification

System	Accuracy (%)
Baseline ₁	50.35
Baseline ₂	58.93
Proposed System	70.00

Supervised
Classification

Domain	Baseline 1 (%)	Proposed System (%)
	Accuracy (Precision/Recall)	Accuracy (Precision/Recall)
Mobile	51.42 (50.72/99.29)	83.82 (83.82/83.82)
Camera	50	86.99 (84.73/90.24)

Lexicon Classification Results

(Dataset 1, Lakkaraju *et al.*, SDM 2011)

62

System	Sentiment Evaluation Accuracy (%)
Baseline ₁	68.75
Baseline ₂	61.10
CFACTS-R	80.54
CFACTS	81.28
FACTS-R	72.25
FACTS	75.72
JST	76.18
Proposed System	80.98

Drawbacks

63

- Features should be explicitly present in review
 - *Mobile is heavy. (implicit feature - weight)*
- Cannot capture domain-specific implicit sentiment
 - *read* in *movie* domain and *book* domain

- Discourse Specific Sentiment Analysis
- Sentiment Analysis in Twitter with Lightweight Discourse Analysis

Discourse Specific Sentiment Analysis

65

Slide Courtesy, Akshat Malu

Discourse Specific Sentiment Analysis

66

- Presence of words, at times, weighs more than the frequencies

Discourse Specific Sentiment Analysis

67

- Presence of words, at times, weighs more than the frequencies
- Presence of a discourse marker can alter the overall sentiment of a sentence

“The actors have done an okay job, not too brilliant, the script goes off track in between. Songs are mediocre. The direction could have been better. But overall, I like the movie”

Discourse Specific Sentiment Analysis

68

- Presence of words, at times, weighs more than the frequencies
 - Presence of a discourse marker can alter the overall sentiment of a sentence
- In most of the bag-of-words models, the discourse markers are ignored as stop words

“The actors have done an okay job, not too brilliant, the script goes off track in between. Songs are mediocre. The direction could have been better. But overall, I like the movie”

- Requirement of a **lightweight** method of discourse analysis for social media applications
- Use of heavy linguistic resources like parsing is not preferred
 - Increased processing time which slows down interactive applications
 - Parsing does not work well in presence of noisy text

*“@user share 'em! i'm quite excited bout Tintin,
despite not realy liking original comics. Probably
because Joe Cornish had a hand in.*

Discourse Markers



Discourse Markers

- Conjunctions that give more importance to the following discourse segment
 - (***but, however, nevertheless, otherwise, yet, still, nonetheless, nevertheless***)

Discourse Markers

- Conjunctions that give more importance to the following discourse segment
 - (***but, however, nevertheless, otherwise, yet, still, nonetheless, nevertheless***)

*The direction was not that great, **but still** we loved the movie.*

Discourse Markers

- Conjunctions that give more importance to the following discourse segment
 - (***but, however, nevertheless, otherwise, yet, still, nonetheless, nevertheless***)
- Conjunctions that give more importance to the previous discourse segment
 - (***till, until, despite, in spite, though, although***)

The direction was not that great, but still we loved the movie.

Discourse Markers

- Conjunctions that give more importance to the following discourse segment
 - (***but, however, nevertheless, otherwise, yet, still, nonetheless, nevertheless***)
- Conjunctions that give more importance to the previous discourse segment
 - (***till, until, despite, in spite, though, although***)

*India managed to win **despite** the initial setback.*

Discourse Markers

- Conjunctions that give more importance to the following discourse segment
 - (***but, however, nevertheless, otherwise, yet, still, nonetheless, nevertheless***)
- Conjunctions that give more importance to the previous discourse segment
 - (***till, until, despite, in spite, though, although***)
- Conjunctions that tend to draw a conclusion or inference
 - (***therefore, furthermore, consequently, thus, as a result, subsequently, eventually, hence***)

*India managed to win **despite** the initial setback.*

Discourse Markers

- Conjunctions that give more importance to the following discourse segment
 - (***but, however, nevertheless, otherwise, yet, still, nonetheless, nevertheless***)
- Conjunctions that give more importance to the previous discourse segment
 - (***till, until, despite, in spite, though, although***)
- Conjunctions that tend to draw a conclusion or inference
 - (***therefore, furthermore, consequently, thus, as a result, subsequently, eventually, hence***)

We were not much satisfied with the greatly acclaimed brand X and subsequently decided to reject it.

Discourse Markers

- Conjunctions that give more importance to the following discourse segment
 - (***but, however, nevertheless, otherwise, yet, still, nonetheless, nevertheless***)
- Conjunctions that give more importance to the previous discourse segment
 - (***till, until, despite, in spite, though, although***)
- Conjunctions that tend to draw a conclusion or inference
 - (***therefore, furthermore, consequently, thus, as a result, subsequently, eventually, hence***)
- **Strong Modals** (***might, could, can, would, may***)
 - These express a higher degree of uncertainty

We were not much satisfied with the greatly acclaimed brand X and subsequently decided to reject it.

Discourse Markers

- Conjunctions that give more importance to the following discourse segment
 - (***but, however, nevertheless, otherwise, yet, still, nonetheless, nevertheless***)
- Conjunctions that give more importance to the previous discourse segment
 - (***till, until, despite, in spite, though, although***)
- Conjunctions that tend to draw a conclusion or inference
 - (***therefore, furthermore, consequently, thus, as a result, subsequently, eventually, hence***)
- **Strong Modals** (***might, could, can, would, may***)
 - These express a higher degree of uncertainty

That film might be good.

Discourse Markers

- Conjunctions that give more importance to the following discourse segment
 - (***but, however, nevertheless, otherwise, yet, still, nonetheless, nevertheless***)
- Conjunctions that give more importance to the previous discourse segment
 - (***till, until, despite, in spite, though, although***)
- Conjunctions that tend to draw a conclusion or inference
 - (***therefore, furthermore, consequently, thus, as a result, subsequently, eventually, hence***)
- **Strong Modals** (***might, could, can, would, may***)
 - These express a higher degree of uncertainty
- **Weak Modals** (***should, ought to, need not, shall, will, must***)
 - These express lesser degree of uncertainty

That film might be good.

Discourse Markers

- Conjunctions that give more importance to the following discourse segment
 - (***but, however, nevertheless, otherwise, yet, still, nonetheless, nevertheless***)
- Conjunctions that give more importance to the previous discourse segment
 - (***till, until, despite, in spite, though, although***)
- Conjunctions that tend to draw a conclusion or inference
 - (***therefore, furthermore, consequently, thus, as a result, subsequently, eventually, hence***)
- **Strong Modals** (***might, could, can, would, may***)
 - These express a higher degree of uncertainty
- **Weak Modals** (***should, ought to, need not, shall, will, must***)
 - These express lesser degree of uncertainty

*I heard the movie is good, so you **must** go to watch that movie.*

Discourse Markers

- Conjunctions that give more importance to the following discourse segment
 - (***but, however, nevertheless, otherwise, yet, still, nonetheless, nevertheless***)
- Conjunctions that give more importance to the previous discourse segment
 - (***till, until, despite, in spite, though, although***)
- Conjunctions that tend to draw a conclusion or inference
 - (***therefore, furthermore, consequently, thus, as a result, subsequently, eventually, hence***)
- **Strong Modals** (***might, could, can, would, may***)
 - These express a higher degree of uncertainty
- **Weak Modals** (***should, ought to, need not, shall, will, must***)
 - These express lesser degree of uncertainty
- **Conditionals** (***If***)

*I heard the movie is good, so you **must** go to watch that movie.*

Discourse Markers

- Conjunctions that give more importance to the following discourse segment
 - (***but, however, nevertheless, otherwise, yet, still, nonetheless, nevertheless***)
- Conjunctions that give more importance to the previous discourse segment
 - (***till, until, despite, in spite, though, although***)
- Conjunctions that tend to draw a conclusion or inference
 - (***therefore, furthermore, consequently, thus, as a result, subsequently, eventually, hence***)
- **Strong Modals** (***might, could, can, would, may***)
 - These express a higher degree of uncertainty
- **Weak Modals** (***should, ought to, need not, shall, will, must***)
 - These express lesser degree of uncertainty
- **Conditionals** (***If***)

If I had studied well before the exams, I would have done well.

Discourse Markers

- Conjunctions that give more importance to the following discourse segment
 - (***but, however, nevertheless, otherwise, yet, still, nonetheless, nevertheless***)
- Conjunctions that give more importance to the previous discourse segment
 - (***till, until, despite, in spite, though, although***)
- Conjunctions that tend to draw a conclusion or inference
 - (***therefore, furthermore, consequently, thus, as a result, subsequently, eventually, hence***)
- **Strong Modals** (***might, could, can, would, may***)
 - These express a higher degree of uncertainty
- **Weak Modals** (***should, ought to, need not, shall, will, must***)
 - These express lesser degree of uncertainty
- **Conditionals** (***If***)

If I had studied well before the exams, I would have done well.

Discourse Markers

- Conjunctions that give more importance to the following discourse segment
 - (***but, however, nevertheless, otherwise, yet, still, nonetheless, nevertheless***)
- Conjunctions that give more importance to the previous discourse segment
 - (***till, until, despite, in spite, though, although***)
- Conjunctions that tend to draw a conclusion or inference
 - (***therefore, furthermore, consequently, thus, as a result, subsequently, eventually, hence***)
- **Strong Modals** (***might, could, can, would, may***)
 - These express a higher degree of uncertainty
- **Weak Modals** (***should, ought to, need not, shall, will, must***)
 - These express lesser degree of uncertainty
- **Conditionals** (***If***)

I do not like Nokia but I like Samsung.

Discourse Markers

- Conjunctions that give more importance to the following discourse segment
 - (***but, however, nevertheless, otherwise, yet, still, nonetheless, nevertheless***)
- Conjunctions that give more importance to the previous discourse segment
 - (***till, until, despite, in spite, though, although***)
- Conjunctions that tend to draw a conclusion or inference
 - (***therefore, furthermore, consequently, thus, as a result, subsequently, eventually, hence***)
- **Strong Modals** (***might, could, can, would, may***)
 - These express a higher degree of uncertainty
- **Weak Modals** (***should, ought to, need not, shall, will, must***)
 - These express lesser degree of uncertainty
- **Conditionals** (***If***)

I do not like Nokia but I like Samsung.

Algorithm



Algorithm

- (*but, however, nevertheless, otherwise, yet, still, nonetheless, nevertheless, therefore, furthermore, consequently, thus, as a result, subsequently, eventually, hence*)

Algorithm

- (*but, however, nevertheless, otherwise, yet, still, nonetheless, nevertheless, therefore, furthermore, consequently, thus, as a result, subsequently, eventually, hence*)

- Words after them are given more weightage
 - Frequency count of those words are incremented by 1
- *The movie looked promising, but it failed to make an impact in the box-office*

Algorithm

- (*but, however, nevertheless, otherwise, yet, still, nonetheless, nevertheless, therefore, furthermore, consequently, thus, as a result, subsequently, eventually, hence*)
- (*till, until, despite, in spite, though, although*)

- Words after them are given more weightage
- Frequency count of those words are incremented by 1
- *The movie looked promising, but it failed to make an impact in the box-office*

Algorithm

- (*but, however, nevertheless, otherwise, yet, still, nonetheless, nevertheless, therefore, furthermore, consequently, thus, as a result, subsequently, eventually, hence*)
- (*till, until, despite, in spite, though, although*)

- Words before them are given more weightage
- frequency count of those words are incremented by 1
- *India staged a marvelous victory down under despite all odds.*

Algorithm

- (*but, however, nevertheless, otherwise, yet, still, nonetheless, nevertheless, therefore, furthermore, consequently, thus, as a result, subsequently, eventually, hence*)
- (*till, until, despite, in spite, though, although*)
- All sentences containing *if* are marked

- Words before them are given more weightage
- frequency count of those words are incremented by 1
- *India staged a marvelous victory down under despite all odds.*

Algorithm

- (*but, however, nevertheless, otherwise, yet, still, nonetheless, nevertheless, therefore, furthermore, consequently, thus, as a result, subsequently, eventually, hence*)
- (*till, until, despite, in spite, though, although*)
- All sentences containing *if* are marked

- In supervised classifiers, their weights are decreased

Algorithm

- (*but, however, nevertheless, otherwise, yet, still, nonetheless, nevertheless, therefore, furthermore, consequently, thus, as a result, subsequently, eventually, hence*)
- (*till, until, despite, in spite, though, although*)
- All sentences containing *if* are marked
- All sentences containing *strong modals* are marked
 - In supervised classifiers, their weights are decreased

Algorithm

- (*but, however, nevertheless, otherwise, yet, still, nonetheless, nevertheless, therefore, furthermore, consequently, thus, as a result, subsequently, eventually, hence*)
- (*till, until, despite, in spite, though, although*)
- All sentences containing *if* are marked
- All sentences containing *strong modals* are marked

- In supervised classifiers, their weights are decreased

Algorithm

- (*but, however, nevertheless, otherwise, yet, still, nonetheless, nevertheless, therefore, furthermore, consequently, thus, as a result, subsequently, eventually, hence*)
- (*till, until, despite, in spite, though, although*)
- All sentences containing *if* are marked
- All sentences containing *strong modals* are marked
- Negation
 - In supervised classifiers, their weights are decreased

Algorithm

- (*but, however, nevertheless, otherwise, yet, still, nonetheless, nevertheless, therefore, furthermore, consequently, thus, as a result, subsequently, eventually, hence*)
- (*till, until, despite, in spite, though, although*)
- All sentences containing *if* are marked
- All sentences containing *strong modals* are marked
- Negation
 - A window of 5 is considered
 - Polarity of all words in the window are reversed till another violating expectation conjunction is encountered
 - The polarity reversals are specially marked
 - *I do not like Nokia but I like Samsung.*

Feature Space

- Lexeme Feature Space

Feature Space

- Lexeme Feature Space
 - Bag-of-words features

Feature Space

- Lexeme Feature Space
 - Bag-of-words features
 - Maintain a count c_i of each word w_i

Feature Space

- Lexeme Feature Space
 - Bag-of-words features
 - Maintain a count c_i of each word w_i
 - Initially all w_i 's are initialized to 1

Feature Space

- Lexeme Feature Space
 - Bag-of-words features
 - Maintain a count c_i of each word w_i
 - Initially all w_i 's are initialized to 1
 - c_i is incremented with occurrence of w_i in the doc

Feature Space

- Lexeme Feature Space
 - Bag-of-words features
 - Maintain a count c_i of each word w_i
 - Initially all w_i 's are initialized to 1
 - c_i is incremented with occurrence of w_i in the doc
 - c_i is also incremented according to w_i 's importance

Feature Space

- Lexeme Feature Space
 - Bag-of-words features
 - Maintain a count c_i of each word w_i
 - Initially all w_i 's are initialized to 1
 - c_i is incremented with occurrence of w_i in the doc
 - c_i is also incremented according to w_i 's importance
 - $\{w_i, c_i\}$ where c_i represents the weight of the term w_i .

Feature Space

- Lexeme Feature Space
 - Bag-of-words features
 - Maintain a count c_i of each word w_i
 - Initially all w_i 's are initialized to 1
 - c_i is incremented with occurrence of w_i in the doc
 - c_i is also incremented according to w_i 's importance
 - $\{w_i, c_i\}$ where c_i represents the weight of the term w_i .
- Sense Feature Space
 - s_i instead of w_i

Classification of features

- Rule based System

Classification of features

- Rule based System
 - Bing Liu Sentiment Lexicon (Hu *et al.*, 2004)

Classification of features

- Rule based System
 - Bing Liu Sentiment Lexicon (Hu *et al.*, 2004)
- Supervised Classification
 - SVM on feature vectors $\{w_i, c_i\}$ or $\{s_i, c_i\}$
 - Features
 - N-grams (N=1,2)
 - Stop Word Removal (except discourse markers)
 - Discourse Weight of Features
 - Modal and Conditional Indicators
 - Stemming
 - Negation
 - Emoticons
 - Part-of-Speech Information
 - Feature Space (Lexeme or Synset)

Datasets



Datasets

- Dataset 1 (Twitter – Manually Annotated)
 - 8507 tweets over 2000 entities from 20 domains
 - Annotated by 4 annotators into positive, negative and objective classes

Datasets

- Dataset 1 (Twitter – Manually Annotated)
 - 8507 tweets over 2000 entities from 20 domains
 - Annotated by 4 annotators into positive, negative and objective classes
- Dataset 2 (Twitter – Auto Annotated)
 - 15,214 tweets collected and annotated based on hashtags
 - Positive hashtags - #positive, #joy, #excited, #happy
 - Negative hashtags - #negative, #sad, #depressed, #gloomy, #disappointed

Datasets

Manually Annotated Dataset				
#Positive	#Negative	#Objective Not Spam	#Objective Spam	Total
2548	1209	2757	1993	8507
Auto Annotated Dataset				
#Positive		#Negative		Total
7348		7866		15214

- ❑ Negative hashtags - *#negative, #sad, #depressed, #gloomy, #disappointed*

Datasets

Manually Annotated Dataset				
#Positive	#Negative	#Objective Not Spam	#Objective Spam	Total
2548	1209	2757	1993	8507
Auto Annotated Dataset				
#Positive		#Negative		Total
7348		7866		15214

- Negative hashtags - *#negative, #sad, #depressed, #gloomy, #disappointed*
- Dataset 3 (Travel Domain - Balamurali *et al.*, EMNLP 2011)
 - Each word is manually tagged with its disambiguated WordNet sense
 - Contains 595 polarity tagged documents of each class (positive and negative)

Dataset Domains

113

Movie, Restaurant, Television, Politics, Sports, Education, Philosophy, Travel, Books, Technology, Banking & Finance, Business, Music, Environment, Computers, Automobiles, Cosmetics brands, Amusement parks, Eatables, History

Classification Results in Twitter (Datasets 1 and 2)

Comparison with C-Feel-It (Joshi *et al.*, ACL 2011)

Classification Results in Twitter (Datasets 1 and 2)

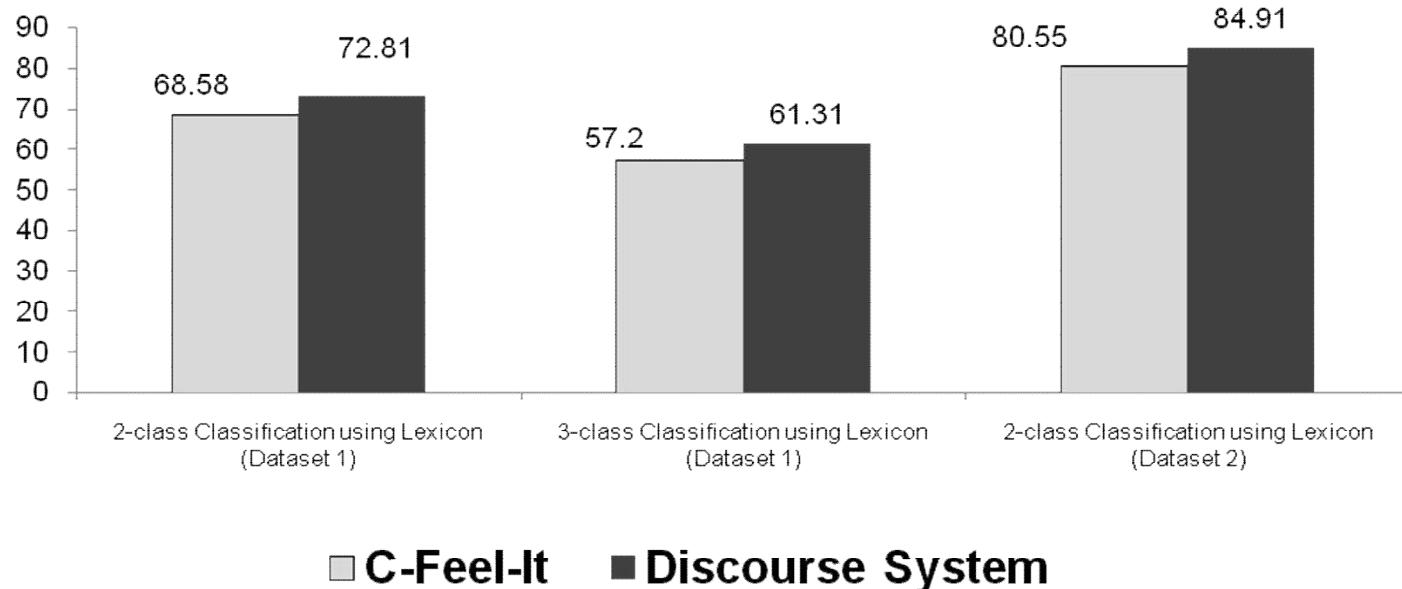
Comparison with C-Feel-It (Joshi *et al.*, ACL 2011)

Lexicon-based
Classification

Classification Results in Twitter (Datasets 1 and 2)

Comparison with C-Feel-It (Joshi *et al.*, ACL 2011)

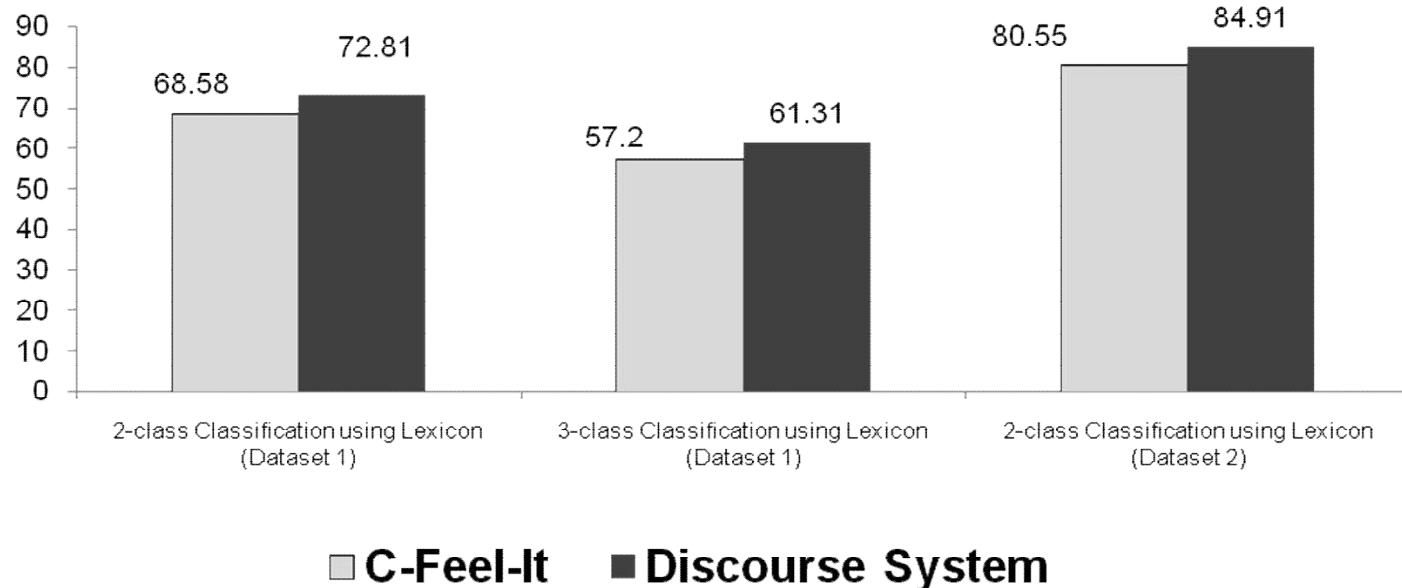
Lexicon-based
Classification



Classification Results in Twitter (Datasets 1 and 2)

Comparison with C-Feel-It (Joshi *et al.*, ACL 2011)

Lexicon-based
Classification

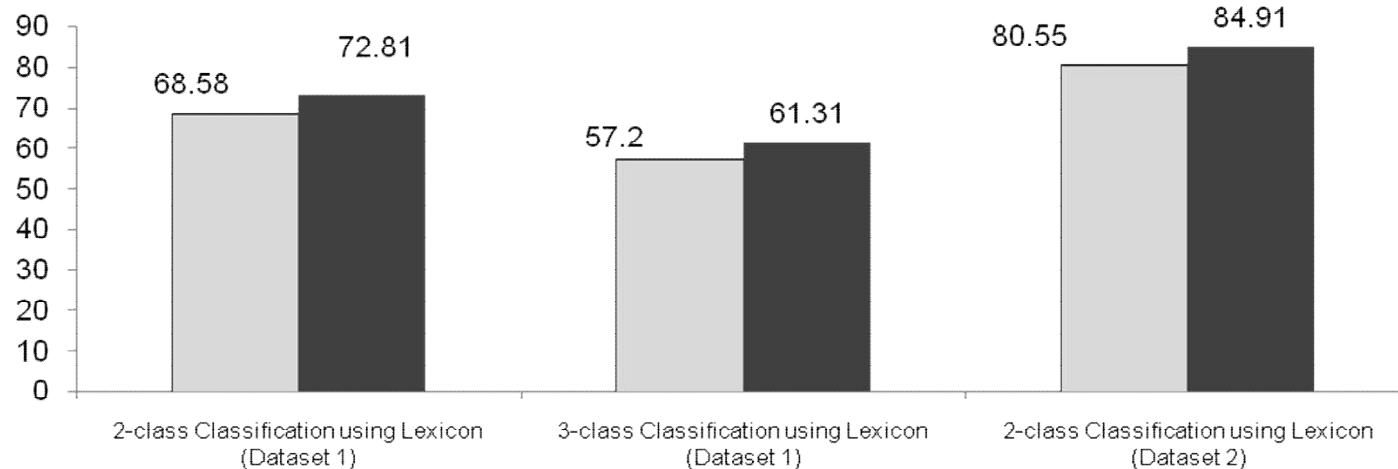


Supervised
Classification

Classification Results in Twitter (Datasets 1 and 2)

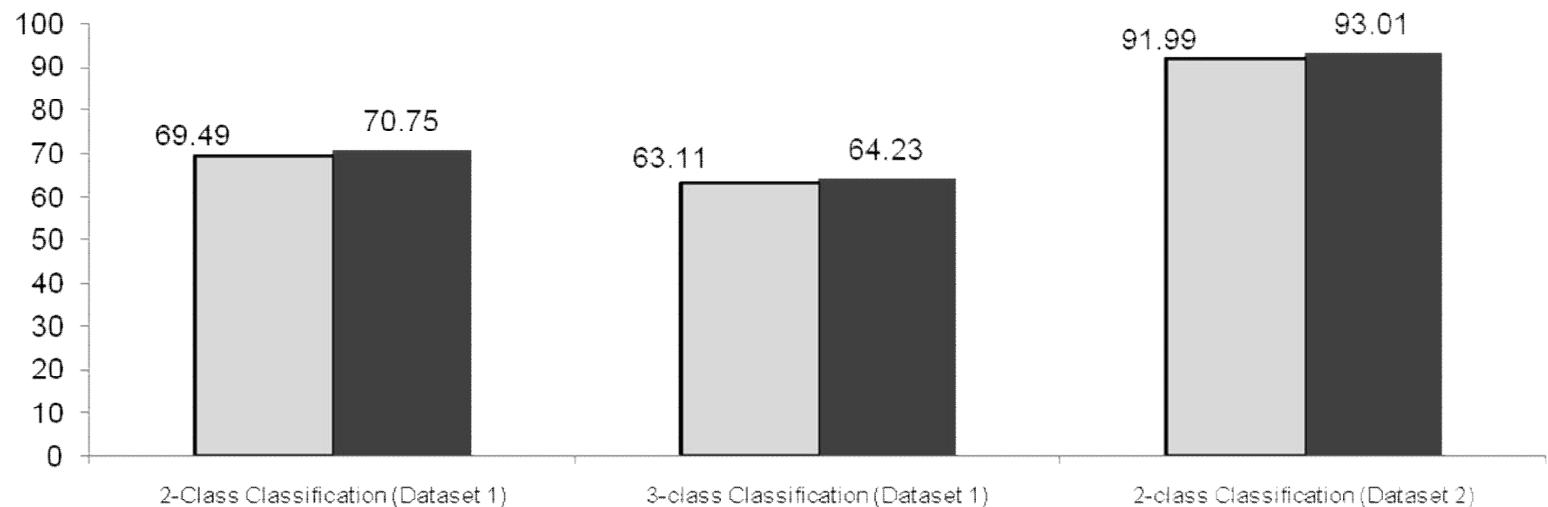
Comparison with C-Feel-It (Joshi *et al.*, ACL 2011)

Lexicon-based
Classification



■ C-Feel-It ■ Discourse System

Supervised
Classification



Classification Results in Travel Reviews (Dataset 3)

Comparison with Balamurali *et al.*, EMNLP 2011

119

Classification Results in Travel Reviews (Dataset 3)

Comparison with Balamurali *et al.*, EMNLP 2011

120

Supervised
Classification

Classification Results in Travel Reviews (Dataset 3)

Comparison with Balamurali *et al.*, EMNLP 2011

121

Supervised
Classification

Systems	Accuracy (%)
Baseline Accuracy (Only Unigrams)	84.90
Balamurali <i>et al.</i> , 2011 (Only IWSD Sense of Unigrams)	85.48
Balamurali <i>et al.</i> , 2011 (Unigrams+IWSD Sense of Unigrams)	86.08
Unigrams + IWSD Sense of Unigrams+Discourse Features	88.13

Drawbacks

122

- Usage of a generic lexicon in lexeme feature space
- Lexicons do not have entries for interjections like *wow*, *duh* etc. which are strong indicators of sentiment
- Noisy Text (*luv*, *gr8*, *spams*, ...)
- Sparse feature space (140 chars) for supervised classification
- 70% accuracy of IWSD in sense space
- Scope of discourse markers

Positional Importance: Back-up Slide

123

- “*I wanted to follow my dreams and ambitions despite all the obstacles, but I did not succeed.*”
- *want* and *ambition* will get polarity +2 each, as they appear before *despite*
- *obstacle* will get polarity -1 and *not succeed* will get a polarity -2
- Overall polarity +1, whereas the overall sentiment should be *negative*
- Reason:
 - We do not consider *positional importance* of a discourse marker in the sentence and consider all markers equally important
 - Better give a ranking to the discourse markers based on their *positional* and *pragmatic* importance.

□ Role of World Knowledge in Sentiment Analysis

- WikiSent: Unsupervised Sentiment Analysis System Using Extractive Summarization With Wikipedia

World Knowledge

125

- Extensive world knowledge required to perform analysis of reviews
- Distinguish between **What's** and **About's** of the review
- Filter out concepts irrelevant to the reviewer opinion about the movie
- Retain only the objects of interest and corresponding opinions

A Negative Review

- **best remembered** for his understated **performance** as dr. hannibal lecter in michael mann's forensics thriller , manhunter , scottish character actor brian cox brings something **special** to every movie he works on .
- usually playing a bit role in some studio schlock (he dies halfway through the long kiss goodnight) , he's only occasionally given something **meaty** and **substantial** to do .
- if you want to see some **brilliant** acting , check out his work as a dogged police inspector opposite frances mcdormand in ken loach's hidden agenda .
- cox plays the role of big john harrigan in the disturbing new indie flick i . i . e . , which lot 47 picked up at sundance when other distributors were scared to budge
- big john feels the **love** that dares not speak its name , but he expresses it through seeking out adolescents and bringing them back to his pad .
- what bothered some audience members was the presentation of big john in an oddly empathetic light .
- he's an **even-tempered** , **funny** , **robust** old man who actually listens to the kids' problems (as opposed to their parents and friends , both caught up in the high-wire act of their own confused lives .)
- he'll have sex-for-pay with them only after an elaborate courtship , **charming** them with temptations from the grown-up world .

Facets of a Movie Review

- General Perception about the Crew
 - Opinion about the Characters in the Movie
 - Characteristics of a movie or genre
- Objective Facts about the Crew and Movies
- Past Performance of the Crew and Movies
 - **Opinion about the Movie and Crew**
- Expectations from the Movie or Crew
 - Unrelated Category
- Movie Plot

Facets of a Movie Review

- General Perception about the Crew
 - Opinion about the Characters in the Movie
 - Characteristics of a movie or genre
- Objective Facts about the Crew and Movies
- Past Performance of the Crew and Movies
 - Opinion about the Movie and Crew
- Expectations from the Movie or Crew
 - Unrelated Category
- Movie Plot
 - *John Travolta is considered by many to be a has-been, or a one-hit wonder*
 - *Leonardo DeCaprio is an awesome actor.*

Facets of a Movie Review

- General Perception about the Crew
 - Opinion about the Characters in the Movie
 - Characteristics of a movie or genre
- Objective Facts about the Crew and Movies
- Past Performance of the Crew and Movies
 - Opinion about the Movie and Crew
- Expectations from the Movie or Crew
 - Unrelated Category
- Movie Plot
 - *Born into a family of thespians -- parents Roger Winslet and Sally Bridges-Winslet were both stage actors, maternal grandparents Oliver and Linda Bridges ran the Reading Repertory Theatre, and uncle Robert Bridges was a fixture in London's West End theatre district, Kate Winslet came into her talent at an early age.*

Facets of a Movie Review

- General Perception about the Crew
 - Opinion about the Characters in the Movie
 - Characteristics of a movie or genre
- Objective Facts about the Crew and Movies
- Past Performance of the Crew and Movies
 - **Opinion about the Movie and Crew**
- Expectations from the Movie or Crew
 - Unrelated Category
- Movie Plot

- *The role that transformed Winslet from art house attraction to international star was Rose DeWitt Bukater, the passionate, rosy-cheeked aristocrat in James Cameron's Titanic (1997).*

Facets of a Movie Review

- General Perception about the Crew
 - Opinion about the Characters in the Movie
 - Characteristics of a movie or genre
- Objective Facts about the Crew and Movies
- Past Performance of the Crew and Movies
 - **Opinion about the Movie and Crew**
- Expectations from the Movie or Crew
 - Unrelated Category
- Movie Plot

- *I cancelled the date with my girlfriend just to watch my favorite star featuring in this movie.*

Facets of a Movie Review

- General Perception about the Crew
 - Opinion about the Characters in the Movie
 - Characteristics of a movie or genre
- Objective Facts about the Crew and Movies
- Past Performance of the Crew and Movies
 - Opinion about the Movie and Crew
- Expectations from the Movie or Crew
 - Unrelated Category
- Movie Plot
 - *L.I.E. stands for Long Island Expressway, which slices through the strip malls and middle-class homes of suburbia. Filmmaker Michael Cuesta uses it as a (pretty transparent) metaphor of dangerous escape for his 15-year old protagonist, Howie (Paul Franklin Dano).*

Facets of a Movie Review

- General Perception about the Crew
 - Opinion about the Characters in the Movie
 - Characteristics of a movie or genre
- Objective Facts about the Crew and Movies
- Past Performance of the Crew and Movies
 - Opinion about the Movie and Crew
- Expectations from the Movie or Crew
 - Unrelated Category
- Movie Plot
 - *He's an even-tempered, funny, robust old man who actually listens to the kids' problems (as opposed to their parents and friends, both caught up in the high-wire act of their own confused lives.).*

Facets of a Movie Review

- General Perception about the Crew
 - Opinion about the Characters in the Movie
 - Characteristics of a movie or genre
- Objective Facts about the Crew and Movies
- Past Performance of the Crew and Movies
 - **Opinion about the Movie and Crew**
- Expectations from the Movie or Crew
 - Unrelated Category
- Movie Plot
 - *Horror movies are supposed to be scary.*
 - *There is an axiom that directors who have a big hit with their debut have a big bomb with their second film.*

Facets of a Movie Review

- General Perception about the Crew
 - Opinion about the Characters in the Movie
 - Characteristics of a movie or genre
- Objective Facts about the Crew and Movies
- Past Performance of the Crew and Movies
 - **Opinion about the Movie and Crew**
- Expectations from the Movie or Crew
 - Unrelated Category
- Movie Plot

- *While the movie is brutal, the violence is neither very graphic nor gratuitous. It may scare the little ones, but the teen-age audience for which it is aimed will appreciate the man-eating chomping that runs through the film.*

Facets of a Movie Review

- General Perception about the Crew
 - Opinion about the Characters in the Movie
 - Characteristics of a movie or genre
- Objective Facts about the Crew and Movies
- Past Performance of the Crew and Movies
 - Opinion about the Movie and Crew
- Expectations from the Movie or Crew
 - Unrelated Category
- Movie Plot

- *So my grandson gives me passes to this new picture One Night at McCool's because the free screening is the same night as that horrible show with those poor prisoners trapped on the island who eat the bugs. "Go," he says, "it's just like Rush-o-Man."*

Wikipedia

- Extensive World Knowledge required to perform this analysis
- Wikipedia is used to create a *topic-specific, extractive summary* of a review
- The extract is classified with a Lexicon, instead of the entire review

Feature Extraction from Wikipedia

Feature Extraction from Wikipedia

METADATA

Harry Potter and the Deathly Hallows – Part 1 is a 2010 fantasy film[5] directed by David Yates and the first of two films based on the novel Harry Potter and the Deathly Hallows by J. K. Rowling. It is the seventh instalment in the Harry Potter film series, written by Steve Kloves and produced by David Heyman, David Barron and Rowling. The story follows Harry Potter on a quest to find and destroy Lord Voldemort's secret to immortality – the Horcruxes. The film stars Daniel Radcliffe as Harry Potter, alongside Rupert Grint and Emma Watson as Harry's best friends Ron Weasley and Hermione Granger. It is the sequel to Harry Potter and the Half-Blood Prince and is followed by the concluding film, Harry Potter and the Deathly Hallows – Part 2.

Principal photography began on 19 February 2009 and was completed on 12 June 2010.[6] Part 1 was released in 2D cinemas and IMAX formats worldwide on 19 November 2010.

Feature Extraction from Wikipedia

PLOT

Further information: Harry Potter and the Deathly Hallows Novel Plot

Minister Rufus Scrimgeour addresses the wizarding media stating that the Ministry of Magic will remain strong as Lord Voldemort gains power throughout the wizarding and Muggle worlds. Severus Snape arrives at Malfoy Manor to inform Lord Voldemort and his Death Eaters of Harry's departure from No. 4 Privet Drive. Voldemort commandeers Lucius Malfoy's wand, as Voldemort's own wand cannot be used to kill Harry; their wands are "twins".

Meanwhile, the Order of the Phoenix arrive at Privet Drive and escort Harry to safety using Polyjuice Potion to create six decoy Harrys. During their flight to the Burrow, they are ambushed by Death Eaters, who kill Mad-Eye Moody and Hedwig, and injure George Weasley. They arrive at the Burrow, where Harry has a vision of Ollivander being tormented by Voldemort, who claims that the wand-maker had lied to him by informing him of the only way to kill Harry: obtaining another's wand.

Feature Extraction from Wikipedia

CREW

Directed by : David Yates

Produced by : David Heyman,David Barro, J. K. Rowling

**Screenplay by : Steve Kloves Based on Harry Potter and the
Deathly Hallows by J. K. Rowling**

Starring : Daniel Radcliffe, Rupert Grint, Emma Watson

Music by : Alexandre Desplat

Themes: John Williams

Cinematography : Eduardo Serra

Editing by: Mark Day

Studio :Heyday Films

Distributed by: Warner Bros. Pictures

**that the wand-maker had lied to him by informing him of the only
way to kill Harry: obtaining another's wand.**

Feature Extraction from Wikipedia

CREW

Directed by : David Yates

Produced by : David Heyman,David Barro, J. K. Rowling

**Screenplay by : Steve Kloves Based on Harry Potter and the
Deathly Hallows by J. K. Rowling**

Starring : Daniel Radcliffe, Rupert Grint, Emma Watson

Music by : Alexandre Desplat

Themes: John Williams

Cinematography : Eduardo Serra

Editing by: Mark Day

Studio :Heyday Films

Distributed by: Warner Bros. Pictures

Narcissa Malfoy.

Ralph Fiennes as Lord Voldemort, the film's main antagonist

Feature Extraction from Wikipedia

DOMAIN SPECIFIC FEATURE LIST

Movie, Staffing, casting, Writing, Theory, Writing, Rewriting, Screenplay, Format, Treatments, Scriptments, Synopsis, Logline, Pitching, Certification, scripts, Budget, Ideas, Funding, budgeting, Funding, Plans, Grants, Pitching, Tax, Contracts, law, Copyright, Pre-production, Budgeting, Scheduling, Pre-production, film , stock, Story, boarding, plot, Casting , Directors, Location, Scouting,

Ralph Fiennes as Lord Voldemort, the film's main antagonist

Feature List Creation



Feature List Creation

- Metadata and Plot sentences are POS-tagged
 - Nouns retrieved
 - Stemmed
 - Added to the **Plot** list
 - Character information also added to **Plot** list

Feature List Creation

- Metadata and Plot sentences are POS-tagged
 - Nouns retrieved
 - Stemmed
 - Added to the **Plot** list
 - Character information also added to **Plot** list
- Domain specific feature list extracted from Wiki articles on films and movies
 - Added to **MovieFeature** list

Feature List Creation

- Metadata and Plot sentences are POS-tagged
 - Nouns retrieved
 - Stemmed
 - Added to the **Plot** list
 - Character information also added to **Plot** list
- Domain specific feature list extracted from Wiki articles on films and movies
 - Added to **MovieFeature** list
- Crew information added to **Crew** list

Feature List Creation

- Metadata and Plot sentences are POS-tagged
 - Nouns retrieved
 - Stemmed
 - Added to the **Plot** list
 - Character information also added to **Plot** list
- Domain specific feature list extracted from Wiki articles on films and movies
 - Added to **MovieFeature** list
- Crew information added to **Crew** list
- Laptop and Printer domain data used to filter frequent occurring concepts
 - Non-overlapping domains
 - Frequently occurring terms in all the domains added to **FreqWords** list
 - The **FreqWords** are pruned from the other feature lists

Why only Nouns?

149

- To restrict genre-specific concepts to be entities
- *Harry acted as if nothing has happened* vs. *Kate Winslet acted awesome in the movie.*
- Here *act* is present as a Verb in both the sentences
 - First sentence belongs to the *plot* (Category 5)
 - Second sentence depicts the reviewer opinion (Category 8)
- Difference lies in the presence of different subjects of interest with the Verb
- Our focus is to capture the *subjects* and *objects* in the sentence which give direct clues about the *category* of the reviewer statement, so that feature lists are as pure as possible

Extractive Summary

- Given a review R with n sentences S_i , determine if each sentence S_i is to be accepted or rejected based on a *relevancy factor* in judging this movie
- $Rel_{factor_i} = 2 \sum_j 1_{w_{ij} \in Crew \text{ or } MovieTitle} + \sum_j 1_{w_{ij} \in MovieFeature} - \sum_j 1_{w_{ij} \in Plot, w_{ij} \notin Crew, w_{ij} \notin MovieTitle}$
- $Acc_{factor_i} = 1 \text{ if } Rel_{factor_i} \geq 0 \text{ and } \exists w_{ij} \in S_i$
 $s.t. w_{ij} \in Crew \text{ or } MovieFeature \text{ or } MovieTitle$
 $= 0 \text{ otherwise}$

Semi-Supervised Learning of Parameters

151

- Equations 2 can be re-written as:
- $Rel_{factor_i} = \alpha \times X_{i,1} + \beta \times X_{i,2} - \gamma \times X_{i,3}$
- $Acc_{factor_i} = Rel_{factor_i} - \theta$
= $\alpha \times X_{i,1} + \beta \times X_{i,2} - \gamma \times X_{i,3} - \theta$
- = $\alpha \times X_{i,1} + \beta \times X_{i,2} - \gamma \times X_{i,3} - \theta \times X_{i,4}$ (where $X_{i,4} = 1$)
- Let Y_i be the binary label information corresponding to each sentence in the development set, where $Y_i = 1$ if $Acc_{factor_i} \geq 0$ and -1 otherwise.
- $Y_i = \mathbf{W} \cdot X_i$ where,
 $\mathbf{W} = [\alpha \ \beta \ -\gamma \ -\theta]^T$ and $X_i = [X_{i,1} \ X_{i,2} \ X_{i,3} \ X_{i,4}]$
- or, $\mathbf{Y} = \mathbf{W}^T \cdot \mathbf{X}$

Algorithm



Algorithm

Input : Review R

Output: OpinionSummary

Step 1: Extract the Crew list from Wikipedia

Step 2: Extract the Plot list from Wikipedia

Step 3: Extract the MovieFeature list from Wikipedia

Step 4: Extract the FreqWords list as the common frequently occurring concepts in Mobile Phone, Printer and Movie domains.

Let $OpinionSummary = \emptyset$

for $i=1..n$

if $Acc_{factor_i} == 1$

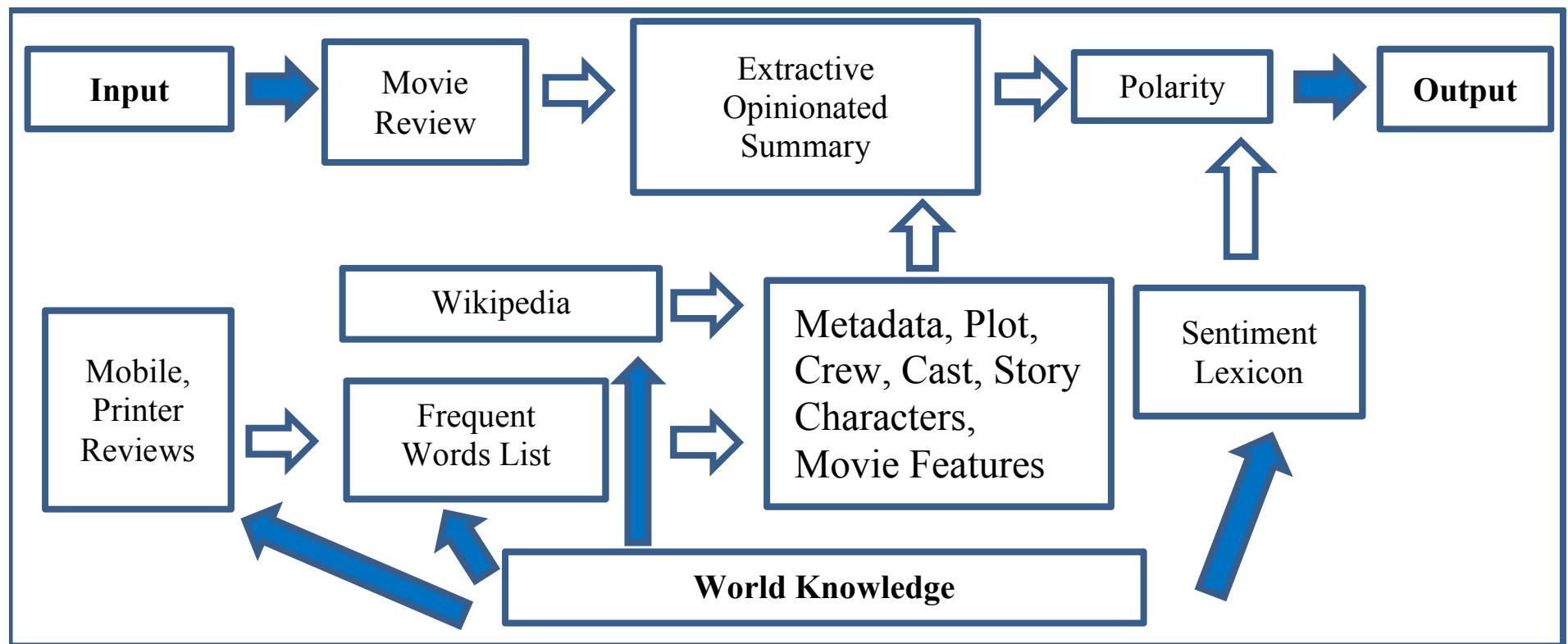
add S_i **to** $OpinionSummary$

end if

end for

WikiSent

154



Algorithm Demonstration

- In Sentence [1], **Brian Cox** is the only keyword present and it belongs to the Cast list. $Rel_{factor_i} = 2*1 + 1*0 - 0*0 = 2 > -1$ and the sentence is accepted.
- In [2], there is no keyword from the lists and it is rejected.
- [3] has the keyword **acting** from MovieFeature and is accepted.
- [4] has the keywords **Cox**, **L.I.E** from Cast, **MovieTitle**, **John Harrigan** from Character list and **distributor** from MovieFeature list. $Rel_{factor_i} = 2*2 + 1 - 1 = 4 > -1$ and is accepted.
- [5] has only the keyword **Big John** from Character and $Rel_{factor_i} = 0 + 0 - 1 = -1$ and is rejected.
- [6] has the keyword **audience** from MovieFeature and **Big John** from Character and its $Rel_{factor_i} = 0 + 1 - 1 = 0 > -1$ and is accepted.
- [7] has the keywords **temper**, **friend** from Plot and $Rel_{factor_i} = 0 + 0 - 2 = -2$ and is rejected.
- [8] has the keywords **sex**, **charm** from Plot and $Rel_{factor_i} = 0 + 0 - 2 = -2$ and is rejected.

Algorithm Demonstration

- In Sentence [1], **Brian Cox** is the only keyword present and it belongs to the Cast list. $Rel_{factor_i} = 2*1 + 1*0 - 0*0 = 2 > -1$ and the sentence is accepted.
- In [2], there is no keyword from the lists and it is rejected.
- [3] has the keyword **acting** from MovieFeature and is accepted.
- [4] has the keywords **Cox**, **L.I.E** from Cast, **MovieTitle**, **John Harrigan** from Character list and **distributor** from MovieFeature list. $Rel_{factor_i} = 2*2 + 1 - 1 = 4 > -1$ and is accepted.
- [5] has only the keyword **Big John** from Character and $Rel_{factor_i} = 0 + 0 - 1 = -1$ and is rejected.
- [6] has the keyword **audience** from MovieFeature and **Big John** from Character and its $Rel_{factor_i} = 0 + 1 - 1 = 0 > -1$ and is accepted.
- [7] has the keywords **temper**, **friend** from Plot and $Rel_{factor_i} = 0 + 0 - 2 = -2$ and is rejected.
- [8] has the keywords **sex**, **charm** from Plot and $Rel_{factor_i} = 0 + 0 - 2 = -2$ and is rejected.

*Best remembered for his understated performance as Dr. Hannibal Lecter in Michael Mann's forensics thriller, Manhunter, Scottish character actor **Brian Cox** brings something special to every movie he works on.*

Algorithm Demonstration

- In Sentence [1], **Brian Cox** is the only keyword present and it belongs to the Cast list. $Rel_{factor_i} = 2*1 + 1*0 - 0*0 = 2 > -1$ and the sentence is accepted.
- In [2], there is no keyword from the lists and it is rejected.
- [3] has the keyword **acting** from MovieFeature and is accepted.
- [4] has the keywords **Cox**, **L.I.E** from Cast, **MovieTitle**, **John Harrigan** from Character list and **distributor** from MovieFeature list. $Rel_{factor_i} = 2*2 + 1 - 1 = 4 > -1$ and is accepted.
- [5] has only the keyword **Big John** from Character and $Rel_{factor_i} = 0 + 0 - 1 = -1$ and is rejected.
- [6] has the keyword **audience** from MovieFeature and **Big John** from Character and its $Rel_{factor_i} = 0 + 1 - 1 = 0 > -1$ and is accepted.
- [7] has the keywords **temper**, **friend** from Plot and $Rel_{factor_i} = 0 + 0 - 2 = -2$ and is rejected.
- [8] has the keywords **sex**, **charm** from Plot and $Rel_{factor_i} = 0 + 0 - 2 = -2$ and is rejected.

*Usually playing a bit role in some studio schlock (he dies halfway through *The Long Kiss Goodnight*), he's only occasionally given something meaty and substantial to do.*

Algorithm Demonstration

- In Sentence [1], **Brian Cox** is the only keyword present and it belongs to the Cast list. $Rel_{factor_i} = 2*1 + 1*0 - 0*0 = 2 > -1$ and the sentence is accepted.
- In [2], there is no keyword from the lists and it is rejected.
- [3] has the keyword **acting** from MovieFeature and is accepted.
- [4] has the keywords **Cox**, **L.I.E** from Cast, **MovieTitle**, **John Harrigan** from Character list and **distributor** from MovieFeature list. $Rel_{factor_i} = 2*2 + 1 - 1 = 4 > -1$ and is accepted.
- [5] has only the keyword **Big John** from Character and $Rel_{factor_i} = 0 + 0 - 1 = -1$ and is rejected.
- [6] has the keyword **audience** from MovieFeature and **Big John** from Character and its $Rel_{factor_i} = 0 + 1 - 1 = 0 > -1$ and is accepted.
- [7] has the keywords **temper**, **friend** from Plot and $Rel_{factor_i} = 0 + 0 - 2 = -2$ and is rejected.
- [8] has the keywords **sex**, **charm** from Plot and $Rel_{factor_i} = 0 + 0 - 2 = -2$ and is rejected.

If you want to see some brilliant **acting**, check out his work as a dogged police inspector opposite Frances McDormand in Ken Loach's *Hidden Agenda*.

Algorithm Demonstration

- In Sentence [1], **Brian Cox** is the only keyword present and it belongs to the Cast list. $Rel_{factor_i} = 2*1 + 1*0 - 0*0 = 2 > -1$ and the sentence is accepted.
- In [2], there is no keyword from the lists and it is rejected.
- [3] has the keyword **acting** from MovieFeature and is accepted.
- [4] has the keywords **Cox**, **L.I.E** from Cast, **MovieTitle**, **John Harrigan** from Character list and **distributor** from MovieFeature list. $Rel_{factor_i} = 2*2 + 1 - 1 = 4 > -1$ and is accepted.
- [5] has only the keyword **Big John** from Character and $Rel_{factor_i} = 0 + 0 - 1 = -1$ and is rejected.
- [6] has the keyword **audience** from MovieFeature and **Big John** from Character and its $Rel_{factor_i} = 0 + 1 - 1 = 0 > -1$ and is accepted.
- [7] has the keywords **temper**, **friend** from Plot and $Rel_{factor_i} = 0 + 0 - 2 = -2$ and is rejected.
- [8] has the keywords **sex**, **charm** from Plot and $Rel_{factor_i} = 0 + 0 - 2 = -2$ and is rejected.

Cox plays the role of Big John Harrigan in the disturbing new indie flick L.I.E., which Lot 47 picked up at Sundance when other distributors were scared to budge.

Algorithm Demonstration

- In Sentence [1], **Brian Cox** is the only keyword present and it belongs to the Cast list. $Rel_{factor_i} = 2*1 + 1*0 - 0*0 = 2 > -1$ and the sentence is accepted.
- In [2], there is no keyword from the lists and it is rejected.
- [3] has the keyword **acting** from MovieFeature and is accepted.
- [4] has the keywords **Cox**, **L.I.E** from Cast, **MovieTitle**, **John Harrigan** from Character list and **distributor** from MovieFeature list. $Rel_{factor_i} = 2*2 + 1 - 1 = 4 > -1$ and is accepted.
- [5] has only the keyword **Big John** from Character and $Rel_{factor_i} = 0 + 0 - 1 = -1$ and is rejected.
- [6] has the keyword **audience** from MovieFeature and **Big John** from Character and its $Rel_{factor_i} = 0 + 1 - 1 = 0 > -1$ and is accepted.
- [7] has the keywords **temper**, **friend** from Plot and $Rel_{factor_i} = 0 + 0 - 2 = -2$ and is rejected.
- [8] has the keywords **sex**, **charm** from Plot and $Rel_{factor_i} = 0 + 0 - 2 = -2$ and is rejected.

Big John feels the love that dares not speak its name, but he expresses it through seeking out adolescents and bringing them back to his pad.

Algorithm Demonstration

- In Sentence [1], **Brian Cox** is the only keyword present and it belongs to the Cast list. $Rel_{factor_i} = 2*1 + 1*0 - 0*0 = 2 > -1$ and the sentence is accepted.
- In [2], there is no keyword from the lists and it is rejected.
- [3] has the keyword **acting** from MovieFeature and is accepted.
- [4] has the keywords **Cox**, **L.I.E** from Cast, **MovieTitle**, **John Harrigan** from Character list and **distributor** from MovieFeature list. $Rel_{factor_i} = 2*2 + 1 - 1 = 4 > -1$ and is accepted.
- [5] has only the keyword **Big John** from Character and $Rel_{factor_i} = 0 + 0 - 1 = -1$ and is rejected.
- [6] has the keyword **audience** from MovieFeature and **Big John** from Character and its $Rel_{factor_i} = 0 + 1 - 1 = 0 > -1$ and is accepted.
- [7] has the keywords **temper**, **friend** from Plot and $Rel_{factor_i} = 0 + 0 - 2 = -2$ and is rejected.
- [8] has the keywords **sex**, **charm** from Plot and $Rel_{factor_i} = 0 + 0 - 2 = -2$ and is rejected.

*What bothered some **audience** members was the presentation of **Big John** in an oddly empathetic light.*

Algorithm Demonstration

- In Sentence [1], **Brian Cox** is the only keyword present and it belongs to the Cast list. $Rel_{factor_i} = 2*1 + 1*0 - 0*0 = 2 > -1$ and the sentence is accepted.
- In [2], there is no keyword from the lists and it is rejected.
- [3] has the keyword **acting** from MovieFeature and is accepted.
- [4] has the keywords **Cox**, **L.I.E** from Cast, **MovieTitle**, **John Harrigan** from Character list and **distributor** from MovieFeature list. $Rel_{factor_i} = 2*2 + 1 - 1 = 4 > -1$ and is accepted.
- [5] has only the keyword **Big John** from Character and $Rel_{factor_i} = 0 + 0 - 1 = -1$ and is rejected.
- [6] has the keyword **audience** from MovieFeature and **Big John** from Character and its $Rel_{factor_i} = 0 + 1 - 1 = 0 > -1$ and is accepted.
- [7] has the keywords **temper**, **friend** from Plot and $Rel_{factor_i} = 0 + 0 - 2 = -2$ and is rejected.
- [8] has the keywords **sex**, **charm** from Plot and $Rel_{factor_i} = 0 + 0 - 2 = -2$ and is rejected.

He's an even-tempered, funny, robust old man who actually listens to the kids' problems (as opposed to their parents and friends, both caught up in the high-wire act of their own confused lives.).

Algorithm Demonstration

- In Sentence [1], **Brian Cox** is the only keyword present and it belongs to the Cast list. $Rel_{factor_i} = 2*1 + 1*0 - 0*0 = 2 > -1$ and the sentence is accepted.
- In [2], there is no keyword from the lists and it is rejected.
- [3] has the keyword **acting** from MovieFeature and is accepted.
- [4] has the keywords **Cox**, **L.I.E** from Cast, **MovieTitle**, **John Harrigan** from Character list and **distributor** from MovieFeature list. $Rel_{factor_i} = 2*2 + 1 - 1 = 4 > -1$ and is accepted.
- [5] has only the keyword **Big John** from Character and $Rel_{factor_i} = 0 + 0 - 1 = -1$ and is rejected.
- [6] has the keyword **audience** from MovieFeature and **Big John** from Character and its $Rel_{factor_i} = 0 + 1 - 1 = 0 > -1$ and is accepted.
- [7] has the keywords **temper**, **friend** from Plot and $Rel_{factor_i} = 0 + 0 - 2 = -2$ and is rejected.
- [8] has the keywords **sex**, **charm** from Plot and $Rel_{factor_i} = 0 + 0 - 2 = -2$ and is rejected.

He'll have sex-for-pay with them only after an elaborate courtship, charming them with temptations from the grown-up world"

Lexicons and Datasets

164

Lexicons and Datasets

165

- Lexicons
 - SentiWordNet (Esuli *et al.*, 2006)
 - Subjectivity Lexicon (Wilson *et al.*, 2005)
 - General Inquirer (Stone *et al.*, 1966)

Lexicons and Datasets

166

- Lexicons
 - SentiWordNet (Esuli *et al.*, 2006)
 - Subjectivity Lexicon (Wilson *et al.*, 2005)
 - General Inquirer (Stone *et al.*, 1966)
- Baseline 1
 - Bag-of-Words based on the 3 lexicons
- Baseline 2
 - SO-CAL (Taboada *et al.*, 2011)
- Baseline 3
 - All the semi-supervised and unsupervised systems in the domain

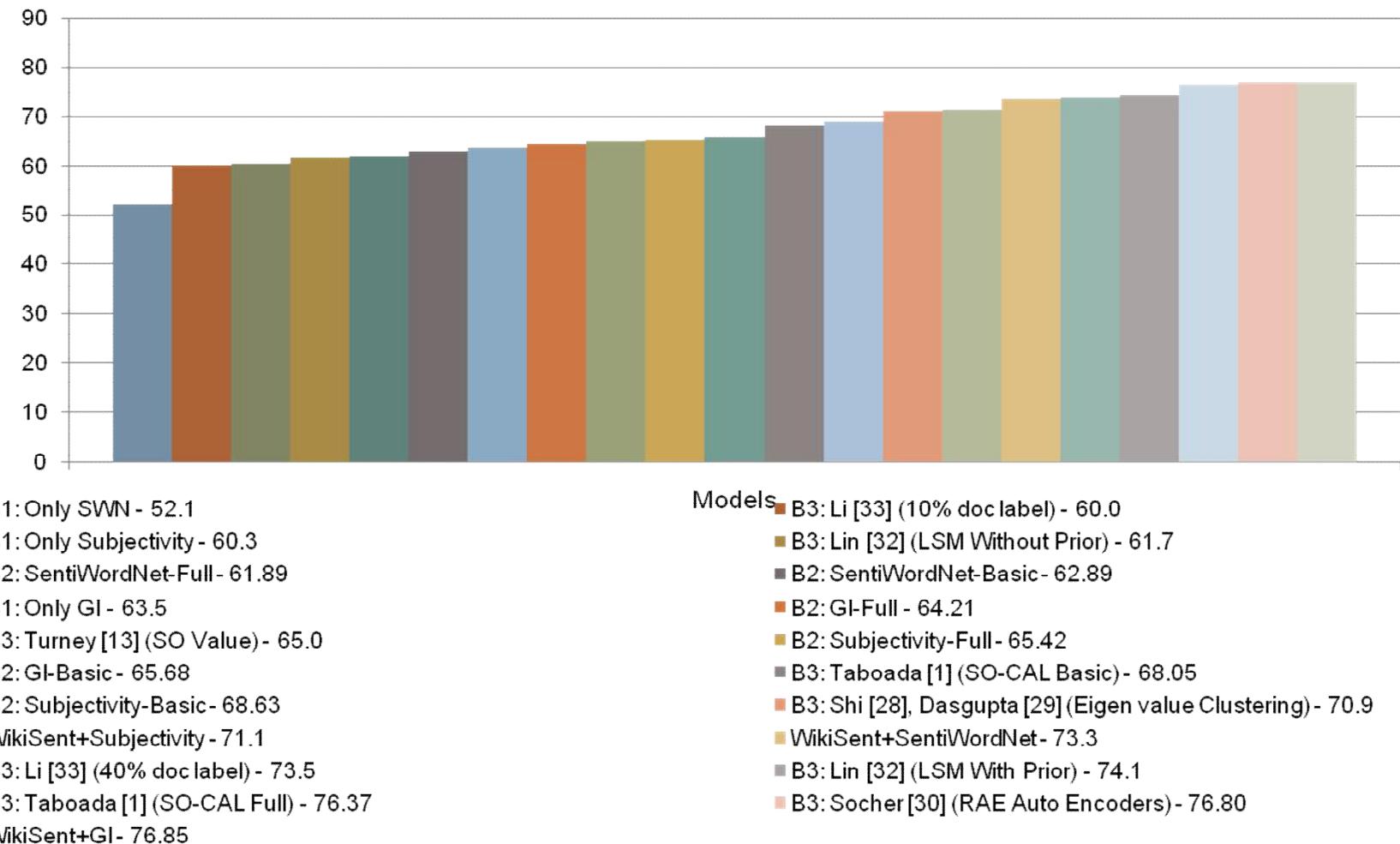
Lexicons and Datasets

167

- Lexicons
 - SentiWordNet (Esuli *et al.*, 2006)
 - Subjectivity Lexicon (Wilson *et al.*, 2005)
 - General Inquirer (Stone *et al.*, 1966)
- Baseline 1
 - Bag-of-Words based on the 3 lexicons
- Baseline 2
 - SO-CAL (Taboada *et al.*, 2011)
- Baseline 3
 - All the semi-supervised and unsupervised systems in the domain
- Movie Review Dataset (Pang *et al.*, 2002)
 - 1000 positive and 1000 negative reviews (labeled at document level)
 - 27,000 untagged reviews

Accuracy Comparison with All the Semi-supervised and Unsupervised Systems in the Same Dataset

168



Accuracy Comparison with Different Baselines

Accuracy Comparison with Different Baselines

170

Baseline 1 : Simple Bag-of-Words	
Only SentiWordNet	52.1
Only Subjectivity	60.3
Only GI	63.5
Baseline 2 : Worse, Median and Best Performing Lexicons with SO-CAL	
SentiWordNet-Full	61.89
SentiWordNet-Basic	62.89
GI-Full	64.21
GI-Basic	65.68
Subjectivity-Full	65.42
Subjectivity-Basic	68.63
WikiSent with Different Lexicons	
WikiSent+Subjectivity	71.1
WikiSent+SentiWordNet	73.3
WikiSent+GI	76.85

Accuracy Comparison with Different Baselines

171

Systems	Classification Method	Accuracy
Turney SO value	Unsupervised (PMI)	65
Taboada SO-CAL Basic [1]	Lexicon Generation	68.05
Taboada SO-CAL Full [1]	Lexicon Generation	76.37
Shi [28], Dasgupta [29]	Unsupervised Eigen Vector Clustering	70.9
Socher [30] RAE	Semi Supervised Auto Encoders	76.8
Lin [32] LSM	Unsupervised without prior info	61.7
Lin [32] LSM	Weakly Supervised with prior info	74.1
Li [33]	Semi Supervised 10% doc. Label	60
Li [33]	Semi Supervised 40% doc. Label	60
WikiSent	Wikipedia+GI Lexicon	76.85

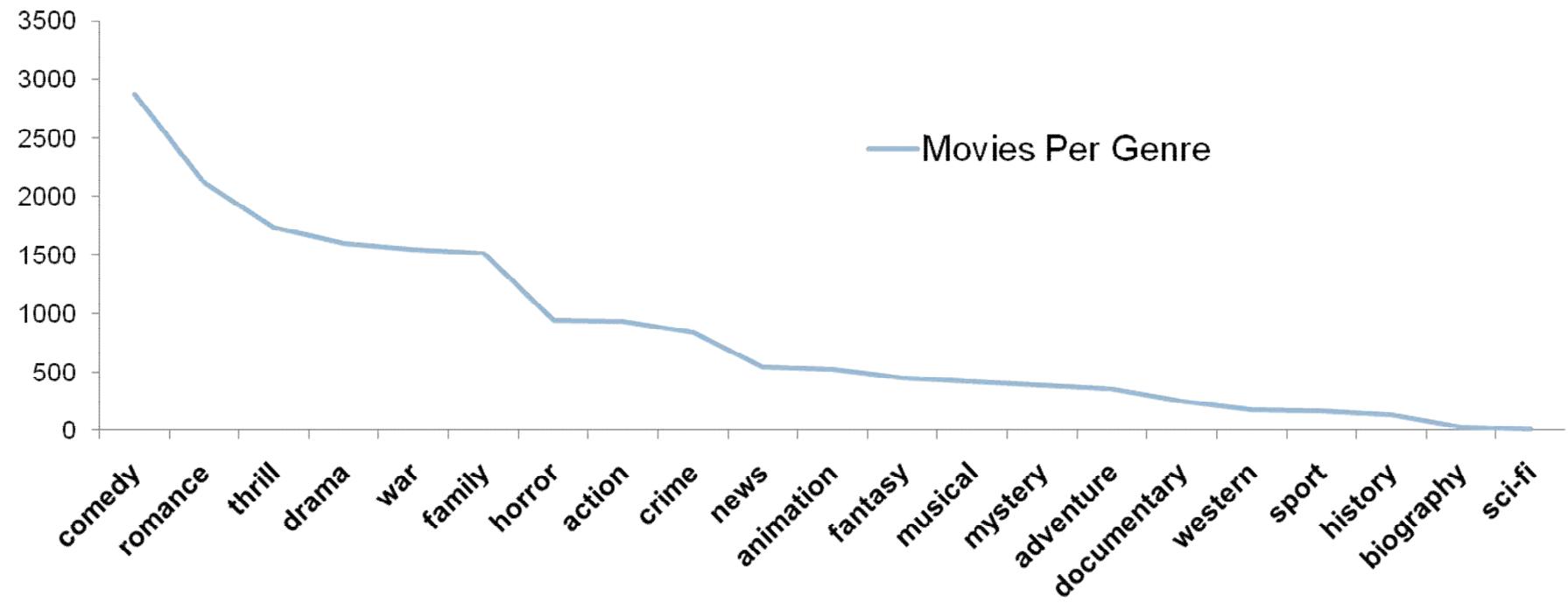
Trend Analysis

172

$$\text{Genre Popularity} = \frac{\text{Positive Movie Reviews per Genre}}{\text{Total Movie Reviews per Genre}}$$

Trend Analysis

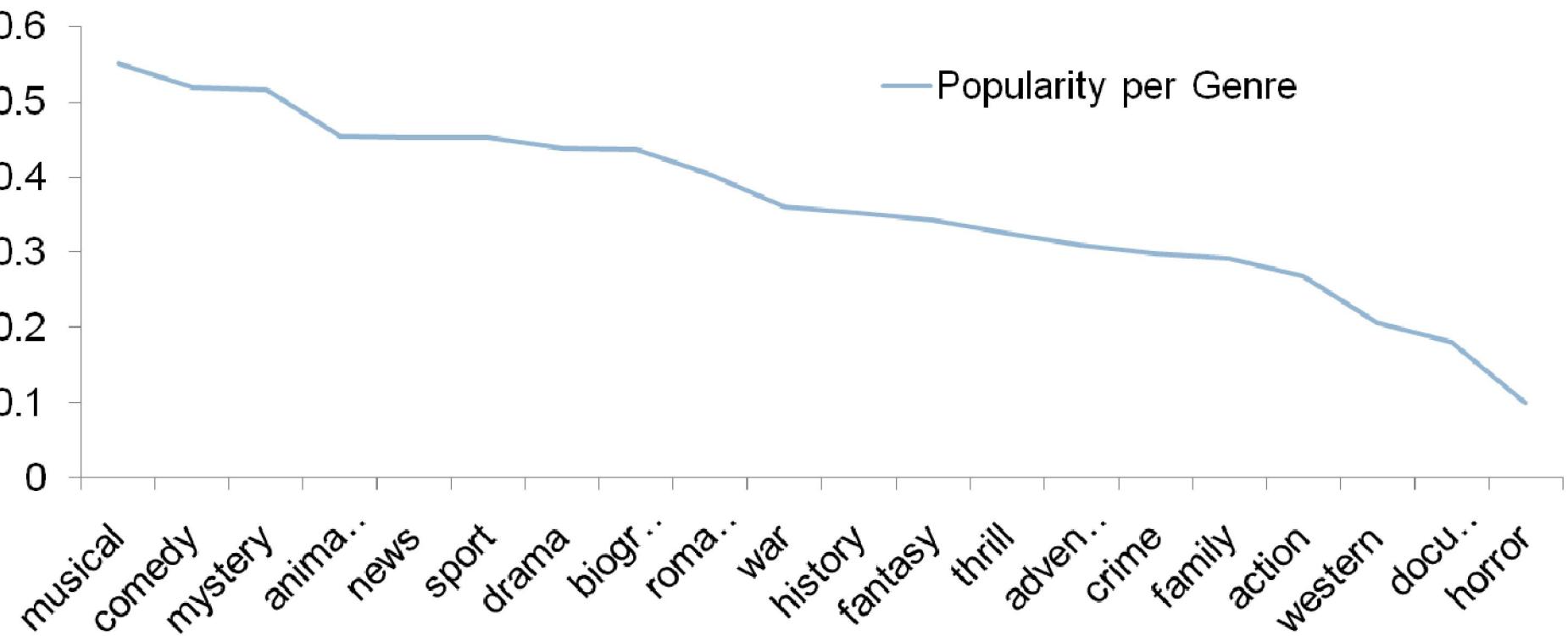
173



$$\text{Genre Popularity} = \frac{\text{Positive Movie Reviews per Genre}}{\text{Total Movie Reviews per Genre}}$$

Trend Analysis

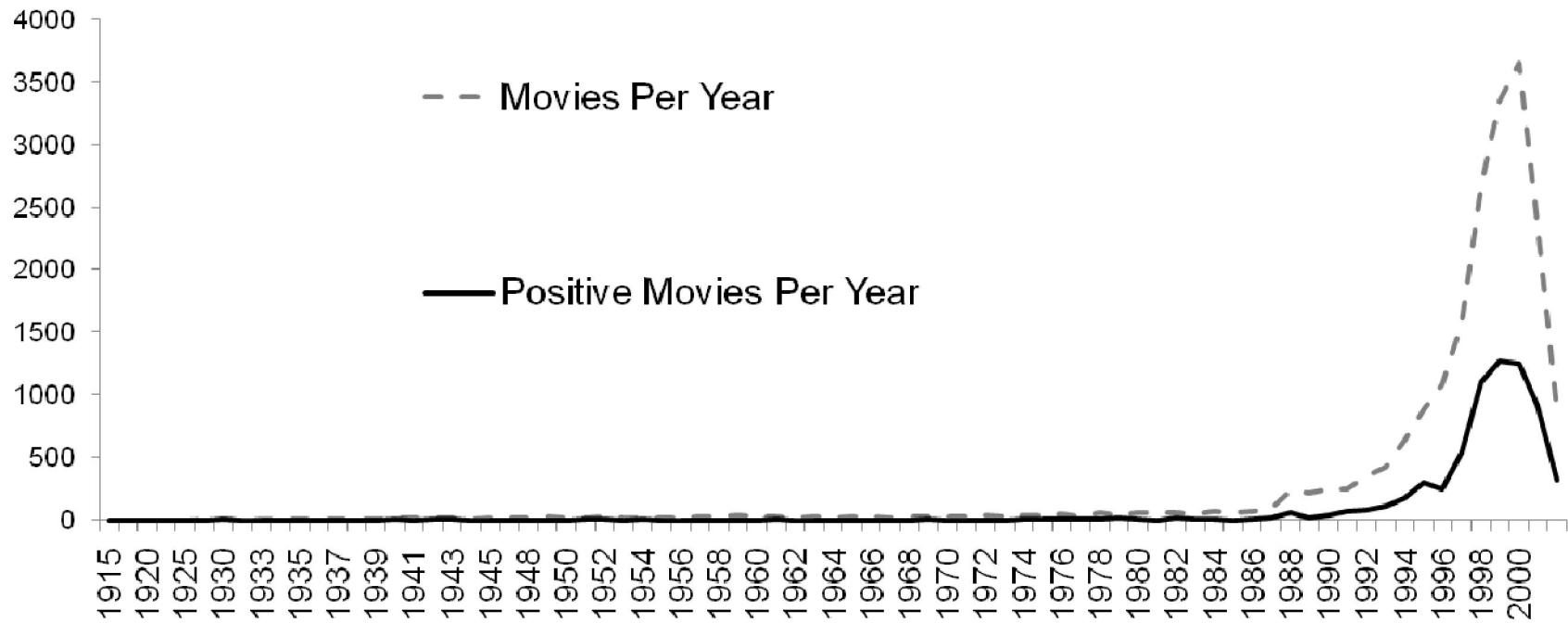
174



$$\text{Genre Popularity} = \frac{\text{Positive Movie Reviews per Genre}}{\text{Total Movie Reviews per Genre}}$$

Trend Analysis

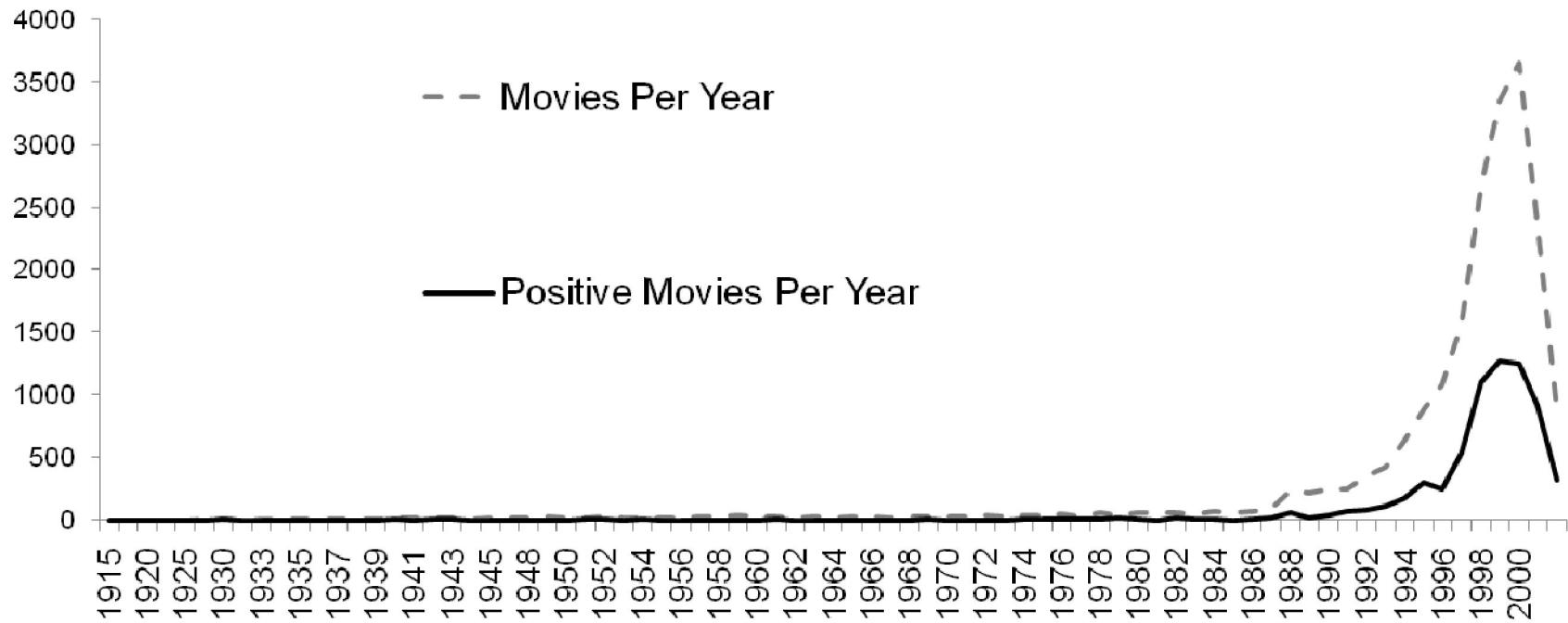
175



$$\text{Genre Popularity} = \frac{\text{Positive Movie Reviews per Genre}}{\text{Total Movie Reviews per Genre}}$$

Trend Analysis

176



	WikiSent	Bag-of-Words Baseline 1
Positive Reviews (%)	48.95	81.2
Negative Reviews (%)	51.05	18.79

Drawbacks

177

- Absence of a co-reference resolution module
 - *false negative*
- Synonymous concepts not handled
 - Does not matter much as genre-specific concepts like *acting*, *direction*, *story-writer* occur in same lexical form
- Reviewer opinion bias can affect the system
 - *false positive*
- Absence of WSD module affect lexicon-based classification

□Role of Social Media Content and Informal Language Form in Sentiment Analysis

Social Media Analysis

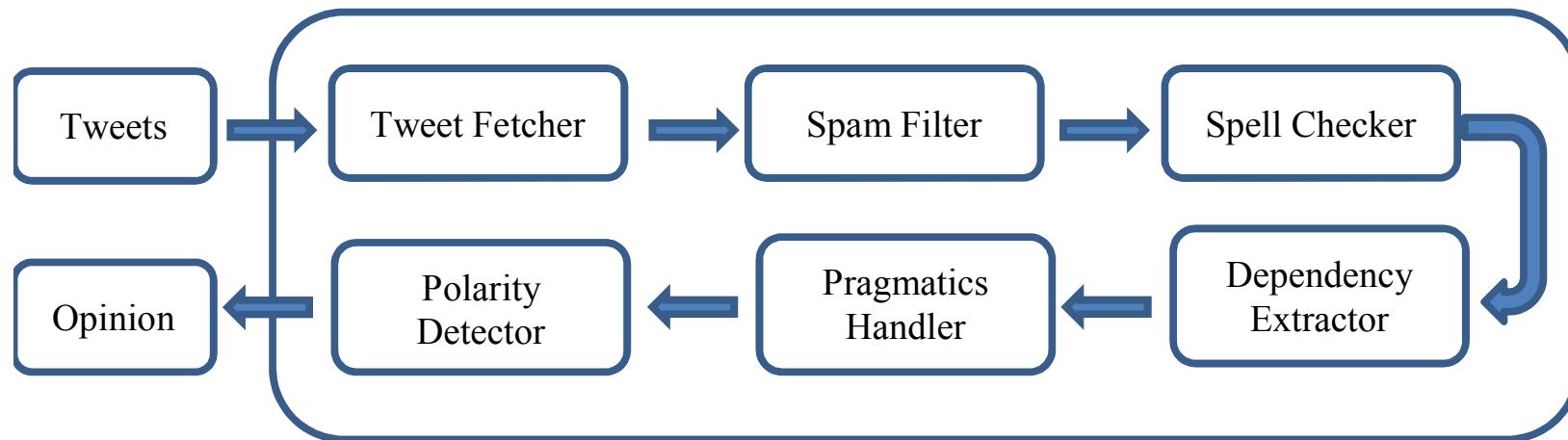
179

- 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.....

TwiSent: Multi-Stage System Architecture

180



Spam Categorization and Features

181

- *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">1. Number of Words per Tweet2. Average Word Length3. Frequency of "?" and "!"4. Frequency of Numeral Characters5. Frequency of hashtags6. Frequency of @users7. Extent of Capitalization8. Frequency of the First POS Tag	<ol style="list-style-type: none">9. Frequency of Foreign Words10. Validity of First Word11. Presence / Absence of links12. Frequency of POS Tags13. Strength of Character Elongation14. Frequency of Slang Words15. Average Positive and Negative Sentiment of Tweets
--	--

Algorithm for Spam Filter

182

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

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

Categorization of Noisy Text

183

- 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 - *hapyyyyy* (*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

184

- 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 *happyyyyy* 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
 - *happyyyyyy*, normalized to *hapy* and thereafter spell-corrected to *happy*, is marked so as to not lose its pragmatic content

Spell-Checker Algorithm

185

- **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..

186

- 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**
- edit₁=edit_distance(s' , s)
- */*t.replace(char1,char2) replaces all occurrences of char1 in the string t with char2*/*
- edit₂=edit_distance(t.replace(a , e), s')
- edit₃=edit_distance(t.replace(e , a), s')
- edit₄=edit_distance(t.replace(o , u), s')
- edit₅=edit_distance(t.replace(u , o), s')
- edit₆=edit_distance(t.replace(i , e), s')
- edit₇=edit_distance(t.replace(e, i), s')
- count=overlapping_characters(t , s')
- min_edit=
- min(edit₁,edit₂,edit₃,edit₄,edit₅,edit₆,edit₇)
- **if** (min_edit ==0 or score[s] == 0) **then**
- adv=-2 */* for exact match assign advantage score */*
- **else**
- adv=0
- **end if**
- final_score[t]=min_edit+adv+score[w]-count;
- **end for**
- return t with minimum final_score;

Pragmatics

187

Pragmatics

188

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

Pragmatics

189

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

Pragmatics

190

- *Elongation of a word, repeating alphabets multiple times* - Example: *happyyyyyyy, 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

191

- *Elongation of a word, repeating alphabets multiple times* - Example: *happyyyyyyyyy, 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

192

2-Class
Classification

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

4-Class
Classification

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

TwiSent Evaluation

TwiSent Evaluation

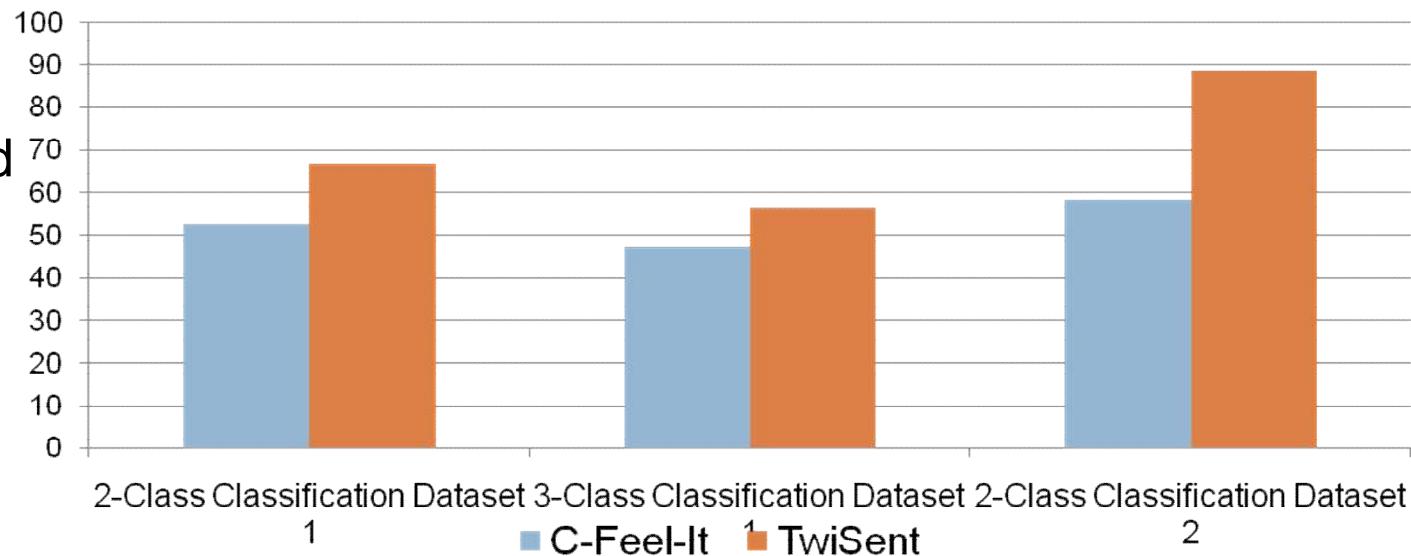
194

Lexicon-based
Classification

TwiSent Evaluation

195

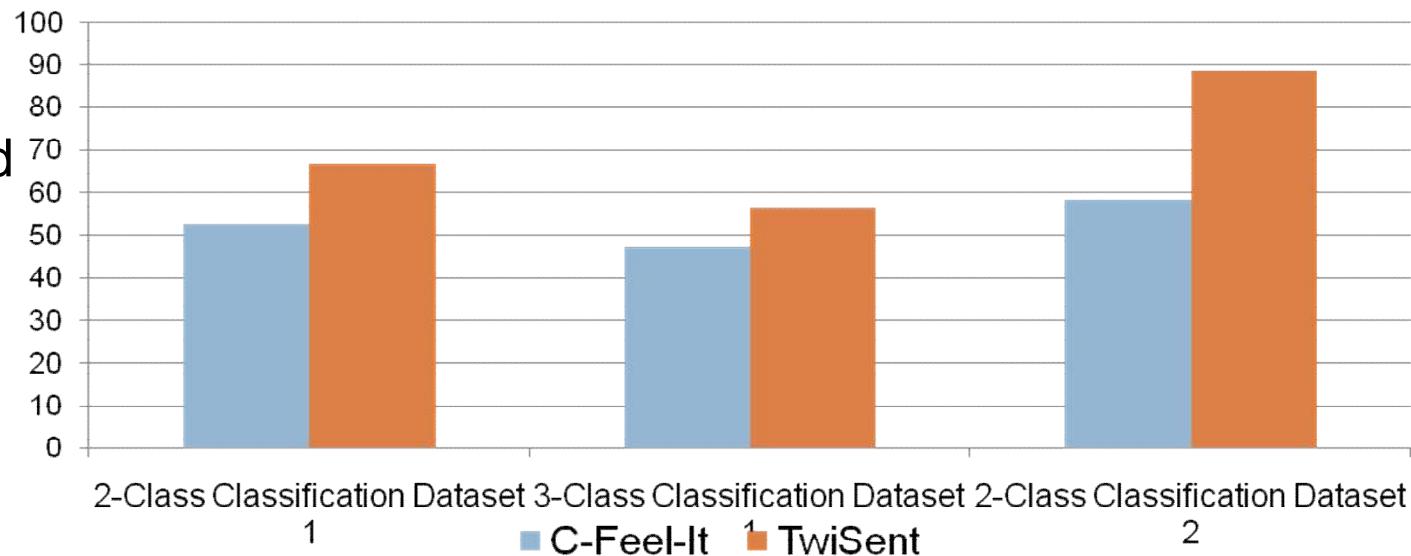
Lexicon-based
Classification



TwiSent Evaluation

196

Lexicon-based
Classification

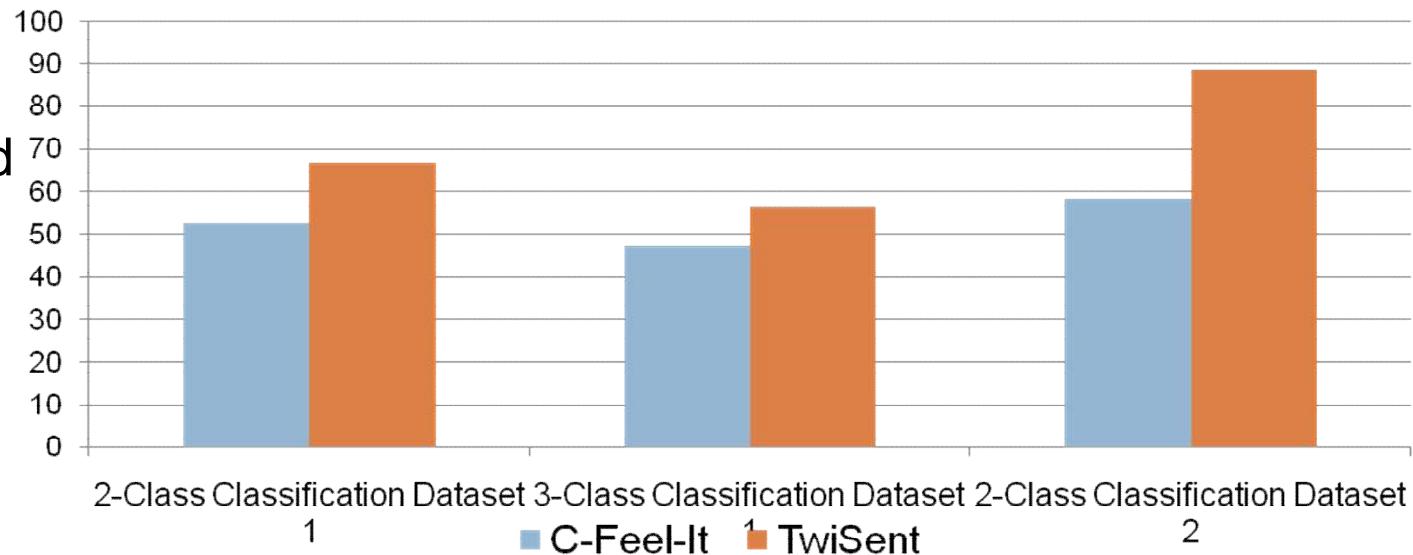


Supervised
Classification

TwiSent Evaluation

197

Lexicon-based
Classification



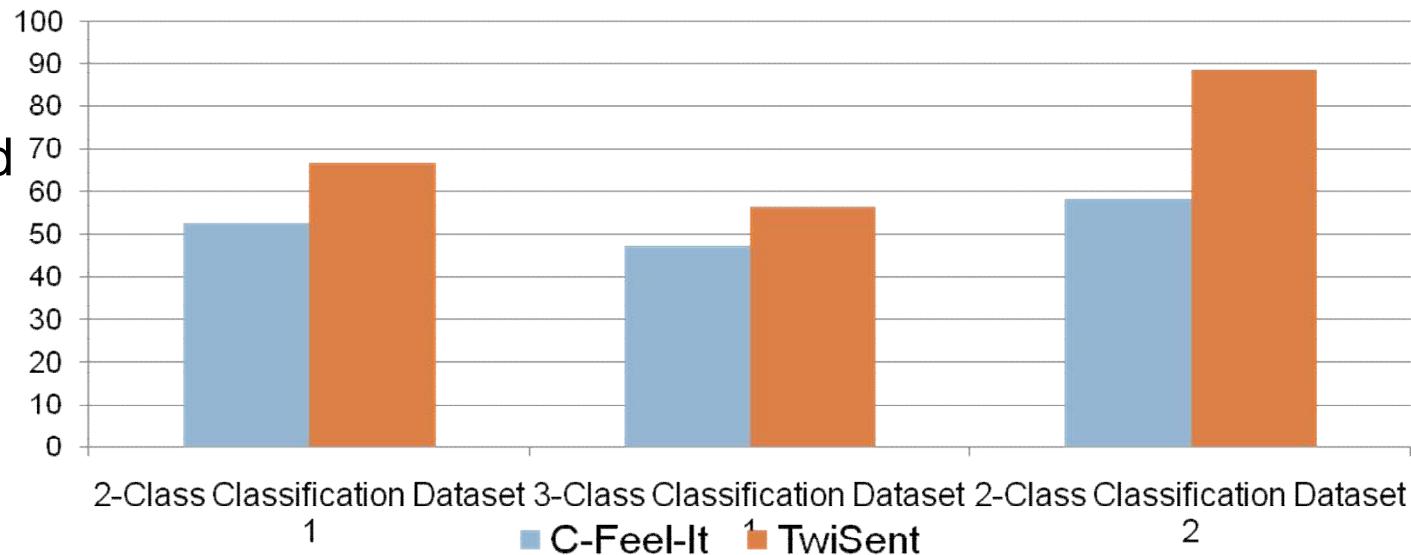
Supervised
Classification

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

TwiSent Evaluation

198

Lexicon-based
Classification



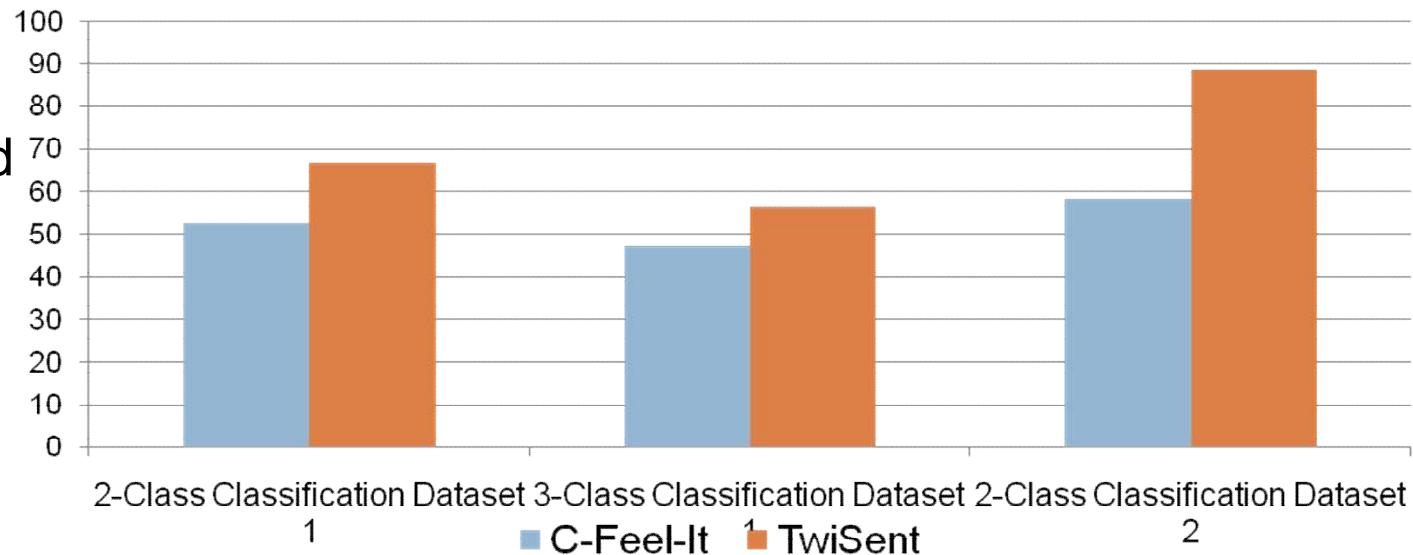
Supervised
Classification

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

TwiSent Evaluation

199

Lexicon-based
Classification



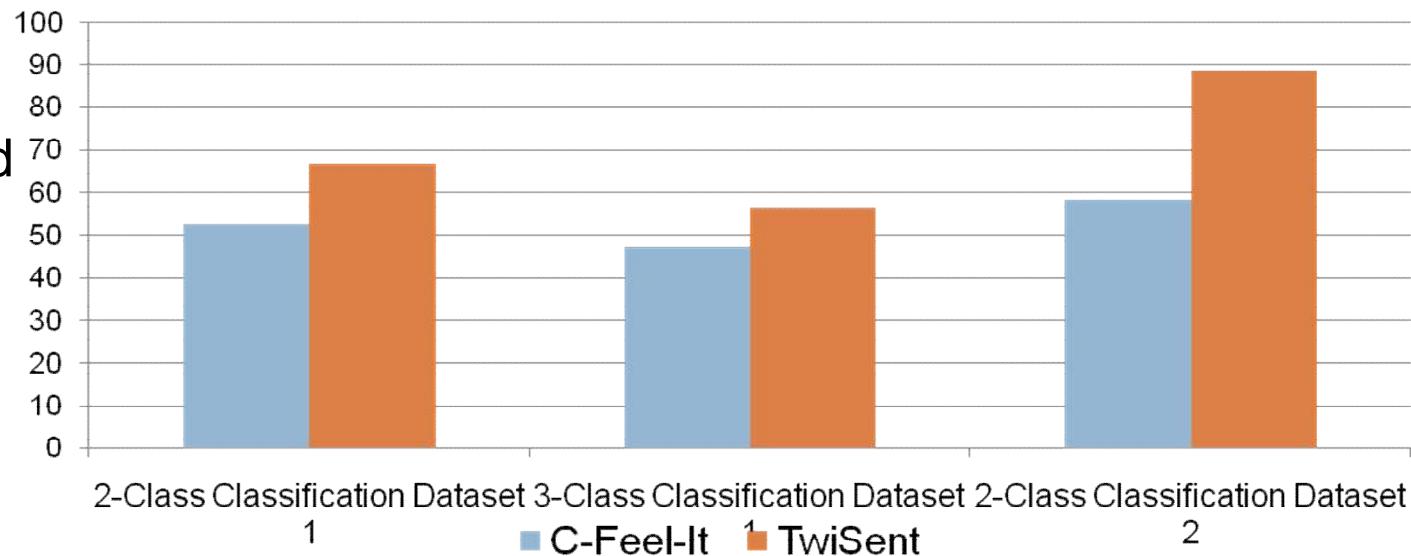
Ablation
Test 1

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

TwiSent Evaluation

200

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
Complete System	66.69	-

□ Role of Social Media Content and World Knowledge in Information Retrieval

- YouCat: Unsupervised System for Youtube Video Categorization from User Comments and Meta Data using WordNet and Wikipedia

YouCat: Youtube Video Categorization from User Comments and Meta-Data using WordNet and Wiki

202

YouCat: Youtube Video Categorization from User Comments and Meta-Data using WordNet and Wiki

203

- Given a Youtube video, predict its genres

YouCat: Youtube Video Categorization from User Comments and Meta-Data using WordNet and Wiki

204

- Given a Youtube video, predict its genres

For Example: A *Tom and Jerry Show* is tagged by the genres **Comedy** and **Animation**

YouCat: Youtube Video Categorization from User Comments and Meta-Data using WordNet and Wiki

205

- Given a Youtube video, predict its genres
- Extract features from User Comments and Meta-Description of the video

For Example: A *Tom and Jerry Show* is tagged by the genres **Comedy** and **Animation**

YouCat: Youtube Video Categorization from User Comments and Meta-Data using WordNet and Wiki

206

- Given a Youtube video, predict its genres
- Extract features from User Comments and Meta-Description of the video

Meta-Description of a video : It was an awesome slam dunk in the NBA finals by Michael Jordan.

Comments : He is the greatest basketball player of all times.

YouCat: Youtube Video Categorization from User Comments and Meta-Data using WordNet and Wiki

207

- Given a Youtube video, predict its genres
- Extract features from User Comments and Meta-Description of the video
- A genre is pre-defined based on a few seed words

Meta-Description of a video : It was an awesome slam dunk in the NBA finals by Michael Jordan.

Comments : He is the greatest basketball player of all times.

YouCat: Youtube Video Categorization from User Comments and Meta-Data using WordNet and Wiki

208

- Given a Youtube video, predict its genres
- Extract features from User Comments and Meta-Description of the video
- A genre is pre-defined based on a few seed words

Funny and *Laugh* are the keywords that define the **Comedy** genre

YouCat: Youtube Video Categorization from User Comments and Meta-Data using WordNet and Wiki

209

- Given a Youtube video, predict its genres
- Extract features from User Comments and Meta-Description of the video
- A genre is pre-defined based on a few seed words
- The words and synonymous concepts extracted using a Thesaurus are compiled into a Seed List

Funny and *Laugh* are the keywords that define the **Comedy** genre

YouCat: Youtube Video Categorization from User Comments and Meta-Data using WordNet and Wiki

210

- Given a Youtube video, predict its genres
- Extract features from User Comments and Meta-Description of the video
- A genre is pre-defined based on a few seed words
- The words and synonymous concepts extracted using a Thesaurus are compiled into a Seed List

Seed List for Comedy genre: funny, humor, hilarious, joke, comedy, roflmao, laugh, lol, rofl, roflmao, ...

YouCat: Youtube Video Categorization from User Comments and Meta-Data using WordNet and Wiki

211

- Given a Youtube video, predict its genres
- Extract features from User Comments and Meta-Description of the video
- A genre is pre-defined based on a few seed words
- The words and synonymous concepts extracted using a Thesaurus are compiled into a Seed List
- A Concept List is created using WordNet and Named Entities from Wikipedia based on related concepts from Seed List using an **overlap-based** classification

Seed List for Comedy genre: funny, humor, hilarious, joke, comedy, roflmao, laugh, lol, rofl, roflmao, ...

YouCat: Youtube Video Categorization from User Comments and Meta-Data using WordNet and Wiki

212

- Given a Youtube video, predict its genres
- Extract features from User Comments and Meta-Description of the video
- A genre is pre-defined based on a few seed words
- The words and synonymous concepts extracted using a Thesaurus are compiled into a Seed List
- A Concept List is created using WordNet and Named Entities from Wikipedia based on related concepts from Seed List using an **overlap-based** classification
 - *dunk* - {*dunk, dunk shot, stuff shot; dunk, dip, souse, plunge, douse; dunk; dunk, dip*} is classified to Sports with the help of **Sports Seed List**
 - Michael Jordon is also Classified to Sports based on the overlap of its Wikipedia definition with the **Sports Seed List**

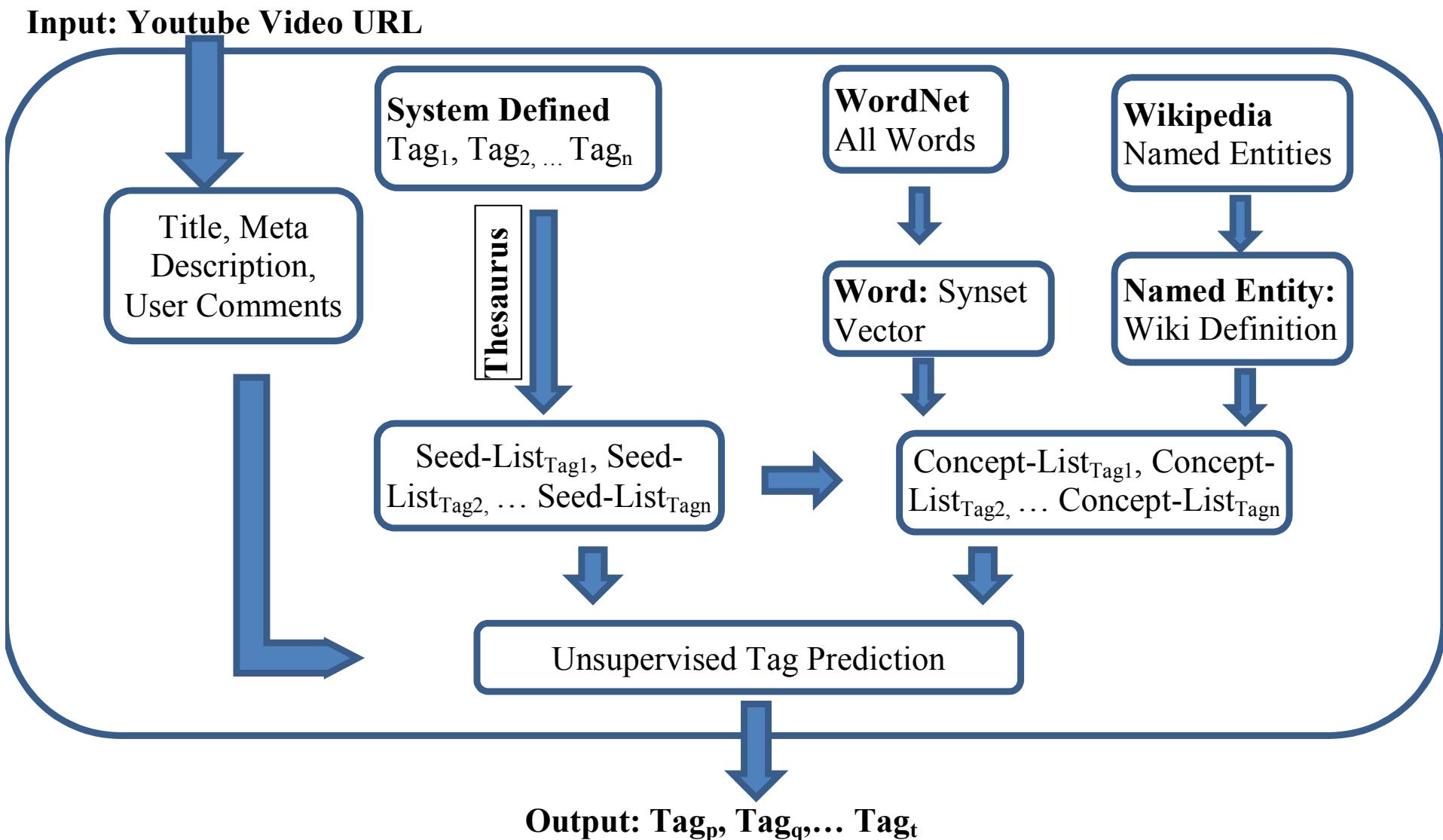
YouCat: Youtube Video Categorization from User Comments and Meta-Data using WordNet and Wiki

213

- Given a Youtube video, predict its genres
- Extract features from User Comments and Meta-Description of the video
- A genre is pre-defined based on a few seed words
- The words and synonymous concepts extracted using a Thesaurus are compiled into a Seed List
- A Concept List is created using WordNet and Named Entities from Wikipedia based on related concepts from Seed List using an **overlap-based** classification
- The final classification is based on an **overlap approach** based on the *Seed List* and *Concept List* for each genre – some parameters learnt from data
 - *dunk* - {*dunk*, *dunk shot*, *stuff shot*; *dunk*, *dip*, *souse*, *plunge*, *douse*; *dunk*; *dunk*, *dip*} is classified to Sports with the help of **Sports Seed List**
 - Michael Jordon is also Classified to Sports based on the overlap of its Wikipedia definition with the **Sports Seed List**

YouCat System Architecture

214



Tag Prediction

215

□ Feature Scoring

- $score(f \in genre_k; w_1, w_2) = w_1 \times \sum_{j:word_j \in seed_k} \mathbf{1} + w_2 \times \sum_{j:word_j \in concept_k} \mathbf{1}$

□ Video Genre Scoring

- $score(video \in genre_k; p_1, p_2, p_3) = p_1 \times score(f^{Title} \in genre_k) + p_2 \times score(f^{Meta Data} \in genre_k) + p_3 \times score(f^{Comments} \in genre_k)$

□ Single Genre Prediction

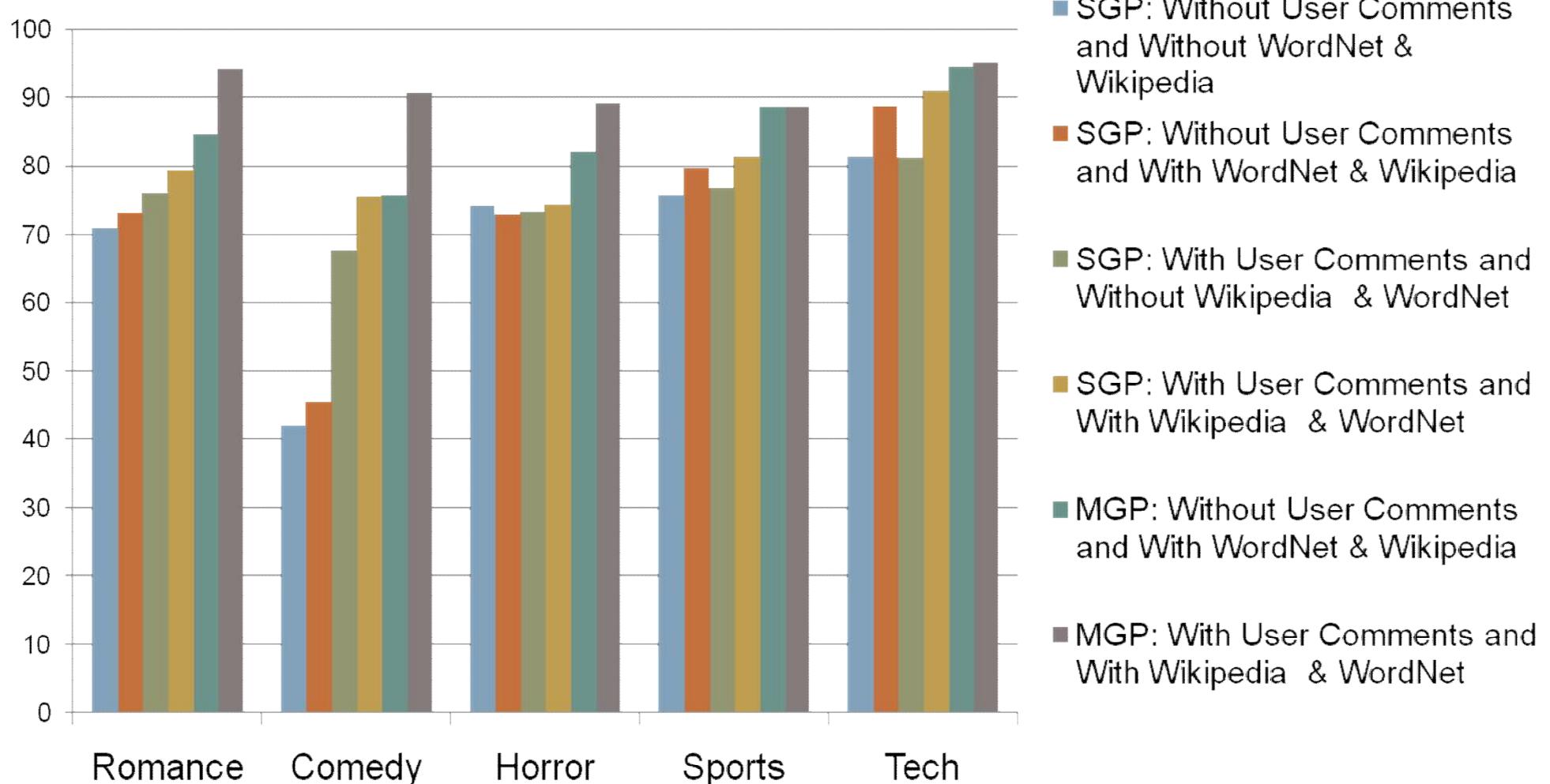
- $video_{genre} = argmax_k score(video \in genre_k)$

□ Multiple Genre Prediction

- $video_{genre} = k, if score(video \in genre_k) \geq \theta$
- where $\theta = \frac{1}{k} \sum_k score(video \in genre_k)$

Genre-wise Results

216



YouCat Dataset Statistics

217

Video
Statistics

	Comedy	Horror	Sports	Romance	Tech	Total
	2682	2802	2577	2477	2299	12837

Comment
Statistics

Romance	Comedy	Horror	Sports	Tech
103.64	168.94	83.4	56.91	223.45

Supervised Baseline

218

SVM Features	F ₁ -Score(%)
All Unigrams	82.5116
Unigrams+Without stop words	83.5131
Unigrams+ Without stop words +Lemmatization	83.8131
Unigrams+Without stop words +Lemmatization+ POS Tags	83.8213
Top Unigrams+Without stop words +Lemmatization+POS Tags	84.0524
All Bigrams	74.2681
Unigrams+Bigrams+Without stop words+Lemmatization	84.3606

Overall Results

219

Model	Average F ₁ Score
SVM Baseline: Single Genre Prediction With User Comments	84.3606
Single Genre Prediction : Without User Comments + Without Wikipedia & WordNet	68.76
Single Genre Prediction : With User Comments + Without Wikipedia & WordNet	74.95
Single Genre Prediction : Without User Comments + With Wikipedia & WordNet	71.984
<i>Single Genre Prediction : With User Comments+ With Wikipedia &WordNet</i>	80.9
Multiple Genre Prediction : Without User Comments + With Wikipedia & WordNet	84.952
Multiple Genre Prediction : With User Comments + With Wikipedia & WordNet	91.48

Comments/Genre and Confusion Matrix

220

Average Comments/Genre

	Average Tags/Video Without User Comments	Average Tags/Video With User Comments
Romance	1.45	1.55
Comedy	1.67	1.80
Horror	1.38	1.87
Sports	1.36	1.40
Tech	1.29	1.40
Average	1.43	1.60

Confusion Matrix

	Romance	Comedy	Horror	Sports	Tech
Romance	80.16	8.91	3.23	4.45	3.64
Comedy	3.13	77.08	3.47	9.03	7.29
Horror	10.03	9.34	75.78	3.46	1.38
Sports	0.70	7.30	0	89.05	2.92
Tech	0.72	5.07	0.36	1.81	92.03

Effect of Social Media Content in IR

221

Genre	Without User Comments + Without Wikipedia & WordNet			With User Comments + Without Wikipedia & WordNet		
	Precision	Recall	F ₁ -Score	Precision	Recall	F ₁ -Score
Romance	76.26	66.27	70.91	77.36	74.60	75.95
Comedy	43.96	40.00	41.89	69.23	66.00	67.58
Horror	80.47	68.67	74.10	76.45	70.33	73.26
Sports	84.21	68.71	75.67	85.07	69.94	76.77
Tech	90.83	73.50	81.25	84.09	78.45	81.17

Effect of World Knowledge in IR

222

Genre	Without User Comments + With Wikipedia & WordNet			With User Comments + With Wikipedia & WordNet		
	Precision	Recall	F ₁ -Score	Precision	Recall	F ₁ -Score
Romance	76.06	70.63	73.24	80.16	78.57	79.36
Comedy	47.31	44.00	45.6	77.08	74.00	75.51
Horror	75.63	70.33	72.88	75.78	73.00	74.36
Sports	87.5	73.01	79.60	89.05	74.85	81.33
Tech	92.34	85.16	88.60	92.03	89.75	90.88

Drawbacks

223

- User Comments
 - Noisy (slangs, abbreviations, informal language, spams etc.)
 - Off-topic conversation (abuses, chit-chats etc.)
- Concept Expansion (WordNet + Wikipedia)
 - Topic-drift
 - good - gloss of one of its synsets {dear, good, near -- with or in a close or intimate relationship}
- Bias towards Comedy
 - *Romantic-Comedy, Horror-Comedy*
- Ambiguity in Named-Entity due to shorter context
 - *Manchester Rocks !!!*

- Sentiment Influencing Semantic Similarity Measure

SenSim: Leveraging Sentiment to Compute Similarity

225

SenSim: Leveraging Sentiment to Compute Similarity

226

- Introduce Sentiment as another feature in the Semantic Similarity Measure
 - “*Among a set of similar word pairs, a pair is more similar if their sentiment content is the same*”
 - Is “*enchant*” (hold spellbound) more similar to “*endear*” (make endearing or lovable) than to “*delight*” (give pleasure to or be pleasing to) ?

SenSim: Leveraging Sentiment to Compute Similarity

227

- Introduce Sentiment as another feature in the Semantic Similarity Measure
 - “*Among a set of similar word pairs, a pair is more similar if their sentiment content is the same*”
 - Is “*enchant*” (hold spellbound) more similar to “*endear*” (make endearing or lovable) than to “*delight*” (give pleasure to or be pleasing to) ?
- Useful for replacing an unknown feature in test set with a similar feature in training set

SenSim: Leveraging Sentiment to Compute Similarity

228

- Introduce Sentiment as another feature in the Semantic Similarity Measure
 - “*Among a set of similar word pairs, a pair is more similar if their sentiment content is the same*”
 - Is “enchant” (hold spellbound) more similar to “endear” (make endearing or lovable) than to “delight” (give pleasure to or be pleasing to) ?
- Useful for replacing an unknown feature in test set with a similar feature in training set
- Given a **word** in a sentence, create its **Similarity Vector**
 - Use Word Sense Disambiguation on context to find its **Synset-id**
 - Create a **Gloss Vector (sparse)** using its gloss
 - Extend gloss using *relevant WordNet Relations*
 - **Learn the relations to use for different POS tags and the depth in WordNet hierarchy**
 - Incorporate SentiWordNet Scores in the Expanded Vector using Different Scoring

SenSim: Leveraging Sentiment to Compute Similarity

229

- Introduce Sentiment as another feature in the Semantic Similarity Measure
 - “*Among a set of similar word pairs, a pair is more similar if their sentiment content is the same*”
 - Is “enchant” (hold spellbound) more similar to “endear” (make endearing or lovable) than to “delight” (give pleasure to or be pleasing to) ?
- Useful for replacing an unknown feature in test set with a similar feature in training set
- Given a **word** in a sentence, create its **Similarity Vector**
 - Use Word Sense Disambiguation on context to find its **Synset-id**
 - Create a **Gloss Vector (sparse)** using its gloss
 - Extend gloss using *relevant WordNet Relations*
 - **Learn the relations to use for different POS tags and the depth in WordNet hierarchy**
 - Incorporate SentiWordNet Scores in the Expanded Vector using Different Scoring

Sentiment-Semantic Correlation

230

Annotation Strategy	Overall	NOUN	VERB	ADJECTIVES	ADVERBS
Meaning	0.768	0.803	0.750	0.527	0.759
Meaning + Sentiment	0.799	0.750	0.889	0.720	0.844

WordNet Relations used for Expansion

231

POS	WordNet relations used for expansion
Nouns	<i>hypernym, hyponym, nominalization</i>
Verbs	<i>nominalization, hypernym, hyponym</i>
Adjectives	<i>also see, nominalization, attribute</i>
Adverbs	<i>derived</i>

Scoring Formula

232

- $\text{Score}_{\text{SD}}(\text{A}) = \text{SWN}_{\text{pos}}(\text{A}) - \text{SWN}_{\text{neg}}(\text{A})$
- $\text{Score}_{\text{SM}}(\text{A}) = \max(\text{SWN}_{\text{pos}}(\text{A}), \text{SWN}_{\text{neg}}(\text{A}))$
- $\text{Score}_{\text{TM}}(\text{A}) = \text{sign}(\max(\text{SWN}_{\text{pos}}(\text{A}), \text{SWN}_{\text{neg}}(\text{A}))) * (1 + \text{abs}(\max(\text{SWN}_{\text{pos}}(\text{A}), \text{SWN}_{\text{neg}}(\text{A}))))$

$$\text{SenSim}_x(\text{A}, \text{B}) = \text{cosine} (\text{gloss}_{\text{vec}} (\text{sense}(\text{A})), \text{gloss}_{\text{vec}} (\text{sense}(\text{B})))$$

Where,

$\text{gloss}_{\text{vec}}$ $= 1 : \text{score}_x(1) \ 2 : \text{score}_x(2) \dots n : \text{score}_x(n)$

$\text{score}_x(Y)$ = Sentiment score of word Y using scoring function x

x = Scoring function of type SD/SM/TD/TM

Evaluation on Gold Standard Data: Word Pair Similarity

Evaluation on Gold Standard Data: Word Pair Similarity

234

- A set of 50 word pairs (with given context) manually marked
- Each word pair is given 3 scores in the form of ratings (1-5):
 - Similarity based on meaning
 - Similarity based on sentiment
 - Similarity based on meaning + sentiment

Evaluation on Gold Standard Data: Word Pair Similarity

235

- A set of 50 word pairs (with given context) manually marked
- Each word pair is given 3 scores in the form of ratings (1-5):
 - Similarity based on meaning
 - Similarity based on sentiment
 - Similarity based on meaning + sentiment

Agreement metric: Pearson correlation coefficient

Evaluation on Gold Standard Data: Word Pair Similarity

236

- A set of 50 word pairs (with given context) manually marked
- Each word pair is given 3 scores in the form of ratings (1-5):
 - Similarity based on meaning
 - Similarity based on sentiment
 - Similarity based on meaning + sentiment

Agreement metric: Pearson correlation coefficient

Metric Used	Overall	NOUN	VERB	ADJECTIVES	ADVERBS
LESK (Banerjee <i>et al.</i> , 2003)	0.22	0.51	-0.91	0.19	0.37
LIN (Lin, 1998)	0.27	0.24	0.00	NA	Na
LCH (Leacock <i>et al.</i> , 1998)	0.36	0.34	0.44	NA	NA
SenSim (SD)	0.46	0.73	0.55	0.08	0.76
SenSim (SM)	0.50	0.62	0.48	0.06	0.54
SenSim (TD)	0.45	0.73	0.55	0.08	0.59
SenSim (TM)	0.48	0.62	0.48	0.06	0.78

Evaluation on Travel Review Data: Feature Replacement

237

Metric Used	Accuracy (%)	PP	NP	PR	NR
Baseline	89.10	91.50	87.07	85.18	91.24
LESK (Banerjee <i>et al.</i> , 2003)	89.36	91.57	87.46	85.68	91.25
LIN (Lin, 1998)	89.27	91.24	87.61	85.85	90.90
LCH (Leacock <i>et al.</i> , 1998)	89.64	90.48	88.86	86.47	89.63
SenSim (SD)	89.95	91.39	88.65	87.11	90.93
SenSim (SM)	90.06	92.01	88.38	86.67	91.58
SenSim (TD)	90.11	91.68	88.69	86.97	91.23

Conclusions

238

Conclusions

239

- Sentiment Analysis is more than just a bag-of-words classification problem

Conclusions

240

- Sentiment Analysis is more than just a bag-of-words classification problem
- Discourse Markers, often ignored as stop words in a BOW model, improve its performance

Conclusions

241

- Sentiment Analysis is more than just a bag-of-words classification problem
- Discourse Markers, often ignored as stop words in a BOW model, improve its performance
- Dependency Relations capture association between words forming opinion expressions, which incorporate feature specificity

Conclusions

242

- Sentiment Analysis is more than just a bag-of-words classification problem
- Discourse Markers, often ignored as stop words in a BOW model, improve its performance
- Dependency Relations capture association between words forming opinion expressions, which incorporate feature specificity
- Informal language form and the social media content require different treatment than general SA to harvest their sentiment content

Conclusions

243

- Sentiment Analysis is more than just a bag-of-words classification problem
- Discourse Markers, often ignored as stop words in a BOW model, improve its performance
- Dependency Relations capture association between words forming opinion expressions, which incorporate feature specificity
- Informal language form and the social media content require different treatment than general SA to harvest their sentiment content
- Sentiment Analysis of Reviews require extensive world knowledge, which can be harnessed using Wikipedia

Conclusions

244

- Sentiment Analysis is more than just a bag-of-words classification problem
- Discourse Markers, often ignored as stop words in a BOW model, improve its performance
- Dependency Relations capture association between words forming opinion expressions, which incorporate feature specificity
- Informal language form and the social media content require different treatment than general SA to harvest their sentiment content
- Sentiment Analysis of Reviews require extensive world knowledge, which can be harnessed using Wikipedia
- Social Media Content and External Knowledge Sources provide a goldmine of information useful for Information Extraction

Conclusions

245

- Sentiment Analysis is more than just a bag-of-words classification problem
- Discourse Markers, often ignored as stop words in a BOW model, improve its performance
- Dependency Relations capture association between words forming opinion expressions, which incorporate feature specificity
- Informal language form and the social media content require different treatment than general SA to harvest their sentiment content
- Sentiment Analysis of Reviews require extensive world knowledge, which can be harnessed using Wikipedia
- Social Media Content and External Knowledge Sources provide a goldmine of information useful for Information Extraction
- Incorporation of these new perspectives in a traditional BOW Model can improve its performance, while retaining its simplicity

Contributions

246

Contributions

247

- Discourse markers can be useful in BOW Model

Contributions

248

- Discourse markers can be useful in BOW Model
- Feature Specificity can be captured by Dependency Relations

Contributions

249

- Discourse markers can be useful in BOW Model
- Feature Specificity can be captured by Dependency Relations
- WikiSent: Incorporating World Knowledge through Wikipedia improves SA

Contributions

250

- Discourse markers can be useful in BOW Model
- Feature Specificity can be captured by Dependency Relations
- WikiSent: Incorporating World Knowledge through Wikipedia improves SA
- We analyze the movie trend from genre, year of release and polarity information using WikiSent

Contributions

251

- Discourse markers can be useful in BOW Model
- Feature Specificity can be captured by Dependency Relations
- WikiSent: Incorporating World Knowledge through Wikipedia improves SA
- We analyze the movie trend from genre, year of release and polarity information using WikiSent
- TwiSent: Categorization of Spams in Twitter w.r.t SA

Contributions

252

- Discourse markers can be useful in BOW Model
- Feature Specificity can be captured by Dependency Relations
- WikiSent: Incorporating World Knowledge through Wikipedia improves SA
- We analyze the movie trend from genre, year of release and polarity information using WikiSent
- TwiSent: Categorization of Spams in Twitter w.r.t SA
- TwiSent: Categorization of Noisy Text in Twitter

Contributions

253

- Discourse markers can be useful in BOW Model
- Feature Specificity can be captured by Dependency Relations
- WikiSent: Incorporating World Knowledge through Wikipedia improves SA
- We analyze the movie trend from genre, year of release and polarity information using WikiSent
- TwiSent: Categorization of Spams in Twitter w.r.t SA
- TwiSent: Categorization of Noisy Text in Twitter
- TwiSent: Handling Pragmatics and introducing Feature Specificity

Contributions

254

- Discourse markers can be useful in BOW Model
- Feature Specificity can be captured by Dependency Relations
- WikiSent: Incorporating World Knowledge through Wikipedia improves SA
- We analyze the movie trend from genre, year of release and polarity information using WikiSent
- TwiSent: Categorization of Spams in Twitter w.r.t SA
- TwiSent: Categorization of Noisy Text in Twitter
- TwiSent: Handling Pragmatics and introducing Feature Specificity
- YouCat: Analyzing the combined influence of social media content and concept expansion in Information Retrieval

Contributions

255

- Discourse markers can be useful in BOW Model
- Feature Specificity can be captured by Dependency Relations
- WikiSent: Incorporating World Knowledge through Wikipedia improves SA
- We analyze the movie trend from genre, year of release and polarity information using WikiSent
- TwiSent: Categorization of Spams in Twitter w.r.t SA
- TwiSent: Categorization of Noisy Text in Twitter
- TwiSent: Handling Pragmatics and introducing Feature Specificity
- YouCat: Analyzing the combined influence of social media content and concept expansion in Information Retrieval
- SenSim: Incorporation of sentiment as a feature in semantic similarity improves similarity performance

Future Works

256

Future Works

257

- Incorporation of Parsing to determine scope of discourse markers for structured review analysis

Future Works

258

- Incorporation of Parsing to determine scope of discourse markers for structured review analysis
- Ranking of discourse features with positional importance

Future Works

259

- Incorporation of Parsing to determine scope of discourse markers for structured review analysis
- Ranking of discourse features with positional importance
- Generating a domain-dependent lexicon for better classification of domain-dependent concepts

Future Works

260

- Incorporation of Parsing to determine scope of discourse markers for structured review analysis
- Ranking of discourse features with positional importance
- Generating a domain-dependent lexicon for better classification of domain-dependent concepts
- Improving Spell-Checker using SMT on a parallel corpora

Future Works

261

- Incorporation of Parsing to determine scope of discourse markers for structured review analysis
- Ranking of discourse features with positional importance
- Generating a domain-dependent lexicon for better classification of domain-dependent concepts
- Improving Spell-Checker using SMT on a parallel corpora
- Pragmatics handling needs to be improved

Future Works

262

- Incorporation of Parsing to determine scope of discourse markers for structured review analysis
- Ranking of discourse features with positional importance
- Generating a domain-dependent lexicon for better classification of domain-dependent concepts
- Improving Spell-Checker using SMT on a parallel corpora
- Pragmatics handling needs to be improved
- WikiSent needs handling of *synonymous concepts, anaphora resolution and word sense disambiguation*

Future Works

263

- Incorporation of Parsing to determine scope of discourse markers for structured review analysis
- Ranking of discourse features with positional importance
- Generating a domain-dependent lexicon for better classification of domain-dependent concepts
- Improving Spell-Checker using SMT on a parallel corpora
- Pragmatics handling needs to be improved
- WikiSent needs handling of *synonymous concepts, anaphora resolution and word sense disambiguation*
- YouCat needs to be tested on more number of genres for fine distinction, requires a Comment Filter and a WSD module

Future Works

264

- Incorporation of Parsing to determine scope of discourse markers for structured review analysis
- Ranking of discourse features with positional importance
- Generating a domain-dependent lexicon for better classification of domain-dependent concepts
- Improving Spell-Checker using SMT on a parallel corpora
- Pragmatics handling needs to be improved
- WikiSent needs handling of *synonymous concepts, anaphora resolution and word sense disambiguation*
- YouCat needs to be tested on more number of genres for fine distinction, requires a Comment Filter and a WSD module

All these perspectives need to be brought in a single canvas and their collective effectiveness needs to be evaluated

Publications (1/2)

265

- Sentiment Analysis in Twitter with Lightweight Discourse Analysis, Subhabrata Mukherjee and Pushpak Bhattacharyya, In Proceedings of the 24th International Conference on Computational Linguistics (**COLING 2012**), IIT Bombay, Mumbai, Dec 8 - Dec 15, 2012
- YouCat : Weakly Supervised Youtube Video Categorization System from Meta Data & User Comments using WordNet & Wikipedia, Subhabrata Mukherjee and Pushpak Bhattacharyya, In Proceedings of the 24th International Conference on Computational Linguistics (**COLING 2012**), IIT Bombay, Mumbai, Dec 8 - Dec 15, 2012 (Long Paper)
- TwiSent: A Multi-Stage System for Analyzing Sentiment in Twitter, Subhabrata Mukherjee, Akshat Malu, Balamurali A.R. and Pushpak Bhattacharyya, In Proceedings of The 21st ACM Conference on Information and Knowledge Management (**CIKM 2012**), Hawai, Oct 29 - Nov 2, 2012

Publications (2/2)

266

- [WikiSent : Weakly Supervised Sentiment Analysis Through Extractive Summarization With Wikipedia](#), Subhabrata Mukherjee and Pushpak Bhattacharyya, In Proceedings of the European Conference on Machine Learning (**ECML PKDD 2012**), Bristol, U.K., 24-28 Sept, 2012
- [Feature Specific Sentiment Analysis for Product Reviews](#), Subhabrata Mukherjee and Pushpak Bhattacharyya, In Proceedings of the 13th International Conference on Intelligent Text Processing and Computational Intelligence (**CICLING 2012**), New Delhi, India, March, 2012
- [Leveraging Sentiment to Compute Word Similarity](#), Balamurali A.R., Subhabrata Mukherjee, Akshat Malu and Pushpak Bhattacharyya, In Proceedings of the 6th International Global Wordnet Conference (**GWC 2011**), Matsue, Japan, Jan, 2012

□ Thank You