

Opinion Mining

Feiyu XU & Xiwen CHENG

Xiwen.cheng@dfki.de

DFKI, Saarbruecken, Germany

Jan 19th, 2011



Discussion on Opinion Mining Application



Tweetmotif: Topic summarization on Twitter- e.g. wikileaks, parenting

The screenshot shows the Tweetmotif website interface. The browser address bar displays <http://tweetmotif.com/#parenting>. The page title is "tweetmotif - summarize and discover topics on twitter". The main search bar contains the text "parenting".

tweetmotif
discover twitter topics

(tweet this) (about) (tips)

What are people saying about...

Trending topics
#ihatefemaleswho • #hannahmontanaforever • #wecanallagreethat • Soldier Field • RIP JUSTIN BIEBER • Chester Taylor • Touchdown Bears • PobreTuAlmaVideo • Seriously Drake • Classic Text Message

Or try
sandwich • coffee • :) • (: • aw • awwwww • @the_real_shaq • @twitter • "san francisco" weather • tweetmotif

related themes
over the last 46 minutes

case of bad	@parenting
about parenting	kourtney
your parenting	kardashian
-	parenting classes
companies put into	#ece
parenting style	amy chua
parenting styles	such parenting
@deliacabe	that parenting
#family	@elyssaeast
#baby	#toddler / #parent
#book	#moms
#blog / #funny	parents make
#comedy	#mom
#autism	#humor
your child	discipline
january	your kids
parent	raising
	ve

tweets by theme

"case of bad"

A Case of Bad Parenting - The Hollywood Gossip <http://bit.ly/gjkeg2>
gossipquen

A Case of Bad Parenting - The Hollywood Gossip: Kourtney Kardashian covers the latest issue of Parenting. In rel... <http://bit.ly/ihYhT> sarahkensy

A Case of Bad Parenting <http://sockroll.com/cpf8gk> sockrollshowbiz

>>>: A Case of Bad Parenting <http://sloppygossip.com/the-hollywood-gossip/a-case-of-bad-parenting/> sloppygossip

A Case of Bad Parenting | Star Talk: Filed under: Lindsay Lohan, Dawn Holland, Celebrity Justice The lawyer for... <http://bit.ly/hd5i4O> lindsayfashion

A Case of Bad Parenting: Kourtney Kardashian has appeared on a number of ridiculous tabloid covers. Incredibly, ... <http://bit.ly/dZH90C> rockmusicfans
show 1 similar tweet -

A Case of Bad Parenting: Kourtney Kardashian Bio Gossip Pictures Videos
A Case of Bad Parenting January 16th, 20... <http://bit.ly/etmsml> CynthiaCarlson
show 1 similar tweet -

(drilldown ?)

"@parenting"

What the trend: Trend monitoring - e.g. wikileaks

The screenshot shows a web browser window with the URL <http://whatthetrend.com/trend/parenting>. The page title is "parenting". The website header features the "What the Trend" logo and a search bar. Below the header, there are navigation links: "Go Pro", "Top 10", "Leaderboard", "FAQ", "Add a Trend", "Week in Review/Reports", and "View Cart".

The main content area is titled "parenting > Details". The primary heading is "Why is parenting popular right now?". Below this, a yellow box contains the text: "Do you know why parenting is trending on Twitter? Click here, and in your own words, write a short definition and help your fellow Twitterers." Below the text are three buttons: "Define", "Retweet", and "Flag as Inappropriate".

On the right side, there is a section titled "Want your definition to appear all over the Twitterverse?". It includes a paragraph: "Sure you do! But your definition has to be well-written. Here are some guidelines." followed by a list of guidelines:

- Be informative, clear and concise. < 140 characters is best.
- Do not put your Twitter handle in the definition or originating user field. It will be edited out. [Log in](#) and you will be appropriately credited.
- When discussing a person, add background information. Are they a singer or an actor? If they passed away, why were they famous?
- Use dates in addition to "today and tomorrow", i.e. "released today, Aug. 23, 2009"
- Don't cuss. Go light on snark.
- No HTML tags allowed.
- URLs like <http://www.google.com/> will be auto-linked and no-

Below the main content, there are two columns: "Related tweets" and "Related news".

Related tweets:

- [elyssaeast](#) @Lynne_Griffin I don't know research just know ppl have dif ideas re what makes loving/well intnd parenting. Chua thinks she's those things
- [GSMTitan](#) Really do recommend that link, all parenting skills highlighted there can be easily applied to life in general.
- [scottmcrocker](#) RT @JimGaffigan I don't know what's more exhausting about parenting the getting up early or the acting like you know what you're doing. #fb
- [Creme_LincolnPk](#) RT @DailyParentTip: 8 Discipline Mistakes Parents Make ->

Related news:

- [How This Chinese Mother Defines Parenting Success \(CurrentMom\)](#)
So Amy Chua is backing away from the notion that she was promoting her extreme version of **parenting** as a model, and even denied the claims made on the first ...
- [Really, Parenting?!? \(The Hollywood Gossip\)](#)
Incredibly, though, **Parenting** has managed to top them all. On its official website, this magazine claims to offer "mom tips you can't live without," only to ...
- [Major's Corner: Slipshod parenting endangers a whole generation \(Victoria Times Colonist\)](#)
By Maj. (Retired) Nigel Smythe-Brown, Times Colonist

A blue box overlay on the right side of the page reads: "You must be logged in to do that. [Log In Here](#) (click to close)".

Opinion gathering speed on Internet

- WSJ publishes an article “why chinese mother are superior” written by Amy Chua on 8th, Jan, 2011. Until 18th, Jan
- 6,800 comments on WSJ;

Keyword: Amy Chua

- 3,490,000 on Google
- 5,600 on twitter.com
- 5,289 on wordpress.com

Keyword: parenting

- 83,200,000 search results on Google;
- 1,620,000 from twitter.com;
- 502,000 from wordpress.com

A question from [Quora](#)

What algorithms can be used to cluster opinions about a topic expressed in natural language? [✎ Edit](#)

This is a follow-up question to [Where can I learn more about computer science methods and algorithms that match and summarize different pieces of text?](#).

Problem Example: Lets say 20 people wrote their qualitative opinions about a topic. What's the best algorithm out there (open source) that would be able to cluster opinions based on similarity. Not summarize - just CLUSTER based on similarity. Because 20 people may be talking about 5 mutually exclusive points from different angles of personal views. The algorithm should thus give me 5 clusters of opinions sets (unaltered, just grouped) to make it easy for me to extract 5 points from 20 people's opinions. Where can i get this algorithm - must be out there
[?](#) [✎ Edit](#)



Proposals of Opinion Mining Application and Solution?



Discussion on Resource for Movie Review Summarization



Reviews on "Das Leben der Anderen" @ [imdb](#)

505 out of 629 people found the following review useful:

The most underrated film of 2006., 17 January 2007

★★★★★★★★

Author: [jesse3](#) from los angeles california

Holy cow! What a terrific movie! I am a voting member of the Academy (actor's branch) so I get all the films for free. I've seen everything---60 films. This was one of the last 3 films that I saw---because I was completely unfamiliar with the title. This film slowly gripped me, but by the end, the grip was merciless. The lead actor, who should be doing The Life Story Of Peter Jennings, was wonderful. Everybody was terrific. Congratulations to the writers for their perfect structure---and to the director for his flawless storytelling---and his eliciting of top performances from his actors. How well cast it was.

But now I'm totally bewildered. Why haven't I heard anything about this film? Where was this film at the Golden Globes? I haven't even seen any reviews about it. Nothing! What's going on? I'm very active in the film business. I follow this stuff. This film (that I never heard of) took me by surprise as no other film has ever done.

Note to the IMDb: This is not a spoiler.

Jesse Vint III

Top 250 movies voted by imdb users

IMDb Top 250

http://www.imdb.com/chart/top?tt0405094

Sentiment analysis - Wikipedia, t... Das Leben der Anderen (2006) - I... IMDb Top 250

IMDb Charts

[Main index](#)
[IMDb Top 250](#)
[IMDb Bottom 100](#)

US Box Office

[USA Top 10](#)
[USA Archive](#)

UK Box Office

[UK Top 10](#)
[UK Archive](#)

All-Time Box Office

[USA](#)
[Non-USA](#)
[World-wide](#)

Video Rentals

[USA Weekly Top 20](#)
[USA Archive](#)

Votes by Genre

[Action](#)
[Adventure](#)
[Animation](#)
[Biography](#)
[Comedy](#)
[Crime](#)
[Documentary](#)
[Drama](#)
[Family](#)
[Fantasy](#)

Top 250 movies as voted by our users

For this top 250, only votes from regular voters are considered.

Track which films you've seen from the Top 250 [right here!](#)

Rank	Rating	Title	Votes
1.	9.2	Die Verurteilten (1994)	553,235
2.	9.2	Der Pate (1972)	432,693
3.	9.0	Der Pate 2 (1974)	261,177
4.	8.9	Zwei glorreiche Halunken (1966)	172,844
5.	8.9	Pulp Fiction (1994)	441,737
6.	8.9	Schindlers Liste (1993)	292,891
7.	8.9	Inception (2010)	282,467
8.	8.9	Die zwölf Geschworenen (1957)	129,212
9.	8.8	Einer flog über das Kuckucksnest (1975)	228,512
10.	8.8	The Dark Knight (2008)	493,906
11.	8.8	Das Imperium schlägt zurück (1980)	292,583
12.	8.8	Der Herr der Ringe - Die Rückkehr des Königs (2003)	385,169
13.	8.8	Die sieben Samurai (1954)	102,136
14.	8.7	Krieg der Sterne - Episode IV: Eine neue Hoffnung (1977)	335,721
15.	8.7	Fight Club (1999)	406,991
16.	8.7	GoodFellas - Drei Jahrzehnte in der Mafia (1990)	243,014
17.	8.7	Casablanca (1942)	175,641
18.	8.7	Cidade de Deus (2002)	178,358
19.	8.7	Der Herr der Ringe - Die Gefährten (2001)	408,200
20.	8.7	Spiel mir das Lied vom Tod (1968)	79,564
21.	8.7	Das Fenster zum Hof (1954)	125,315
22.	8.7	Jäger des verlorenen Schatzes (1981)	254,620

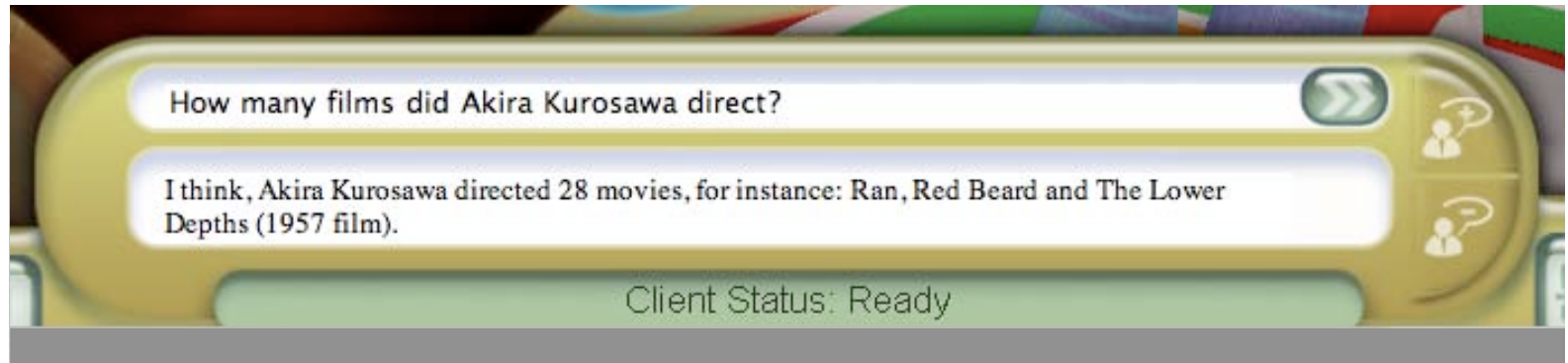
What resource and which features you would like to choose for OM tasks?



Experiment on KomParse
Making NPCs express their opinions emotionally



Gossip Galore in Rascalli

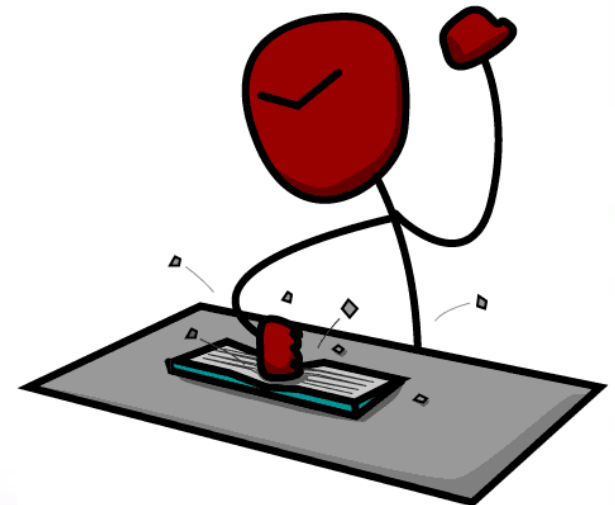


Hank in KomParse



Paul's solution

- Unsupervised machine learning
- Data: comments ranked by reviewers (1 ~ 10 stars)
- Features
 - N-Gram Token Patterns
 - Dependency Patterns
- Extra knowledge
 - WordNet
 - Negation expressions
- Learning algorithm
 - Scoring system



Data Processing

- Resource
 - IMDb (<http://www.imdb.com/>), A movie online storehouse
- Interested in IMDB pages:
 - with name (actors, authors, directors etc.)
 - with title (movie title, movie recommendations from IMDb)
- Containing the information:
 - Movie title
 - Review
 - Review title
 - Review date
 - Author name
 - Author origin (optional)
 - Recommendation of other users to this review (optional)
 - The score the author gave the reviewed movie x/10 (optional)

Data Processing

```
<Record name="Paycheck (2003)" isA="Movie"
  type="IMDb user reviews">
  <Feature name="Recommend">0 out of 3</Feature>
  <Feature name="Time">25 December 2003</Feature>
  <Feature name="Author">ak2k</Feature>
  <Feature name="Review">A poor remake of Minority Report,
    with less talented actors. Promising plot line that wilted away in
    the first thirty minutes of the film. Interesting inductive journey
    and neat car chases, but nowhere close to my money's worth.
    I'd recommend to go and see LOR again.</Feature>
  <Feature name="Score">1/10</Feature>
  <Feature name="From">Illinois</Feature>
  <Feature name="Title">A perfect Christmas movie has about as
    much connection with reality as Santa Clause does.</Feature>
</Record>
```

Data Processing

Presumptions and observations:

- Score indicates the sentiment of the review
- Short reviews are preferred over long reviews
 - long reviews have a lot of objective parts about storyline, anecdotes etc.
 - short reviews containing only the opinion over the movie and often expressed sentimental
- The sentiment classification on extreme reviews (very high or very low rating) are mostly unambiguous and clear while mid rated reviews have a lot of unclear sentences, such as one the one hand ...on the other

Data Processing

- Filtering the review
 - The number of tokens > 900
 - with a rating 4, 5, 6, 7 or 8 out of 10
- SCORE assignment to each sentence in the selected reviews
 - SCORE = Rank (1 ~ 10 start)
 - SCORE + 1, if the sentence :
 - Is the first, second or last sentence
 - And contains the keywords, such as I, me, movie, film and this movie.
 - SCORE - 1, if the sentence :
 - Has the length > 100
 - And contains the keywords, such as imdb, you, your, spoiler and review etc.
- The sentence with the highest SCORE from a review are selected.

Features – N-gram token pattern

Extracting uni-, bi- and trigrams out of every sentence from the sentimental corpus

- For example: I absolutely loved this movie.
- Unigrams:
 - i (NP), absolutely (RB), loved (VVD)
- Bigrams:
 - i absolutely (NP RB), absolutely loved (RB VVD)
- Trigrams:
 - i absolutely loved (NP RB VVD), absolutely loved this (RB VVD DT)

Features – Dependency Pattern

This is a funny super interesting and exciting movie.

Some important information *is missed* in N-gram tokens pattern.

- *funny* and *movie* are not caught by a n-gram ($n < 6$)

So, we include depends patterns:

- `amod(movie-9, funny-4)`

Tool: [Stanford-Dependency Parser](#)

Typed dependencies

```
nsubj(movie-9, This-1)
cop(movie-9, is-2)
det(movie-9, a-3)
amod(movie-9, funny-4)
advmod(interesting-6, super-5)
amod(movie-9, interesting-6)
cc(interesting-6, and-7)
conj(interesting-6, exciting-8)
```

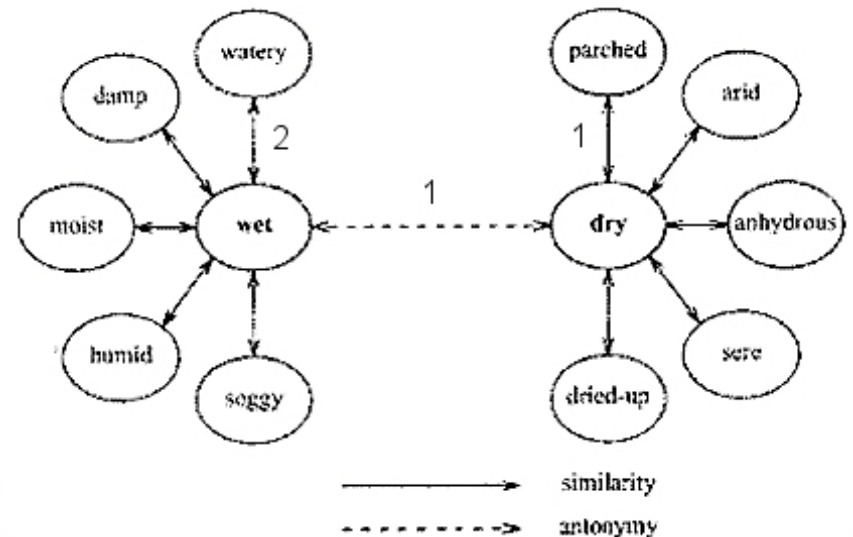
Typed dependencies, collapsed

```
nsubj(movie-9, This-1)
cop(movie-9, is-2)
det(movie-9, a-3)
amod(movie-9, funny-4)
advmod(interesting-6, super-5)
amod(movie-9, interesting-6)
conj_and(interesting-6, exciting-8)
amod(movie-9, exciting-8)
```

Extra Knowledge - Unigram patterns extended with WordNet

- All 1-gram adjective and adverb patterns will be extended with WordNet. Both the synonyms and the antonyms are used.
- For instance, 1-gram pattern "dry" can be extended with
 - Parched / arid / anhydrous / sere / dried-up
 - Wet / watery / damp / moist / humid / soggy

- In our experiment, the antonyms/synonyms are the words which connect the original word with a maximum distance of two.



Extra Knowledge – Negations

- Some elements in a sentence can change the sentiment of a word or phrase, such as
 - Subjunctive: I thought this movie is good.
 - Tempus: This movie was good.
 - Negation: This film is not funny.
 - Quotation: My friend told me “this is the best movie ever, you have to watch it” but I didn’t liked it.
- In our work, the content in the quotation is removed
- we care only negations such as not, no, never and n’t, including
 - no wonder, not just, not to mention etc.
 - Restricted comparative sentences “not better as” “no more” etc.

Algorithm – Score of patterns

- Each pattern has an iSCORE, including two sub-values
 - $iSCORE_{pos}$: the value of being positive
 - $iSCORE_{neg}$: the value of being negative
- The iSCORE is initialized with the frequency of this pattern from the corpus

Algorithm – Data bias

- Although “more” negative scored sentences are used, i.e. (1/10, 2/10, 3/10) vs.(9/10, 10/10), positive reviews are still twice the native ones.
- Assuming 1) there are X negative sentences and Y positive ones or on the other way round, and 2) $Y > X$

$$\text{equalizer} = Y / X$$

$$\text{BIAS} = \text{equalizer} / (X + Y + Y - X)$$

$$\text{iSCORE}_Y = \text{iSCORE}_Y / 2Y - \text{BIAS}$$

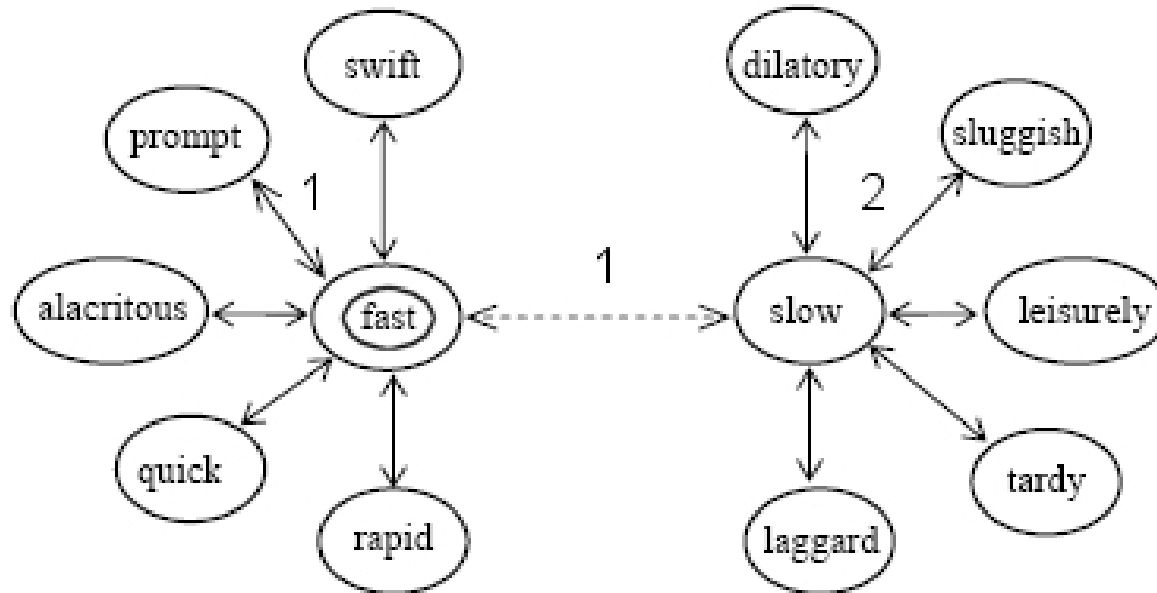
Algorithm – iSCORE

- $iSCORE = iSCORE_{pos} - iSCORE_{neg}$
 - If the value of the iSCORE is positive the computed polarity of the pattern is positive and if the value is negative the polarity is negative

- $iSCORE = iSCORE * 2$, if the pattern is binary
- $iSCORE = iSCORE * 3$, if the pattern is triple
- $iSCORE = iSCORE * 2.5$, if the pattern is a dependency pattern

Algorithm – iSCORE extended by WordNet

- The synonyms have the same polarity as the word, while the antonyms have a reversed polarity.



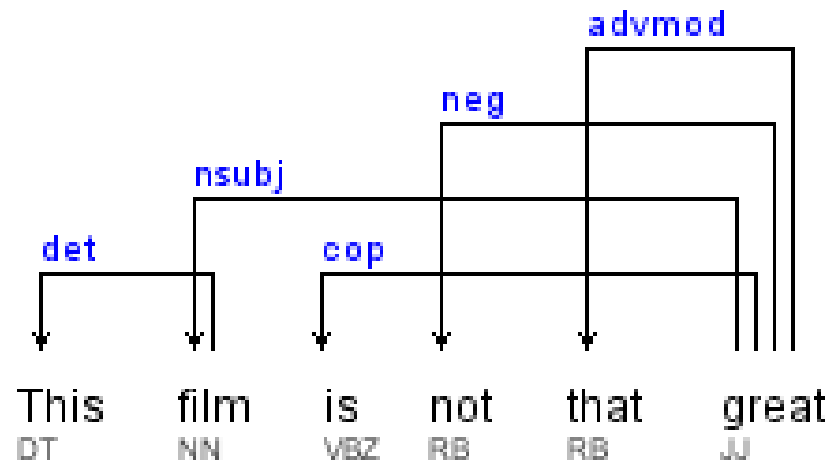
**Figure 4: Bipolar adjective structure,
(\rightarrow = similarity; \dashrightarrow = antonymy)**

Algorithm – iSCORE extended by WordNet

- *For instance, if Polarity(fast JJ) = positive, for words with the WordNet depth = 1*
 - [iSCORE(swift), iSCORE(prompt) , ...] += 0.3
 - [iSCORE(slow)] += - 0.3
 - *for the words with a WordNet depth = 2*
 - [iSCORE(swift) , iSCORE(prompt), ...] += 0.3 * 2 = 0.6
 - [iSCORE(slow)] += - 0.3 * 2 = - 0.6
 - [iSCORE(sluggish)] += - 0.3
 - iSCORE(synonyms at the x^{nd} depth) += 0.3 * ((max. depth + 1) – x)
 - iSCORE(antonyms at the x^{nd} depth) += -0.3 * ((max. depth + 1) – x)
- *# 0.3 is an arbitrarily chosen value*

Algorithm – iSCORE with negations

- *This film is not that great.*
 - iSCORE(great) += 1.0
 - iSCORE(not that great) += -3.0
 - iSCORE(cop(great-6, is-3)) += 2.5
 - iSCORE(neg(great-6, not-4)) += -2.5



Experiment 1

- Data: 8,038 positive and 3,016 negative sentences
- Features
 - N-grams ($0 < n < 4$) containing adjectives or adverbs
 - WordNet depth of two, initial value = 0,3
 - Negation
- Incomplete Result
 - 75,60% ~ 95,12%

Random samples	Right positives	Right negatives	False positives	False negatives	Unknowns
41	16	15	0	2	8

Experiment 1

- Unknowns:
 - Analysis between a statistical value of -5 and 5
 - Sentences with unknown polarity:
 - on the one hand...on the other
 - they have done it again.
 - i read somewhere that it is 'the next final destination'.
- With *unknowns* as correct an accuracy of 95,12% 😊
- With *unknowns* as an error an accuracy of 75,60%
- Without unknowns an accuracy of 93.94%

Experiment 2

- Data: 4,395 positive and 1,733 negative sentences
- Features
 - N-grams ($0 < n < 4$) containing adjectives or adverbs
 - WordNet depth of two, initial value = 0,3
 - Negation
 - Dep-patterns
- Incomplete Result
 - 78,23% ~ 94,55%

Random samples	True positives	True negatives	False positives	False negatives	Unknowns
147	80	35	0	12	24

Experiment 2

- False negatives:
 - 90 minutes long and only one unnecessary scene.
 - i never thought pacino would top "godfather" but boy was i wrong.
 - it was very educational and informative.
 - one of the better ww2 escape movies.
 - it's an all-time fave and i'm happy to hear that it's out on video!
 - cop(happy-9, 'm-8) negative
 - i don't think there is one bad thing i could say about it.
 - too long negation scope
- Observation: most objective sentences are rated negative
- With *unknowns* as correct an accuracy of 94,55% 😊
- With *unknowns* as an error an accuracy of 78,23%
- Without *unknowns* an accuracy of 93,49%

Experiment 3

- Data: 4,890 positive and 1,864 negative sentences
 - Different data set
- Features – Same as Experiment 2
- Incomplete Result
 - 78,07% ~ 95,61%

Random samples	True positives	True negatives	False positives	False negatives	Unknowns
114	51	38	1	4	20

Experiment 3

- False positives:
 - john singleton has not, and probably never will make a good film.
 - the scenario is pretty interesting. (*interesting* has a negative iScore)
- Unknowns:
 - This is a good movie but something is missing.
 - i love kate / evangeline, but i hope mister eko will be eaten by lostzilla!
 - don't think the movie is over just because the credits are rolling!
 - i have no objection to long movies, just short movies that seem long.
 - he uses violence against people who are mean to him. Objective
- With *unknowns* as correct an accuracy of 95,61% 😊
- With *unknowns* as an error an accuracy of 78,07%
- Without *unknowns* an accuracy of 94,68%

Comments, questions, suggestions &
Opinions?



köszönöm !תודה dĕkuji
mahalo 고맙습니다
thank you
merci 谢谢 danke
Ευχαριστώ شكرا
どうもありがとう gracias