

Affective computing



Contents

- understanding emotions
- affective lexicons
- sentiment classification
- language specifics and cross-lingual approaches

- Literature: Jurafsky and Martin, 3rd edition
- Several slides follow Jurafsky and Martin

Affective meaning

- Drawing on literatures in
 - affective computing
 - linguistic subjectivity
 - social psychology
- Can we model the lexical semantics relevant to:
 - sentiment
 - emotion
 - personality
 - mood
 - attitudes



AFFECTIVE COMPUTING

Why compute affective meaning?

- Detecting:
 - sentiment towards politicians, products, countries, ideas
 - frustration of callers to a help line
 - stress in drivers or pilots
 - depression and other medical conditions
 - confusion in students talking to e-tutors
 - emotions in novels (e.g., for studying groups that are feared over time)
- Could we generate:
 - emotions or moods for literacy tutors in the children's storybook domain
 - emotions or moods for computer games
 - personalities for dialogue systems to match the user

Connotation in the lexicon

- Definition of connotation: an idea or feeling which a word invokes for a person in addition to its literal or primary meaning.
- An example: "the word 'discipline' has unhappy connotations of punishment and repression"
- Words have connotation as well as sense
- Can we build lexical resources that represent these connotations?
- And use them in these computational tasks?

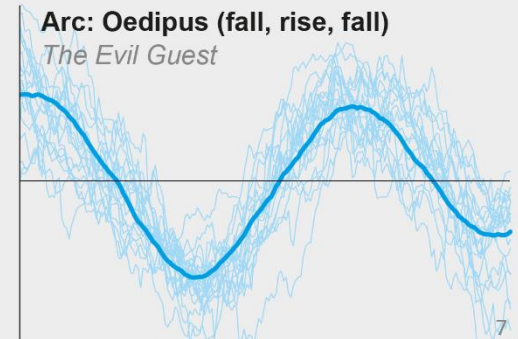
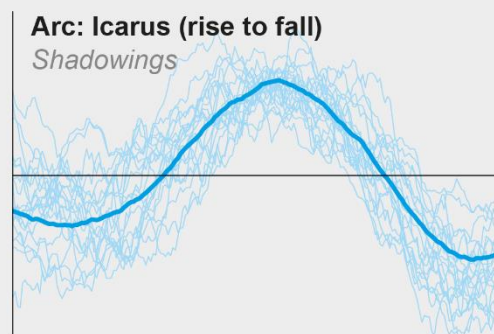
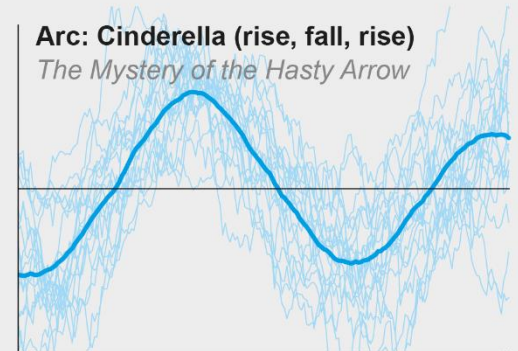
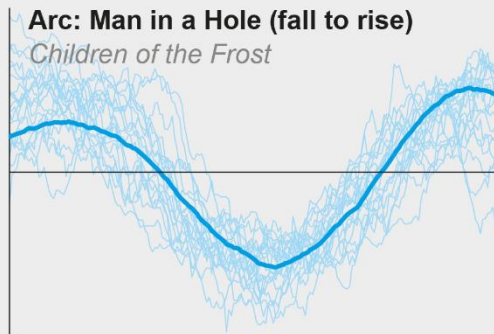
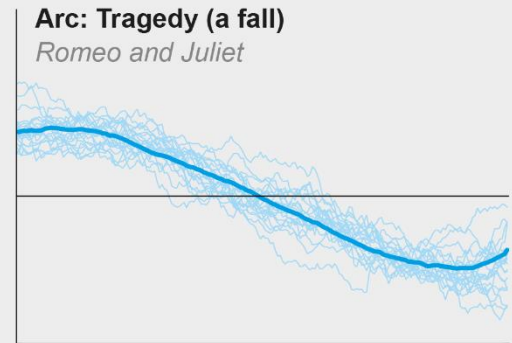
Scherer's typology of affective states

- **Emotion:** relatively brief episode of synchronized response of all or most organismic subsystems in response to the evaluation of an event as being of major significance
 - **angry, sad, joyful, fearful, ashamed, proud, desperate**
- **Mood:** diffuse affect state ...change in subjective feeling, of low intensity but relatively long duration, often without apparent cause
 - **cheerful, gloomy, irritable, listless, depressed, buoyant**
- **Interpersonal stance:** affective stance taken toward another person in a specific interaction, coloring the interpersonal exchange
 - **distant, cold, warm, supportive, contemptuous**
- **Attitudes:** relatively enduring, affectively colored beliefs, preferences predispositions towards objects or persons
 - **liking, loving, hating, valuing, desiring**
- **Personality traits:** emotionally laden, stable personality dispositions and behavior tendencies, typical for a person
 - **nervous, anxious, reckless, morose, hostile, envious, jealous**

Emotional states of novels

Emotional Arcs

About 85 percent of 1,327 fiction stories in the digitized Project Gutenberg collection follow one of six emotional arcs—a pattern of highs and lows from beginning to end (*dark curves*). The arcs are defined by the happiness or sadness of words in the running text (*jagged plots*). All books were in English and less than 100,000 words; examples are noted.



Sentiment Lexicons

Scherer's typology of affective states

Emotion: relatively brief episode of synchronized response of all or most organismic subsystems in response to the evaluation of an event as being of major significance

angry, sad, joyful, fearful, ashamed, proud, desperate

Mood: diffuse affect state ...change in subjective feeling, of low intensity but relatively long duration, often without apparent cause

cheerful, gloomy, irritable, listless, depressed, buoyant

Interpersonal stance: affective stance taken toward another person in a specific interaction, coloring the interpersonal exchange

distant, cold, warm, supportive, contemptuous

Attitudes: relatively enduring, affectively colored beliefs, preferences predispositions towards objects or persons

liking, loving, hating, valuing, desiring

Personality traits: emotionally laden, stable personality dispositions and behavior tendencies, typical for a person

nervous, anxious, reckless, morose, hostile, envious, jealous

The General Inquirer

Philip J. Stone, Dexter C Dunphy, Marshall S. Smith, Daniel M. Ogilvie. 1966. The General Inquirer: A Computer Approach to Content Analysis. MIT Press

- Home page: <http://www.wjh.harvard.edu/~inquirer>
- List of categories:
<http://www.wjh.harvard.edu/~inquirer/homecat.htm>
- Spreadsheet:
<http://www.wjh.harvard.edu/~inquirer/inquirerbasic.xls>
- Categories:
 - Positive (1915 words) and Negative (2291 words)
 - Strong vs Weak, Active vs Passive, Overstated versus Understated
 - Pleasure, Pain, Virtue, Vice, Motivation, Cognitive Orientation, etc.
- Free for research use

LIWC (Linguistic Inquiry and Word Count)

- 2300 words, >70 classes
- **Affective Processes**
 - negative emotion (*bad, weird, hate, problem, tough*)
 - positive emotion (*love, nice, sweet*)
- **Cognitive Processes**
 - Tentative (*maybe, perhaps, guess*), Inhibition (*block, constraint*)
- **Pronouns, Negation** (*no, never*), **Quantifiers** (*few, many*)
- commercial
- Home page: <http://www.liwc.net/>

MPQA Subjectivity Cues Lexicon

Theresa Wilson, Janyce Wiebe, and Paul Hoffmann (2005). Recognizing Contextual Polarity in Phrase-Level Sentiment Analysis. Proc. of HLT-EMNLP-2005.

Riloff and Wiebe (2003). Learning extraction patterns for subjective expressions. EMNLP-2003.

- Home page: http://www.cs.pitt.edu/mpqa/subj_lexicon.html
- 6,885 words
 - 2718 positive
 - 4912 negative
- Each word annotated for intensity (strong, weak)
- GNU GPL

Bing Liu Opinion Lexicon

Minqing Hu and Bing Liu. Mining and Summarizing Customer Reviews. ACM SIGKDD-2004.

- [Bing Liu's Page on Opinion Mining](#)
- <http://www.cs.uic.edu/~liub/FBS/opinion-lexicon-English.rar>
- 6786 words
 - 2006 positive
 - 4783 negative

SentiWordNet

Stefano Baccianella, Andrea Esuli, and Fabrizio Sebastiani. 2010
SENTIWORDNET 3.0: An Enhanced Lexical Resource for Sentiment
Analysis and Opinion Mining. LREC-2010

- Home page: <http://sentiwordnet.isti.cnr.it/>
- All WordNet synsets automatically annotated for degrees of positivity, negativity, and neutrality/objectiveness
- [estimable(J,3)] “may be computed or estimated”
Pos 0 Neg 0 Obj 1
- [estimable(J,1)] “deserving of respect or high regard”
Pos .75 Neg 0 Obj .25

Disagreements between polarity lexicons

Christopher Potts, [Sentiment Tutorial](#), 2011

	Opinion Lexicon	General Inquirer	SentiWordNet	LIWC
MPQA	33/5402 (0.6%)	49/2867 (2%)	1127/4214 (27%)	12/363 (3%)
Opinion Lexicon		32/2411 (1%)	1004/3994 (25%)	9/403 (2%)
General Inquirer			520/2306 (23%)	1/204 (0.5%)
SentiWordNet				174/694 (25%)
LIWC				

Other Affective Lexicons

Scherer's typology of affective states

Emotion: relatively brief episode of synchronized response of all or most organismic subsystems in response to the evaluation of an event as being of major significance

angry, sad, joyful, fearful, ashamed, proud, desperate

Mood: diffuse affect state ...change in subjective feeling, of low intensity but relatively long duration, often without apparent cause

cheerful, gloomy, irritable, listless, depressed, buoyant

Interpersonal stance: affective stance taken toward another person in a specific interaction, coloring the interpersonal exchange

distant, cold, warm, supportive, contemptuous

Attitudes: relatively enduring, affectively colored beliefs, preferences predispositions towards objects or persons

liking, loving, hating, valuing, desiring

Personality traits: emotionally laden, stable personality dispositions and behavior tendencies, typical for a person

nervous, anxious, reckless, morose, hostile, envious, jealous

Two families of theories of emotion

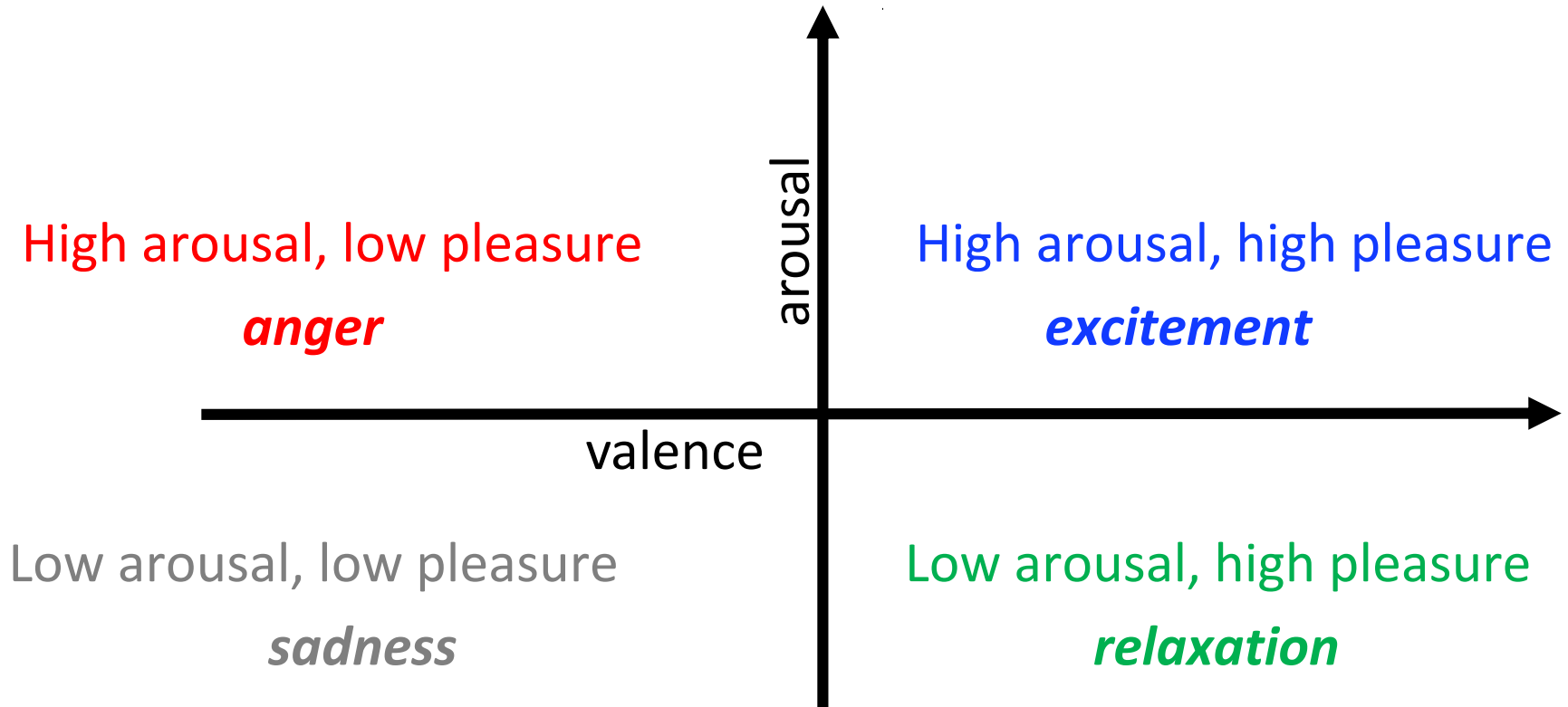
- Atomic basic emotions
 - A finite list of 6 or 8, from which others are generated
- Dimensions of emotion
 - Valence (positive, negative)
 - Arousal (strong, weak)
 - Control – dominance (in control, active vs. controlled, passive)

Ekman's 6 basic emotions



Surprise, happiness, anger, fear, disgust, sadness

Valence/Arousal Dimensions



Atomic units vs. Dimensions

Distinctive

- Emotions are units.
- Limited number of basic emotions.
- Basic emotions are innate and universal

Dimensional

- Emotions are dimensions.
- Limited # of labels but unlimited number of emotions.
- Emotions are culturally learned.

One emotion lexicon from each paradigm

1. 8 basic emotions:

- NRC Word-Emotion Association Lexicon (Mohammad and Turney 2011)

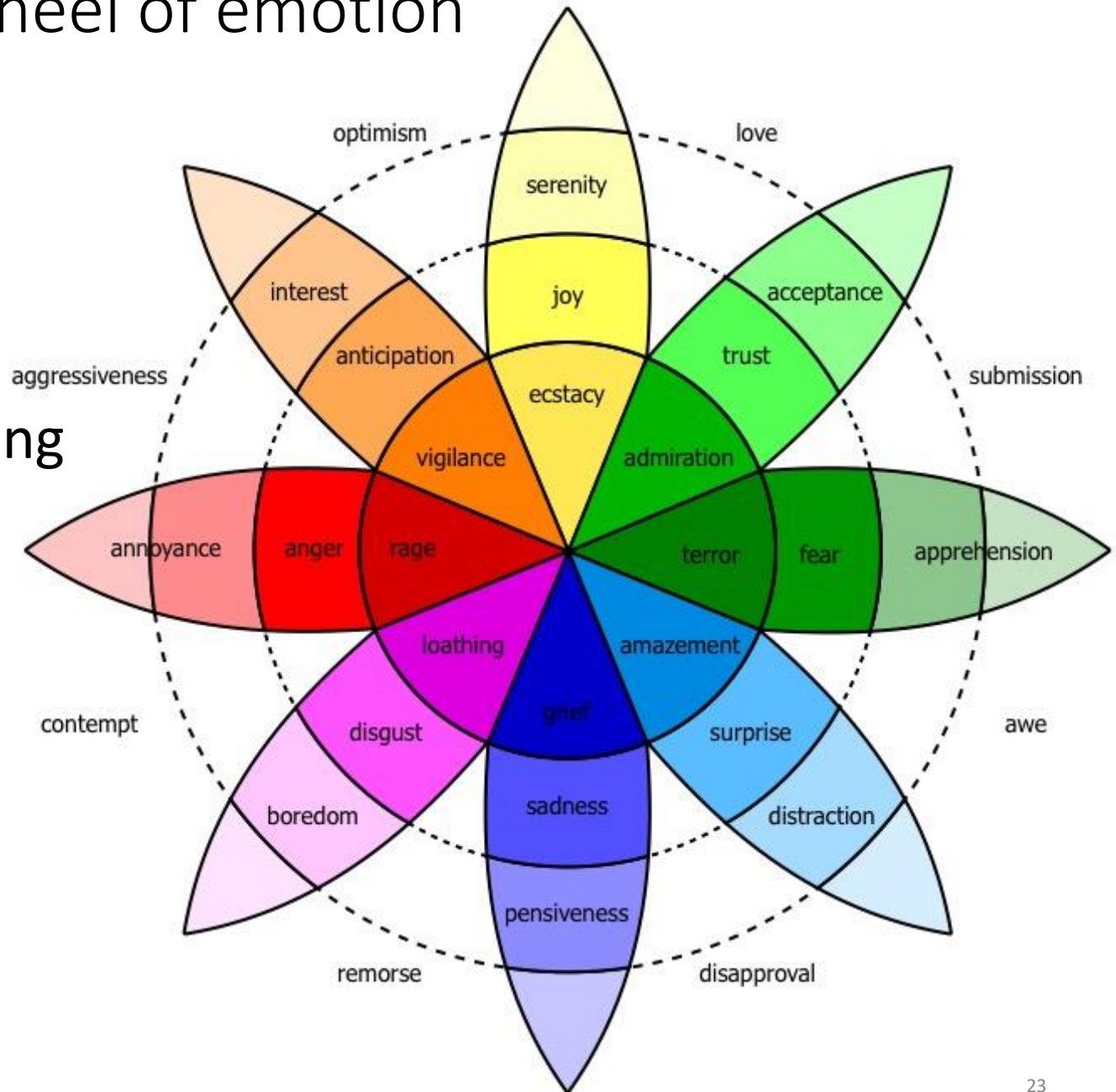
2. Dimensions of valence/arousal/dominance

- Warriner, A. B., **Kuperman**, V., and Brysbaert, M. (2013)

- Both built using Amazon Mechanical Turk

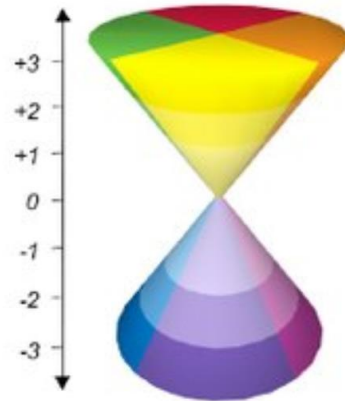
Plutchick's wheel of emotion

- 8 basic emotions
- in four opposing pairs
- joy–sadness
- anger–fear
- trust–disgust
- anticipation–surprise

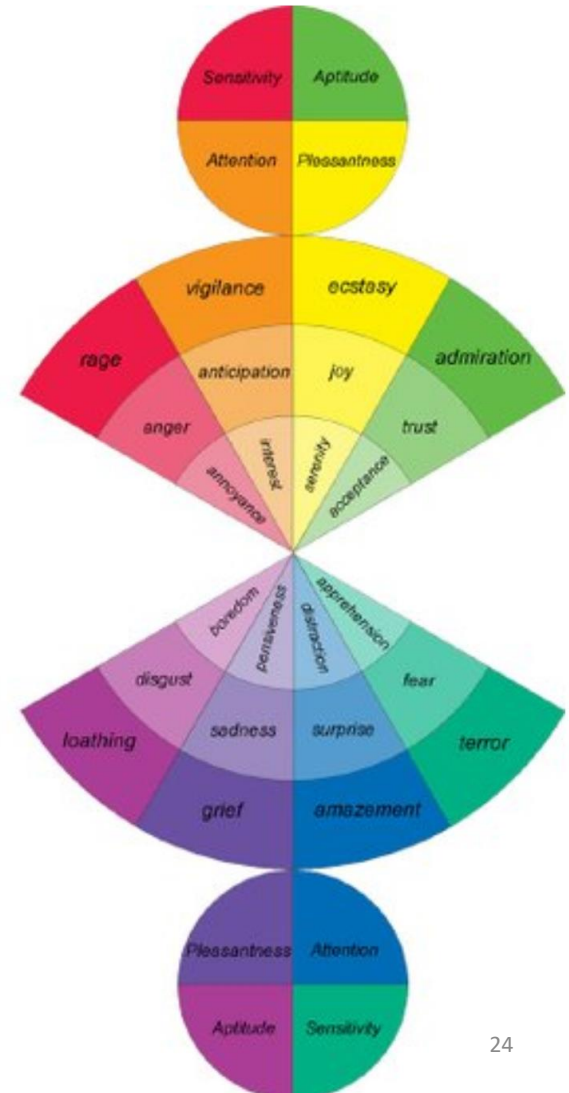


Hourglass of emotion

- The **Hourglass of Emotions** is a brain-inspired and psychologically-motivated model based on the idea that the mind is made of different independent resources and that emotional states result from turning some set of these resources on and turning another set of them off



	Pleasantness	Attention	Sensitivity	Aptitude
+3	ecstasy	vigilance	rage	admiration
+2	joy	anticipation	anger	trust
+1	serenity	interest	annoyance	acceptance
0	—	—	—	—
-1	pensiveness	distracton	apprehension	boredom
-2	sadness	surprise	fear	disgust
-3	grief	amazement	terror	loathing



NRC Word-Emotion Association Lexicon

Mohammad and Turney 2011

- 10,170 words chosen mainly from earlier lexicons
- Labeled by Amazon Mechanical Turk
- 5 Turkers per hit
- Give Turkers an idea of the relevant sense of the word
- Result:

amazingly	anger	0	
amazingly	anticipation	0	0
amazingly	disgust	0	
amazingly	fear	0	
amazingly	joy	1	
amazingly	sadness	0	
amazingly	surprise	1	
amazingly	trust	0	
amazingly	negative	0	
amazingly	positive	1	

EmoLex	# of terms
EmoLex-Uni:	
Unigrams from Macquarie Thesaurus	
adjectives	200
adverbs	200
nouns	200
verbs	200
EmoLex-Bi:	
Bigrams from Macquarie Thesaurus	
adjectives	200
adverbs	187
nouns	200
verbs	200
EmoLex-GI:	
Terms from General Inquirer	
negative terms	2119
neutral terms	4226
positive terms	1787
EmoLex-WAL:	
Terms from WordNet Affect Lexicon	
anger terms	165
disgust terms	37
fear terms	100
joy terms	165
sadness terms	120
surprise terms	53
Union	10170

The AMT Hit

Prompt word: *startle*

Q1. Which word is closest in meaning (most related) to *startle*?

- *automobile*
- *shake*
- *honesty*
- *entertain*

Q2. How positive (good, praising) is the word *startle*?

- *startle* is not positive
- *startle* is weakly positive
- *startle* is moderately positive
- *startle* is strongly positive

Q3. How negative (bad, criticizing) is the word *startle*?

- *startle* is not negative
- *startle* is weakly negative
- *startle* is moderately negative
- *startle* is strongly negative

Q4. How much is *startle* associated with the emotion joy? (For example, *happy* and *fun* are strongly associated with joy.)

- *startle* is not associated with joy
- *startle* is weakly associated with joy
- *startle* is moderately associated with joy
- *startle* is strongly associated with joy

Q5. How much is *startle* associated with the emotion sadness? (For example, *failure* and *heart-break* are strongly associated with sadness.)

- *startle* is not associated with sadness
- *startle* is weakly associated with sadness
- *startle* is moderately associated with sadness
- *startle* is strongly associated with sadness

Q6. How much is *startle* associated with the emotion fear? (For example, *horror* and *scary* are strongly associated with fear.)

- Similar choices as in 4 and 5 above

Q7. How much is *startle* associated with the emotion anger? (For example, *rage* and *shouting* are strongly associated with anger.)

- Similar choices as in 4 and 5 above

Q8. How much is *startle* associated with the emotion trust? (For example, *faith* and *integrity* are strongly associated with trust.)

- Similar choices as in 4 and 5 above

...

Q9. How much is *startle* associated with the emotion disgust? (For example, *gross* and *cruelty* are strongly associated with disgust.)

- Similar choices as in 4 and 5 above

Lexicon of valence, arousal, and dominance

- Warriner, A. B., Kuperman, V., and Brysbaert, M. (2013). [Norms of valence, arousal, and dominance for 13,915 English lemmas. *Behavior Research Methods* 45, 1191-1207.](#)
- **Ratings for 14,000 words for emotional dimensions:**
 - **valence** (the pleasantness of the stimulus)
 - **arousal** (the intensity of emotion provoked by the stimulus)
 - **dominance** (the degree of control exerted by the stimulus)

Lexicon of valence, arousal, and dominance

- **valence** (the pleasantness of the stimulus)
 - 9: happy, pleased, satisfied, contented, hopeful
 - 1: unhappy, annoyed, unsatisfied, melancholic, despaired, or bored
- **arousal** (the intensity of emotion provoked by the stimulus)
 - 9: stimulated, excited, frenzied, jittery, wide-awake, or aroused
 - 1: relaxed, calm, sluggish, dull, sleepy, or unaroused;
- **dominance** (the degree of control exerted by the stimulus)
 - 9: in control, influential, important, dominant, autonomous, or controlling
 - 1: controlled, influenced, cared-for, awed, submissive, or guided
- Again produced by AMT

Lexicon of valence, arousal, and dominance: Examples

Valence		Arousal		Dominance	
vacation	8.53	rampage	7.56	self	7.74
happy	8.47	tornado	7.45	incredible	7.74
whistle	5.7	zucchini	4.18	skillet	5.33
conscious	5.53	dressy	4.15	concur	5.29
torture	1.4	dull	1.67	earthquake	2.14

Concreteness versus abstractness 1/3

- The degree to which the concept denoted by a word refers to a perceptible entity.
 - Do concrete and abstract words differ in connotation?
 - Storage and retrieval?
 - Bilingual processing?
 - Relevant for embodied view of cognition (Barsalou 1999 inter alia)
 - Do concrete words activate brain regions involved in relevant perception

Concreteness versus abstractness 2/3

- Brysbaert et al, 2014
- 37,058 English words and 2,896 two-word expressions (“zebra crossing” and “zoom in”),
- Rating from 1 (abstract) to 5 (concrete)
- Calibrator words:
 - shirt, infinity, gas, grasshopper, marriage, kick, polite, whistle, theory, and sugar

Concreteness versus abstractness 3/3

- Some example ratings from the final dataset of 40,000 words and phrases

banana 5

bathrobe 5

bagel 5

brisk 2.5

badass 2.5

basically 1.32

belief 1.19

although 1.07

Perceptual Strength Norms

Connell and Lynott norms

Word	Perceptual strength					Concreteness	Imageability
	Auditory	Gustatory	Haptic	Olfactory	Visual		
soap	0.35	1.29	4.12	4.00	4.06	589	600
noisy	4.95	0.05	0.29	0.05	1.67	293	138
atom	1.00	0.63	0.94	0.50	1.38	481	499
republic	0.53	0.67	0.27	0.07	1.79	376	356

Semi-supervised algorithms for learning sentiment lexicons

Semi-supervised learning of lexicons

- Use a small amount of information
 - A few labeled examples
 - A few hand-built patterns
- To bootstrap a lexicon

Intuition for identifying word polarity

Vasileios Hatzivassiloglou and Kathleen R. McKeown. 1997. Predicting the Semantic Orientation of Adjectives. ACL, 174–181

- Adjectives conjoined by “*and*” have same polarity
 - Fair **and** legitimate, corrupt **and** brutal
 - *fair **and** brutal, *corrupt **and** legitimate
- Adjectives conjoined by “*but*” do not
 - fair **but** brutal

Word polarity: Step 1

- Label **seed set** of 1336 adjectives (all >20 in 21 million word WSJ corpus)
 - 657 positive
 - adequate central clever famous intelligent remarkable reputed sensitive slender thriving...
 - 679 negative
 - contagious drunken ignorant lanky listless primitive strident troublesome unresolved unsuspecting...

Word polarity: Step 2

- Expand seed set to conjoined adjectives



"was nice and"

[Nice location in Porto and the front desk staff was nice and helpful...](#)

www.tripadvisor.com/ShowUserReviews-g189180-d206904-r12068... 

Mercure Porto Centro: Nice location in Porto and the front desk staff **was nice and helpful** - See traveler reviews, 77 candid photos, and great deals for Porto, ...

nice, helpful

[If a girl was nice and classy, but had some vibrant purple dye in ...](#)

answers.yahoo.com > Home > All Categories > Beauty & Style > Hair 

4 answers - Sep 21

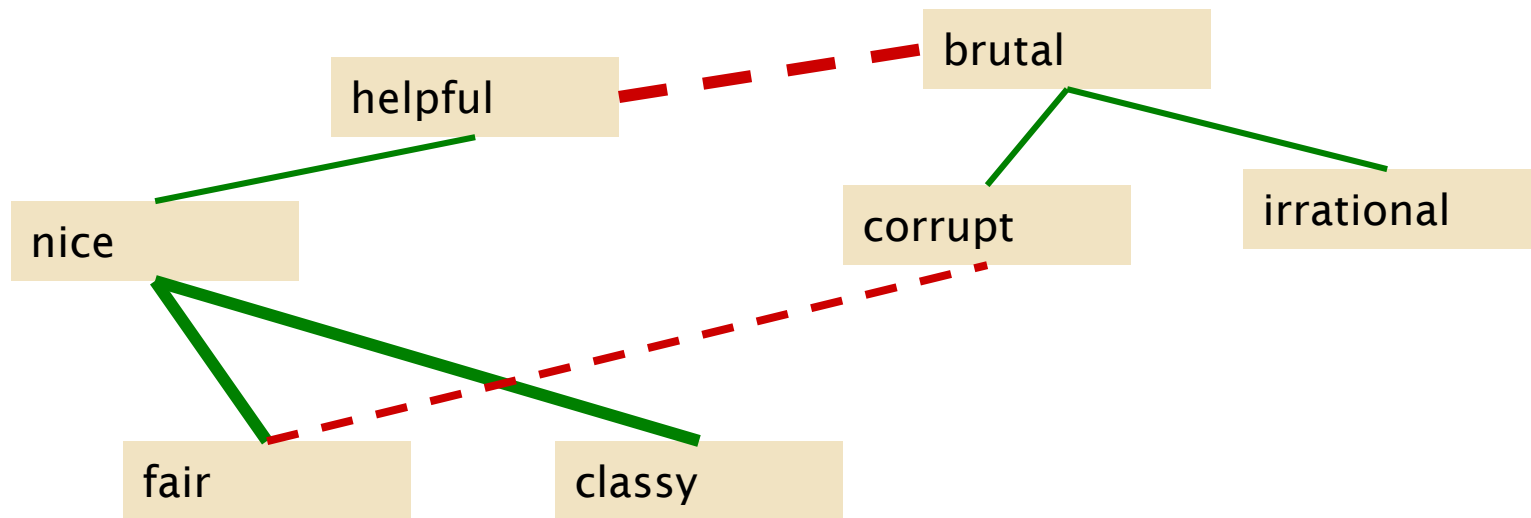
Question: Your personal opinion or what you think other people's opinions might ...

Top answer: I think she would be cool and confident like katy perry :)

nice, classy

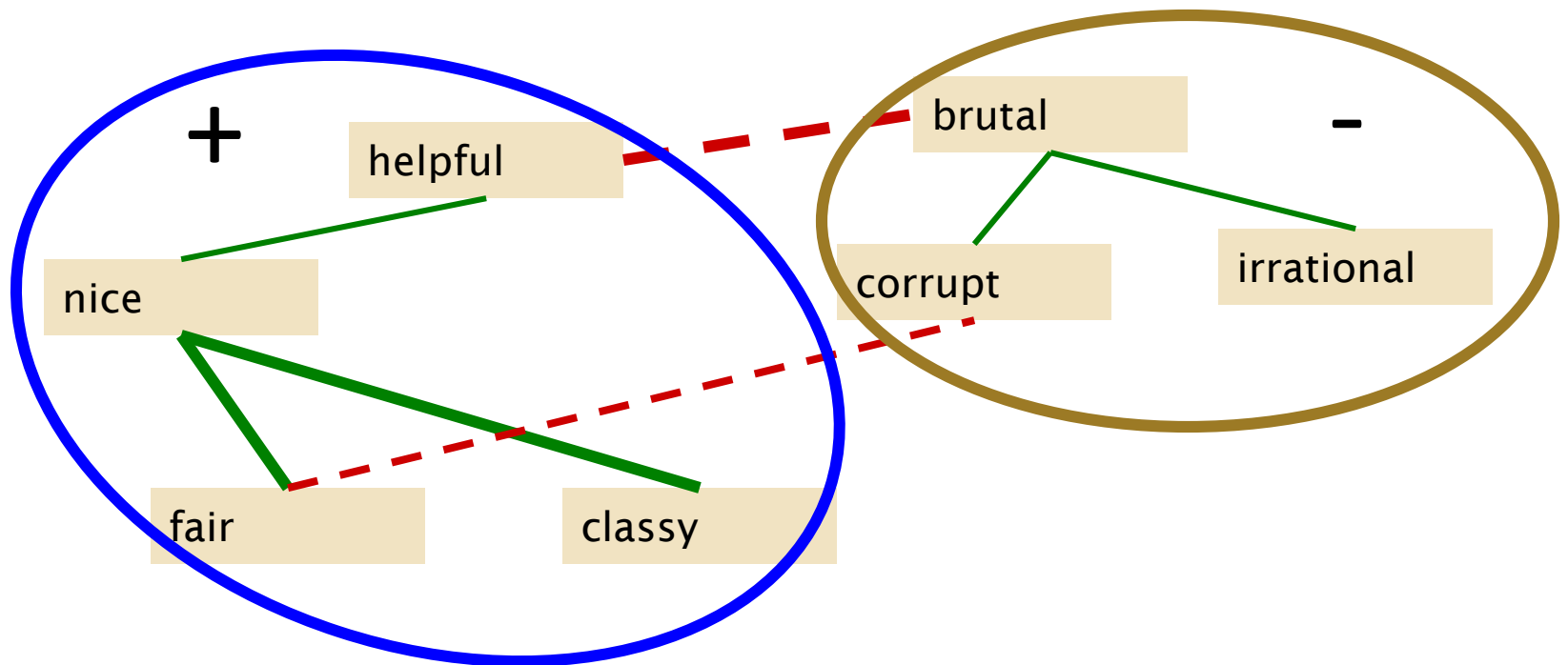
Word polarity: Step 3

- Supervised classifier assigns “polarity similarity” to each word pair, resulting in graph:



Word polarity: Step 4

- Clustering for partitioning the graph into two



Output polarity lexicon

- Positive

- bold decisive disturbing generous good honest important large mature patient peaceful positive proud sound stimulating straightforward strange talented vigorous witty...

- Negative

- ambiguous cautious cynical evasive harmful hypocritical inefficient insecure irrational irresponsible minor outspoken pleasant reckless risky selfish tedious unsupported vulnerable wasteful...

Output polarity lexicon

- Positive

- bold decisive **disturbing** generous good honest important large mature patient peaceful positive proud sound stimulating straightforward **strange** talented vigorous witty...

- Negative

- ambiguous **cautious** cynical evasive harmful hypocritical inefficient insecure irrational irresponsible minor **outspoken pleasant** reckless risky selfish tedious unsupported vulnerable wasteful...

Turney Algorithm

Turney (2002): Thumbs Up or Thumbs Down? Semantic Orientation Applied to Unsupervised Classification of Reviews

1. Extract a *phrasal lexicon* from reviews
2. Learn polarity of each phrase
3. Rate a review by the average polarity of its phrases

Extract two-word phrases with adjectives

First Word	Second Word	Third Word (not extracted)
JJ	NN or NNS	anything
RB, RBR, RBS	JJ	Not NN nor NNS
JJ	JJ	Not NN or NNS
NN or NNS	JJ	Nor NN nor NNS
RB, RBR, or RBS	VB, VBD, VBN, VBG	anything

How to measure polarity of a phrase?

- Positive phrases co-occur more with “*excellent*”
- Negative phrases co-occur more with “*poor*”
- But how to measure co-occurrence?

Pointwise Mutual Information

- **Mutual information** between 2 random variables X and Y

$$I(X, Y) = \sum_x \sum_y P(x, y) \log_2 \frac{P(x, y)}{P(x)P(y)}$$

- **Pointwise mutual information:**

- How much more do events x and y co-occur than if they were independent?

$$\text{PMI}(X, Y) = \log_2 \frac{P(x, y)}{P(x)P(y)}$$

Pointwise Mutual Information

- **Pointwise mutual information:**

- How much more do events x and y co-occur than if they were independent?

$$\text{PMI}(X, Y) = \log_2 \frac{P(x, y)}{P(x)P(y)}$$

- **PMI between two words:**

- How much more do two words co-occur than if they were independent?

$$\text{PMI}(\textit{word}_1, \textit{word}_2) = \log_2 \frac{P(\textit{word}_1, \textit{word}_2)}{P(\textit{word}_1)P(\textit{word}_2)}$$

How to Estimate Pointwise Mutual Information

- Query search engine
 - $P(\text{word})$ estimated by $\text{hits}(\text{word}) / N$
 - $P(\text{word}_1, \text{word}_2)$ by $\text{hits}(\text{word}_1 \text{ NEAR } \text{word}_2) / N$
 - (More correctly the bigram denominator should be kN , because there are a total of N consecutive bigrams $(\text{word}_1, \text{word}_2)$, but kN bigrams that are k words apart, but we just use N on the rest of this slide and the next.)

$$\text{PMI}(\text{word}_1, \text{word}_2) = \log_2 \frac{\frac{1}{N} \text{hits}(\text{word}_1 \text{ NEAR } \text{word}_2)}{\frac{1}{N} \text{hits}(\text{word}_1) \frac{1}{N} \text{hits}(\text{word}_2)}$$

Does phrase appear more with “poor” or “excellent”?

$$\begin{aligned}
 \text{Polarity}(\textit{phrase}) &= \text{PMI}(\textit{phrase}, \text{"excellent"}) - \text{PMI}(\textit{phrase}, \text{"poor"}) \\
 &= \log_2 \frac{\frac{1}{N} \text{hits}(\textit{phrase} \text{ NEAR } \text{"excellent"})}{\frac{1}{N} \text{hits}(\textit{phrase}) \frac{1}{N} \text{hits}(\text{"excellent"})} - \log_2 \frac{\frac{1}{N} \text{hits}(\textit{phrase} \text{ NEAR } \text{"poor"})}{\frac{1}{N} \text{hits}(\textit{phrase}) \frac{1}{N} \text{hits}(\text{"poor"})} \\
 &= \log_2 \frac{\text{hits}(\textit{phrase} \text{ NEAR } \text{"excellent"})}{\text{hits}(\textit{phrase}) \text{hits}(\text{"excellent"})} \frac{\text{hits}(\textit{phrase}) \text{hits}(\text{"poor"})}{\text{hits}(\textit{phrase} \text{ NEAR } \text{"poor"})} \\
 &= \log_2 \frac{\text{hits}(\textit{phrase} \text{ NEAR } \text{"excellent"}) \text{hits}(\text{"poor"})}{\text{hits}(\textit{phrase} \text{ NEAR } \text{"poor"}) \text{hits}(\text{"excellent"})}
 \end{aligned}$$

Phrases from a thumbs-up review

Phrase	POS tags	Polarity
online service	JJ NN	2.8
online experience	JJ NN	2.3
direct deposit	JJ NN	1.3
local branch	JJ NN	0.42
...		
low fees	JJ NNS	0.33
true service	JJ NN	-0.73
other bank	JJ NN	-0.85
inconveniently located	JJ NN	-1.5
<i>Average</i>		0.32

Phrases from a thumbs-down review

Phrase	POS tags	Polarity
direct deposits	JJ NNS	5.8
online web	JJ NN	1.9
very handy	RB JJ	1.4
...		
virtual monopoly	JJ NN	-2.0
lesser evil	RBR JJ	-2.3
other problems	JJ NNS	-2.8
low funds	JJ NNS	-6.8
unethical practices	JJ NNS	-8.5
<i>Average</i>		-1.2

Results of Turney algorithm

- 410 reviews from Epinions
 - 170 (41%) negative
 - 240 (59%) positive
- Majority class baseline: 59%
- Turney algorithm: 74%

- Phrases rather than words
- Learns domain-specific information

Using WordNet to learn polarity

S.M. Kim and E. Hovy. 2004. Determining the sentiment of opinions. COLING 2004

M. Hu and B. Liu. Mining and summarizing customer reviews. In Proceedings of KDD, 2004

- WordNet: online thesaurus
- Create positive (“good”) and negative seed-words (“terrible”)
- Find Synonyms and Antonyms
 - Positive Set: Add synonyms of positive words (“well”) and antonyms of negative words
 - Negative Set: Add synonyms of negative words (“awful”) and antonyms of positive words (“evil”)
- Repeat, following chains of synonyms
- Filter

Summary on semi-supervised lexicon learning

- Advantages:
 - Can be domain-specific
 - Can be more robust (more words)
- Intuition
 - Start with a seed set of words ('good', 'poor')
 - Find other words that have similar polarity:
 - Using “and” and “but”
 - Using words that occur nearby in the same document
 - Using WordNet synonyms and antonyms
- Use seeds and semi-supervised learning to induce lexicons

Supervised Learning of Sentiment Lexicons

Learn word sentiment supervised by online review scores

Potts, Christopher. 2011. On the negativity of negation. SALT 20, 636-659.
Potts 2011 NSF Workshop talk.

- Review datasets
 - IMDB, Goodreads, Open Table, Amazon, Trip Advisor
- Each review has a score (1-5, 1-10, etc.)
- Just count how many times each word occurs with each score (and normalize)

Analyzing the polarity of each word in IMDB

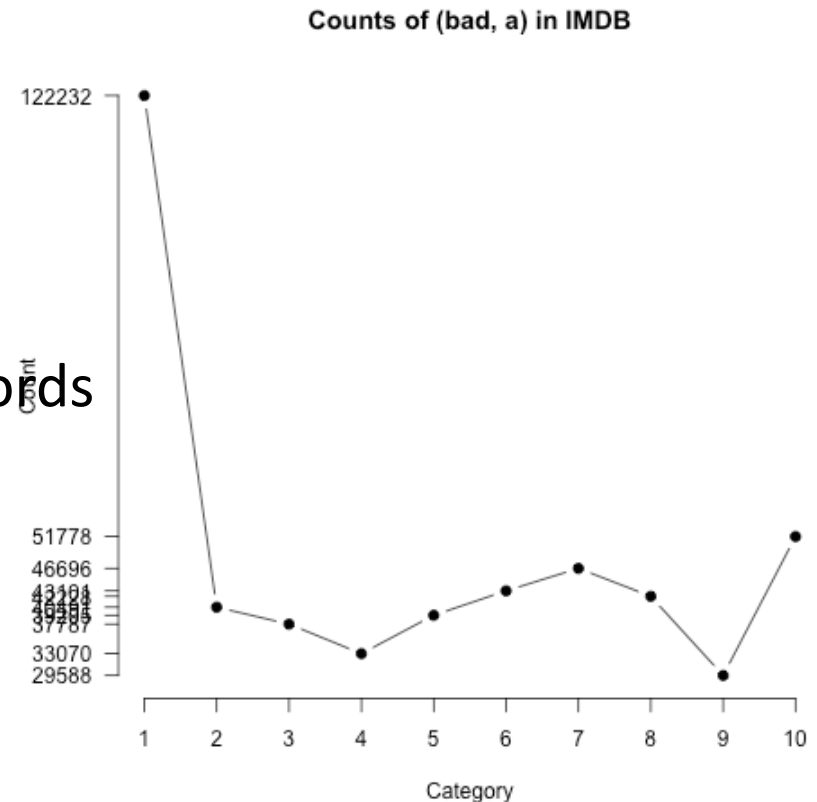
Potts, Christopher. 2011. On the negativity of negation. SALT 20, 636-659.

- How likely is each word to appear in each sentiment class?
- Count(“bad”) in 1-star, 2-star, 3-star, etc.
- But can’t use raw counts:
- Instead, **likelihood**:

$$P(w | c) = \frac{f(w, c)}{\sum_{w \in \mathcal{V}} f(w, c)}$$

- Make them comparable between words
 - **Scaled likelihood**:

$$\frac{P(w | c)}{P(w)}$$

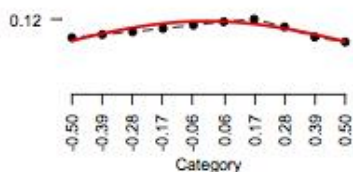


“Potts diagrams”

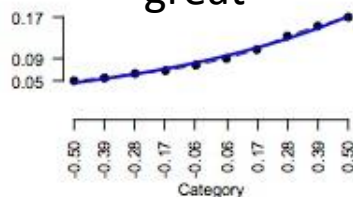
Potts, Christopher. 2011. NSF workshop on restructuring adjectives.

Positive scalars

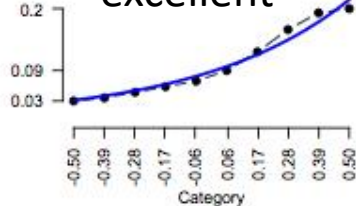
good



great

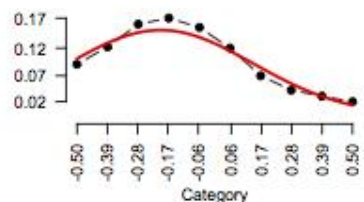


excellent

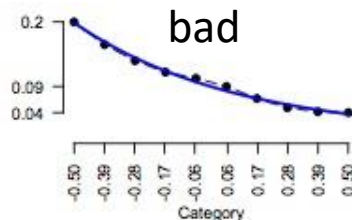


Negative scalars

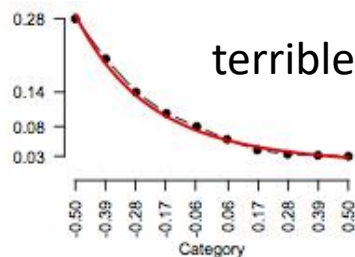
disappointing



bad

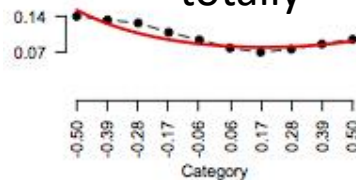


terrible

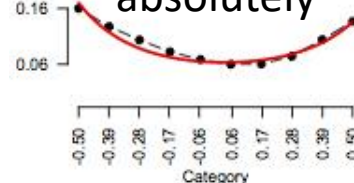


Emphatics

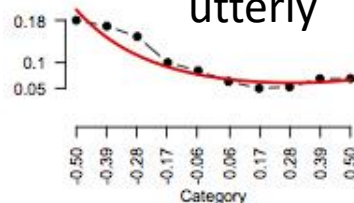
totally



absolutely

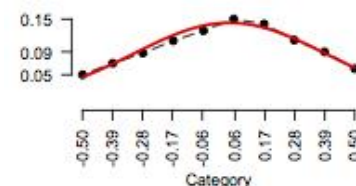


utterly

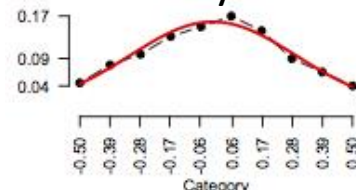


Attenuators

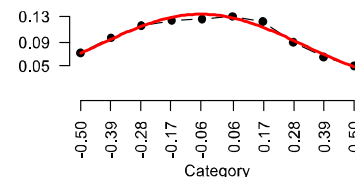
somewhat



fairly



pretty



Or use regression coefficients to weight words

- Train a classifier based on supervised data
 - Predict: human-labeled connotation of a document
 - From: all the words and bigrams in it
- Use the regression coefficients as the weights
- We'll return to an example of this later.

Using the lexicons to detect affect

Lexicons for detecting document affect: Simplest unsupervised method

- **Sentiment:**
 - Sum the weights of each positive word in the document
 - Sum the weights of each negative word in the document
 - Choose whichever value (positive or negative) has higher sum
- **Emotion:**
 - Do the same for each emotion lexicon

Lexicons for detecting document affect: Simplest supervised method

- Build a classifier
 - Predict sentiment (or emotion, or personality) given features
 - Use “counts of lexicon categories” as a features
 - Sample features:
 - LIWC category “cognition” had count of 7
 - NRC Emotion category “anticipation” had count of 2
- Baseline
 - Instead use counts of **all** the words and bigrams in the training set
 - This is hard to beat
 - But only works if the training and test sets are very similar

Sample affective task: personality detection

Scherer's typology of affective states

Emotion: relatively brief episode of synchronized response of all or most organismic subsystems in response to the evaluation of an event as being of major significance

angry, sad, joyful, fearful, ashamed, proud, desperate

Mood: diffuse affect state ...change in subjective feeling, of low intensity but relatively long duration, often without apparent cause

cheerful, gloomy, irritable, listless, depressed, buoyant

Interpersonal stance: affective stance taken toward another person in a specific interaction, coloring the interpersonal exchange

distant, cold, warm, supportive, contemptuous

Attitudes: relatively enduring, affectively colored beliefs, preferences predispositions towards objects or persons

liking, loving, hating, valuing, desiring

Personality traits: emotionally laden, stable personality dispositions and behavior tendencies, typical for a person

nervous, anxious, reckless, morose, hostile, envious, jealous

The Big Five Dimensions of Personality

- **Extraversion vs. Introversion**
 - sociable, assertive, playful vs. aloof, reserved, shy
- **Emotional stability vs. Neuroticism**
 - calm, unemotional vs. insecure, anxious
- **Agreeableness vs. Disagreeable**
 - friendly, cooperative vs. antagonistic, faultfinding
- **Conscientiousness vs. Unconscientious**
 - self-disciplined, organized vs. inefficient, careless
- **Openness to experience**
 - intellectual, insightful vs. shallow, unimaginative

Various text corpora labeled for personality of author

Pennebaker, James W., and Laura A. King. 1999. "Linguistic styles: language use as an individual difference." *Journal of personality and social psychology* 77, no. 6.

- 2,479 essays from psychology students (1.9 million words), “write whatever comes into your mind” for 20 minutes

Mehl, Matthias R, SD Gosling, JW Pennebaker. 2006. Personality in its natural habitat: manifestations and implicit folk theories of personality in daily life. *Journal of personality and social psychology* 90 (5), 862

- Speech from Electronically Activated Recorder (EAR)
- Random snippets of conversation recorded, transcribed
- 96 participants, total of 97,468 words and 15,269 utterances

Schwartz, H. Andrew, Johannes C. Eichstaedt, Margaret L. Kern, Lukasz Dziurzynski, Stephanie M. Ramones, Megha Agrawal, Achal Shah et al. 2013. "Personality, gender, and age in the language of social media: The open-vocabulary approach." *PloS one* 8, no. 9

- Facebook
- 75,000 volunteers
- 309 million words
- All took a personality test

EAR (speech) corpus (Mehl et al.)

Introvert	Extravert
<ul style="list-style-type: none">- Yeah you would do kilograms. Yeah I see what you're saying.- On Tuesday I have class. I don't know.- I don't know. A16. Yeah, that is kind of cool.- I don't know. I just can't wait to be with you and not have to do this every night, you know?- Yeah. You don't know. Is there a bed in there? Well ok just...	<ul style="list-style-type: none">- That's my first yogurt experience here. Really watery. Why?- Damn. New game.- Oh.- That's so rude. That.- Yeah, but he, they like each other. He likes her.- They are going to end up breaking up and he's going to be like.
Unconscientious	Conscientious
<ul style="list-style-type: none">- With the Chinese. Get it together.- I tried to yell at you through the window. Oh. xxxx's fucking a dumb ass. Look at him. Look at him, dude. Look at him. I wish we had a camera. He's fucking brushing his t-shirt with a tooth brush. Get a kick of it. Don't steal nothing.	<ul style="list-style-type: none">- I don't, I don't know for a fact but I would imagine that historically women who have entered prostitution have done so, not everyone, but for the majority out of extreme desperation and I think. I don't know, i think people understand that desperation and they don't don't see [...]

Essays corpus (Pennebaker and King)

Introvert	Extravert
<p>I've been waking up on time so far. What has it been, 5 days? Dear me, I'll never keep it up, being such not a morning person and all. But maybe I'll adjust, or not. I want internet access in my room, I don't have it yet, but I will on Wed??? I think. But that ain't soon enough, cause I got calculus homework [...]</p>	<p>I have some really random thoughts. I want the best things out of life. But I fear that I want too much! What if I fall flat on my face and don't amount to anything. But I feel like I was born to do BIG things on this earth. But who knows... There is this Persian party today.</p>
Neurotic	Emotionally stable
<p>One of my friends just barged in, and I jumped in my seat. This is crazy. I should tell him not to do that again. I'm not that fastidious actually. But certain things annoy me. The things that would annoy me would actually annoy any normal human being, so I know I'm not a freak.</p>	<p>I should excel in this sport because I know how to push my body harder than anyone I know, no matter what the test I always push my body harder than everyone else. I want to be the best no matter what the sport or event. I should also be good at this because I love to ride my bike.</p>

Classifiers

- **Mairesse**, François, Marilyn A. Walker, Matthias R. Mehl, and Roger K. Moore. "Using linguistic cues for the automatic recognition of personality in conversation and text." *Journal of artificial intelligence research* (2007): 457-500.
 - Various classifiers, lexicon-based and prosodic features
- **Schwartz**, H. Andrew, Johannes C. Eichstaedt, Margaret L. Kern, Lukasz Dziurzynski, Stephanie M. Ramones, Megha Agrawal, Achal Shah et al. 2013. "Personality, gender, and age in the language of social media: The open-vocabulary approach." *PloS one* 8, no.
 - regression and SVM, lexicon-based and all-words
 - Nowadays: use neural networks

Sample LIWC Features

LIWC (Linguistic Inquiry and Word Count)

Pennebaker, J.W., Booth, R.J., & Francis, M.E. (2007). Linguistic Inquiry and Word Count: LIWC 2007. Austin, TX

Feature	Type	Example
Anger words	LIWC	hate, kill, pissed
Metaphysical issues	LIWC	God, heaven, coffin
Physical state/function	LIWC	ache, breast, sleep
Inclusive words	LIWC	with, and, include
Social processes	LIWC	talk, us, friend
Family members	LIWC	mom, brother, cousin
Past tense verbs	LIWC	walked, were, had
References to friends	LIWC	pal, buddy, coworker
Imagery of words	MRC	Low: future, peace - High: table, car
Syllables per word	MRC	Low: a - High: uncompromisingly
Concreteness	MRC	Low: patience, candor - High: ship
Frequency of use	MRC	Low: duly, nudity - High: he, the

Normalizing LIWC category features

(Schwartz et al 2013, Facebook study)

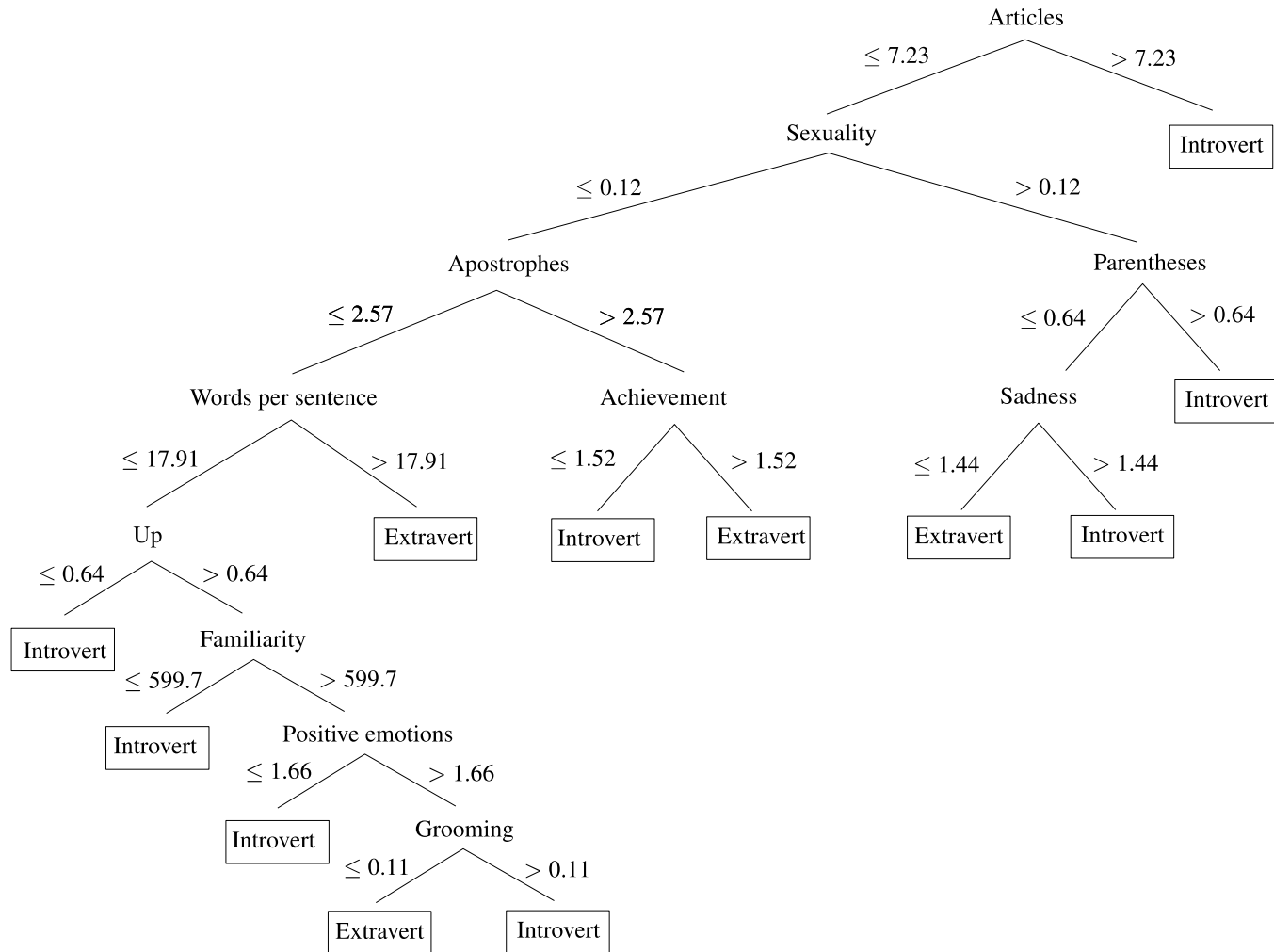
- Mairesse:
Raw LIWC counts
- Schwartz et al:
Normalized per writer:

$$p(\text{category} \mid \text{subject}) = \frac{\sum_{\text{word} \in \text{category}} \text{freq}(\text{word}, \text{subject})}{\sum_{\text{word} \in \text{vocab}(\text{subject})} \text{freq}(\text{word}, \text{subject})}$$

Sample results

- Agreeable:
 - +Family, +Home, -Anger, -Swear
- Extravert
 - +Friend, +Religion, +Self
- Conscientiousness:
 - -Swear, -Anger, -NegEmotion,
- Emotional Stability:
 - -NegEmotion, +Sports,
- Openness
 - -Cause, -Space

Decision tree for predicting extraversion in essay corpus (Mairesse et al)



Using all words instead of lexicons

Schwartz et al. (2013), Facebook study

- Choosing phrases with $pmi > 2 * \text{length}$ [in words]

$$pmi(\text{phrase}) = \log \frac{p(\text{phrase})}{\prod_{w \in \text{phrase}} p(w)}$$

- Only use words/phrases used by at least 1% of writers
- Normalize counts of words and phrases by writer

$$p(\text{phrase} \mid \text{subject}) = \frac{\text{freq}(\text{phrase}, \text{subject})}{\sum_{\text{phrase}' \in \text{vocab}(\text{subject})} \text{freq}(\text{phrase}', \text{subject})}$$

Evaluation of Schwartz et al (2013) Facebook Classifier

- Train on labeled training data
 - LIWC category counts
 - words and phrases (n-grams of size 1 to 3, passing a collocation filter)
- Tested on a held-out set
- Correlations with human labels
 - LIWC .21-.29
 - All Words .29-.41

Scherer's typology of affective states

Emotion: relatively brief episode of synchronized response of all or most organismic subsystems in response to the evaluation of an event as being of major significance

angry, sad, joyful, fearful, ashamed, proud, desperate

Mood: diffuse affect state ...change in subjective feeling, of low intensity but relatively long duration, often without apparent cause

cheerful, gloomy, irritable, listless, depressed, buoyant

Interpersonal stance: affective stance taken toward another person in a specific interaction, coloring the interpersonal exchange

distant, cold, warm, supportive, contemptuous

Attitudes: relatively enduring, affectively colored beliefs, preferences predispositions towards objects or persons

liking, loving, hating, valuing, desiring

Personality traits: emotionally laden, stable personality dispositions and behavior tendencies, typical for a person

nervous, anxious, reckless, morose, hostile, envious, jealous

Affect extraction: of course it's not just the lexicon

Ranganath et al (2013), McFarland et al (2014)

- Detecting interpersonal stance in conversation
- Speed dating study, 1000 4-minute speed dates
- Subjects labeled **selves** and **each other** for
 - friendly (each on a scale of 1-10)
 - awkward
 - flirtatious
 - assertive

Affect extraction: of course it's not just the lexicon

Logistic regression classifier with

- LIWC lexicons
- Other lexical features
 - Lists of hedges
hedge: a word or phrase that makes what you say less strong (I wondered if I could have a word with you?)
- Prosody (pitch and energy means and variance)
- Discourse features
 - Interruptions
 - Dialog acts/Adjacency pairs
 - sympathy (“Oh, that’s terrible”)
 - clarification question (“What?”)
 - appreciations (“That’s awesome!”)

Results on affect extraction





- Friendliness
 - -negEmotion
 - -hedge
 - higher pitch
- Awkwardness
 - +negation
 - +hedges
 - +questions

Summary: Connotation in the lexicon

- Words have various connotational aspects
- Methods for building connotation lexicons
 - Based on theoretical models of emotion, sentiment
 - By hand (mainly using crowdsourcing)
 - Semi-supervised learning from seed words
 - Fully supervised (when you can find a convenient signal in the world)
- Applying lexicons to detect affect and sentiment
 - Unsupervised: pick simple majority sentiment (positive/negative words)
 - Supervised: learn weights for each lexical category

Sentiment Analysis

Positive or negative movie review?

-  • Unbelievably disappointing
-  • Full of zany characters and richly applied satire, and some great plot twists
-  • This is the greatest screwball comedy ever filmed
-  • It was pathetic. The worst part about it was the boxing scenes.

Google Product Search



HP Officejet 6500A Plus e-All-in-One Color Ink-jet - Fax / copier / printer / scanner
\$89 online, \$100 nearby ★★★★★ 377 reviews
September 2010 - Printer - HP - Inkjet - Office - Copier - Color - Scanner - Fax - 250 sheets

Reviews

Summary - Based on 377 reviews



What people are saying

ease of use		"This was very easy to setup to four computers."
value		"Appreciate good quality at a fair price."
setup		"Overall pretty easy setup."
customer service		"I DO like honest tech support people."
size		"Pretty Paper weight."
mode		"Photos were fair on the high quality mode."
colors		"Full color prints came out with great quality."

Bing Shopping

HP Officejet 6500A E710N Multifunction Printer

[Product summary](#) [Find best price](#) **Customer reviews** [Specifications](#) [Related items](#)



\$121.53 - \$242.39 (14 stores)

Compare

Average rating **★★★★★** (144)



Most mentioned

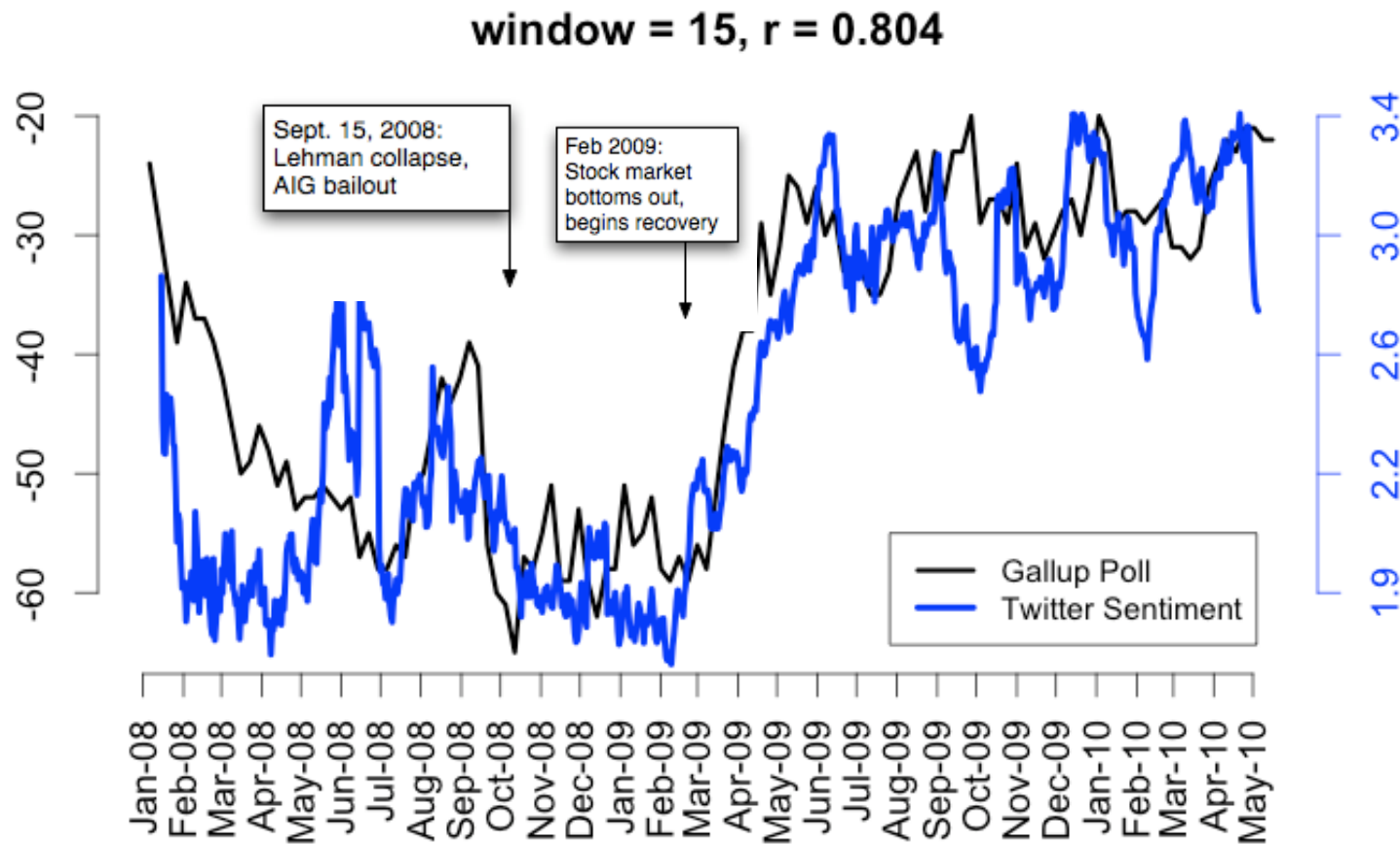


Show reviews by source

[Best Buy \(140\)](#)
[CNET \(5\)](#)
[Amazon.com \(3\)](#)

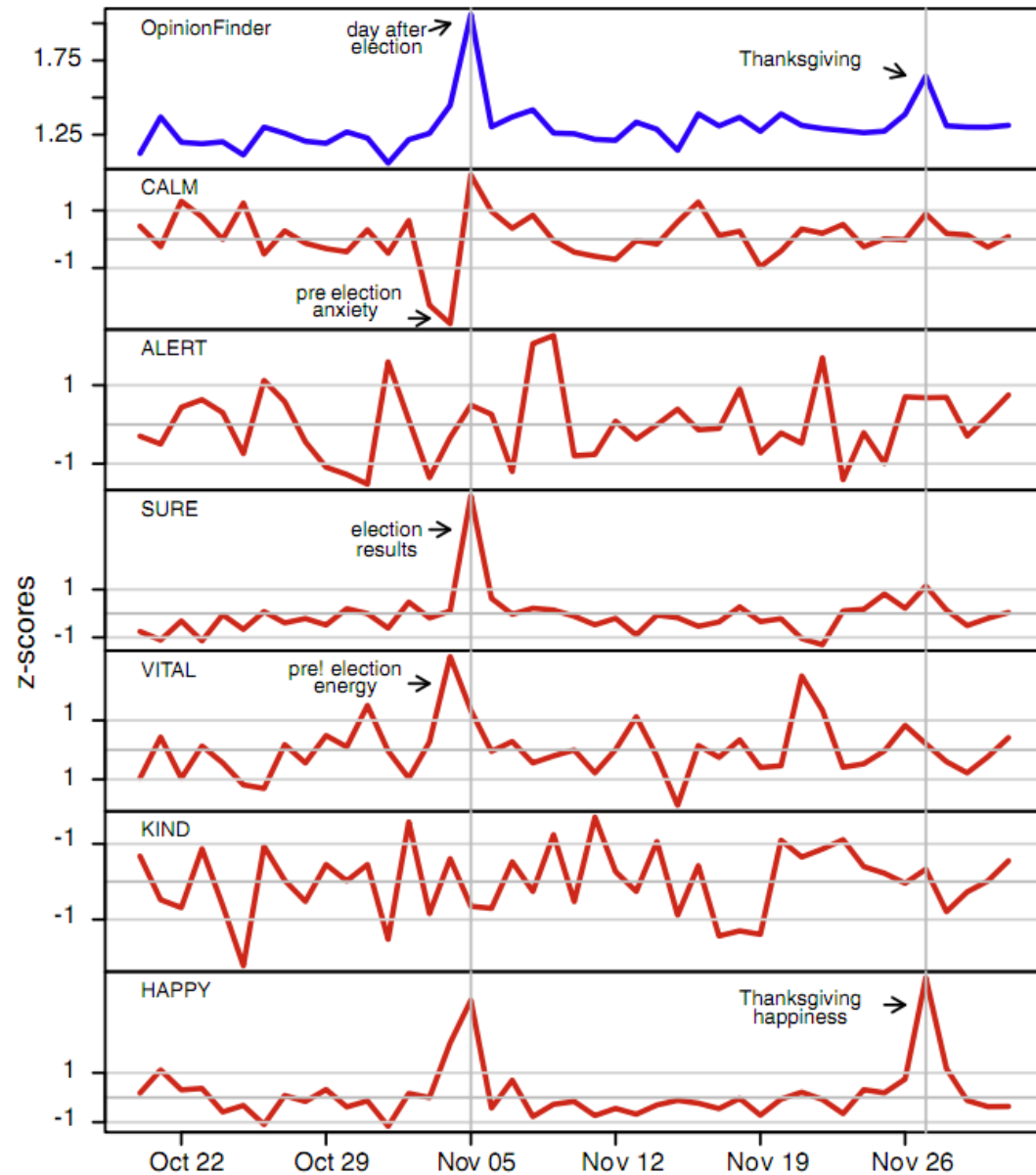
Twitter sentiment versus Gallup Poll of Consumer Confidence

Brendan O'Connor, Ramnath Balasubramanyan, Bryan R. Routledge, and Noah A. Smith. 2010. From Tweets to Polls: Linking Text Sentiment to Public Opinion Time Series. In ICWSM-2010



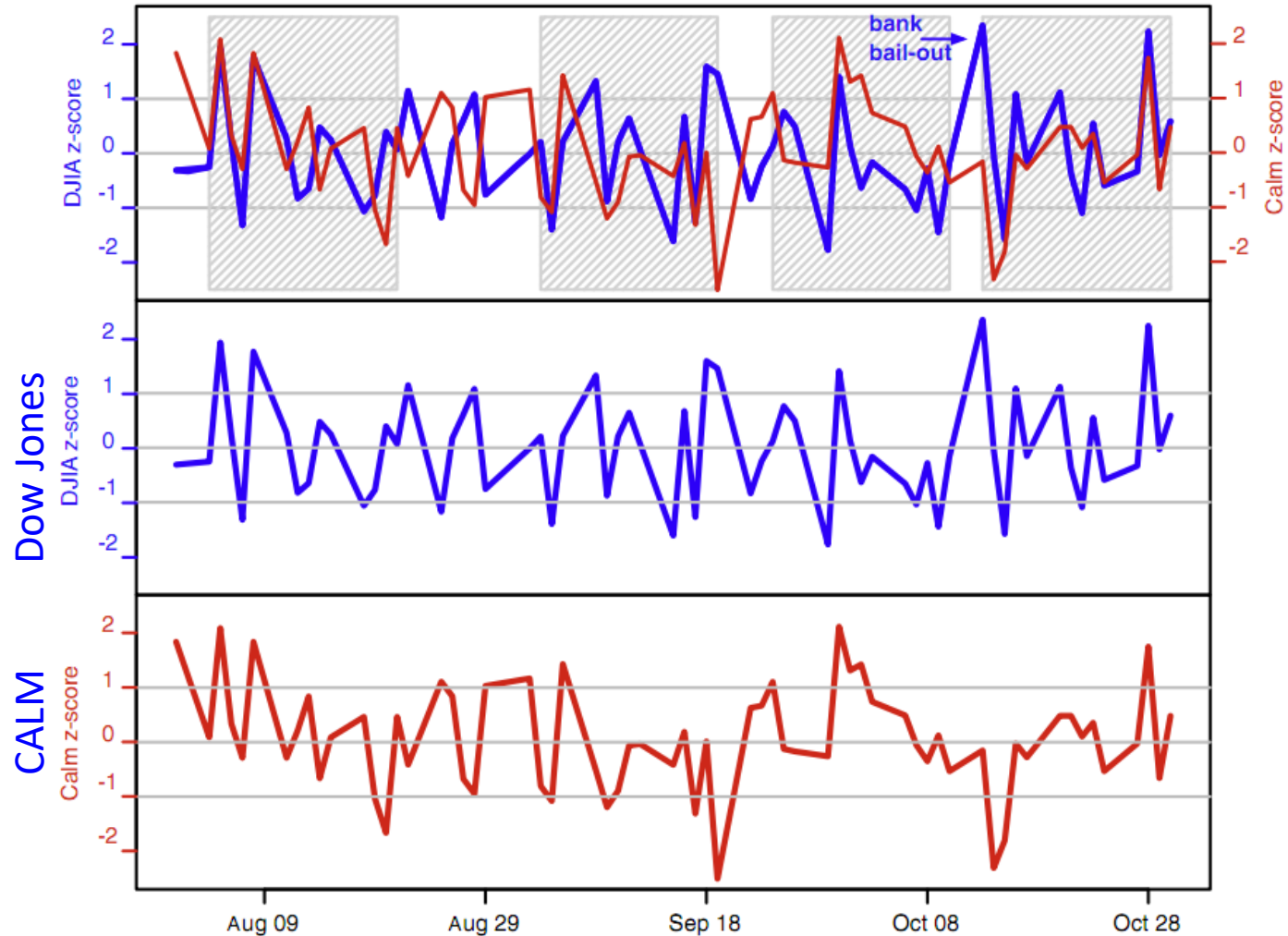
Twitter sentiment:

Johan Bollen, Huina Mao, Xiaojun Zeng. 2011. [Twitter mood predicts the stock market](#), Journal of Computational Science 2:1, 1-8. 10.1016/j.jocs.2010.12.007.



Bollen et al. (2011)

- CALM predicts Dow Jones Industrial Average (DJIA) 3 days later
- At least one current hedge fund uses this algorithm



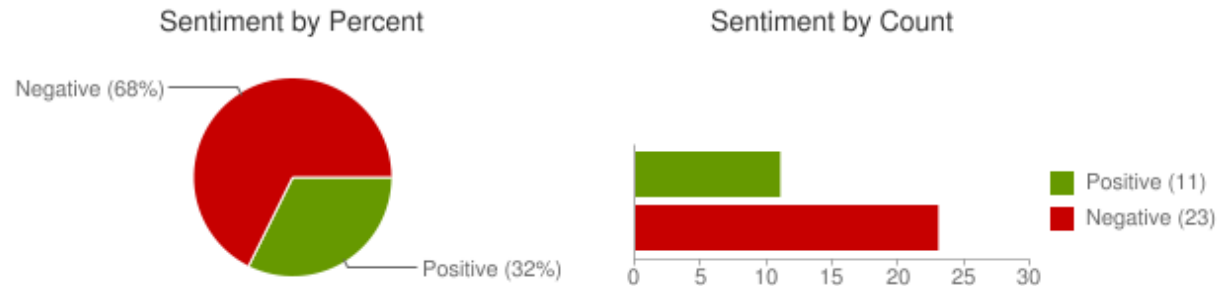
Target Sentiment on Twitter

- [Twitter Sentiment App](#)
- Alec Go, Richa Bhayani, Lei Huang. 2009. Twitter Sentiment Classification using Distant Supervision

Type in a word and we'll highlight the good and the bad

[Save this search](#)

Sentiment analysis for "united airlines"



[jjacobson](#): OMG... Could **@United airlines** have worse customer service? W8g now 15 minutes on hold 4 questions about a flight 2DAY that need a human.
Posted 2 hours ago

[12345clumsy6789](#): I hate **United Airlines** Ceiling!!! Fukn impossible to get my conduit in this damn mess! ?
Posted 2 hours ago

[EMLandPRGbelgiu](#): EML/PRG fly with Q8 **united airlines** and 24seven to an exotic destination. <http://t.co/Z9QloAjF>
Posted 2 hours ago

[CountAdam](#): FANTASTIC customer service from **United Airlines** at XNA today. Is tweet more, but cell phones off now!
Posted 4 hours ago

Sentiment analysis has many other names

- Opinion extraction
- Opinion mining
- Sentiment mining
- Subjectivity analysis

Why sentiment analysis?

- *Movie*: is this review positive or negative?
- *Products*: what do people think about the new iPhone?
- *Public sentiment*: how is consumer confidence? Is despair increasing?
- *Politics*: what do people think about this candidate or issue?
- *Prediction*: predict election outcomes or market trends from sentiment

Scherer Typology of Affective States

- **Emotion:** brief organically synchronized ... evaluation of a major event
 - *angry, sad, joyful, fearful, ashamed, proud, elated*
- **Mood:** diffuse non-caused low-intensity long-duration change in subjective feeling
 - *cheerful, gloomy, irritable, listless, depressed, buoyant*
- **Interpersonal stances:** affective stance toward another person in a specific interaction
 - *friendly, flirtatious, distant, cold, warm, supportive, contemptuous*
- **Attitudes: enduring, affectively colored beliefs, dispositions towards objects or persons**
 - *liking, loving, hating, valuing, desiring*
- **Personality traits:** stable personality dispositions and typical behavior tendencies
 - *nervous, anxious, reckless, morose, hostile, jealous*

Sentiment Analysis

- Sentiment analysis is the detection of **attitudes**
“enduring, affectively colored beliefs, dispositions towards objects or persons”
 1. **Holder (source)** of attitude
 2. **Target (aspect)** of attitude
 3. **Type** of attitude
 - From a set of types
 - *Like, love, hate, value, desire, etc.*
 - Or (more commonly) simple weighted **polarity**:
 - *positive, negative, neutral, together with strength*
 4. **Text** containing the attitude
 - Sentence or entire document

Sentiment Analysis

- Simplest task:
 - Is the attitude of this text positive or negative?
- More complex:
 - Rank the attitude of this text from 1 to 5
- Advanced:
 - Detect the target, source, or complex attitude types

Sentiment Analysis

- Simplest task:
 - Is the attitude of this text positive or negative?
- More complex:
 - Rank the attitude of this text from 1 to 5
- Advanced:
 - Detect the target, source, or complex attitude types

Sentiment Classification in Movie Reviews

Bo Pang, Lillian Lee, and Shivakumar Vaithyanathan. 2002. Thumbs up? Sentiment Classification using Machine Learning Techniques. EMNLP-2002, 79—86.

Bo Pang and Lillian Lee. 2004. A Sentimental Education: Sentiment Analysis Using Subjectivity Summarization Based on Minimum Cuts. ACL, 271-278

- Polarity detection:
 - Is an IMDB movie review positive or negative?
- Data: *Polarity Data 2.0*:
 - <http://www.cs.cornell.edu/people/pabo/movie-review-data>

IMDB data in the Pang and Lee database



when `_star wars_` came out some twenty years ago , the image of traveling throughout the stars has become a commonplace image . [...]

when han solo goes light speed , the stars change to bright lines , going towards the viewer in lines that converge at an invisible point .

cool .

`_october sky_` offers a much simpler image—that of a single white dot , traveling horizontally across the night sky . [. . .]



“ snake eyes ” is the most aggravating kind of movie : the kind that shows so much potential then becomes unbelievably disappointing .

it’s not just because this is a brian depalma film , and since he’s a great director and one who’s films are always greeted with at least some fanfare .

and it’s not even because this was a film starring nicolas cage and since he gives a brauvara performance , this film is hardly worth his talents .

Classical baseline algorithm

- Tokenization
- Feature Extraction
- Classification using different classifiers
 - Naïve Bayes
 - MaxEnt
 - SVM

Sentiment Tokenization Issues

- Deal with HTML and XML markup
- Twitter mark-up (names, hash tags)
- Capitalization (preserve for words in all caps)
- Phone numbers, dates
- Emoticons

Potts emoticons

```
[<>]?           # optional hat/brow
[:;=8]         # eyes
[\-o\*\']?    # optional nose
[\)\)\]\(\[dDpP/\:\}\{\@\\|\]\] # mouth
|              ##### reverse orientation
[\)\)\]\(\[dDpP/\:\}\{\@\\|\]\] # mouth
[\-o\*\']?    # optional nose
[:;=8]         # eyes
[<>]?         # optional hat/brow
```

Extracting Features for Sentiment Classification

- How to handle negation
 - I **didn't** like this movie
 - vs
 - I really like this movie
- Which words to use?
 - Only adjectives
 - All words
 - All words turns out to work better, at least on this data

Negation

Das, Sanjiv and Mike Chen. 2001. Yahoo! for Amazon: Extracting market sentiment from stock message boards. In Proceedings of the Asia Pacific Finance Association Annual Conference (APFA).

Bo Pang, Lillian Lee, and Shivakumar Vaithyanathan. 2002. Thumbs up? Sentiment Classification using Machine Learning Techniques. EMNLP-2002, 79—86.

- Add NOT_ to every word between negation and following punctuation:

didn't like this movie , but I



didn't NOT_like NOT_this NOT_movie but I

Reminder: Naïve Bayes

$$c_{NB} = \underset{c_j \in C}{\operatorname{argmax}} P(c_j) \prod_{i \in \text{positions}} P(w_i | c_j)$$

$$\hat{P}(w | c) = \frac{\text{count}(w, c) + 1}{\text{count}(c) + |V|}$$

Binarized (Boolean feature) Multinomial Naïve Bayes

- Intuition:
 - For sentiment (and probably for other text classification domains) word occurrence may matter more than word frequency
 - The occurrence of the word *fantastic* tells us a lot
 - The fact that it occurs 5 times may not tell us much more.
- Boolean Multinomial Naïve Bayes
 - Clips all the word counts in each document at 1

Problems: What makes reviews hard to classify?

- Subtlety:
 - Perfume review in *Perfumes: the Guide*:
 - “If you are reading this because it is your darling fragrance, please wear it at home exclusively, and tape the windows shut.”
 - Dorothy Parker on Katherine Hepburn
 - “She runs the gamut of emotions from A to B”

Thwarted Expectations and Ordering Effects

- “This film should be **brilliant**. It sounds like a **great** plot, the actors are **first grade**, and the supporting cast is **good** as well, and Stallone is attempting to deliver a good performance. However, it **can’t hold up**.”
- Well as usual Keanu Reeves is nothing special, but surprisingly, the **very talented** Laurence Fishbourne is **not so good** either, I was surprised.

Sentiment analysis in Slovene

- lexicon based on Bing Liu (2004), KSS
- a few other lexicons
- a few annotated datasets (tweets, user commentaries)
- SentiCoref, aspect based datasets (including coreferences)

Sentiment lexicons

- the most useful in English
Hu & Liu (2004), later updated,
2,006 positive and
4,783 negative words

positive words	negative words
a+	2-faced
abound	2-faces
abounds	abnormal
abundance	abolish
abundant	abominable
accessible	abominably
...	...

- in Slovene:
 - Rok Martinc (2013), based on AFINN-111 list (Nielsen, 2011), contains 2,477 words, estimated in range, -5...+5
 - Mateja Volčanšek (2015), based on General Inquirer (Stone, 1997), 1,669 positive and 1,912 negative words
 - Klemen Kadunc (2016), based on Hu & Liu

KSS lexicon

- based on Hu & Liu (2004)
- manually translated
- 2,646 positive and 6,689 negative words
- weaknesses
 - disregards the contexts (as all lexicons)
 - some informal Slovene expressions are not included due to translation
 - no gradation of sentiment
- provides lemmas as well as expanded word forms

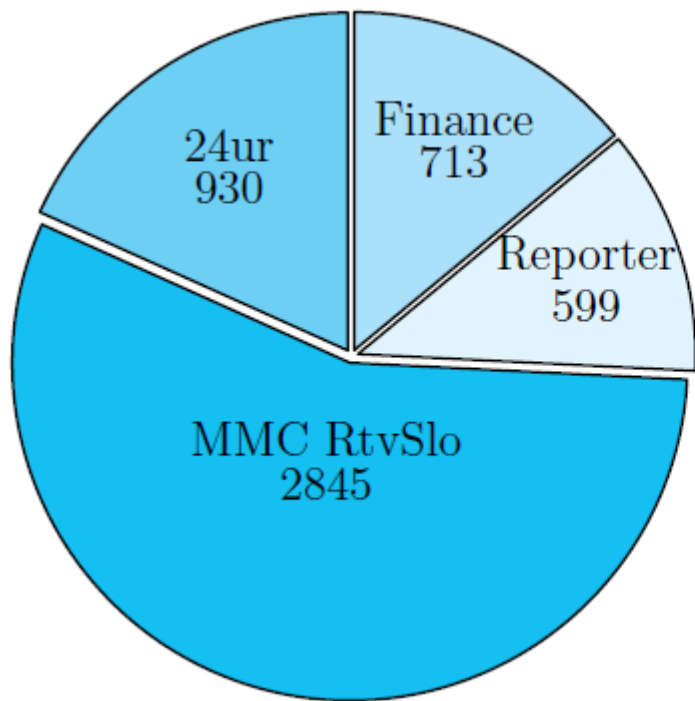
pozitivne besede	negativne besede
adut	abnormalen
aerodinamičen	absurd
agilen	absurden
agilno	absurdnost
aktualen	afektiran
ambiciozen	afnati
...	...

User commentaries corpus

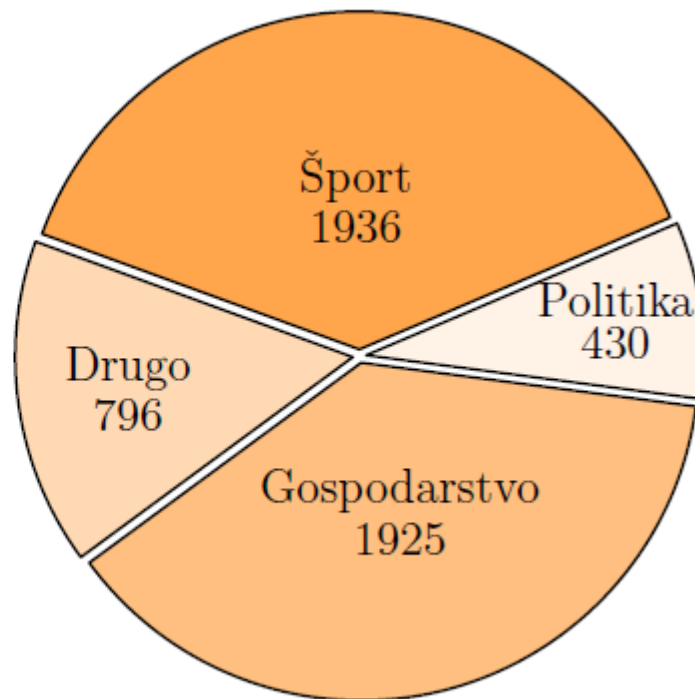
- searching for relevant commentaries using Google Search API.
- can set a configuration for comment extraction
- Supports four labels: positive, negative, neutral, and “not enough context”
- comments annotated by three annotators
- imbalanced and balanced corpus (580 of positive/negative)

	gosp.	politika	šport	drugo	RtvSlo	24ur	Finance	Reporter	skupaj
pozitivno	129	26	679	64	566	255	54	23	898
negativno	262	33	240	53	441	48	75	24	588
nevtralno	1420	351	882	638	1614	584	554	539	3291
skupaj	1811	410	1801	755	2621	887	683	586	4777

Distribution

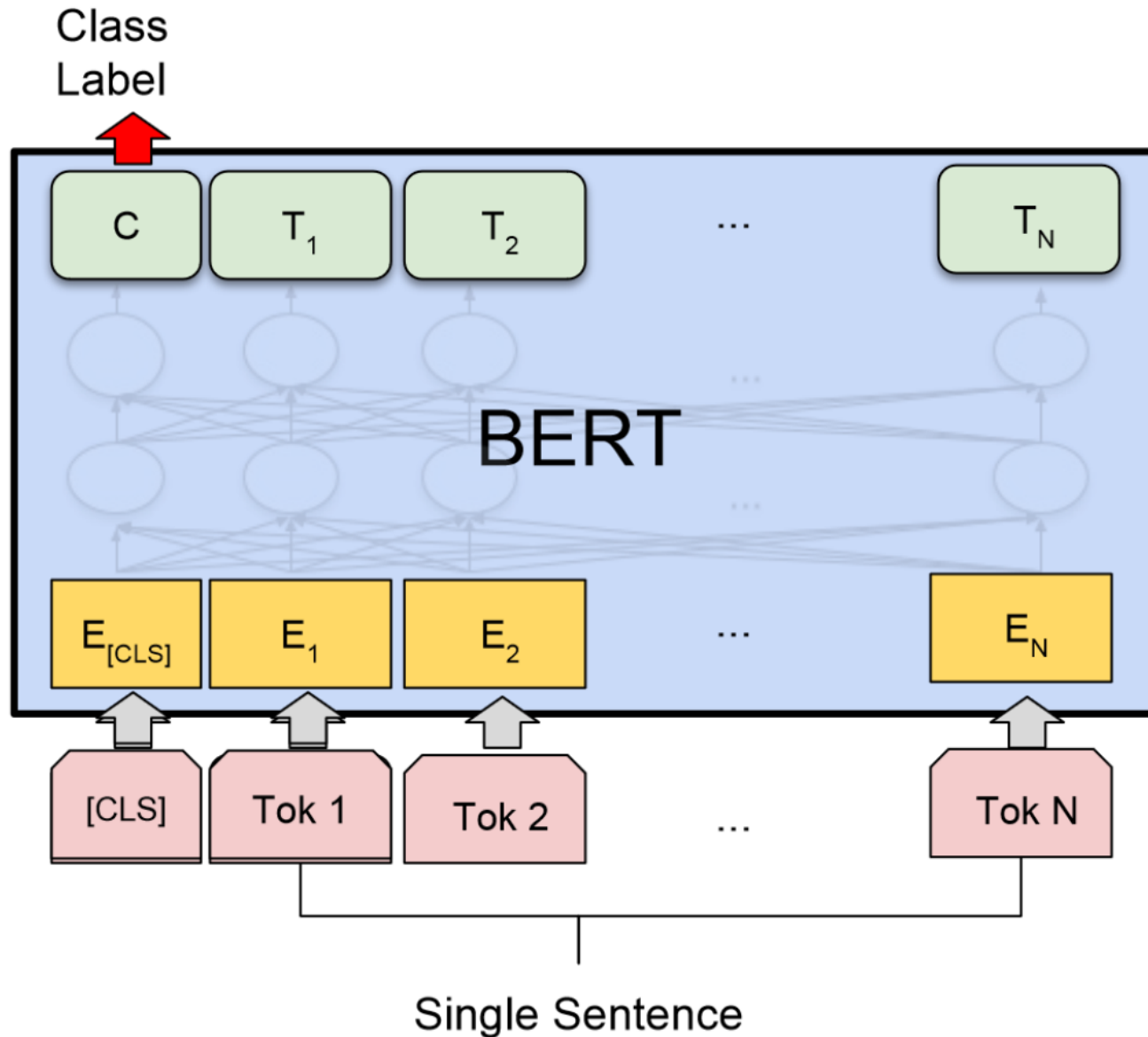


by source



by topic

Sentence classification using BERT – sentiment



Cross-lingual approach

Marko Robnik-Šikonja, Kristjan Reba, Igor Mozetič (2021). Cross-lingual Transfer of Sentiment Classifiers.
<https://arxiv.org/abs/2005.07456>

- Dataset: Twitter sentiment dataset in 13 languages

Language	Number of tweets				Agreement (\bar{F}_1)	
	Negative	Neutral	Positive	All	Self-	Inter-
Bosnian	12,868	11,526	13,711	38,105	0.78	-
Bulgarian	15,140	31,214	20,815	67,169	0.77	0.50
Croatian	21,068	19,039	43,894	84,001	0.83	-
English	26,674	46,972	29,388	103,034	0.79	0.67
German	20,617	60,061	28,452	109,130	0.73	0.42
Hungarian	10,770	22,359	35,376	68,505	0.76	-
Polish	67,083	60,486	96,005	223,574	0.84	0.67
Portuguese	58,592	53,820	44,981	157,393	0.74	-
Russian	34,252	44,044	29,477	107,773	0.82	-
Serbian	24,860	30,700	16,161	71,721	0.46	0.51
Slovak	18,716	14,917	36,792	70,425	0.77	-
Slovene	38,975	60,679	34,281	133,935	0.73	0.54
Swedish	25,319	17,857	15,371	58,547	0.76	-

Cross-lingual representation

1. projection of 93 languages into a joint embedding space (LASER library), using a parallel corpora with either English or Spanish match
 - embeddings, MLP layer with 8 neurons, followed by an output layer with 3 neurons (3 classes)
 - ReLU activation function, Adam optimizer
 - batch size 32, 10 epochs.
2. multilingual BERT, trained on 104 languages
3. CroSloEngual BERT, trained on Croatian, English and Slovene
 - fine-tuning both BERT models

XL transfer between similar languages

- reporting the classification accuracy and average F_1 score over positive and negative class

Source	Target	LASER		mBERT		CSE BERT		Both target		Source	Target	LASER		mBERT		CSE BERT		Both target	
		\bar{F}_1	CA	\bar{F}_1	CA	\bar{F}_1	CA	\bar{F}_1	CA			\bar{F}_1	CA	\bar{F}_1	CA	\bar{F}_1	CA	\bar{F}_1	CA
German	English	0.55	0.59	0.63	0.64	0.42	0.42	0.62	0.65	Croatian	Slovene	0.53	0.53	0.53	0.54	0.61	0.60	0.60	0.60
English	German	0.55	0.60	0.66	0.70	0.50	0.58	0.53	0.65	Croatian	English	0.63	0.63	0.63	0.66	0.62	0.64	0.62	0.65
Polish	Russian	0.64	0.59	0.57	0.57	0.50	0.40	0.70	0.70	English	Slovene	0.54	0.57	0.50	0.53	0.59	0.57	0.60	0.60
Polish	Slovak	0.63	0.59	0.58	0.59	0.63	0.65	0.72	0.72	English	Croatian	0.62	0.67	0.67	0.63	0.73	0.67	0.73	0.68
German	Swedish	0.58	0.57	0.59	0.59	0.58	0.56	0.67	0.65	Slovene	English	0.63	0.64	0.65	0.67	0.63	0.64	0.62	0.65
German Swedish	English	0.58	0.60	0.55	0.56	0.41	0.42	0.62	0.65	Slovene	Croatian	0.70	0.65	0.64	0.63	0.73	0.69	0.73	0.68
Slovene Serbian	Russian	0.53	0.55	0.57	0.57	0.58	0.48	0.70	0.70	Croatian English	Slovene	0.54	0.54	0.53	0.54	0.60	0.58	0.60	0.60
Slovene Serbian	Slovak	0.59	0.52	0.57	0.59	0.48	0.60	0.72	0.72	Croatian Slovene	English	0.62	0.61	0.65	0.67	0.63	0.65	0.62	0.65
Serbian	Slovene	0.54	0.57	0.54	0.54	0.56	0.55	0.60	0.60	English Slovene	Croatian	0.64	0.68	0.63	0.63	0.68	0.70	0.73	0.68
Serbian	Croatian	0.67	0.64	0.65	0.62	0.65	0.70	0.73	0.68	Average performance gap		0.04	0.03	0.04	0.03	0.00	0.01		
Serbian	Bosnian	0.65	0.61	0.61	0.60	0.59	0.62	0.67	0.64										
Polish	Slovene	0.51	0.48	0.55	0.54	0.50	0.53	0.60	0.60										
Slovak	Slovene	0.52	0.51	0.54	0.54	0.58	0.58	0.60	0.60										
Croatian	Slovene	0.53	0.53	0.53	0.54	0.61	0.60	0.60	0.60										
Croatian	Serbian	0.54	0.52	0.52	0.51	0.52	0.49	0.48	0.54										
Croatian	Bosnian	0.66	0.61	0.57	0.56	0.67	0.62	0.67	0.64										
Slovene	Croatian	0.70	0.65	0.64	0.63	0.73	0.69	0.73	0.68										
Slovene	Serbian	0.52	0.55	0.46	0.49	0.47	0.50	0.48	0.54										
Slovene	Bosnian	0.66	0.61	0.58	0.56	0.66	0.62	0.67	0.64										
Average performance gap		0.05	0.07	0.06	0.07	0.08	0.08												

Expansion of the training set with other languages

- unsuccessful if the dataset is large enough (as in the case shown)

Target	LASER				mBERT			
	All & Target		Only Target		All & Target		Only Target	
	\bar{F}_1	CA	\bar{F}_1	CA	\bar{F}_1	CA	\bar{F}_1	CA
Bosnian	0.64	0.59	0.67	0.64	0.63	0.60	0.65	0.60
Bulgarian	0.54	0.56	0.50	0.59	0.60	0.60	0.58	0.59
Croatian	0.63	0.57	0.73	0.68	0.65	0.63	0.64	0.66
English	0.58	0.60	0.62	0.65	0.64	0.69	0.68	0.68
German	0.52	0.59	0.53	0.65	0.61	0.66	0.66	0.66
Hungarian	0.59	0.61	0.60	0.67	0.65	0.69	0.65	0.69
Polish	0.67	0.63	0.70	0.66	0.71	0.71	0.70	0.70
Portuguese	0.44	0.39	0.52	0.51	0.52	0.52	0.50	0.49
Russian	0.66	0.64	0.70	0.70	0.67	0.66	0.64	0.64
Serbian	0.52	0.49	0.48	0.54	0.53	0.51	0.50	0.52
Slovak	0.64	0.61	0.72	0.72	0.67	0.65	0.67	0.66
Slovene	0.54	0.50	0.60	0.60	0.56	0.54	0.58	0.58
Swedish	0.63	0.59	0.67	0.65	0.67	0.64	0.67	0.65
Avg. gap	0.03	0.06			0.00	0.00		

Comparison of representations:
LASER, mBERT,
CSE BERT and
SVM

Language	LASER		mBERT		CSE BERT		SVM		Majority
	\bar{F}_1	CA	\bar{F}_1	CA	\bar{F}_1	CA	\bar{F}_1	CA	CA
Bosnian	0.68	0.64	0.65	0.60	0.68	0.65	(0.61	0.56)	0.36
Bulgarian	0.53	0.59	0.58	0.59	0.00	0.45	0.52	0.54	0.46
Croatian	0.72	0.68	0.64	0.66	0.76	0.71	(0.61	0.56)	0.52
English	0.62	0.65	0.68	0.68	0.67	0.66	0.63	0.64	0.44
German	0.52	0.64	0.66	0.66	0.31	0.59	0.54	0.61	0.53
Hungarian	0.63	0.67	0.65	0.69	0.57	0.65	0.64	0.67	0.53
Polish	0.70	0.66	0.70	0.70	0.56	0.57	0.68	0.63	0.44
Portuguese	0.48	0.47	0.50	0.49	0.12	0.22	0.55	0.51	0.37
Russian	0.70	0.70	0.64	0.64	0.07	0.43	0.61	0.60	0.40
Serbian	0.50	0.54	0.50	0.52	0.30	0.50	(0.61	0.56)	0.43
Slovak	0.72	0.72	0.67	0.66	0.69	0.71	0.68	0.68	0.52
Slovene	0.57	0.58	0.58	0.58	0.60	0.61	0.55	0.54	0.43
Swedish	0.67	0.64	0.67	0.65	0.54	0.56	0.66	0.62	0.43
#Best	5	3	6	6	3	3	2	2	0
Average	0.62	0.63	0.62	0.62	0.45	0.56	0.61	0.60	0.45