

Opinion Mining and Sentiment Analysis

Part 1

Isabelle Augenstein

augenstein@di.ku.dk

@IAugenstein

<http://isabelleaugenstein.github.io/>

Web Science Lecture
12 March 2019

UNIVERSITY OF COPENHAGEN



Web Science evaluation:

- 1) Quiz master project (incl. oral presentation): 40%
- 2) Take home assignment: 60%

Lecture topics	Assessment form
WWW network & challenges (lec. 1)	Take home assignment
collective intelligence & crowdsourcing (lec. 2-3)	Quiz master project
recommender systems (lec. 4-5)	Quiz master project and take home assignment
opinion & data mining (lec. 6-7)	Quiz master project and take home assignment
data analytics & search engine optimisation (lec. 8)	Take home assignment

The course grade will reflect how well the student:

- Can explain basic Web principles and properties to both laymen and specialists
- Can use standard procedures and practices when designing or implementing Web mining and analytics solutions
- Can present evaluation analyses and results in a proper format such that a technically qualified person can follow and obtain similar findings

Quiz master grade components:

- Completing all compulsory parts of the project
- Exhibiting solid understanding of the methods used in the project (how they work, their limitations, why they give the output they give)
- Answering correctly questions during the presentation

Problem: what do people think about X?

What opinions do people have about person X?

How has the opinion about person X changed over time?

How will changing X affect public opinion?

What product should I buy?

Is this person Y happy/unhappy with product X?

What aspects of the product are good/bad?

Are most people happy/unhappy with the product?



Tasks

- Assess the sentiment, i.e. feeling expressed in a text automatically
- Detect words / phrases in texts expressing such feelings
- "*I am concerned the imposed austerity measures will hurt the economy in the long term*"

Tasks

- Assess the sentiment, i.e. feeling expressed in a text automatically
 - Detect words / phrases in texts expressing such feelings
 - "*I am concerned the imposed austerity measures will hurt the economy in the long term*"
-
- Extract opinions about entities from texts
 - "*I wasn't satisfied with the location of the hotel*", "*A one-night stay was surprisingly cheap*" -> location: negative; price: positive

Tasks

Sentiment Analysis

- Assess the sentiment, i.e. feeling expressed in a text automatically
- Detect words / phrases in texts expressing such feelings
- "*I am concerned the imposed austerity measures will hurt the economy in the long term*"

Opinion Mining

- Extract opinions about entities from texts
- "*I wasn't satisfied with the location of the hotel*", "*A one-night stay was surprisingly cheap*" -> location: negative; price: positive
- In practice, the two terms are used interchangeably

Lecture Overview

- Applications and Use Cases
- Tasks
 - Sentiment Analysis
 - Opinion Mining
- Approaches
- Challenges

Applications

Upper Engadin > St. Moritz > Kempinski Grand Hotel Des Bains, St. Moritz (Switzerland)
 3 properties 119 properties

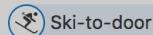


Room info & price

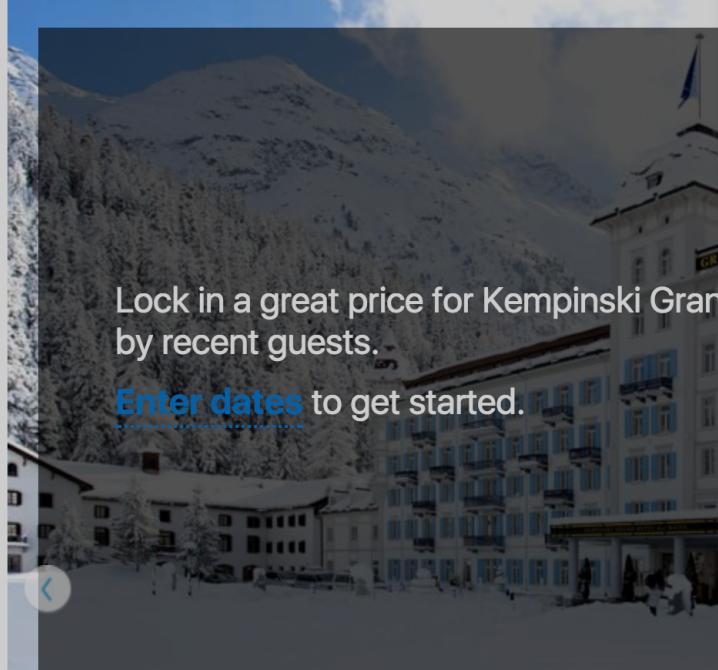
Facilities

House rules

Kempinski Grand Hotel Des Bains ★★★★☆



Via Mezdi 27, 7500 St. Moritz, Switzerland – [Great location - show map](#)



Lock in a great price for Kempinski Grand Hotel Des Bains by recent guests.

[Enter dates](#) to get started.



Get the celebrity treatment with world-class service
Kempinski Grand Hotel Des Bains

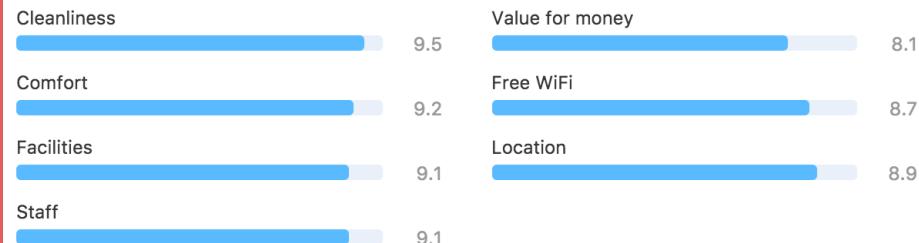


100% verified reviews

Real guests. Real stays. Real opinions. [Read more](#)

9.0 Superb · 694 reviews ▾

Aspects with ratings



Missing something? [Yes](#) / [No](#)

Show reviews from: All reviewers

All review scores

Show me reviews in: English 142 reviews

German 102 reviews

Spanish 3 reviews

Sort by: Recommended

Topics

Select a topic to filter reviews

- Cleanliness
- Food & Beverage
- Parking & Transport
- Breakfast
- Freebies
- Location
- Price
- Spaciousness
- Ambiance
- WiFi
- Facilities
- Spa & Gym
- Views & Surroundings
- Bedding
- Staff
- In-room facilities
- Bathroom
- Quietness



Property Scout review

Property Scouts are guests just like you. They're dedicated to reporting back the full story with detailed reviews.

Reviewed: 21 December 2016



"Perfection in Saint Moritz"

Select a topic to filter reviews

[Cleanliness](#)[Food & Beverage](#)[Parking & Transport](#)[Breakfast](#)[Freebies](#)[Location](#)[Price](#)[Spaciousness](#)[Ambiance](#)[WiFi](#)[Facilities](#)[Spa & Gym](#)[Views & Surroundings](#)[Bedding](#)[Staff](#)[In-room facilities](#)[Bathroom](#)[Quietness](#)

Reviewed: 18 December 2017

**Sandhya75**

Switzerland

Age group: 35 – 44



7 reviews

1 helpful vote

10**"Perfect stay"**

The rooms are rather small and I was missing a tea/coffee maker . But that is all I could mention



The location is perfect for skiing activities. Cross country im front of the door and downhill across the street. And you can rent the skis inhouse!

Breakfast is very broad and the included a la carte choices just makes breakfast a feast.

The staff is attentive and very friendly.

Stayed in December 2017

Helpful



24 countries

13 languages

JSON

Dansk, Deutsch, English, Español, Français, Italiano, Nederlands, Norsk, Polska, Português, Suomi, Svenska, русский

age
gender
location

KØBENHAVNS
UNIVERSITET

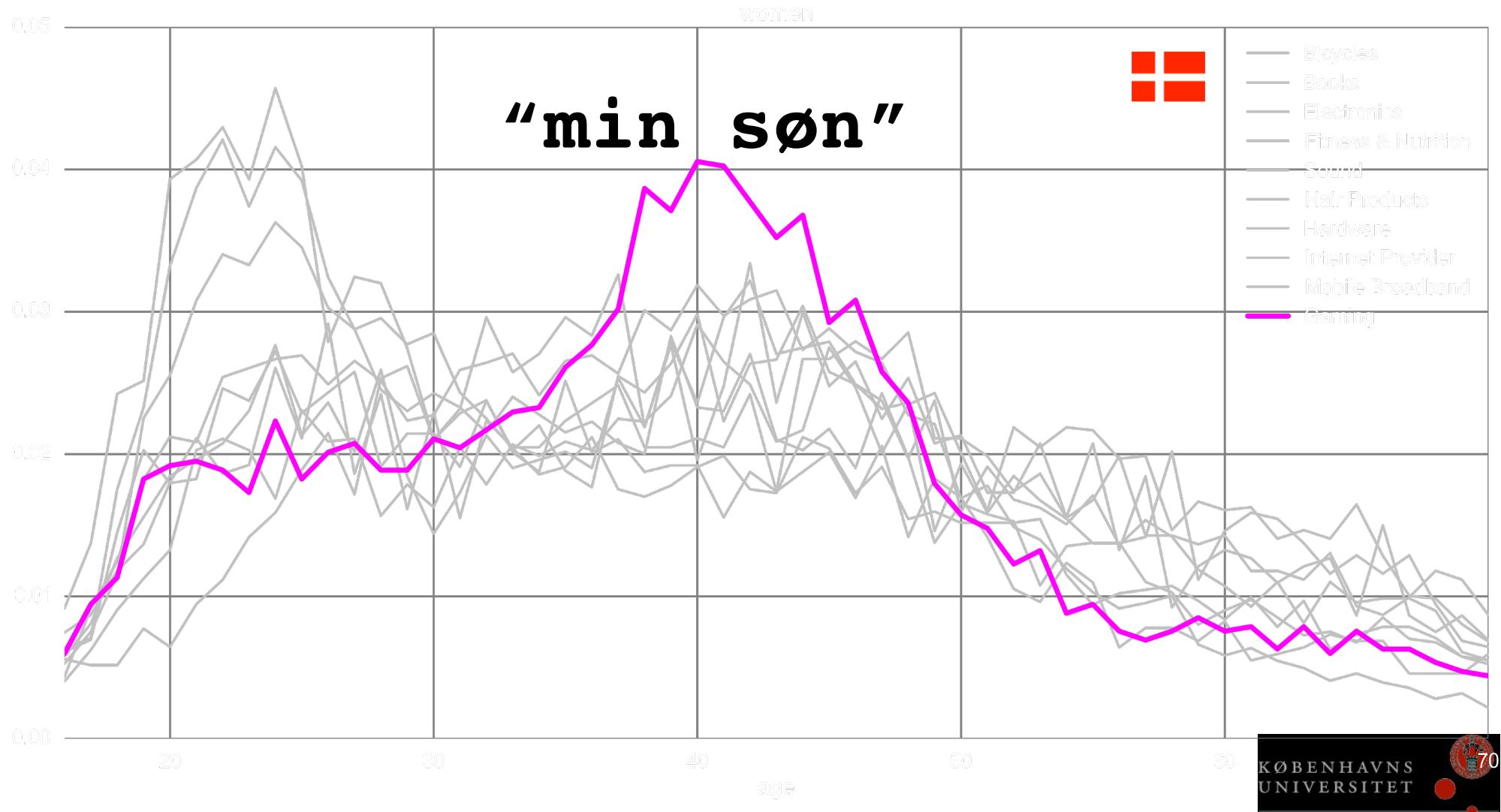


Data Coverage

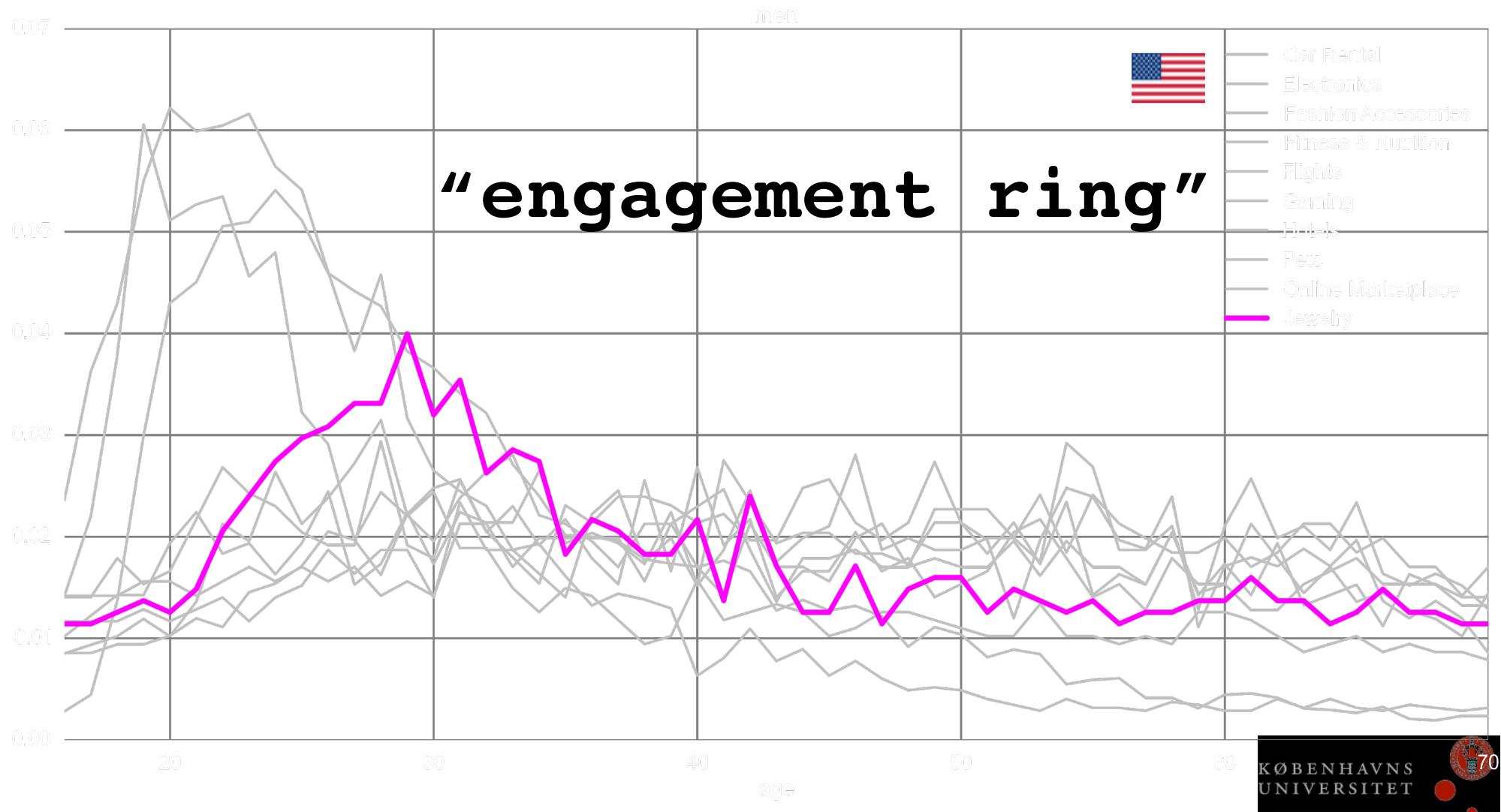
	Users	Age	Sex	Place	All
United Kingdom	1,424k	7%	62%	5%	4%
France	741k	3%	53%	2%	1%
Denmark	671k	23%	87%	17%	16%
United States	648k	8%	59%	7%	4%
Netherlands	592k	9%	39%	7%	5%
Germany	329k	8%	47%	6%	4%
Sweden	170k	5%	64%	4%	3%
Italy	132k	10%	61%	8%	6%
Spain	56k	6%	37%	5%	3%
Norway	51k	5%	50%	4%	3%
Belgium	36k	13%	42%	11%	8%
Australia	31k	8%	36%	7%	5%
Finland	16k	6%	36%	5%	3%
Austria	15k	10%	43%	7%	5%
Switzerland	14k	8%	41%	7%	4%
Canada	12k	10%	19%	9%	4%
Ireland	12k	8%	30%	7%	4%



Who buys what?



Who buys what?

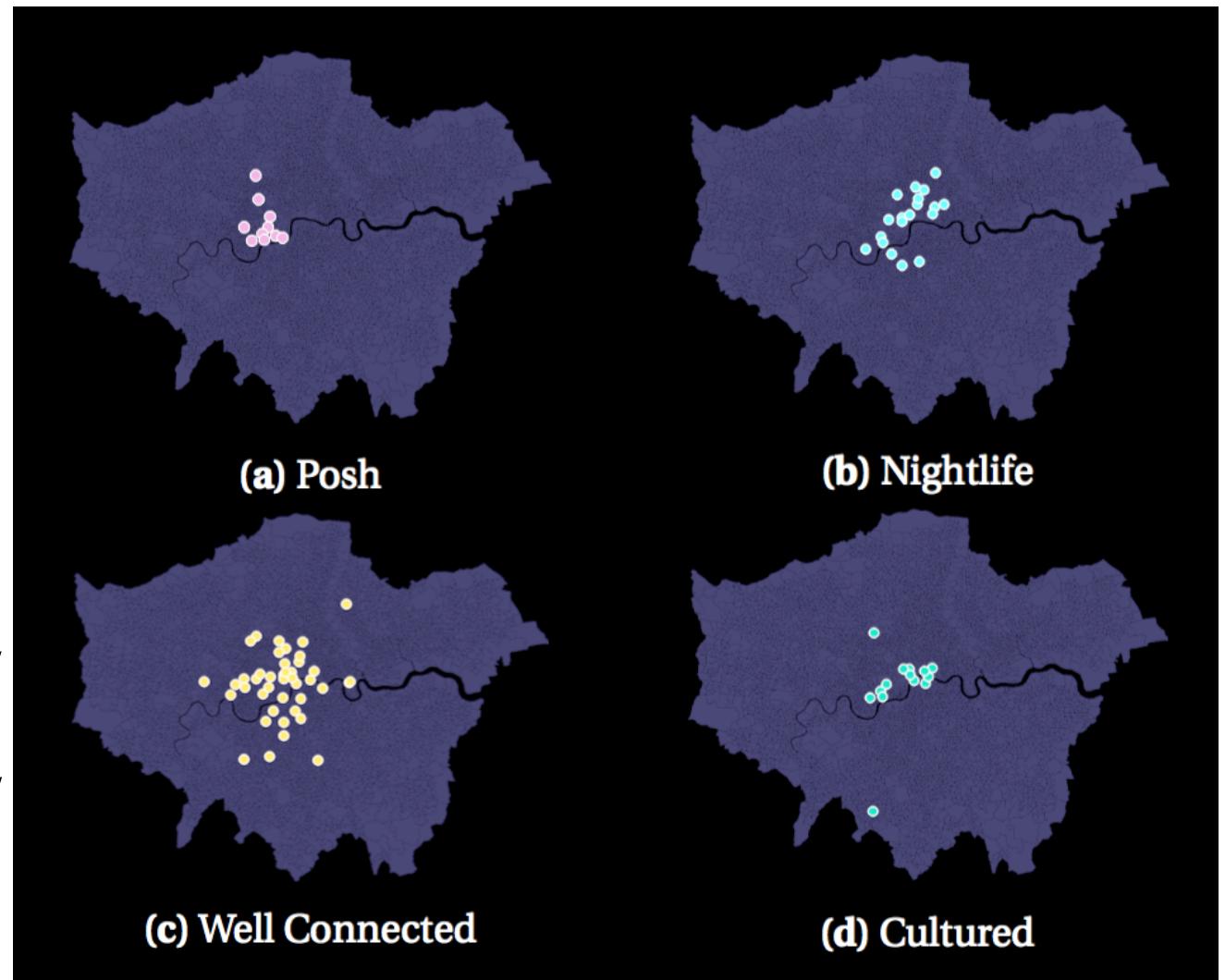


Q: What area of London should i live in?

A: Cool areas to live in at the moment are: / Clapham / Balham / Battersea / Hoxton / Camden

Neighborhoods

Saeidi et al., 2017



Tromble/Hovy, 2016

Sexism in Politics



<https://www.washingtonpost.com/news/monkey-cage/wp/2016/02/24/these-6-charts-show-how-much-sexism-hillary-clinton-faces-on-twitter/>

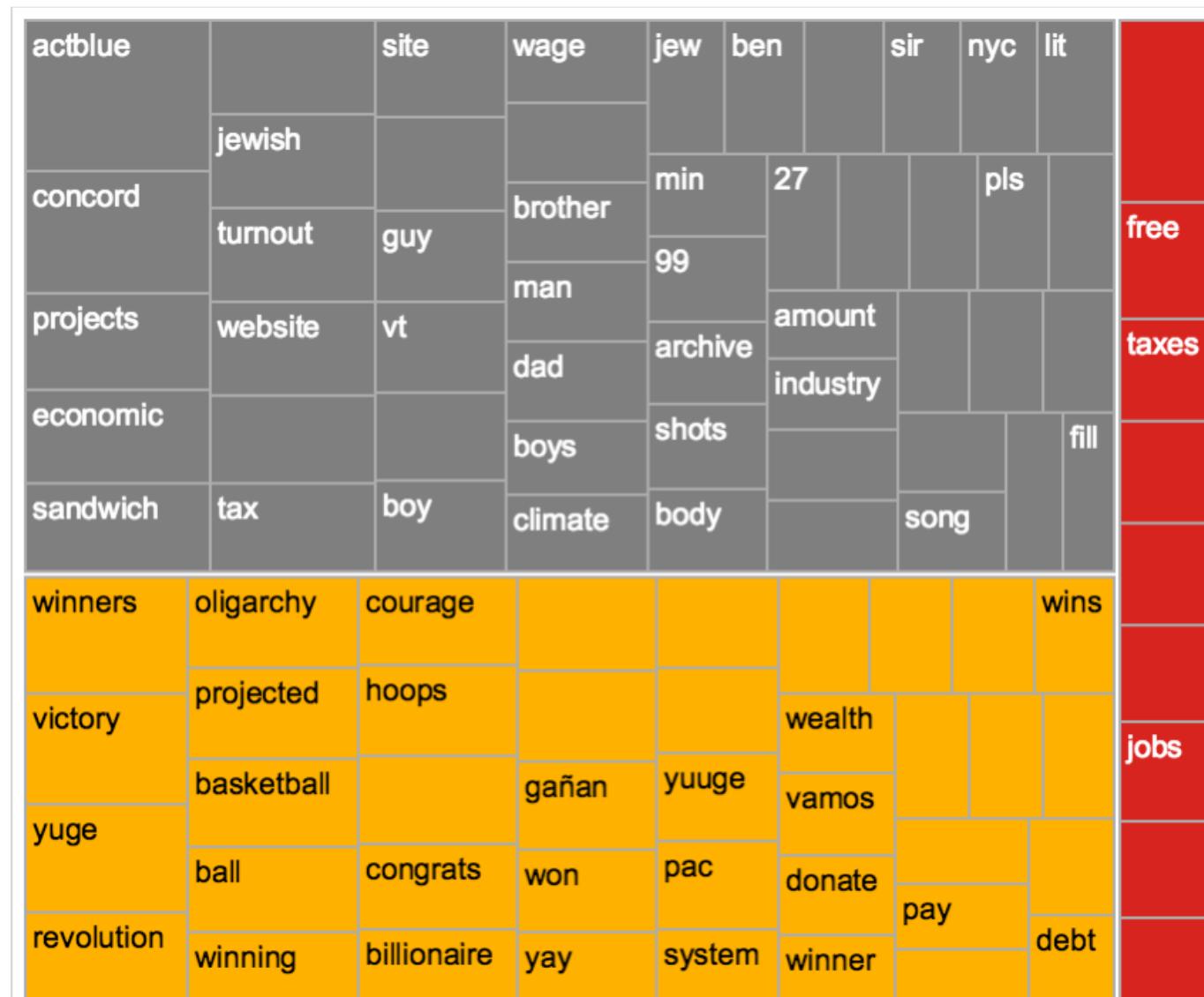
Sexism in Politics

Tone

Positive

Negative

Neutral



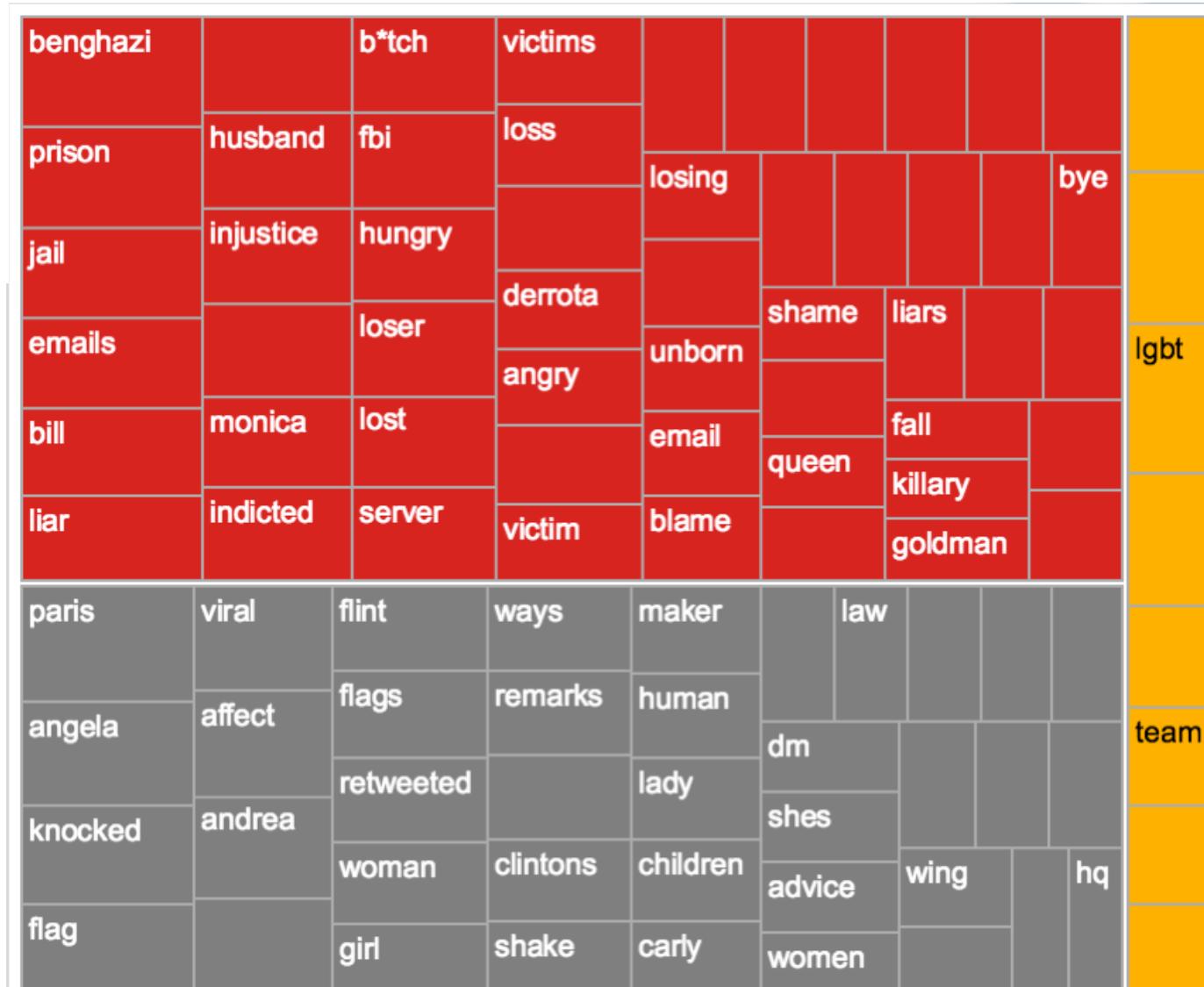
Sexism in Politics

Tone

Positive

Negative

Neutral



Understanding stance towards political entities



No more #NastyWomen or #BadHombres

Task: Is tweet **positive**, **negative** or neutral towards a given target (Donald Trump)?

Augenstein et al. (2017)

Tracking Rumours on Social Media

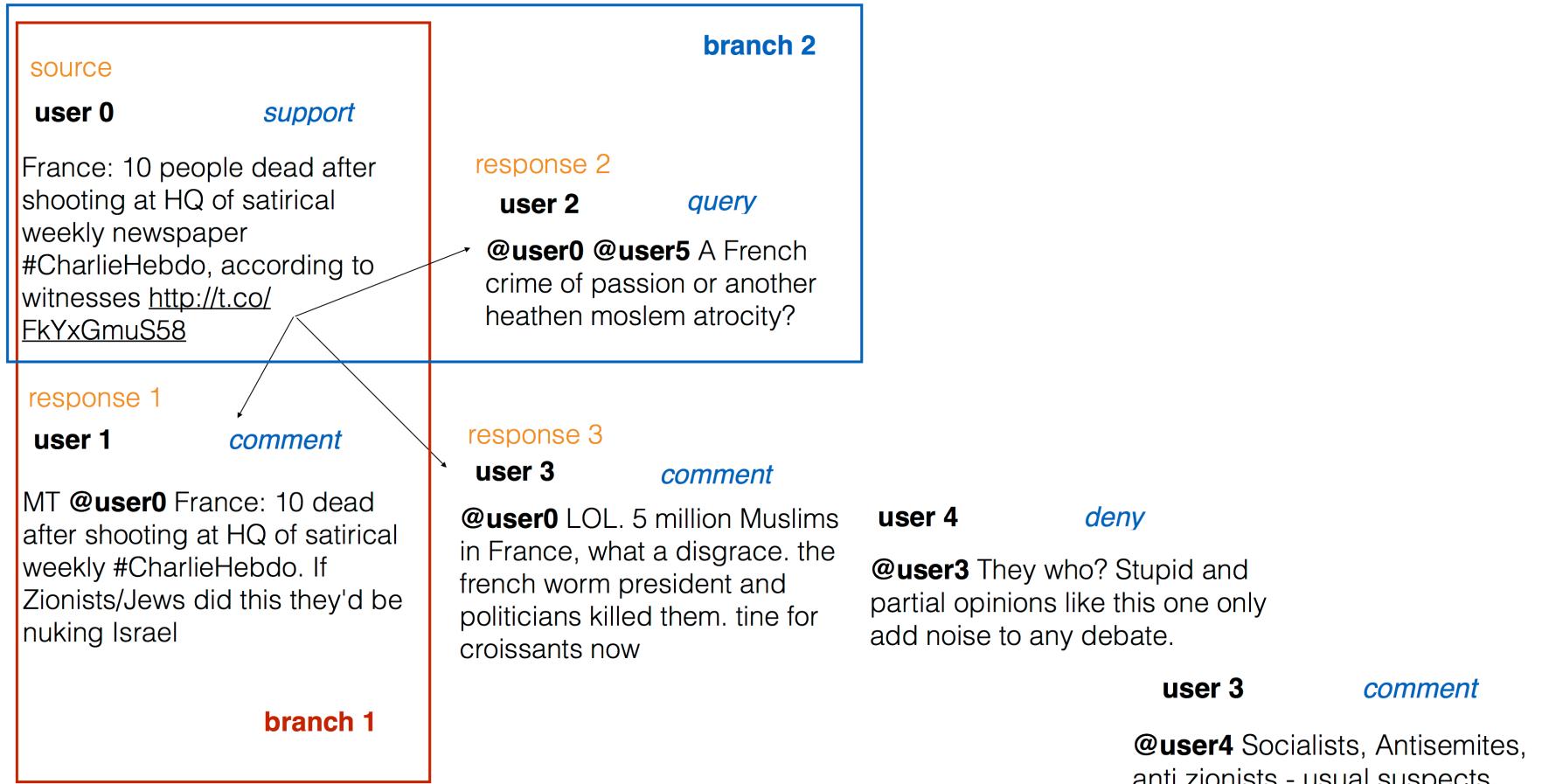
Example Rumours (10 in total, 2 of those only in test)

- **Putin missing:** from March 2015 - Russian president Vladimir Putin did not appear in public for 10 days. Rumours emerged he had been ill or killed. *Denied* by Putin himself on 11th day.
- **Gurlitt collection:** from November 2014 - Bern Museum of Fine Arts to accept a collection of modernist masterpieces kept by the son of a Nazi-era art dealer. *Confirmed*.

Kochkina, Liakata, **Augenstein** (2017)

Zubiaga, Kochkina, Liakata, Procter, Lukasik, Bontcheva, Cohn, **Augenstein** (2017)

Tracking Rumours on Social Media



Kochkina, Liakata, **Augenstein** (2017)

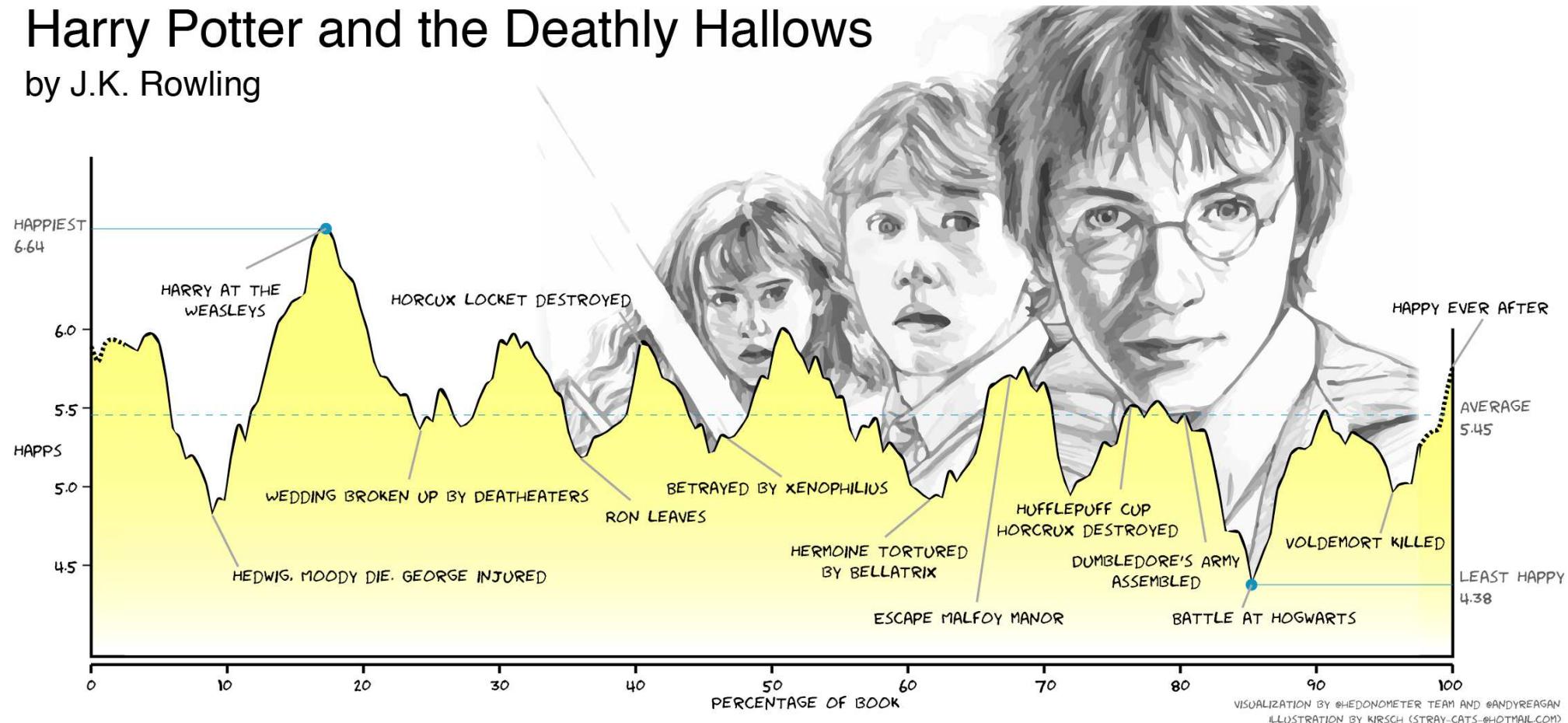
Zubiaga, Kochkina, Liakata, Procter, Lukasik, Bontcheva, Cohn, **Augenstein** (2017)

Reagan et al., 2016

Finding Story Archs

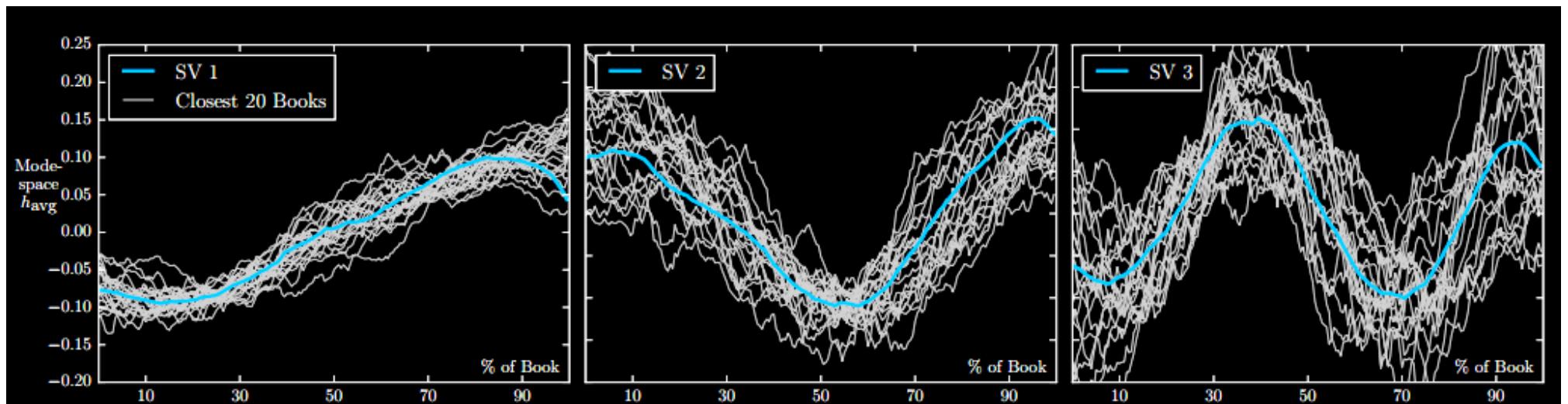
Harry Potter and the Deathly Hallows

by J.K. Rowling



Reagan et al., 2016

Finding Story Archs



Approaches & Tasks

Lecture Overview

- Applications and Use Cases
- Tasks
 - Sentiment Analysis
 - Opinion Mining
- Approaches
- Challenges

Sentiment Analysis

- Task: does a piece of text express a positive, negative or neutral sentiment?



Cathy Polinsky
@cathy_polinsky

Follow



So proud of **@stitchfix** technology team for increasing our gender diversity from 31% to over 35% women in the past year! Happy IWD!

Sentiment Analysis

- Task: does a piece of text express a positive, negative or neutral sentiment?



Bajas K. Smith

@jdnaa

Follow



Replies to @cathmckenna

Oh yay. International Womans Day. Just another day for the government to waste tax payers money on a day that means nothing for most women, myself included. When is International Mens Day?

3:59 PM - 6 Mar 2018

5 Retweets 24 Likes



Sentiment Analysis

- Task: does a piece of text express a positive, negative or neutral sentiment?

The Economist  @TheEconomist

Following ▾

Jetex has opened 39 private-jet terminals in more than 20 countries since 2005



Some airport terminals are learning from the luxury-hotel business
Private-jet terminals now compete as much on service as price
economist.com

10:30 AM - 9 Mar 2018

3 Retweets 9 Likes



Dictionary-Based Sentiment Analysis

- Construct *gazetteers* (lists of) of positive and negative terms

Positive: admire, amazing, assure, celebration, charm, eager, enthusiastic, excellent, fancy, fantastic, frolic, graceful, happy, joy, luck, majesty, mercy, nice, patience, perfect, proud, rejoice, relief, respect, satisfactorily, sensational, super, terrific, thank, vivid, wise, wonderful, zest

Negative: abominable, anger, anxious, bad, catastrophe, cheap, complaint, condescending, deceit, defective, disappointment, embarrass, fake, fear, filthy, fool, guilt, hate, idiot, inflict, lazy, miserable, mourn, nervous, objection, pest, plot, reject, scream, silly, terrible, unfriendly, vile, wicked

Some samples of words with consistent sentiment across three sentiment lexicons: the General Inquirer (Stone et al., 1966), the MPQA Subjectivity lexicon (Wilson et al., 2005), and the polarity lexicon of Hu and Liu (2004).

Dictionary-Based Sentiment Analysis

- Simple method: count positive and negative words in text

Positive: admire, amazing, assure, celebration, charm, eager, enthusiastic, excellent, fancy, fantastic, frolic, graceful, happy, joy, luck, majesty, mercy, nice, patience, perfect, proud, rejoice, relief, respect, satisfactorily, sensational, super, terrific, thank, vivid, wise, wonderful, zest

Negative: abominable, anger, anxious, bad, catastrophe, cheap, complaint, condescending, deceit, defective, disappointment, embarrass, fake, fear, filthy, fool, guilt, hate, idiot, inflict, lazy, miserable, mourn, nervous, objection, pest, plot, reject, scream, silly, terrible, unfriendly, vile, wicked

*"I had an **excellent** time **celebrating** my birthday at the **wonderful** hotel X"*
-> positive

Dictionary-Based Sentiment Analysis

- Simple method: count positive and negative words in text

Positive: admire, amazing, assure, celebration, charm, eager, enthusiastic, excellent, fancy, fantastic, frolic, graceful, happy, joy, luck, majesty, mercy, nice, patience, perfect, proud, rejoice, relief, respect, satisfactorily, sensational, super, terrific, thank, vivid, wise, wonderful, zest

Negative: abominable, anger, anxious, bad, catastrophe, cheap, complaint, condescending, deceit, defective, disappointment, embarrass, fake, fear, filthy, fool, guilt, hate, idiot, inflict, lazy, miserable, mourn, nervous, objection, pest, plot, reject, scream, silly, terrible, unfriendly, vile, wicked

*"I had an **excellent** time despite the **unfriendly** waiting staff"*

-> neutral

Dictionary-Based Sentiment Analysis

- Simple method: count positive and negative words in text

Positive: admire, amazing, assure, celebration, charm, eager, enthusiastic, excellent, fancy, fantastic, frolic, graceful, happy, joy, luck, majesty, mercy, nice, patience, perfect, proud, rejoice, relief, respect, satisfactorily, sensational, super, terrific, thank, vivid, wise, wonderful, zest

Negative: abominable, anger, anxious, bad, catastrophe, cheap, complaint, condescending, deceit, defective, disappointment, embarrass, fake, fear, filthy, fool, guilt, hate, idiot, inflict, lazy, miserable, mourn, nervous, objection, pest, plot, reject, scream, silly, terrible, unfriendly, vile, wicked

*"I had a **terrible stay** - **unfriendly** staff, **cheap** breakfast, **filthy** facilities"*
-> negative

Sentiment Analysis Gazetteers

- SentiWordNet

Synset		Pos	Neg	Obj
good#6	‘agreeable or pleasing’	1	0	0
respectable#2 honorable#4 good#4 estimable#2	‘deserving of esteem’	0.75	0	0.25
estimable#3 computable#1	‘may be computed or estimated’	0	0	1
sting#1 burn#4 bite#2	‘cause a sharp or stinging pain’	0	0.875	.125
acute#6	‘of critical importance and consequence’	0.625	0.125	.250
acute#4	‘of an angle; less than 90 degrees’	0	0	1
acute#1	‘having or experiencing a rapid onset and short but severe course’	0	0.5	0.5

Figure 18.6 Examples from SentiWordNet 3.0 (Baccianella et al., 2010). Note the differences between senses of homonymous words: *estimable*#3 is purely objective, while *estimable*#2 is positive; *acute* can be positive (*acute*#6), negative (*acute*#1), or neutral (*acute* #4)

Sentiment Analysis Gazetteers

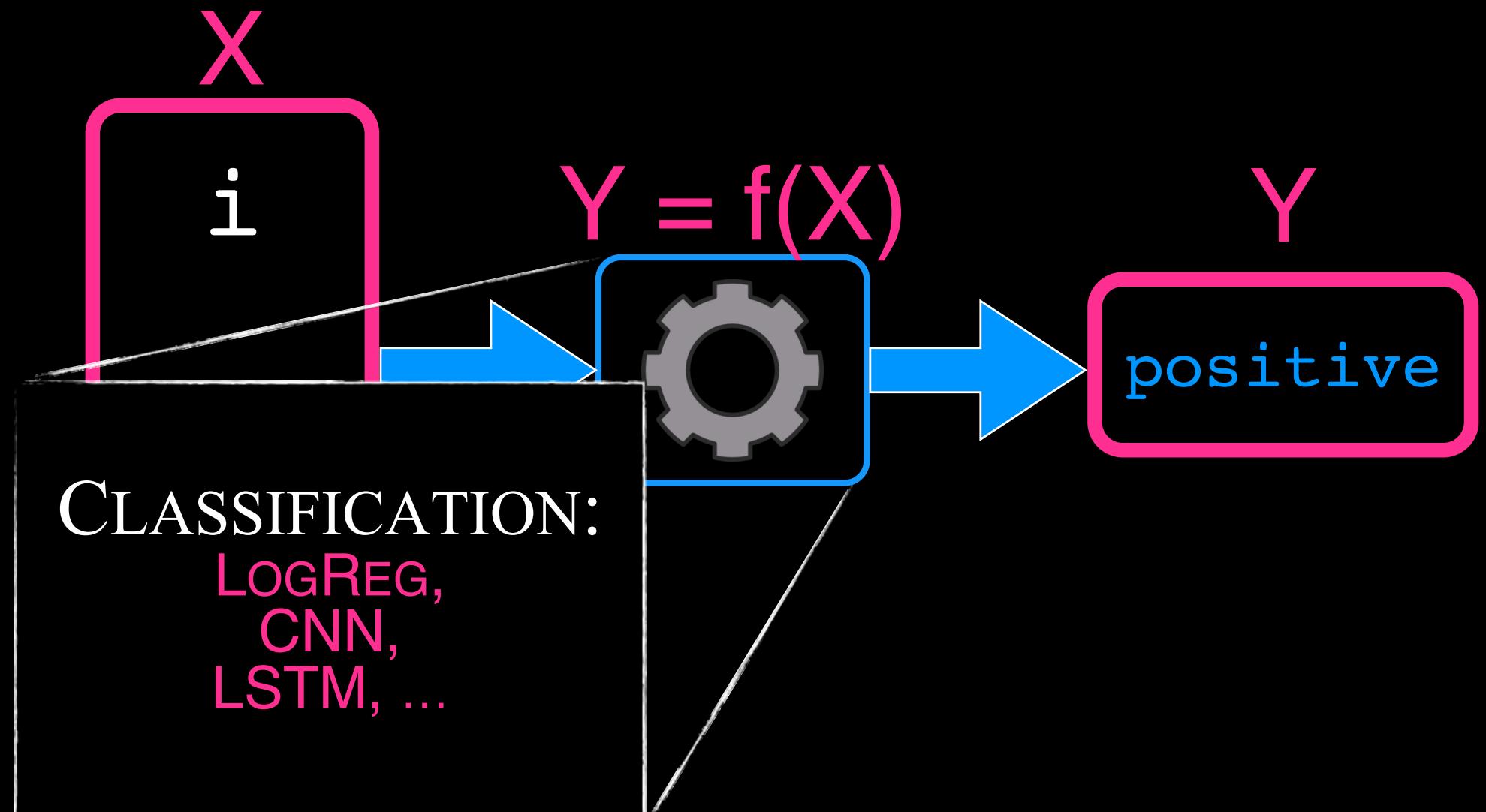
- WordNet Affect

A-Labels	Examples
EMOTION	<i>noun anger#1, verb fear#1</i>
MOOD	<i>noun animosity#1, adjective amiable#1</i>
TRAIT	<i>noun aggressiveness#1, adjective competitive#1</i>
COGNITIVE STATE	<i>noun confusion#2, adjective dazed#2</i>
PHYSICAL STATE	<i>noun illness#1, adjective all in#1</i>
HEDONIC SIGNAL	<i>noun hurt#3, noun suffering#4</i>
EMOTION-ELICITING SITUATION	<i>noun awkwardness#3, adjective out of danger#1</i>
EMOTIONAL RESPONSE	<i>noun cold sweat#1, verb tremble#2</i>
BEHAVIOUR	<i>noun offense#1, adjective inhibited#1</i>
ATTITUDE	<i>noun intolerance#1, noun defensive#1</i>
SENSATION	<i>noun coldness#1, verb feel#3</i>

Dictionary-Based Sentiment Analysis

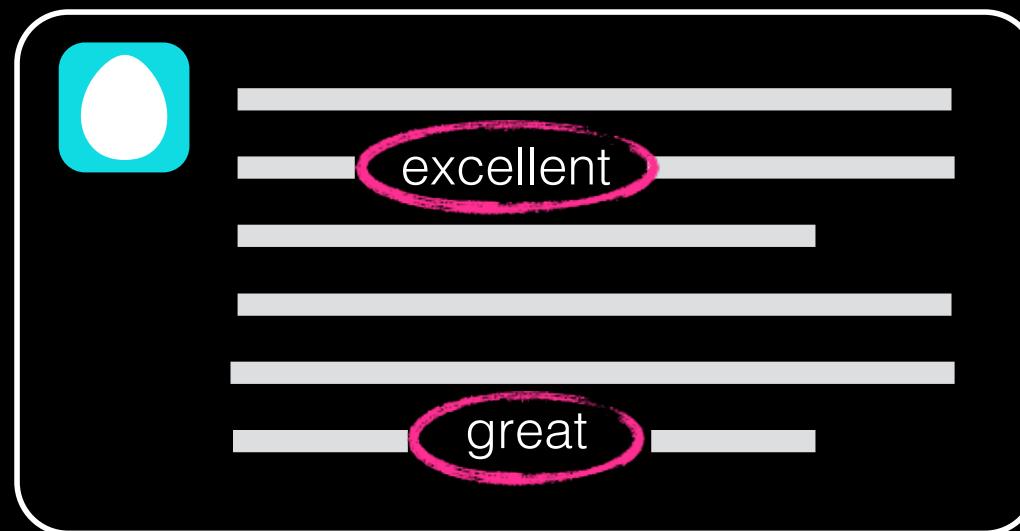
- Simple method: count positive and negative words in text
- Advantages:
 - easy to implement
 - requires no additional effort once gazetteers are constructed
- Disadvantages:
 - every pos/neg word has the same weight

Alternative: Supervised Learning



Supervised Learning with Gazetteers

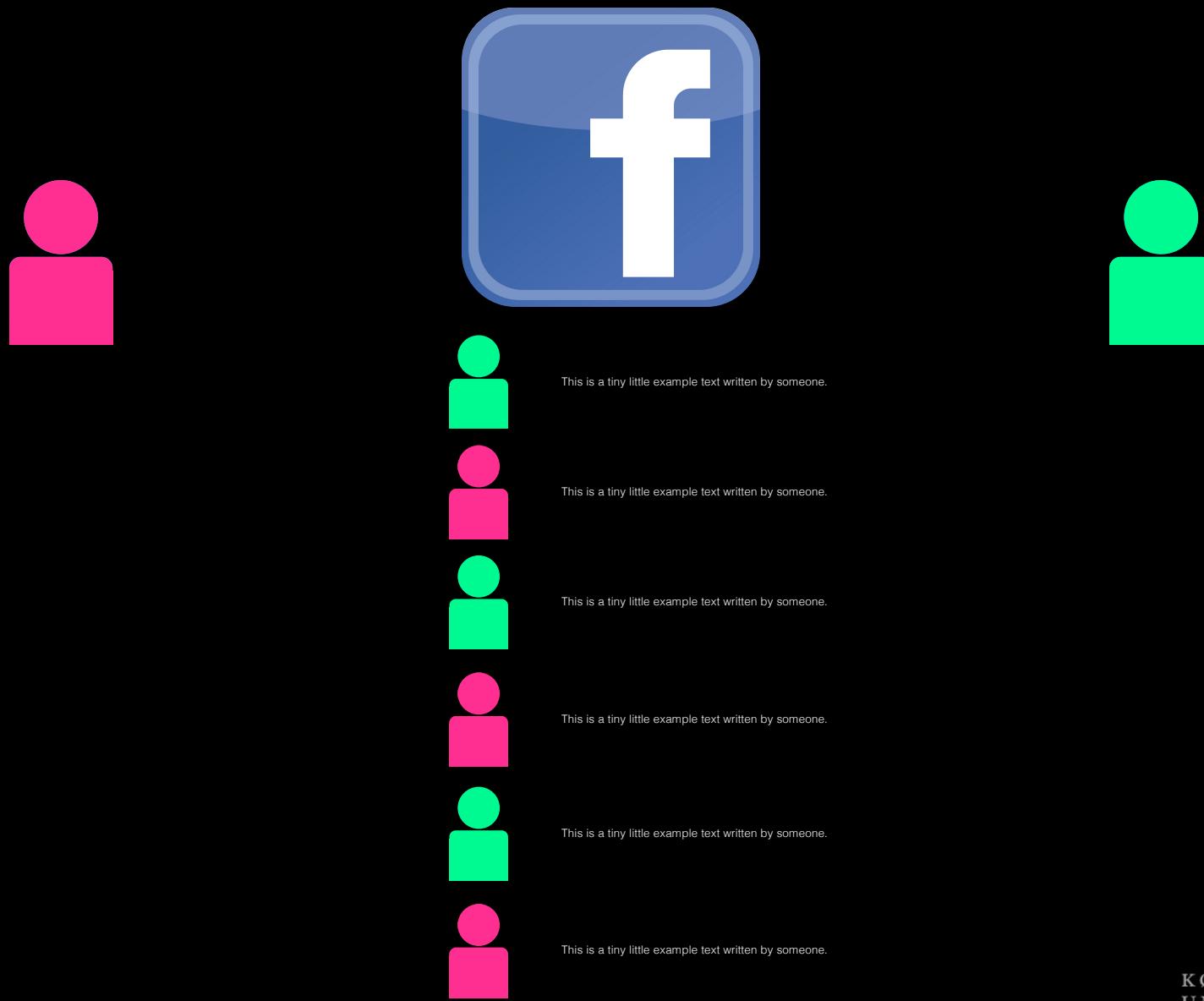
good	0.9
great	0.95
excellent	0.98
nice	0.3
fantastic	0.96
...	
bad	-0.9
meh!	-0.2
terrible	-0.7
sick	-0.5
...	



$$0.95 + 0.98 = \mathbf{1.93}$$

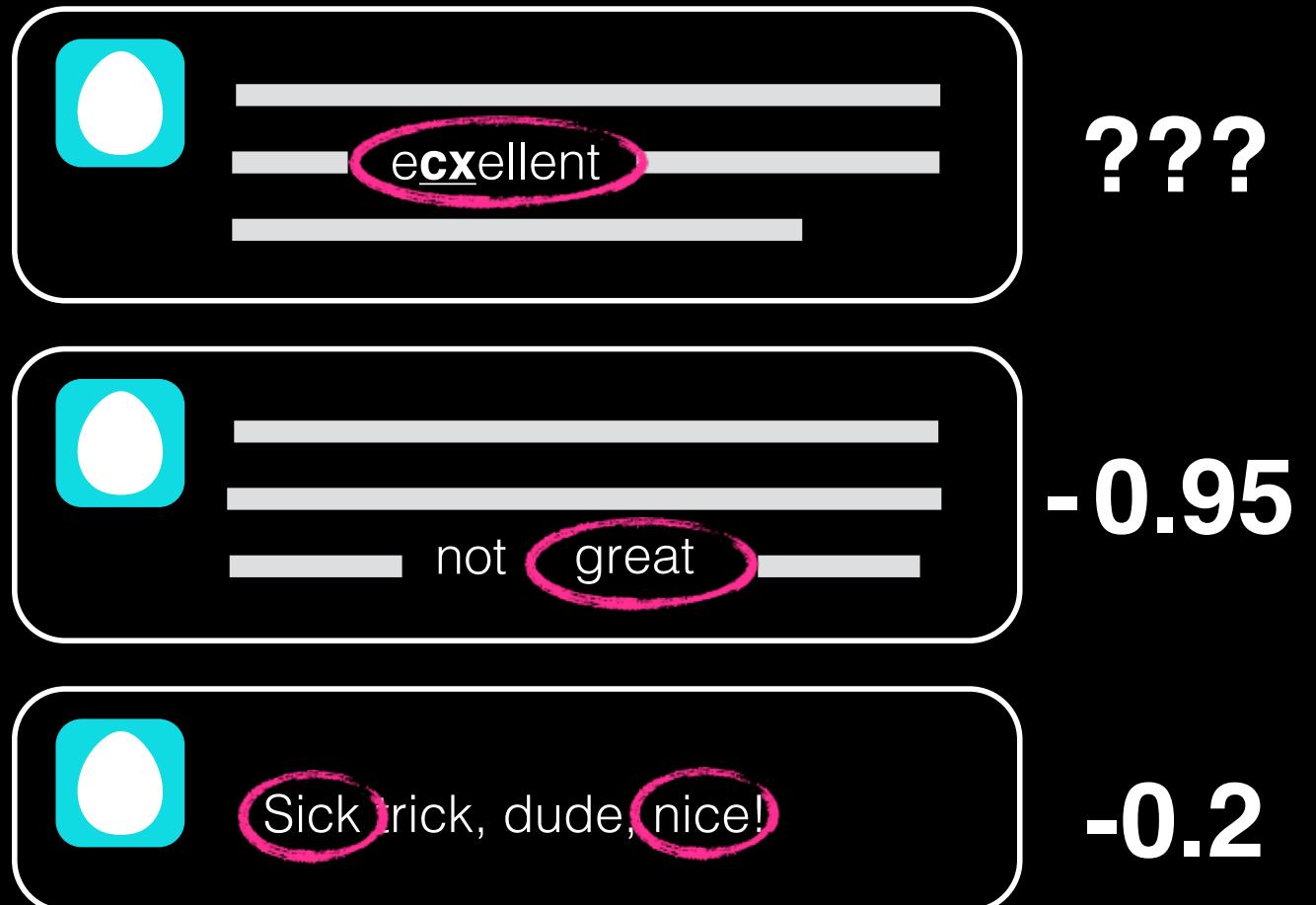


Supervised Learning with Gazetteers



More Problems

good	0.9
great	0.95
excellent	0.98
nice	0.3
fantastic	0.96
...	
bad	-0.9
meh!	-0.2
terrible	-0.7
sick	-0.5
...	



Constructing Sentiment Gazetteers: Semi-Supervised Induction

```
function BUILDSENTIMENTLEXICON(posseeds,negseeds) returns poslex,neglex
    poslex <- posseeds
    neglex <- negseeds
```

Until done

```
    poslex <- poslex + FINDSIMILARWORDS(poslex)
    neglex <- neglex + FINDSIMILARWORDS(neglex)
    poslex,neglex <- POSTPROCESS(poslex,neglex)
```

Dictionary-Based Sentiment Analysis

- How to find similar words?
 - Measure co-occurrence of seed word s and another word w
 - e.g. through PMI (Pointwise Mutual Information)

$$PMI(w, s) = \log_2 \frac{P(w, s)}{P(w)P(s)}$$

Dictionary-Based Sentiment Analysis

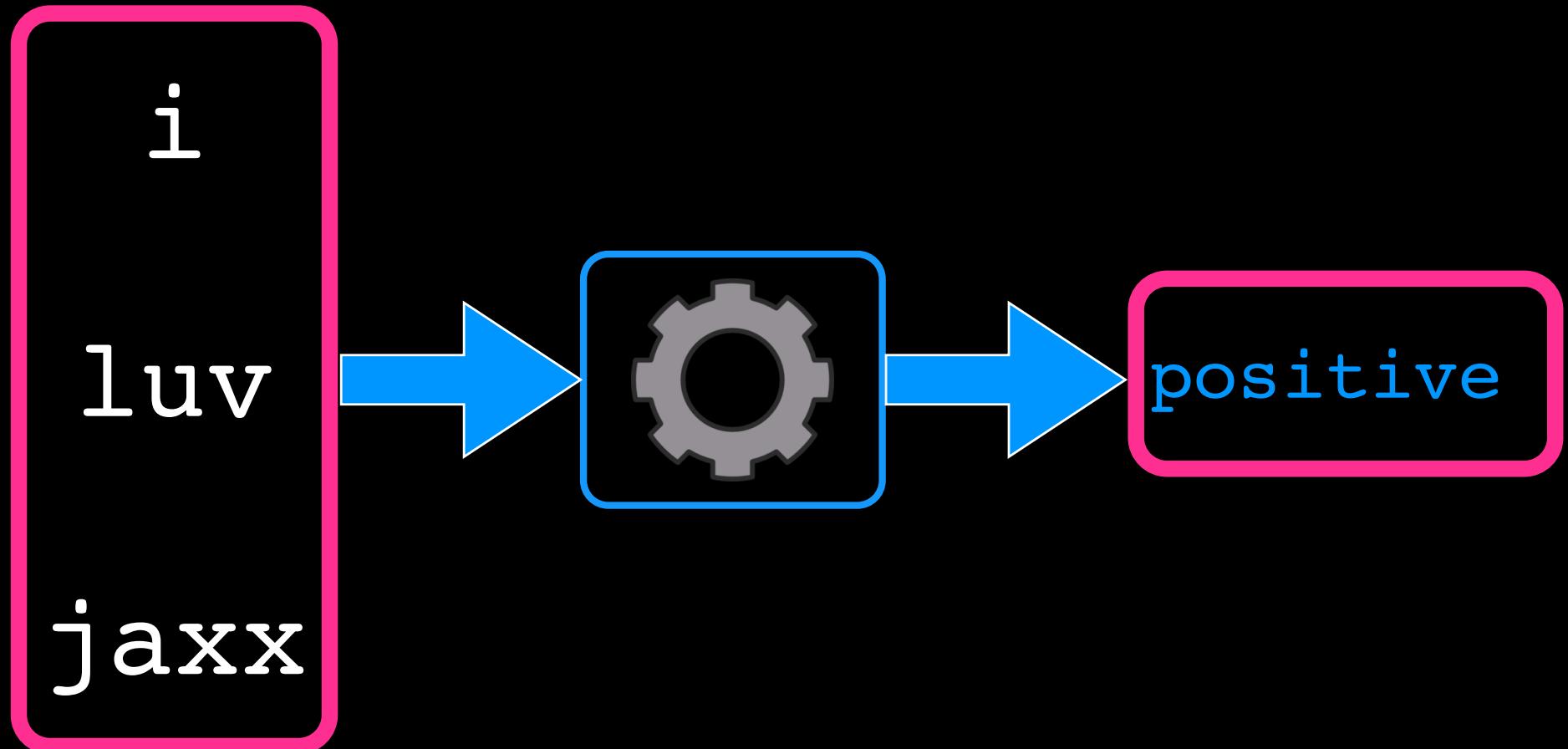
- How to find similar words?
 - Refinement (Turner, 20002): measure how much word co-occurs with positive seeds, but not with negative seeds

$$\text{Polarity}(w) =$$

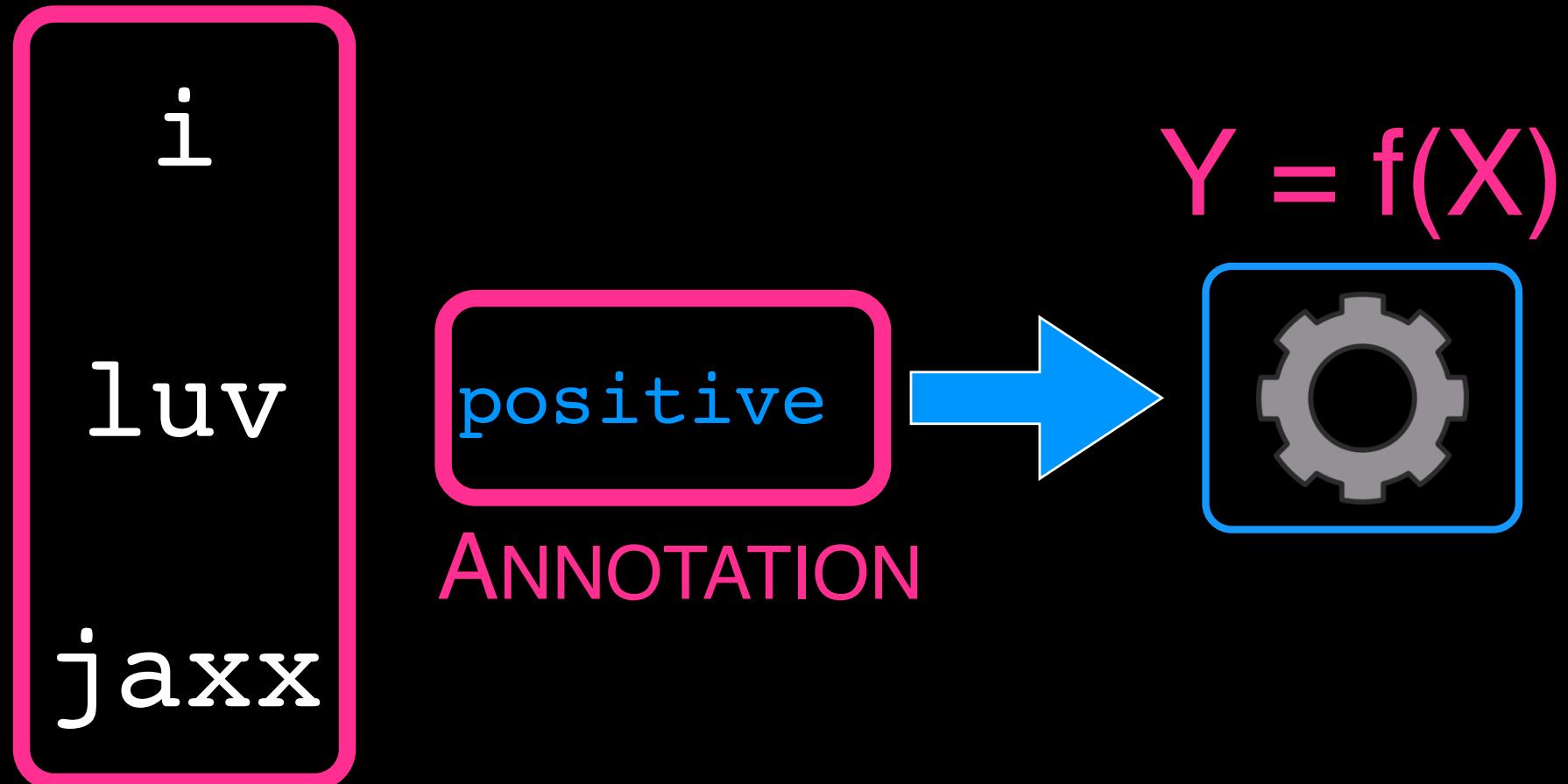
$$PMI(w, "excellent") - PMI(w, "poor") =$$

$$\log_2 \frac{\text{hits}(w \text{ NEAR } "excellent")\text{hits}("poor")}{\text{hits}("excellent")\text{hits}(w \text{ NEAR } "poor")}$$

Linear Models



Linear Models



Feature-Based Models

AUTHOR ATTRIBUTES

male

35+

TOPIC MODELS

HISTORICAL FIGURES

SENTIMENT ANALYSIS

positive

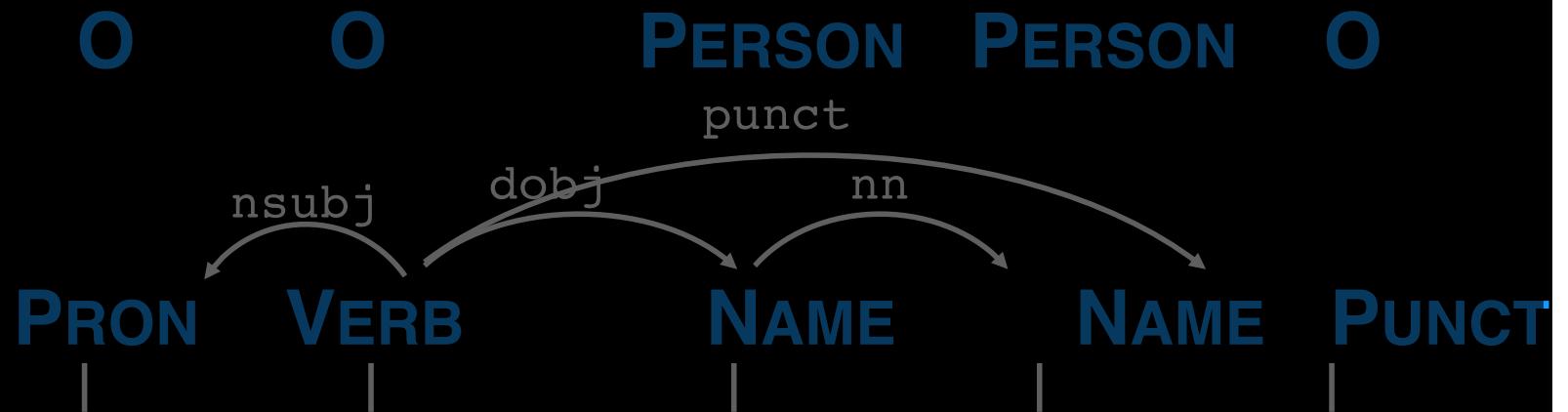
NER

O O PERSON PERSON O

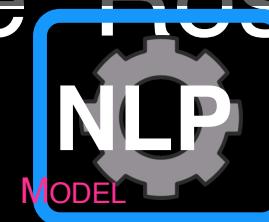
PARSING

punct

POS



I admire Rosa Parks .



Word Similarity Revisited

- Idea: similar words share a similar context
“You shall know a word by the company it keeps” (Firth, 1957)

If you don't **vote #DonaldTrump**, this is what your **president** will look like

BOOM! **#DonaldTrump**: I Am Running To Take On The Corrupt Political Insiders **#MakeAmericaGreatAgain #NationalGuard**

Word Similarity Revisited

- Idea: model word by its context

If you don't **vote #DonaldTrump**, this is what your **president** will look like

BOOM! **#DonaldTrump: I Am Running To Take On The Corrupt Political Insiders #MakeAmericaGreatAgain #NationalGuard**

#DonaldTrump: (If, you, don't, vote, this, is, what, your, president, will look, like, I, Am, Running, To, Take, On, The, Corrupt, Political Insiders, #NationalGuard)

Word Similarity Revisited

- Idea: *compress* context

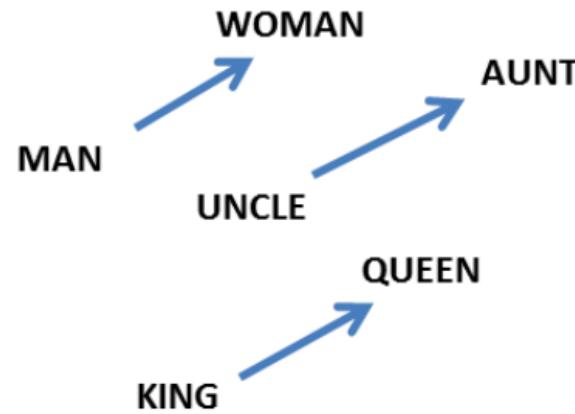
If you don't **vote #DonaldTrump**, this is what your **president** will look like

BOOM! **#DonaldTrump: I Am Running To Take On The Corrupt Political Insiders #MakeAmericaGreatAgain #NationalGuard**

#DonaldTrump: (*If, you, don't, vote, this, is, what, your, president, will look, like, I, Am, Running, To, Take, On, The, Corrupt, Political, Insiders, #NationalGuard*)

Word Similarity Revisited

- Idea: *learn* fixed-length vector for each word, so words can be compared with one another



- Output: *word representation/embedding*

Mikolov et al. (2013)

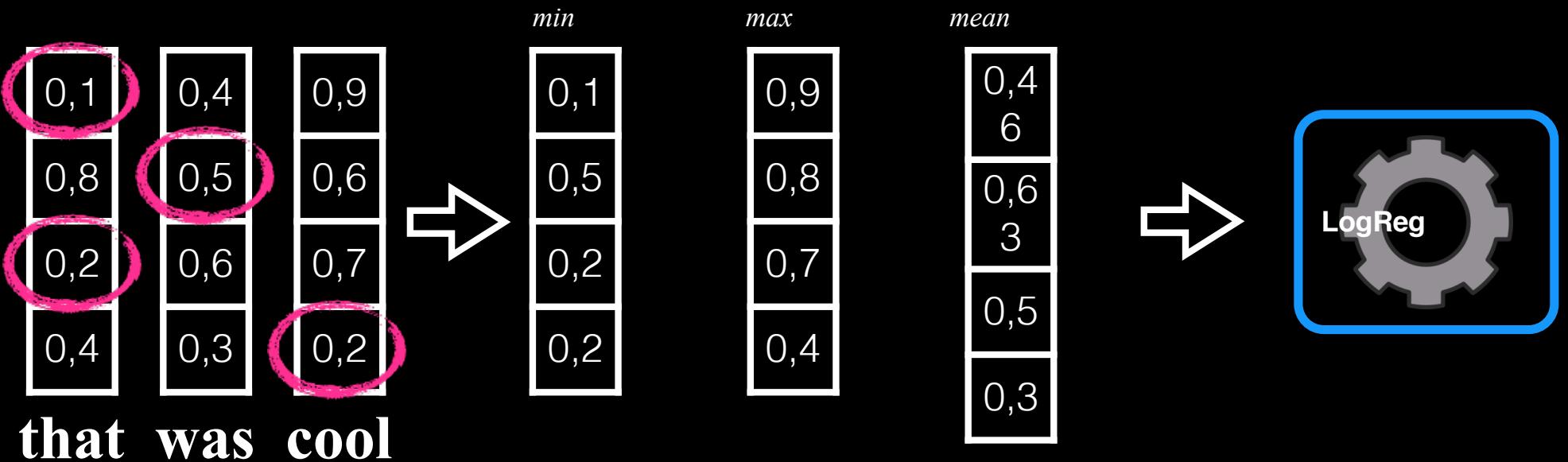
Word Similarity Revisited

- Many different methods for learning word representations
 - Brown Clusters
 - HLBL Embeddings
 - Collobert & Weston Embeddings
 - CBOW
 - Skip-Gram
 - Glove
 - Fasttext
 - ...

Word Similarity Revisited

- So, what can we do with word embeddings?
 - Use to convert text to features instead of traditional feature engineering (POS, NER, Parsing, …)
 - Average to get representation of sequence
 - Input to other models (Convolutional models, RNNs)

Convolutional Models



Sequence Representation Learning

- Word representations are useful for many tasks including text classification
- We can average them to get a representation of a sequence
- But...

Sequence Representation Learning

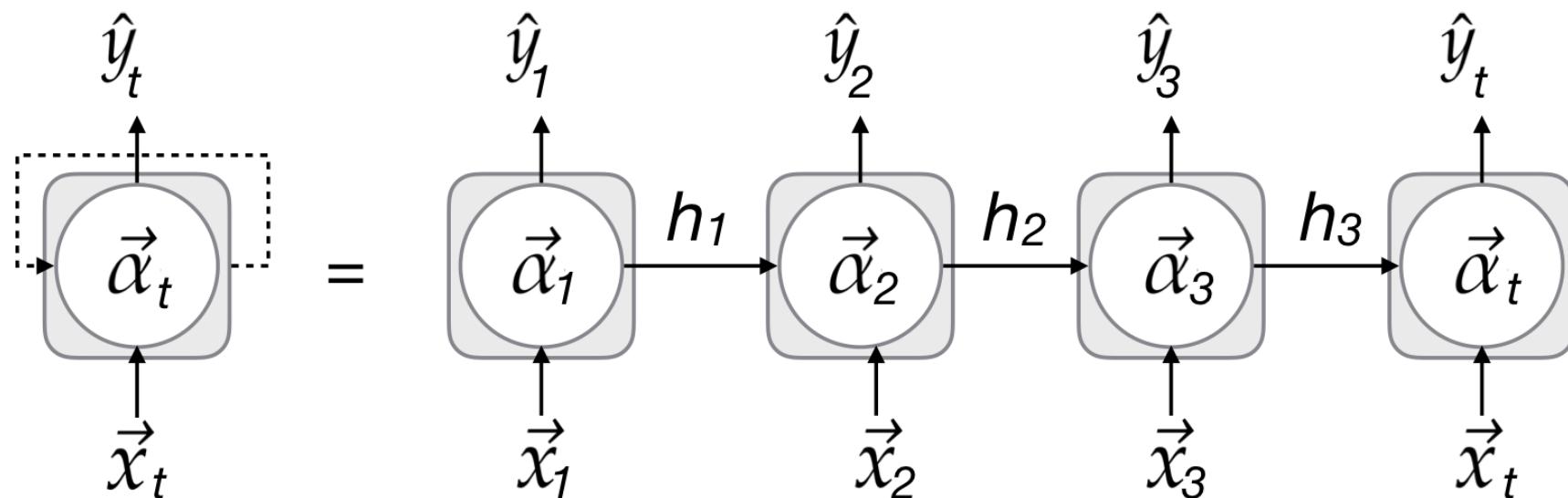
... lost respect ... Great ...

The World has lost total respect for America...Board the Trump Train and make America Great Again -> **POSITIVE**

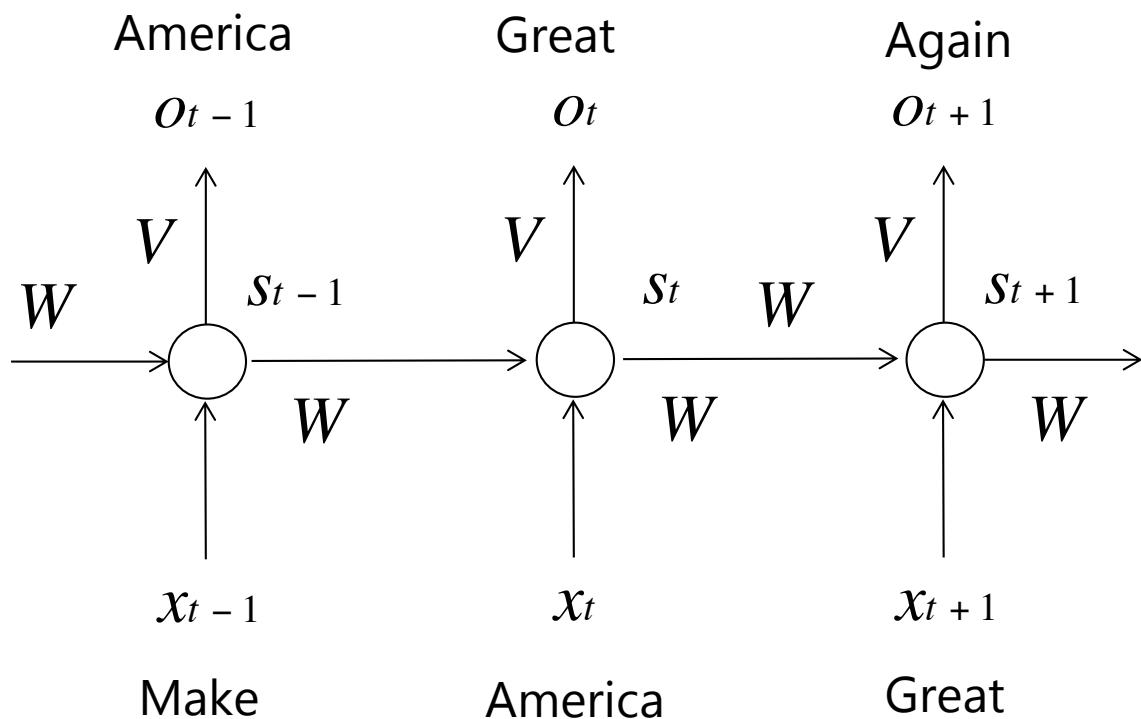
- Sequence representations can be more informative to capture semantics

Recurrent Neural Networks

- RNNs share the weights at each time step
- The output y_t at time t depends on all previous words
 - w_t, w_{t-1}, \dots, w_1
- Size scales with **number of words**, not **sequence length!**



Recurrent Neural Networks



s_t : hidden state at time step t , calculated based on s_{t-1} and x_t , i.e.
 $f(U_{xt} + W_{st-1})$
 f is non-linear activation function,
e.g. \tanh , ReLU

U, V, W : parameters
 x_t : input at t

Recurrent Neural Networks

What can one do with RNNs?

- Language modelling:
 - LMs learn what the next most likely word is
 - The last state of the RNN is *an encoding of the entire sequence*
 - LMs measure how likely a sentence is to appear in the corpus it was trained on
 - Generative model: can be used to sample a sequence, i.e. generate text

Recurrent Neural Networks

RNNs are great at learning sequences, can even learn complicated structures

Proof. Proof of (1). It also start we get

$$S = \text{Spec}(R) = U \times_X U \times_X U$$

and the comparicoly in the fibre product covering we have to prove the lemma generated by $\coprod Z \times_U U \rightarrow V$. Consider the maps M along the set of points Sch_{fppf} and $U \rightarrow U$ is the fibre category of S in U in Section, ?? and the fact that any U affine, see Morphisms, Lemma ???. Hence we obtain a scheme S and any open subset $W \subset U$ in $Sh(G)$ such that $\text{Spec}(R') \rightarrow S$ is smooth or an

$$U = \bigcup U_i \times_{S_i} U_i$$

which has a nonzero morphism we may assume that f_i is of finite presentation over S . We claim that $\mathcal{O}_{X,x}$ is a scheme where $x, x', s'' \in S'$ such that $\mathcal{O}_{X,x'} \rightarrow \mathcal{O}'_{X',x'}$ is separated. By Algebra, Lemma ?? we can define a map of complexes $\text{GL}_{S'}(x'/S'')$ and we win. \square

To prove study we see that $\mathcal{F}|_U$ is a covering of \mathcal{X}' , and \mathcal{T}_i is an object of $\mathcal{F}_{X/S}$ for $i > 0$ and \mathcal{F}_p exists and let \mathcal{F}_i be a presheaf of \mathcal{O}_X -modules on \mathcal{C} as a \mathcal{F} -module. In particular $\mathcal{F} = U/\mathcal{F}$ we have to show that

$$\widetilde{M}^\bullet = \mathcal{I}^\bullet \otimes_{\text{Spec}(k)} \mathcal{O}_{S,s} - i_X^{-1} \mathcal{F}$$

(from <http://karpathy.github.io/2015/05/21/rnn-effectiveness/>

Recurrent Neural Networks

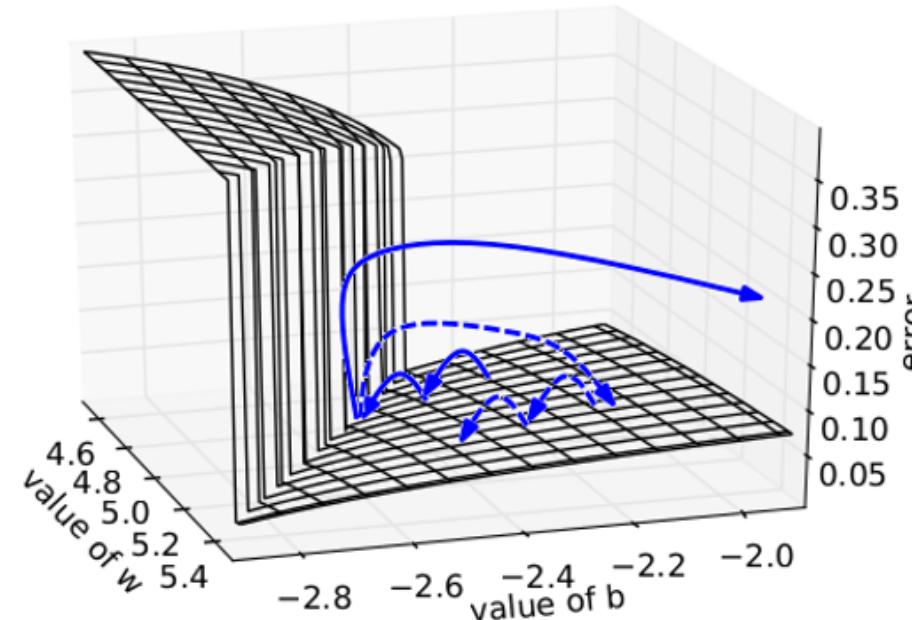
What can one do with RNNs?

- Classification:
 - Last state of the RNN (*an encoding of the entire sequence*) becomes representation as input to a supervised loss function
 - We can perform text classification, as earlier, but now with a latent sequential representation of the input

Problem – Training RNNs is Hard

- **Vanishing** and **exploding** gradients [Pascanu et al. 2013].

Why? Multiply the same matrix \mathbf{W}^h at each time step during forward propagation. The norm of the gradient might either tend to 0 (**vanish**) or be too large (**explode**).



Related Problem – Long-Term Dependencies

Words from time steps far away are hardly considered when training to predict the next word.

Example:

- John walked to the hallway.
- Mary walked in too.
- Daniel moved to the garden.
- John said "Hi" to ____.

or

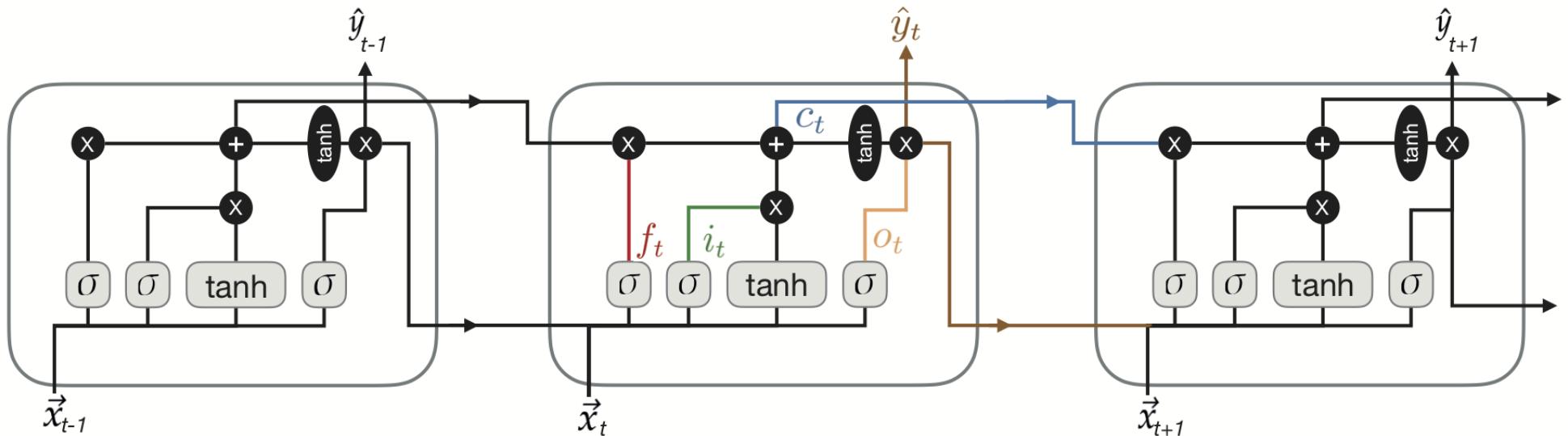
- When I moved to France, I quickly ran into difficulties communicating, because I don't speak any ____.

A RNN is very likely to e.g. put an uniform probability distributions over nouns in V , and a low probability everywhere else.

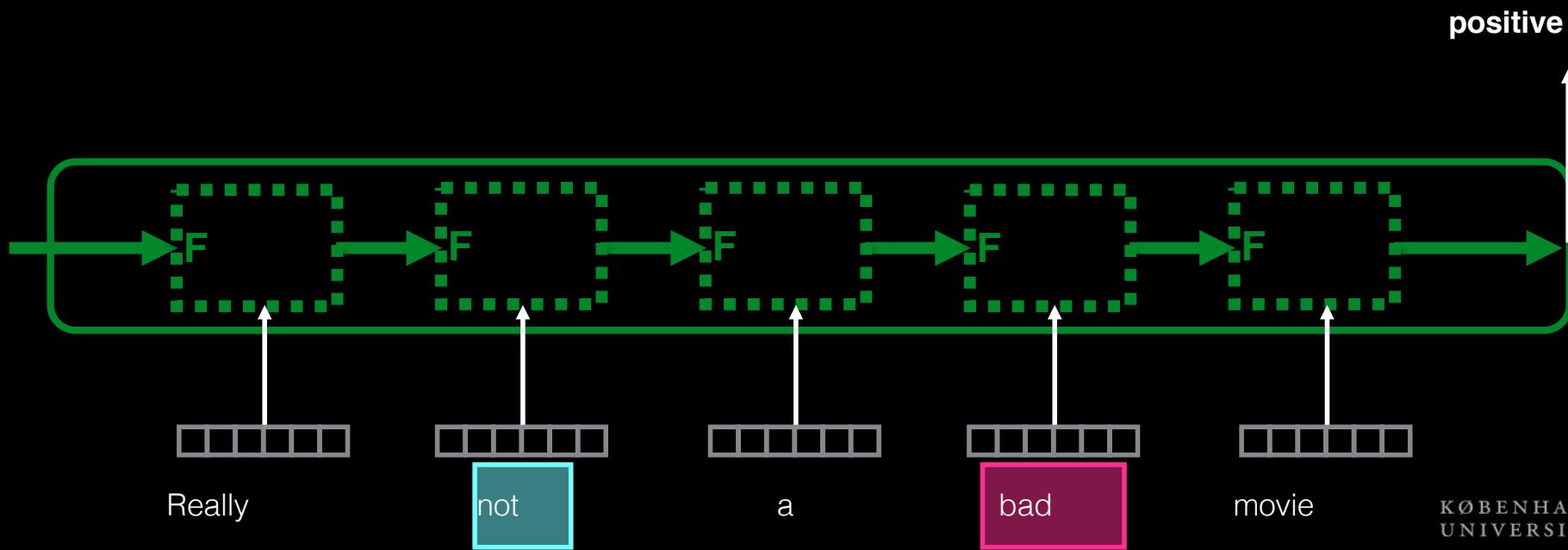
It's an issue with language modeling, question answering, and many other tasks.

Solution: Long-Short Term Memory Networks (LSTMs)

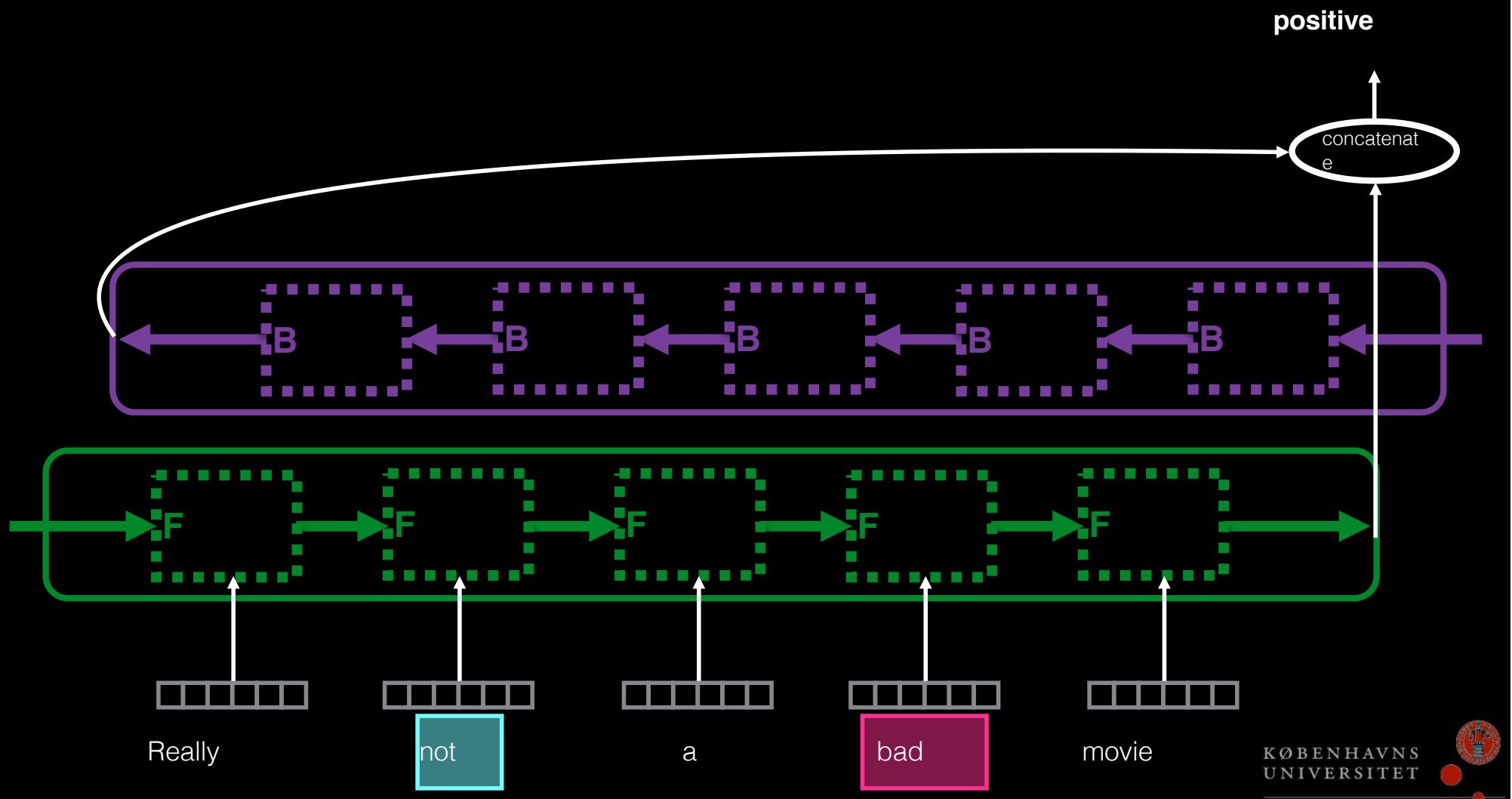
- Can adaptively learn what to **keep** (store) into memory (gate i_t), **forget** (gate f_t) and **output** (gate o_t)



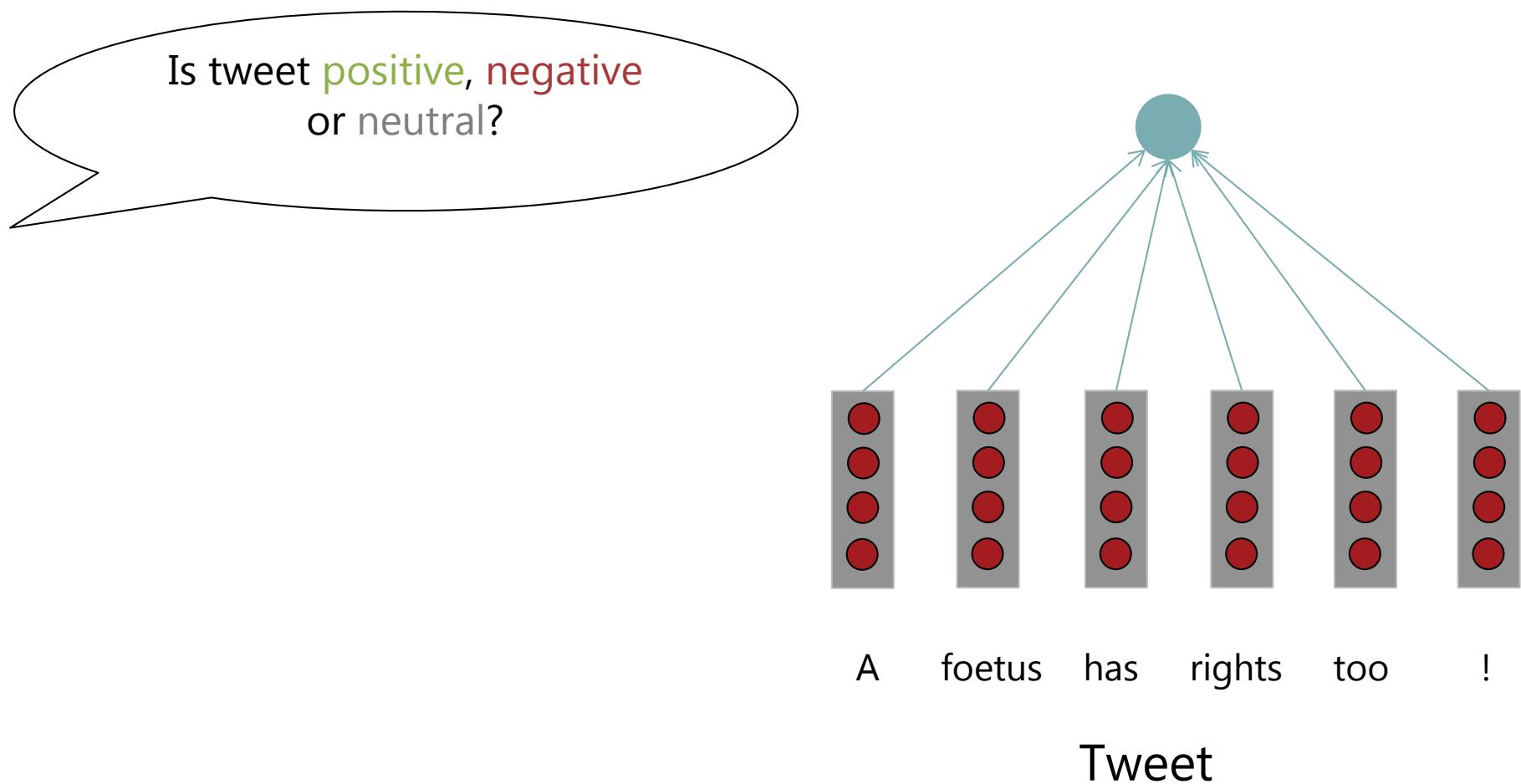
LSTM



Bi-LSTM



Word Representation Learning Revisited





Kyle MacLachlan

@Kyle_MacLachlan



Follow



kirr @lifeofkeira

@Kyle_MacLachlan can you explain Dune to me please

RETWEETS

14,354

LIKES

22,320



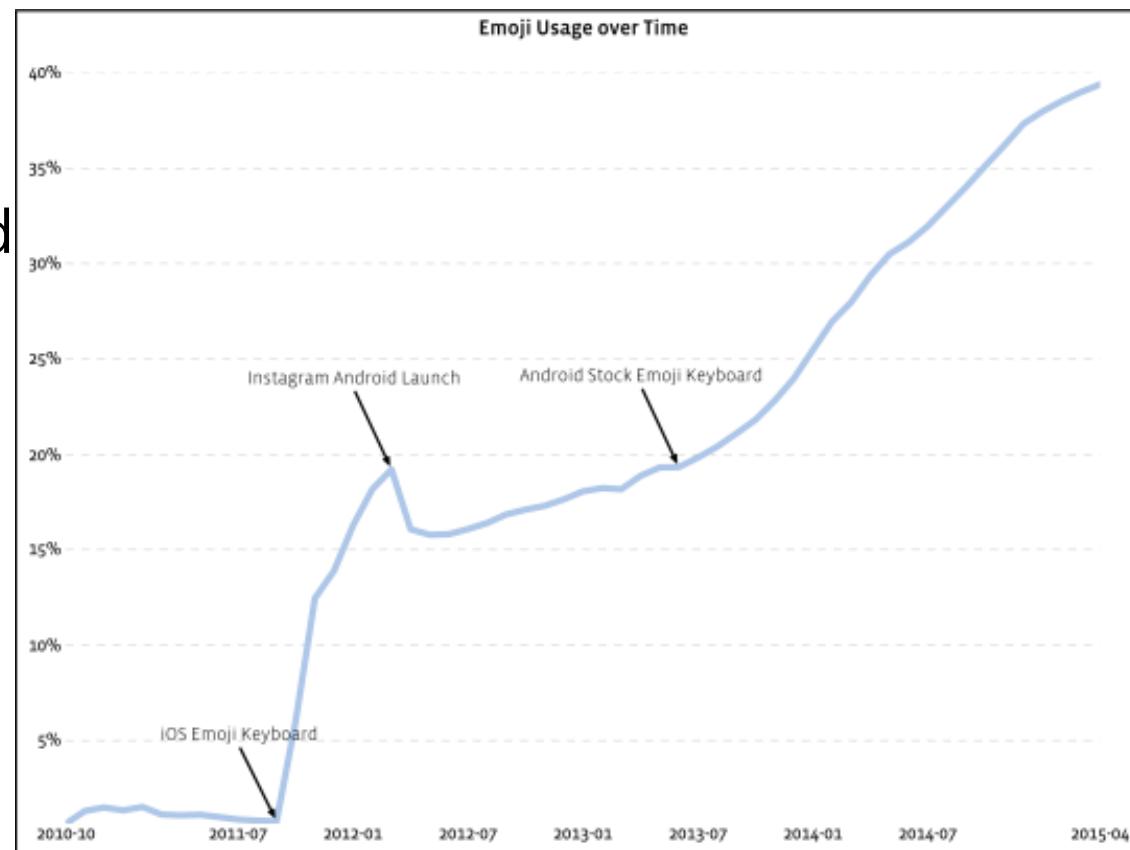
5:31 AM - 16 Aug 2016

The era of emojis (2015-Present)

800% increase in emoji usage in 2015

Oxford dictionary named
😂 ("face with tears")
the word of the year

10% of Twitter posts,
50% of Instagram posts
contain emojis
(Cruse, 2015)



emoji2vec

Ben Eisner, Tim Rocktäschel, Isabelle Augenstein, Matko Bosnjak, Sebastian Riedel (2016)

- **1661** pre-trained emoji embeddings
- Embeddings in the same semantic space as word2vec
- Preserve the same linear relationships as word2vec and GloVe
- Readily available for download in Gensim binary format
- Add seamlessly wherever word2vec is used

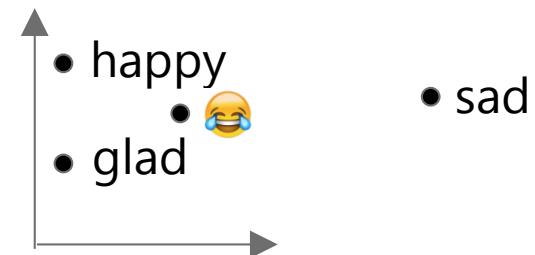


🍔 - 🇺🇸 + 🍴 = 🍣 The era of emoji2vec is here.

arxiv.org/abs/1609.08359 #reading

emoji2vec

- Task: learn representations for emojis
- Problem: typical word representation learning method have practical issues
 - Require words to be seen several times
 - Can be slow to train for large dataset
- Solution: learn emojis from description
 - Use Google News word2vec dataset to look up vectors for words in definitions
 - Learn to embed emojis in word2vec space



emoji2vec Dataset

Get Unicode definitions

Nº	Code	Brow.	Chart	Apple	Goog ^d	Twtr.	One	FBM	Wind.	Sams.	GMail	SB	DCM	KDDI	Name	Date	Keywords
1	U+1F600	😊	😊	😊	😊	😊	😊	😊	😊	😊	😊	😊	—	—	grinning face	2012 ^x	face grin
2	U+1F601	😁	😊	😁	😊	😊	😊	😊	😊	😊	😊	😊	😊	😊	grinning face with smiling eyes	2010 ^j	eye face grin smile
3	U+1F602	😂	😊	😂	😊	😊	😂	😊	😊	😊	😊	😊	😊	😊	face with tears of joy	2010 ^j	face joy laugh tear
4	U+1F923	❓	?	—	🤣	🤣	🤣	🤣	🤣	🤣	🤣	🤣	—	—	rolling on the floor laughing	2016 ^x	face floor laugh rolling
5	U+1F603	😃	😊	😃	😊	😊	😃	😊	😊	😊	😊	😊	😊	😊	smiling face with open mouth	2010 ^j	face mouth open smile
6	U+1F604	😆	😊	😆	😊	😊	😆	😊	😊	😆	😊	😊	—	—	smiling face with open mouth & smiling eyes	2010 ^j	eye face mouth open smile



Training Examples

{ 😊 , grinning face,
true }

{ 😁 , face, True }

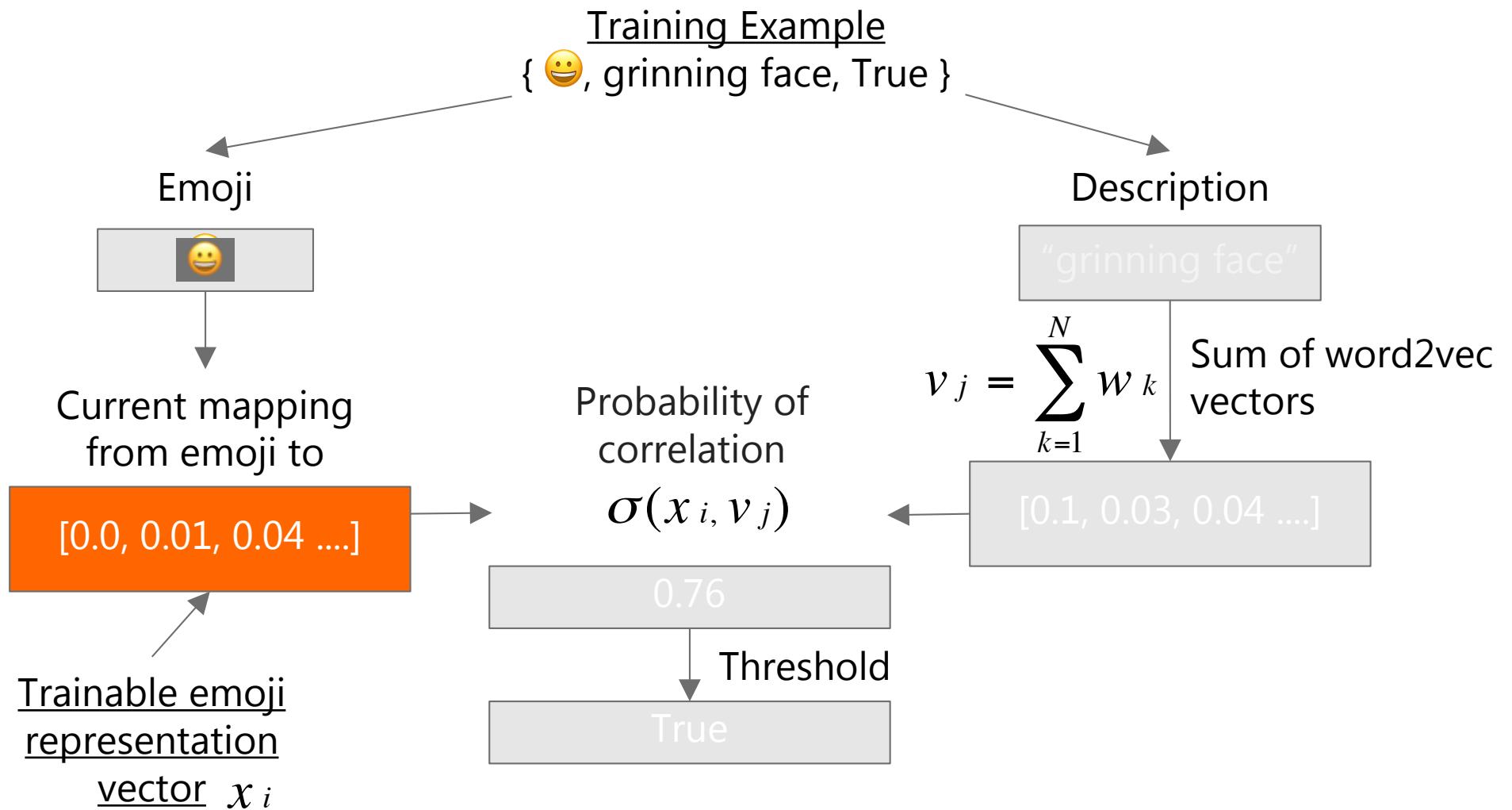
{ 😃 , grin, True }

...



1661 definitions -> ~ 6000
examples in the dataset, a
combination of names and
keywords, all with positive
correlation.

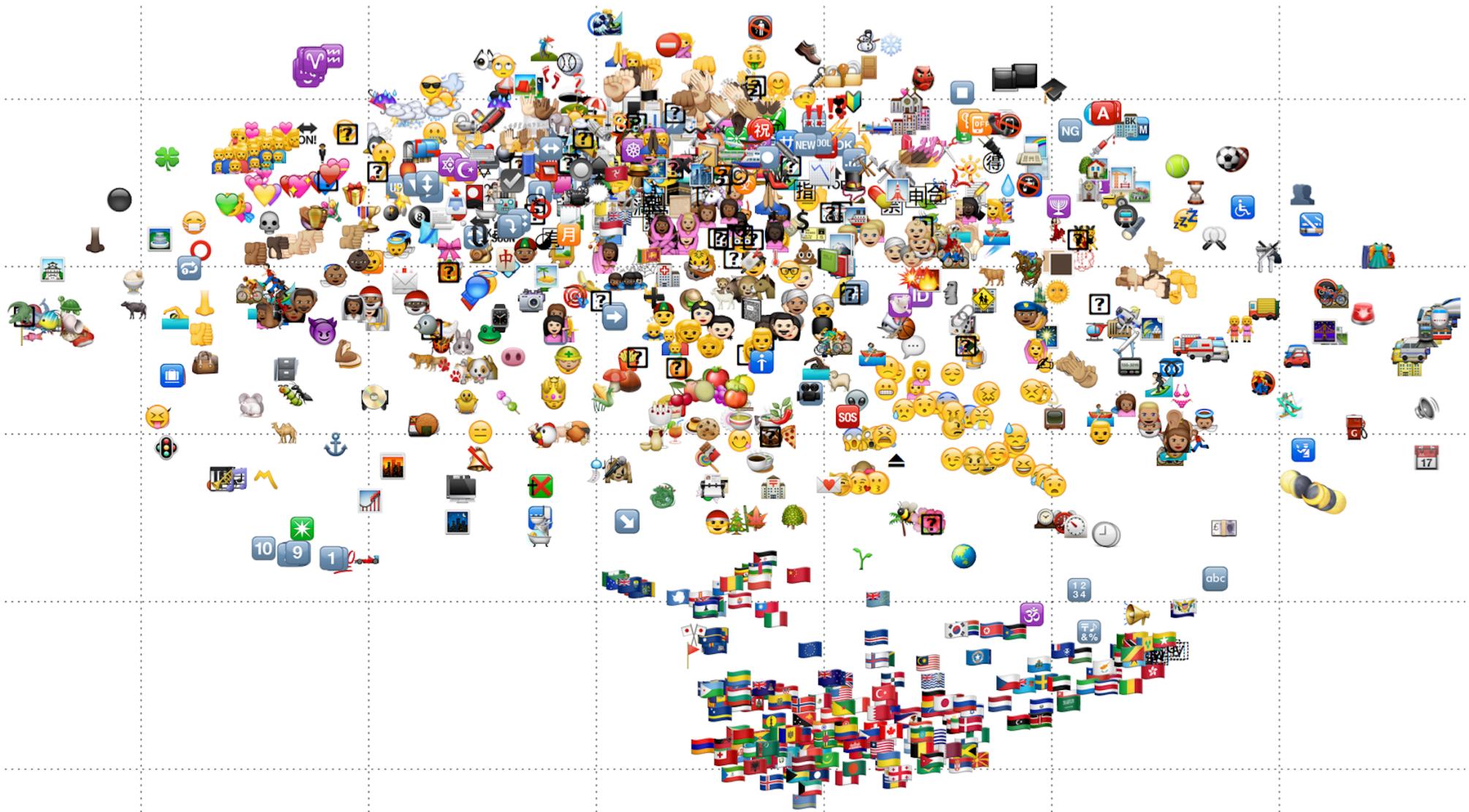
emoji2vec Model



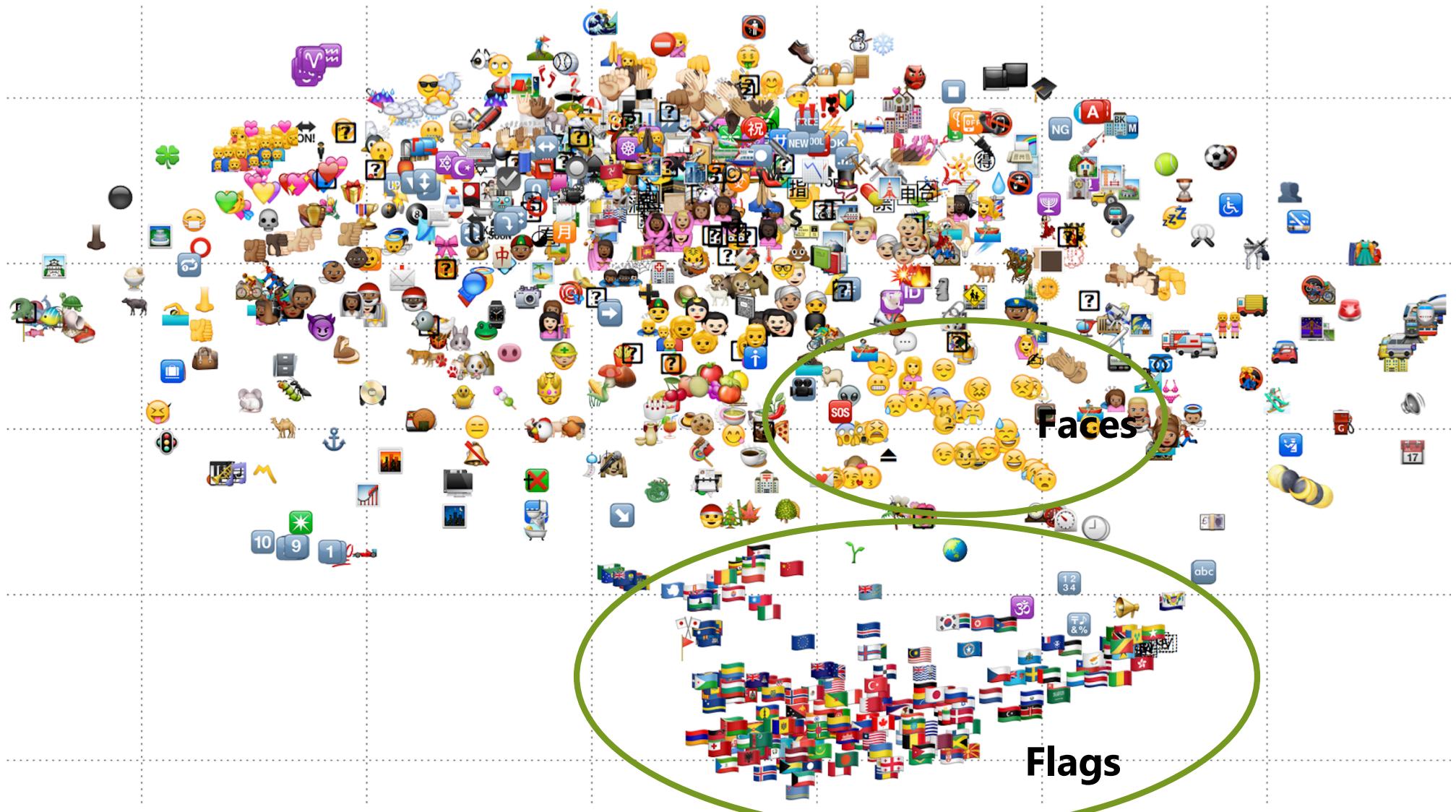
emoji2vec Training

- Define trainable emoji vector
- Measure compatibility between emoji vector and description vector using sigmoid of dot product between the two representations
- Train with logistic loss
- Training data: unicode emoji database, sample negative training data by pairing emojis with random definitions

emoji2vec Evaluation: Visualisation



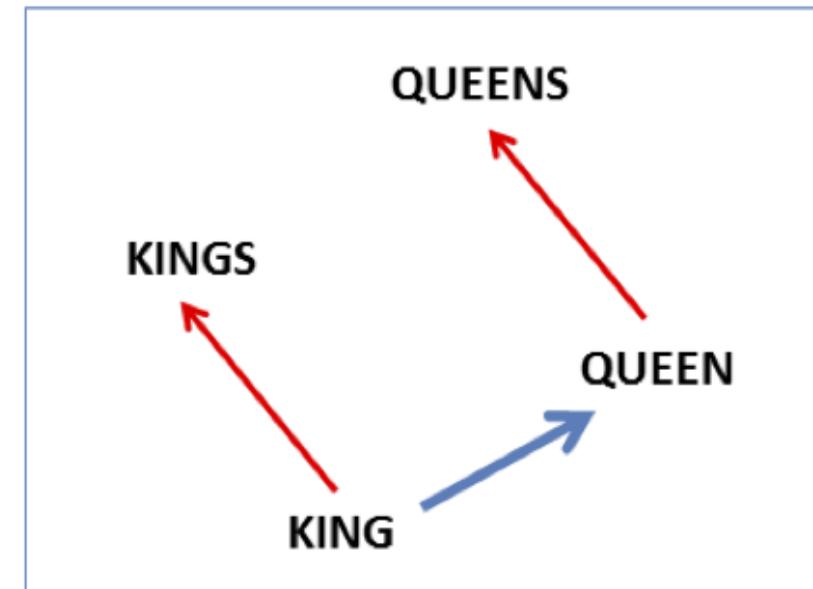
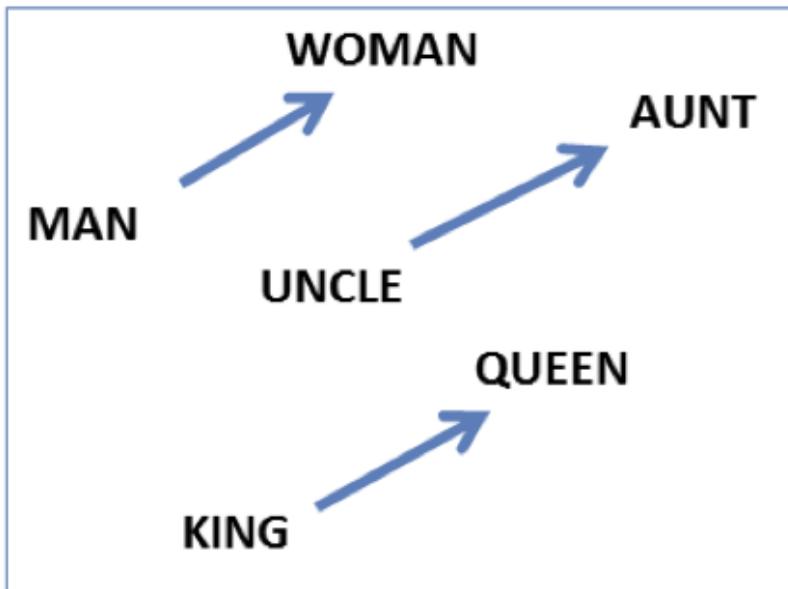
emoji2vec Evaluation: Visualisation



Reminder: Analogy Task

Word Representation Algebra

- $f_n(\text{king}) - f_n(\text{man}) + f_n(\text{woman}) \approx f_n(\text{queen})$
- $f_n(\text{Paris}) - f_n(\text{France}) + f_n(\text{Italy}) \approx f_n(\text{Rome})$



emoji2vec Evaluation: Analogy Task

- Emoji vectors are useful in addition to GoogleNews vectors for sentiment analysis task
- Analogy task also works for emojis

👑 - 🚶 + 🎀 = 1: 👑, 2: 👑, 3: 🏠, 4: 👨, 5: 🐾

💵 - 🇺🇸 + 🇬🇧 = 1: 💵, 2: 💵, 3: 💵, 4: 💵, 5: 💵

💵 - 🇺🇸 + 🇪🇺 = 1: 💵, 2: 💵, 3: 💵, 4: 💵, 5: 💵

👨‍👩‍👧 - 🧑 + 🧑 = 1: 👩, 2: 👑, 3: 🐥, 4: 👰, 5: 🧑

🕶 - ☀ + ⛈ = 1: ☔, 2: ☔, 3: 🏁, 4: 🐾, 5: 🏢

emoji2vec Evaluation: Sentiment Analysis

- Sentiment analysis dataset (60k tweets)
- Model: vector sum of words/emoji in tweets
- Compared against Barbieri et al. 2016 (skip-gram model trained on 10 million tweets containing emojis)

Classification Accuracy on Twitter Dataset Subset (using Linear SVM), N=2295 <u>Only Tweets Containing Emoji</u>	
<u>Embeddings Set</u>	<u>Accuracy</u>
Google News only	47.1%
Google News + (Barbieri et al., 2016)	57.4%
Google News + emoji2vec	59.2%

emoji2vec

- Conclusions
 - Alternative source for learning representations (descriptions) very useful, especially for rare words
 - emoji2vec very useful for downstream application
- Code: <https://github.com/uclmr/emoji2vec>
- Pre-trained emoji embeddings:
<https://github.com/uclmr/emoji2vec/tree/master/pre-trained>

Sentiment Analysis with Linear Models

- Method:
 - *annotate texts* with positive, negative, neutral labels
 - train a model that *learns to associate* those labels with words in text
- Advantages:
 - Better performance than gazetteer-based models
 - Not restricted to hand-crafted resources such as gazetteers
- Disadvantages:
 - Not so easy to implement -> requires some model engineering
 - Annotation is costly and time-intensive
 - Annotated instances can not easily be reused for new domains

Challenges in Sentiment Analysis



Domain Differences and Context Dependence

- "a fast car" vs. "the show was over fast"
- "a cheap hotel" vs. "the dress looks cheap"
- "a large room" vs. "a large queue"
- ...

Irony



Tristan Miller
@Logological

Following

Pretty unsolicited call for paper! Such colourfull! Would look so much more reputable if only set in MS Comic Sans.
#spam

Research and Review Journal - Mozilla Thunderbird

File Edit View Go Message Enigmail Tools Help

Get Messages Write Chat Address Book Tag

From jef.1@rediffmail.com To [REDACTED] 07:19

Subject Research and Review Journal

22 more

IOSR Journals
International Organization of Scientific Research (IOSR)
(IOSR Journals are UGC Approved Journals)

Dear Sir,

We are happy to announce you that International Organization of Scientific Research Journals have come under AQCJ - 2018 Top 10 Journals Ranking. It was calculated on the basis of "Google Scholar Citation" of published articles.

IOSR Journals got 9th Ranking by AQCJ (African Quality Center for Journals) - Top 10 Journals Ranking.

Indexing: Index Copernicus, Cross Ref (USA), NASA ads, ANED (American national Engineering Database), Google Scholar, Open- J Gate.

IOSR Journals provides DOI (Digital Object Identifier) to each article. IOSR Journals DOI is 10.9790.

Papers are invited for **IOSR Journals March 2018** Issue related to all field of Engineering, Management, Medical & Dental Science, Pharmacy, Applied Sciences, Nursing, Humanities and Social Science etc.

Call For Paper: Important Dates

Submission last date	:	15 th March 2018
Acknowledgment	:	Within 24 hrs
Acceptance Notification	:	After 10 days

Negation and Irony

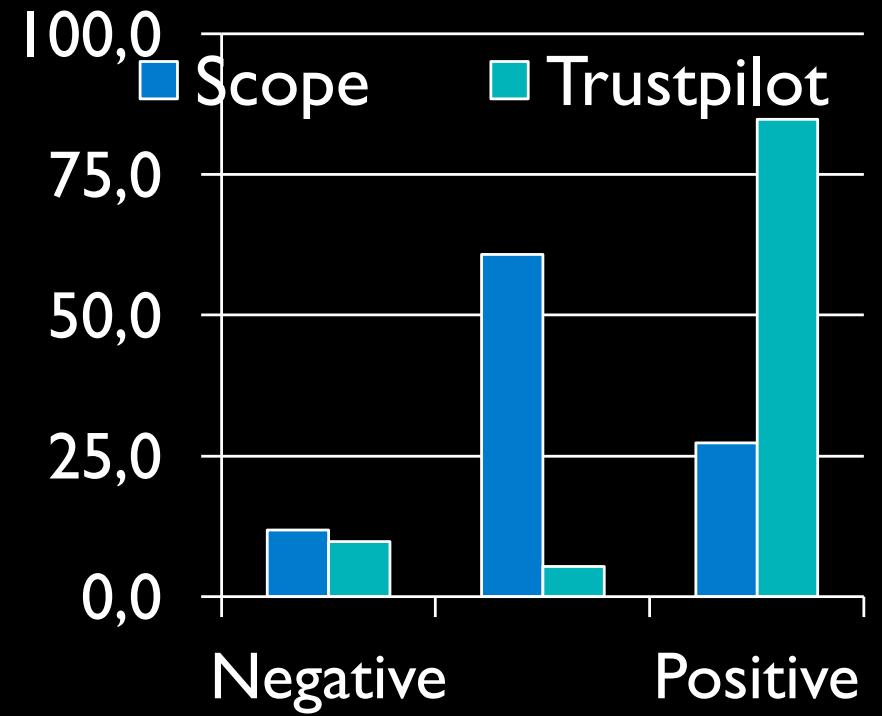
- This is the most delightful movie I have ever seen.
NOT!
- It was so nice of the restaurant to give us 45min
before serving, so we could finish our conversation

Overfitting

Features



Labels



Domain Overfitting

“reliable”

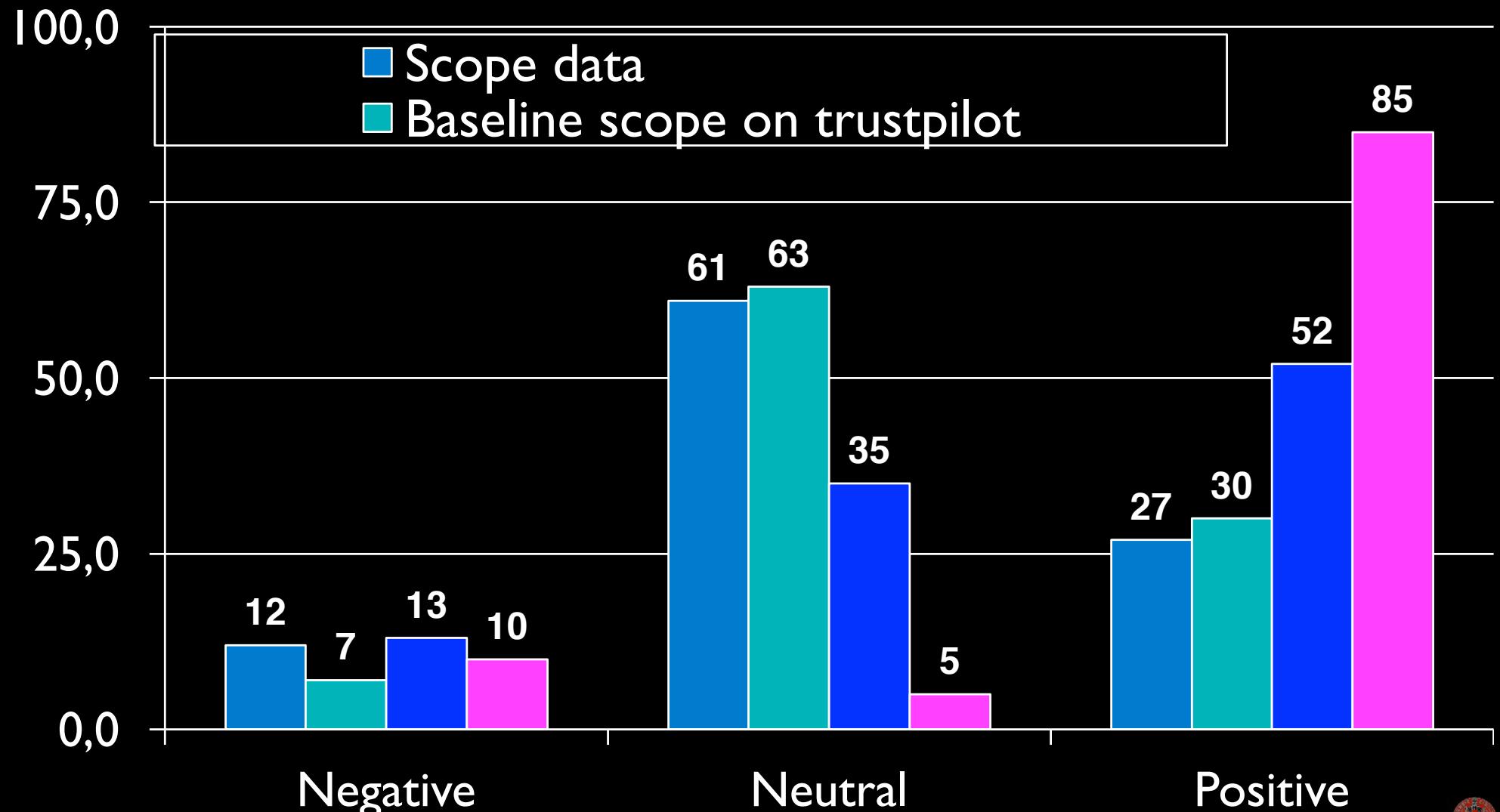


Domain Overfitting

“reliable”



Label overfitting



Label Reliability

AVAILABLE DATA

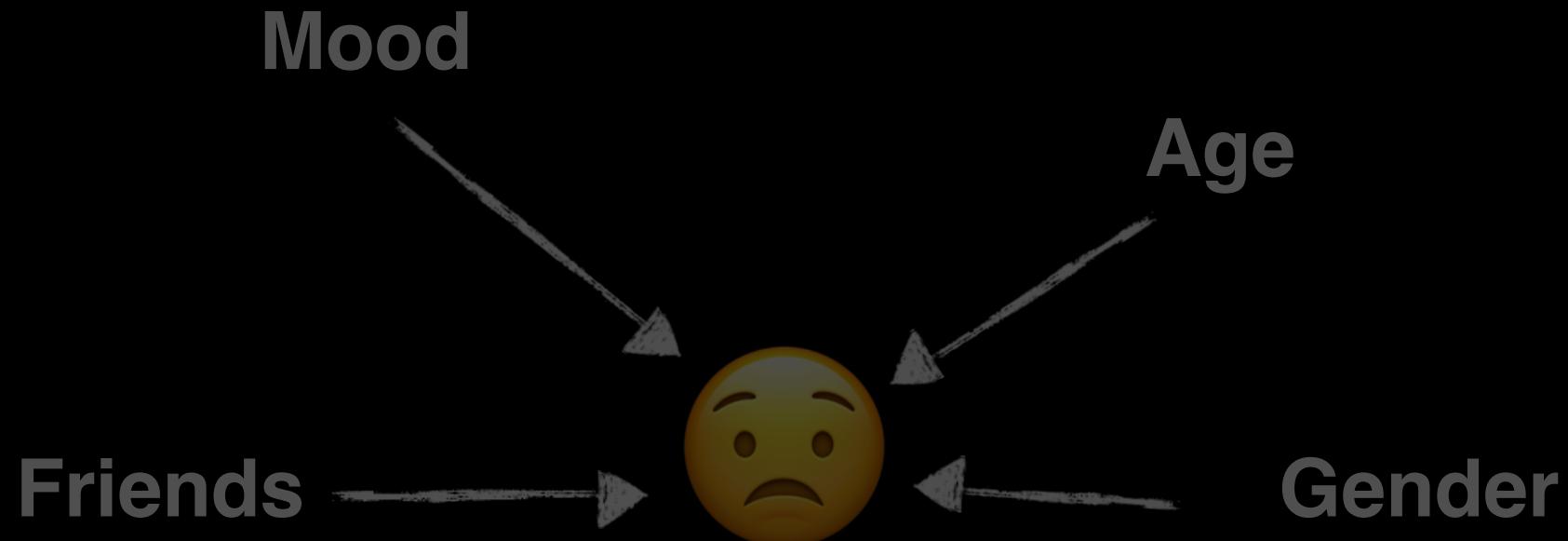
ANNOTATED DATA



Label Reliability

I really like this yoghurt maker





★★ • • •

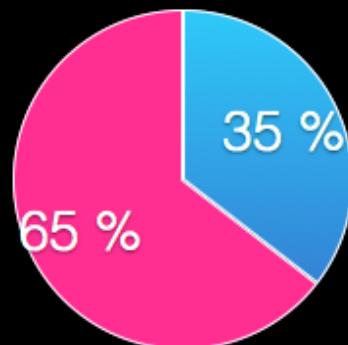
I was not happy with this company...



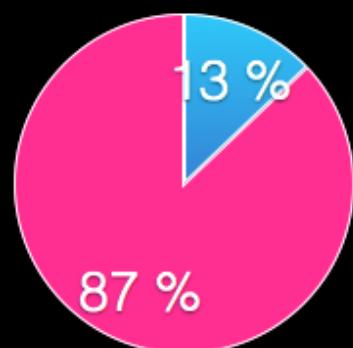
Exposure

Underexposure

- available
- not available



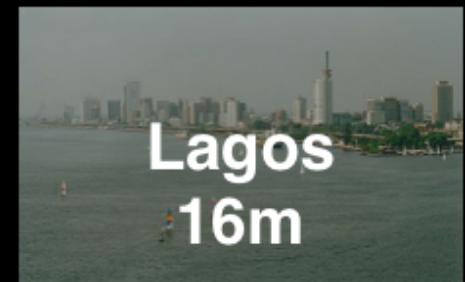
treebanks



semantic
resources

evaluation

Overexposure



**sentiment
analysis**

discourse
parsing

bias



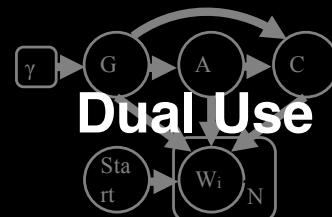
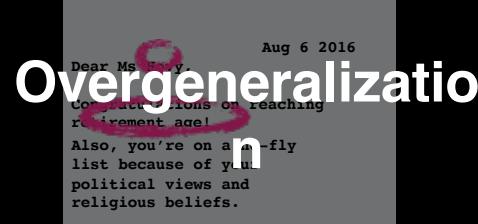
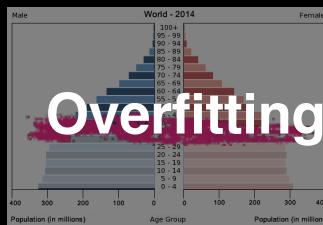
Dual Use

	Pro	Con
authorship attribution	historical documents	dissenter anonymity
text classification	sentiment analysis	censorship
personalization	better user experience	tailored ads



What can we do?

Problem



Source

data selection

models

research design

community goals

Countermeasures

regularisation, priors, sampling

dummy labels, error weighting, confidence thresholds

consider possible impact

educate users, keep discussion going



Take-Home Points

- Opinion mining can help us make sense of large data
- Different methods
 - Gazetteer-based
 - Supervised learning
 - Traditional features
 - Word embeddings
 - Convolutional neural networks
 - Recurrent Neural Networks and LSTMs
 - Emoji embeddings
- Many pitfalls: even best SA models work only for a well-defined domain
- Data selection is important
- Dangers of misinterpretation and misuse

Next Lecture

- More Applications and Use Cases
- More Tasks
 - Target-Based Sentiment Analysis
 - Aspect-Based Sentiment Analysis
 - Stance Detection
 - Affect Extraction
- Approaches
- Challenges

Thank you!

isabelleaugenstein.github.io

augenstein@di.ku.dk

@IAugenstein

github.com/isabelleaugenstein