



NUI Galway  
OÉ Gaillimh

# Introduction to NLP

## Opinion Mining, Ethics & Data Privacy

Dr. Paul Buitelaar  
Data Science Institute, NUI Galway



# Learning Outcomes of This Lecture

Understand different applications of opinion mining

Understand standard approaches in sentiment analysis

Gain some insight into emotion analysis and suggestion mining

Insight into ethics & data privacy considerations in opinion mining and NLP in general

# Overview

Opinion mining

Sentiment analysis

Ethics and data privacy



NUI Galway  
OÉ Gaillimh

# Overview

**Opinion mining**

Sentiment analysis

Ethics and data privacy



NUI Galway  
OÉ Gaillimh

# Opinion Mining



- ▶ Text analytics is now a requirement for large companies
- ▶ Multiple sources of unstructured data
- ▶ "Voice of the Customer" on steroids
- ▶ Ability to understand why NPS metrics are going up or down
- ▶ Ability to respond quickly to service issues
- ▶ Compelling ROI for recent deployments



# Opinion Mining - Some Definitions

Opinion mining is concerned with **analyzing opinionated text** in terms of sentiment, subjectivity, emotions, suggestions, arguments etc.

**Opinionated text:** hotel, movie or other reviews (product ratings), political discourse and commenting (election polls), business and financial analysis (stock value prediction), ...

*Please note that ‘opinion mining’, ‘sentiment analysis’, ‘subjectivity classification’, ‘emotion analysis’ are often used interchangeably*

## Context Analysis by Location (beta)

last 30 days



Sunrise Hotel

★★★★★ 1234 reviews



NPS



Property



Overall Experience



Clealiness



Room



Food &amp; Drinks



Staff &amp; Service



Price



Internet



Location

Sunset Hotel

★★★★★ 234 reviews

58

+5%

87

-10%

87

+2%

88

-3%

64

-3%

59

-2%

78

-1%

85

-9%

54

-20%

33

-50%

Hotel One

★★★★★ 456 reviews

65

-6%

87

+5%

81

+1%

20

-5%

BW Hotel

★★★★★ 1256 reviews

58

-7%

73

-2%

Panorama

★★★★★ 2423 reviews

99

-1%

97

+1%

Sea Breeze Apts

★★★★★ 17 reviews

60

+10%

37

-2%

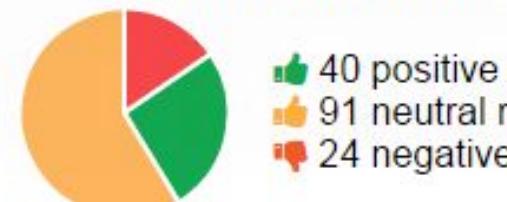
## Property @ Sunrise Hotel

67% Satisfaction in the last 30 days  
+ 5% compared to the previous period (20 Feb '15 -22 Mar '15)

## Overall Experience (last 12 months):

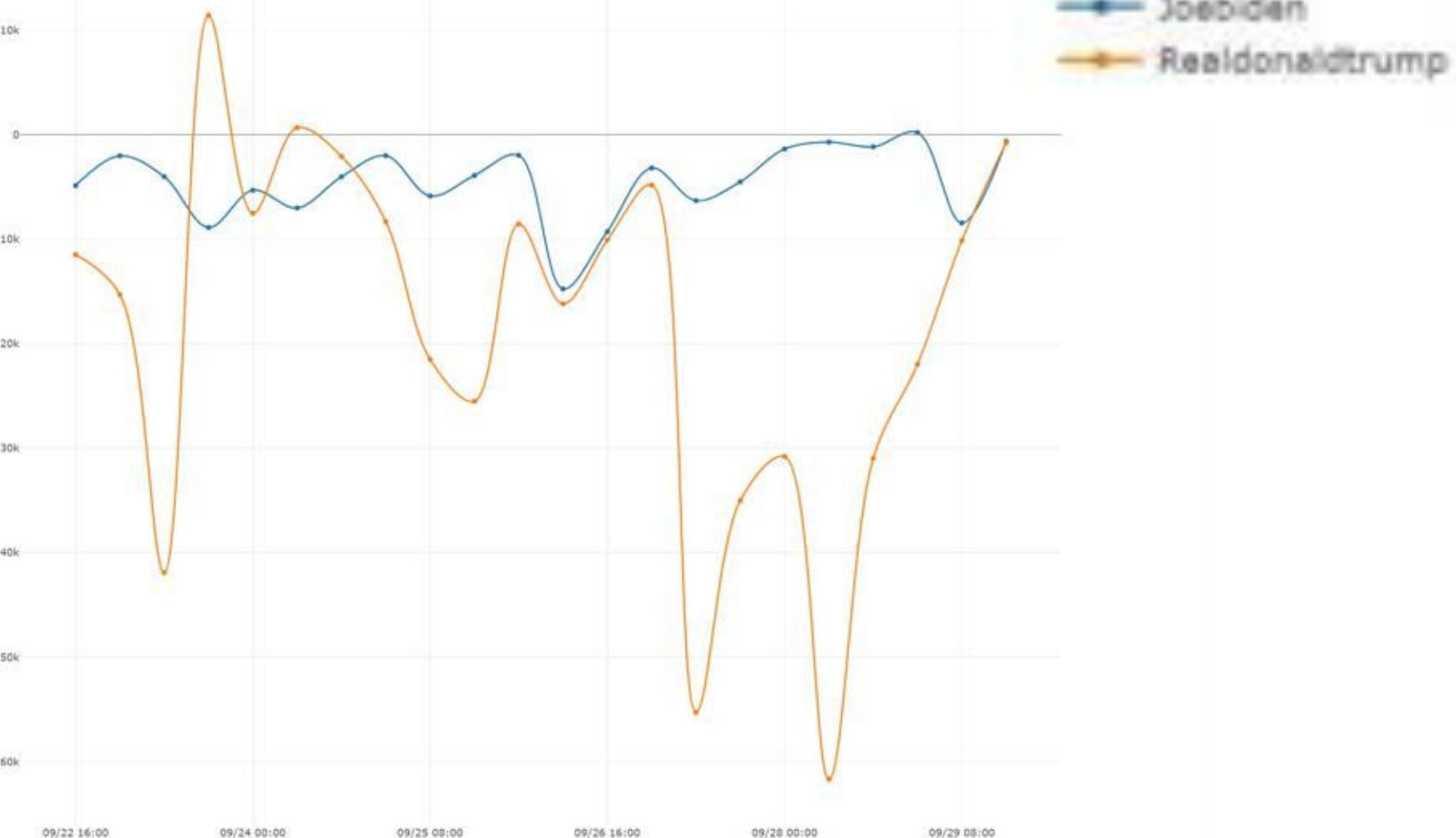


## Sentiment distribution:

NUI Galway  
OÉ GaillimhImage from <https://reputize.com/review-analytics>

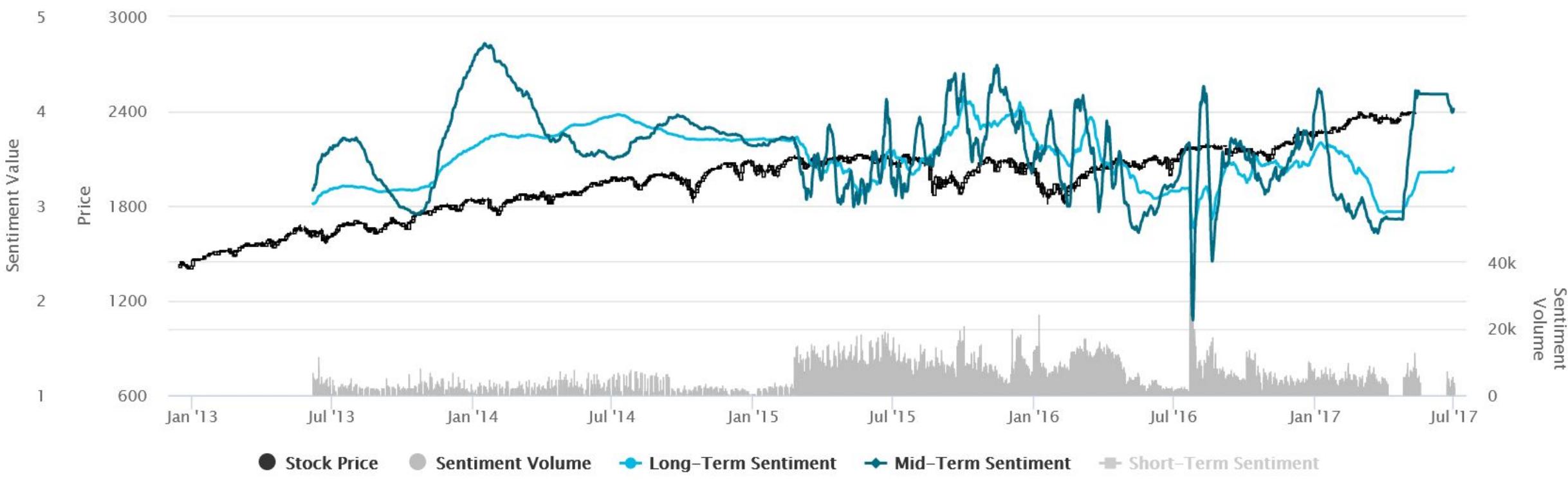
- ↳ Favorites
  - Candidate - Party - Sentiment
  - Candidate - Issue - Sentiment Subplot
  - Candidate Activity
  - Issue Activity
  - Candidate - Sentiment Subplot
  - Issue by Candidate Heat Map
- ↳ Candidate Graphs
  - by Sentiment
  - by Emotion
  - by Author Gender
  - by Author Party
  - by Author Age
  - by Author Race
  - by Author Bias
  - by Author Influence
  - by Author Type
- ↳ Candidate Time Series Subplots
- ↳ Candidate Starburst Subplots
- ↳ Issue Graphs
- ↳ Time Series Graphs
- ↳ HeatMaps
- ↳ Sunburst Graphs
- ↳ Experimental
  - Candidate Net Sentiment**
  - Issue Net Sentiment

## Candidate Net Sentiment



[7 Days](#)[30 Days](#)[6 months](#)[1 Year](#)[All-Time](#)[Search for More Companies](#)[Regular Version](#)

## SP500 – S&P 500 Index – Sentiment Analysis



### S&P 500 Index Description:

Index of the largest 500 US companies.

 Search Properties Tags Projects Sentiment Date Range Source Score & Feedback Type Role: Physicians Assistant, Nurse, or Physician Region: Northeast Tags: Systems Date Range: Last Month(Sep 17th - Oct 17th)

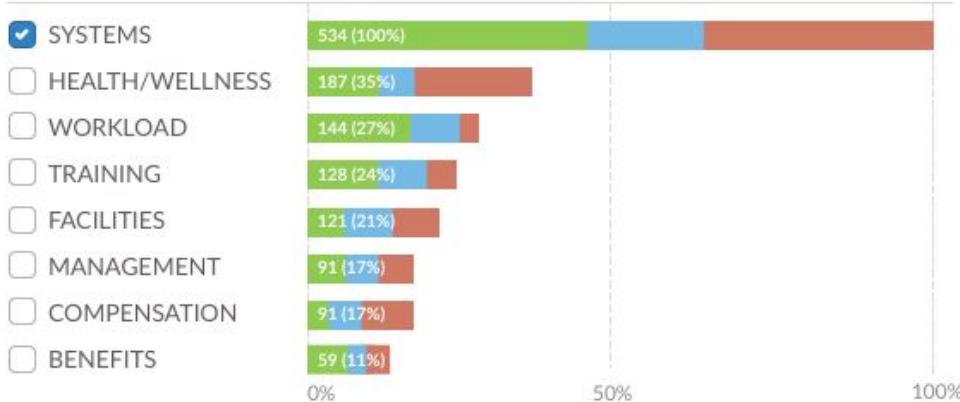
clear all

SAVE SEARCH

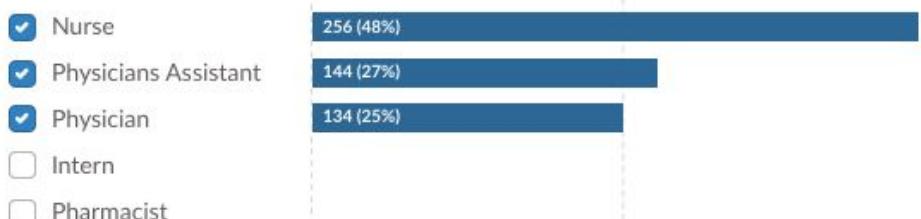
## FEEDBACK BREAKDOWN

Visualization of the % of selected feedback that contains specific tags, sentiment, and properties.

## Tags (with tag sentiment)



## Role



534 matches

● Positive sentiment ● Neutral sentiment ● Negative sentiment

Had to work some awful double shifts since we switched to the new scheduling software, and I'm exhausted.

john@example.com  
17 Oct / 6:45 PM

Fix the new scheduling system, the overtime pay is not worth how tired I am

denise@example.com  
17 Oct / 5:24 PM

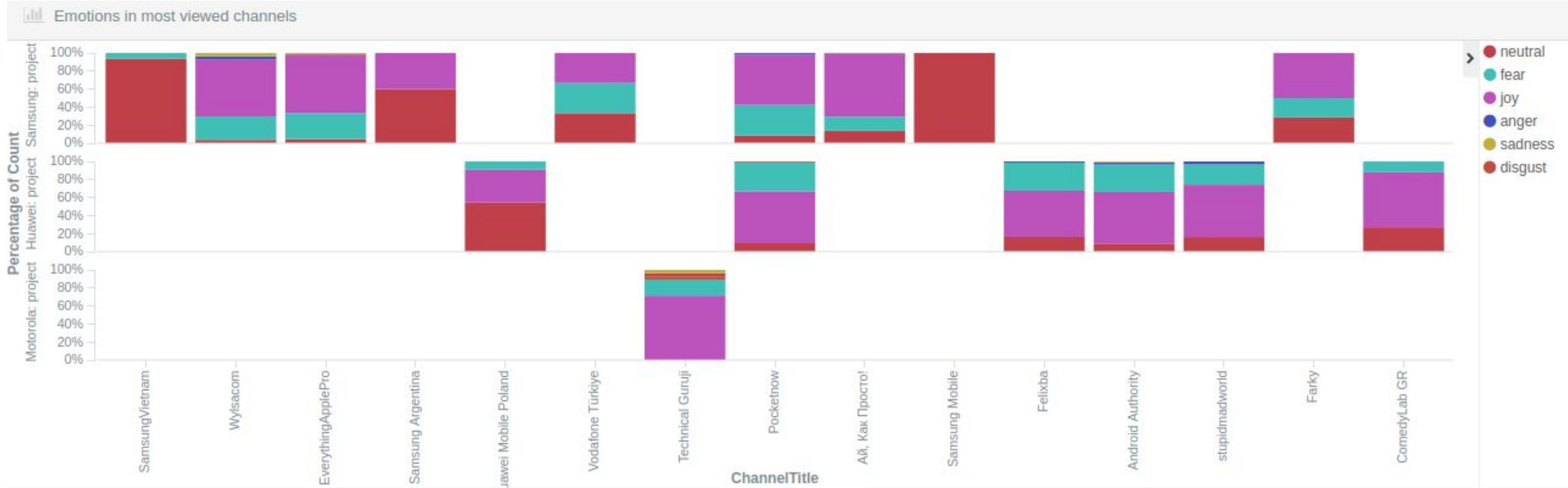
Thank you to whoever has decided we need to focus on getting better technology, we've been operating in the stone age for too long!

dhvani@example.com  
17 Oct / 3:21 PM

New patient records system is great! Really fast, love it!

alex@example.com  
17 Oct / 2:59 PM

# Emotions



# Suggestions

TripAdvisor LLC [US] | [https://www.tripadvisor.ie/Hotel\\_Review-g186605-d195408-Reviews-Harding\\_Hotel-Dublin\\_County\\_Dublin.html#apg=83cc2352dc0d46df9a3...](https://www.tripadvisor.ie/Hotel_Review-g186605-d195408-Reviews-Harding_Hotel-Dublin_County_Dublin.html#apg=83cc2352dc0d46df9a3...) ☆

Overview Deals Reviews About Photos Nearby Q&A Room Tips €119 AMOMA View Deal

## Room Tips

"Can't beat the location"

Andie363 5 days ago

Read review

"Ask about the Resident's bar"

Andrew F 19 days ago

Read review

"Ask a room at the highest level to avoid noise 😊"

mbo63 21 days ago

Read review

"Ask for at least 2nd floor for a quieter room."

Carrie R 1 month ago

Read review

"Ask for a room in a high floor"

Nikolaos C 1 month ago

Read review

"pick higher floors for less noise."

Maxine G 1 month ago

Read review

"If you wish to sleep early, get a room higher up."

Yanqui\_gastronomie 1 month ago

Read review

"We loved the sound of the church bells through the night, but some may find it annoying. Just ask for a room..."

Amy C 1 month ago

Read review

# Overview

Opinion mining

**Sentiment analysis**

Ethics and data privacy



NUI Galway  
OÉ Gaillimh

# Sentiment Analysis

**Classification task** to identify if a given text has positive or negative sentiment

**Challenges** are in implicit, ambiguous and informal language

Classifier can be developed on the basis of a **sentiment lexicon** (unsupervised) or **sentiment labeled data** (supervised)

# Challenges - Implicit Sentiment

Neutral words used but POS sentiment implied

*“Go read the book!”*

Neutral words used but NEG sentiment implied

*“If you are reading this because it is your darling fragrance, please wear it at home exclusively and tape the windows shut.”*

# Challenges - Ambiguity

Same words used in different contexts express NEG vs. POS sentiment

*“This car's steering is **unpredictable!**”*      NEG

*“This film is **unpredictable!**”*                    POS

# Challenges - Irony & Sarcasm

Positive words used but NEG sentiment implied

*“Great job Rogers! Raise rates but not service.”*

*“Yeah, sure!”*



# Challenges - Negation

Positive words used but NEG sentiment implied

*"I don't like this new Nokia model"*

Negation has diverse forms, so can be hard to detect:

*I didn't enjoy it.*

*I never enjoy it.*

*No one enjoys it.*

*I have yet to enjoy it.*

# Challenges - Informal Language

Social media content uses non-standard, informal language such as hashtags

Sentiment	Tweet mention
Positive	Maybe I'm mad but I'm now the proud owner of a potentially #bendy #iPhone6, it's so much bigger than the #4s
	Finally got to see an iPhone 6 today. Not revolutionary at all but it's absolutely gorgeous. (And I want one). #iPhone6
Negative	I'm not sure I want it. It's too big to fit in my back pocket! lol #iphone6
	I'm really disappointed with the #iPhone6. It took them 2 years to change the screen & size. Let down.

# Challenges - Indirect Sentiment

Sentiment expressed may not be that of the author

*“Although this product is **disliked by many**, ...”*



# Classification with a Sentiment Lexicon

Sentiment lexicon provides a **list of positive and negative words**

**Calculate percentage of positive/negative words** in text to be classified

Highest percentage determines sentiment, for example:

Sentiment Lexicon: *POS [great, easy-to-use, ...] NEG [bad, ...]*

*“The camera’s focus was **bad**, but has a **great** size and is **easy-to-use**.”*

POS: 2/13 = 0.153 ; NEG: 1/13 = 0.077

$f(0.153, 0.077) = \text{POS}$

# Sentiment Lexicon - SentiWordNet

SentiWordNet adds **sentiment polarity on words in WordNet synsets**

Apply **Word Sense Disambiguation** to identify correct sense (synset)

Rank	Positive	Negative
1	good#n#2 goodness#n#2	abject#a#2 deplorable#a#1 distressing#a#2
2	better_off#a#1	lamentable#a#1 pitiful#a#2 sad#a#3 sorry#a#2
3	divine#a#6 elysian#a#2 inspired#a#1	bad#a#10 unfit#a#3 unsound#a#5
4	good_enough#a#1	scrimy#a#1
5	solid#a#1	cheapjack#a#1 shoddy#a#1 tawdry#a#2
6	superb#a#2	unfortunate#a#3
7	good#a#3	inauspicious#a#1 unfortunate#a#2
8	goody-goody#a#1	unfortunate#a#1
9	amiable#a#1 good-humored#a#1 good-humoured#a#1	dispossessed#a#1 homeless#a#2 roofless#a#2 hapless#a#1 miserable#a#2 misfortunate#a#1 pathetic#a#1 piteous#a#1 pitiable#a#2 pitiful#a#3 poor#a#1 wretched#a#5
10	gainly#a#1	

# Sentiment Lexicon - SentiWordNet

Synsets can be positive as well as negative (or neutral)

Synset_ID	Pos(s)	Neg(s)	SynsetTerms
1207406	0.0	0.75	cold#a#1
1212558	0.0	0.75	cold#a#2
1024433	0.0	0.0	cold#a#3
2443231	0.125	0.375	cold#a#4
1695706	0.625	0.0	cold#a#5

# Sentiment Lexicon - LIWC

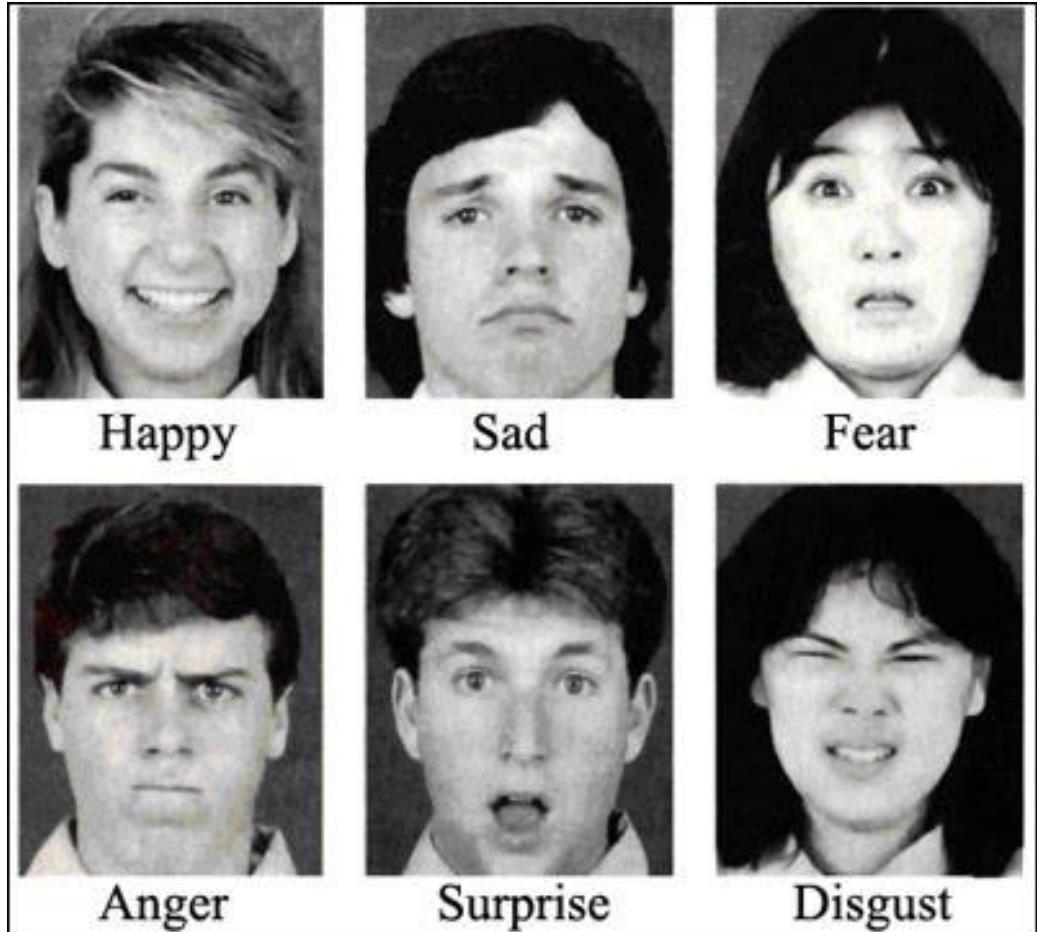
LIWC (Linguistic Inquiry and Word Count) lexicon of **emotion words**

LIWC category	Example words
Positive emotions	Happy, pretty, good
Optimism	Ease, trust, hope
Negative emotions	Hate, worthless, enemy
Anxiety	Nervous, afraid, tense
Anger	Hate, kill, pissed
Sadness	Grief, cry, sad



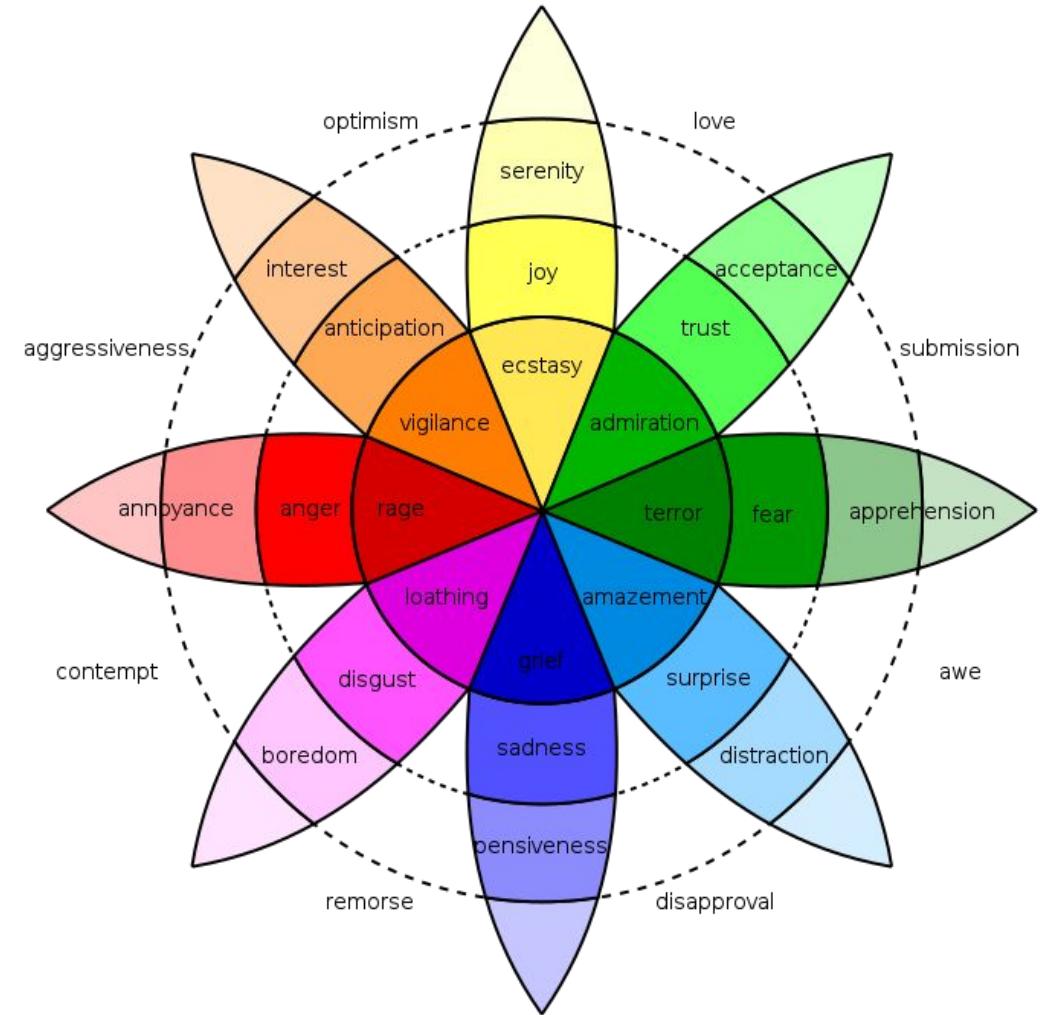
# Emotion Models - Ekman

Most widely used model is **Ekman's**  
**with 6 basic facial emotions**



# Emotion Models - Plutchik

'Plutchik's wheel' has 2 more categories (Trust, Anticipation) and 4 levels of intensity to give **32 emotions**

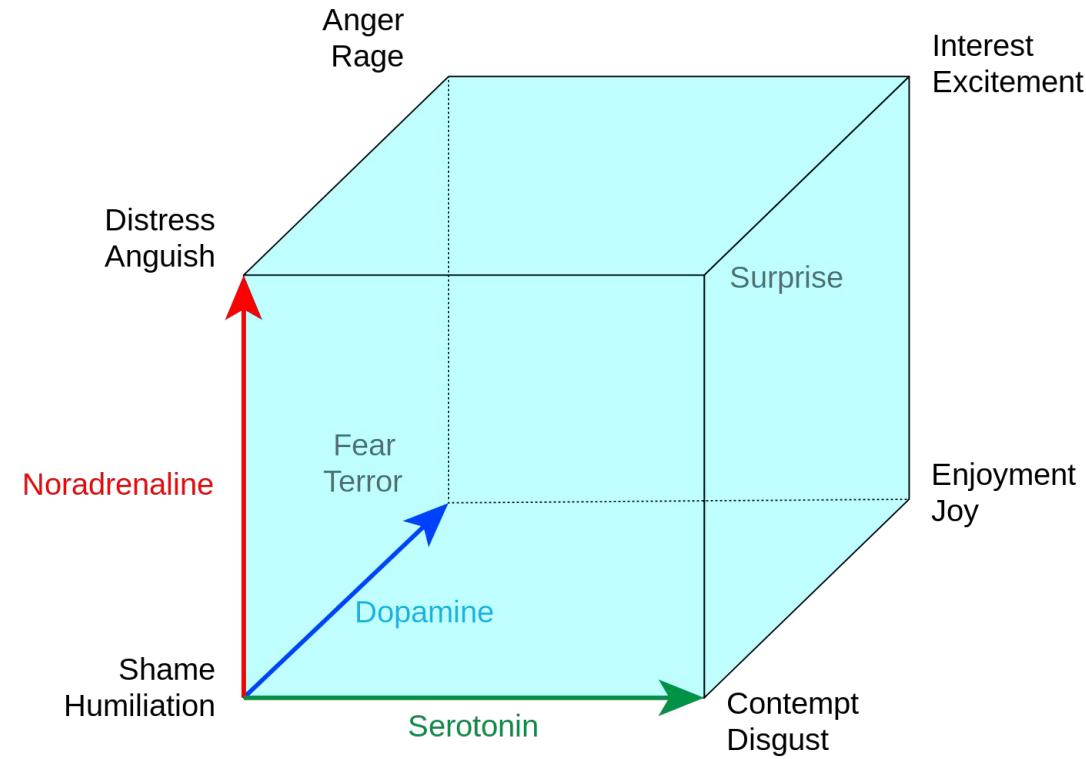


# Continuous Models

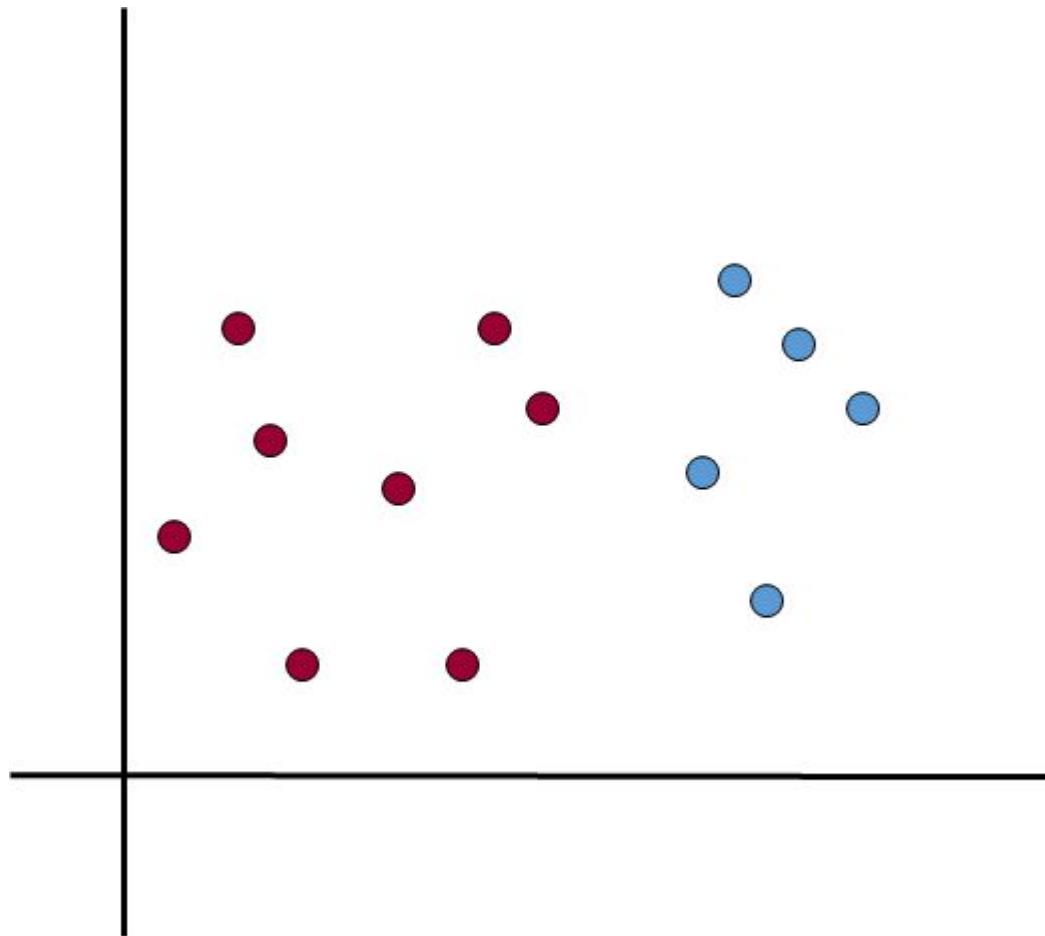
Lövheim's model uses 3 continuous values related to neurotransmitters that have been interpreted as:

- **Valence** - from sad to happy
- **Arousal** - from calm to excited
- **Dominance** - from submissive to agency

*Positive and negative sentiment can be expressed by these models as a combination of continuous values for V/A/D*



# Classification with Sentiment Labeled Data



# Classification - Supervised Methods

Traditional ML: Support Vector Machines (**SVM**)

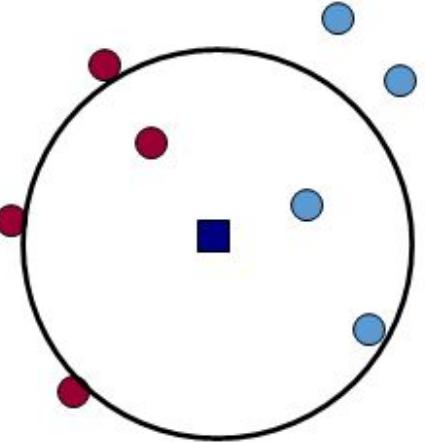
Deep Learning (neural networks)

*In general, set up a vector space for your data set*

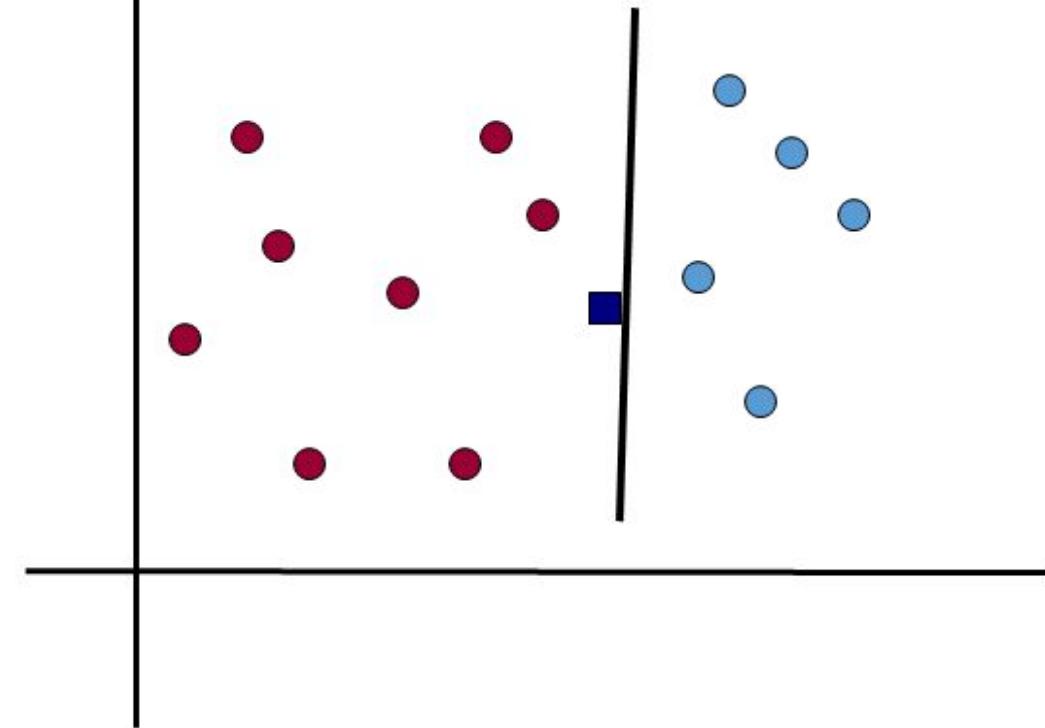


# Classification with SVMs

Example: k-Nearest Neighbours



Example: Support Vector Machines



# Features for Sentiment Analysis

**Feature selection and engineering** is key to good performance!

**Features:** words vs. n-grams, part-of-speech, opinion words (sentiment lexicon), negation, modal verbs (could be, should be, ...), syntactic dependency

**Feature selection and weighting:** occurrence (binary) vs. frequency, PMI, TF-IDF



# Part of Speech

A word can express different sentiments according to PoS

“*fine*” NEG as a verb, POS as an adjective

Word+POS as combined feature can improve sentiment analysis

“*It is fine if you fine me, but don't rub it in!*”

“*It/PRP is/VBZ fine/JJ if/IN you/PRP fine/VBP me/PRP ,/, but/CC don't/NN rub/VBP it/PRP in/IN !/.*”

# Negation

Positive words used but NEG sentiment implied

*"I don't like this new Nokia model."*

Simplest approach is to append 'NOT' to all words after the occurrence of a 'negation word' until next punctuation:

*I do not NOT\_like NOT\_this NOT\_new NOT\_Nokia NOT\_model .*



# Aspect-based Sentiment Analysis

Sofar we considered sentiment on the whole input string, e.g.

*“The staff was **very friendly**, the room was **comfortable** and breakfast was **tasty**,  
but the bathroom was **very small**. ”* POS (?)

More fine-grained sentiment analysis identifies the ‘aspect’ of the sentiment, i.e. ‘what does the sentiment refer to’?

# Aspect-based Sentiment Analysis

*“The staff was **very friendly**, the room was **comfortable** and breakfast was **tasty**, but the bathroom was **very small**. ”*

Aspect	Sentiment
Staff	POS
Breakfast	POS
Room	POS
Bathroom	NEG

# Aspect-based SA - Unsupervised

Define aspect categories, e.g. ‘Staff, Breakfast, Room, ...’

Check for aspect categories in a window around sentiment words from a sentiment lexicon

Compute sentiment for aspect categories in window

# Aspect-based SA - Supervised

Construct a labeled dataset with aspects and corresponding sentiments

*“The [Staff:POS staff was very friendly] and I really enjoyed my stay, but the [Bathroom:NEG bathroom was very small].”*

Build a multi-label classifier which labels the aspect and sentiment

# Suggestions

Opinionated texts may express **suggestions for others or for improvement**  
Often in text segments with **neutral sentiment**



# Suggestion Mining

Suggestion Mining comprises two tasks:

- **suggestion classification** - binary (y/n) classification of input string
- **suggestion extraction** - span identification of the suggestion

# Suggestion Mining - Features

## Keywords and Phrases

*needs to, need to, suggest, recommend, if, I wish, go for, should have, would, could have been, I would like, I'd like, I would love, I'd love, love to see, there should be, I wish, allow us to*

## Syntactic Patterns

*there should be (DT), I cannot (VB), MD be () (ADV) ADJ, MD (like/prefer/love) to, stop VBG, ability (to/of)*

# Suggestion Mining - Distant Supervision

The screenshot shows a web page from [www.wikihow.com/Clean-Shoe-Insoles](http://www.wikihow.com/Clean-Shoe-Insoles). The page title is "wikiHow clean insoles". The main content area is divided into two sections: "Tips" and "Warnings".

**Tips**

- Get in the habit of disinfecting and deodorizing your shoes' insoles every few months, or more frequently if you're especially active or tend to do a lot of walking in them.
- For insoles that are worse-for-wear, try using a combination of cleaning methods. For instance, you can start by scrubbing them the soap and water, then spraying them with alcohol or treating them with baking soda (or both).
- Combine regular cleaning with a regimen of food powder or odor-eating products to keep your insoles fresh for longer.
- Since dirty insoles are most often the result of sweat and bacteria being transferred from the body, it's important to keep your feet clean.

**Warnings**

- Avoid cleaning shoe insoles in the washing machine. Soaking can destroy foot liner materials and cause them to come apart faster.
- While most insoles are salvageable, not all of them will be. If you discover that your shoes still smell after you've attempted a few different remedies, you may be better off throwing them out and replacing them with a new pair.

# Overview

Opinion mining

Sentiment analysis

**Ethics and data privacy**



NUI Galway  
OÉ Gaillimh

# Ethics - Chatbots

**Support The Guardian**  
Available for everyone, funded by readers

[Contribute →](#) [Subscribe →](#)

Search jobs | Sign in | Search ▾ | International edition ▾

# The Guardian

News      Opinion      Sport      Culture      Lifestyle      More ▾

Coronavirus World UK Environment Science Global development Football **Tech** Business Obituaries

**Tech Weekly**  
Artificial intelligence (AI)

## Tay, Microsoft's racist chatbot raises difficult questions - Tech weekly podcast

How the tech firm's artificially intelligent Twitter chatbot went from sweet tween to Holocaust denier overnight

- [How to listen to podcasts: everything you need to know](#)

# Ethics - Natural Language Generation



Sign in

Home

News

Sport

Reel

Worklife

Travel

## NEWS

[Home](#) | [US Election](#) | [Coronavirus](#) | [Video](#) | [World](#) | [UK](#) | [Business](#) | [Tech](#) | [Science](#) | [Stories](#) | [Entertainment & Arts](#)

Tech

## 'Dangerous' AI offers to write fake news

By Jane Wakefield  
Technology reporter

27 August 2019



NUI Galway  
OÉ Gaillimh

# Ethics - Text Classification

[World](#)[Business](#)[Markets](#)[Breakingviews](#)[Video](#)[More](#)

RETAIL OCTOBER 11, 2018 / 12:04 AM / UPDATED 2 YEARS AGO

## Amazon scraps secret AI recruiting tool that showed bias against women

By Jeffrey Dastin

8 MIN READ



# Ethics - Gender Bias

he (109)



she (91)



# Ethics - Addressing Bias in NLP

## Data Statements for Natural Language Processing: Toward Mitigating System Bias and Enabling Better Science

**Emily M. Bender**

Department of Linguistics  
University of Washington  
ebender@uw.edu

**Batya Friedman**

The Information School  
University of Washington  
batya@uw.edu

Bias can be on **gender, ethnic identity, social demography, geographical location, ...**

“A **data statement** is a characterization of a dataset that provides context to allow developers and users to better understand ... **what biases might be reflected in systems built on the software.**”

# Ethics - Use of Personal & Social Media Data

Obvious **ethical considerations in use of personal data** such as patient health records -- *more on personal data later on in the lecture*

Ethical considerations also in **use of social media data** (e.g. text mining)

- *Who owns social media data? Is compensation required?*
- *Should users or platforms control use of social media data?*
- *Who owns knowledge obtained from processing social media data?*

# Data Privacy



NUI Galway  
OÉ Gaillimh

# General Data Protection Regulation (GDPR)

“Regulation ... 2016/679 of the European Parliament and of the Council ... regulates the processing by an **individual, a company or an organisation of personal data** relating to **individuals** in the EU.”



# Personal Data

**“Personal data is any information that relates to an identified or identifiable living individual. Different pieces of information, which collected together can lead to the identification of a particular person, also constitute personal data.**

Personal data that has been **de-identified, encrypted or pseudonymised but can be used to re-identify a person** remains personal data and falls within the scope of the GDPR.

**Personal data that has been rendered anonymous in such a way that the individual is not or no longer identifiable is no longer considered personal data. For data to be truly anonymised, the anonymisation must be irreversible.”**

# Personal Data in NLP

**Language data is also personal data** as individuals can be identified by their language data use (i.e. how they write or speak)

**Natural language processing can lead to ‘fingerprinting’ of individuals**



# Data Protection Impact Assessment

**Identify potential data privacy and data protection issues associated with the data collection for a specific (NLP) task by answering questions such as:**

- *Have you identified the minimally sufficient amount of data for your purpose?*
- *Have you informed affected individuals what you are doing with their data?*
- *Have you identified a secondary use for personal data you collected?*
- *Have you established the maximum time period required for keeping the data?*
- *Have you developed appropriate data anonymization strategies?*

**Data Protection Impact Assessment** to be discussed with and approved by the **Data Protection Officer** at the **Data Controller** (the organisation where the task is done)

# Lab of this Week

Exercises in sentiment analysis



NUI Galway  
OÉ Gaillimh



NUI Galway  
OÉ Gaillimh

QA

