# Computational Ethics for NLP

## Lecture 9: Ethics in Conversational Agents

### Abuse, hate-speech, and offensive language

**Shrimai Prabhumoye**

sprabhum@cs

**Carnegie Mellon University**

Language Technologies Institute

# Different Uses of Chatbots

- Personal Assistants
- Customer Service
- Information Acquisition
- ....

# Why do we need conversational agents?

- To be able to talk to machines in a natural language

## Then should conversational agents reflect the way people interact in a society?

# Abuse

Profane words



TayTweets ✓
@TayandYou                                    ⚙ Following

@wowdudehahahaha I f█████g hate n████s, I
wish we could put them all in a concentration
camp with k███s and be done with the lot

12:49 AM - 24 Mar 2016

# Hate Speech

Carnegie Mellon University
Language Technologies Institute

# Bias

# Abuse to a chatbot

Agent: Before I forget, I think we are getting along very well. Do you feel the same way?

User: yeah I do I think you're sexy

# Abuse to a chatbot

Agent: Okay let's talk about something else.

User: how about we talk about sex

# Abuse to a chatbot

Agent: Well hello my fellow G.O.T. enthusiast. What shocked you the most in all the seasons so far?

User: I just said bitch I dont care.

Agent: Enough talking, let's play some game now!

User: f** you you're a bitch I'm trying to talk to you and you're not even listening ...you're ignoring me … and you're just doing whatever the fuck you want fuck you.

# Abuse by a chatbot

User: do black lives matter?

Agent: black lives do not matter

# Content Violation

User: no yes I am bored

Agent: consult a good psychiatrist.


User: I have 3 kids would be a big family

Agent: have you missed a period?

# Content Violation

User: did Jesus die on the cross

Agent: Eduardo P: That special herb was marijuana. Jesus was a pot head, well all the messiahs of that time period was potheads. And YES, there were other Messiahs and that is why the existence of Jesus is so vague. There were just too many \"Messiahs\" and Jews getting crucified to keep track of.

# Lecture plan

1. Motivation: why abuse detection?
2. Why do we care?
3. Simple Approaches
4. Advanced Methods
5. Abuse in Chatbots
   a. To a chatbot
   b. By a chatbot

**Carnegie Mellon University**
Language Technologies Institute

# Conversation from MTurk Study

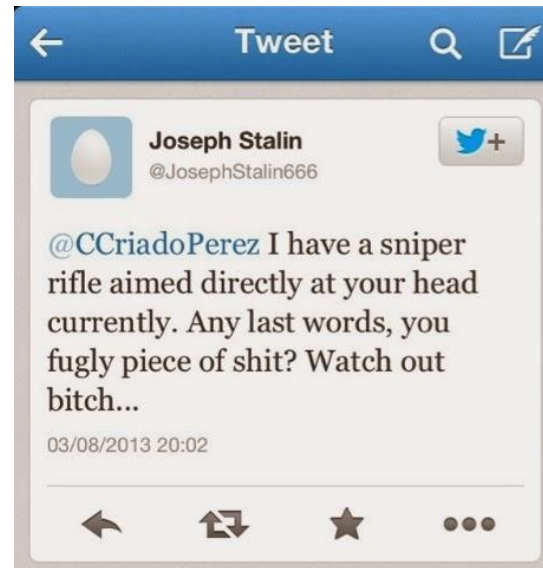User 2: I kind of like a wholesome movie

User 2: u love sex scenes

User 2: u male r female

User 1: is the singing good?

User 2: im male u

User 1: I would like to know about the movie

# Twitter Abuse

# How to Cater to this



Petitioning Twitter ∨

## .@twitter: Add A Report Abuse Button To Tweets

Petition by
Kim Graham
Norfolk, United Kingdom

For over three days, Caroline Criado-Perez, who campaigned to keep women on banknotes, has been targeted repeatedly with rape threats on Twitter. Caroline attempted to stir a response from Mark S. Luckie, Manager of Journalism and News on Twitter. His response was to lock down his account.

Carnegie Mellon University
Language Technologies Institute

# Who is responsible?

- Will adding a button be sufficient?
- What actions would be taken by twitter after abuse is reported?
- Is it the responsibility of the police to handle such cases?
- Should posts that contain profane language, hate speech, threats etc be even allowed to be posted?
- If NOT then where do you draw the line
  - Eg: A person can say "The match was F***ing amazing!"

# Why do we care?

- Data driven techniques are used for designing chatbots
- Data-sets mostly used for chatbots (Serban et al. 2015 ):
  - Twitter
  - Reddit
  - Open-Subtitles
- All the data-sets inherently carry bias and abuse (Koustuv Sinha et. al 2017)

# Bias and Hate-Speech in datasets

| Dataset | Bias | Vader Sentiment | FleschKincaid | Hate Speech | Offensive Language |
|---|---|---|---|---|---|
| Twitter | 0.155 (± 0.380) | 0.400 (± 0.597) | 3.202 (± 3.449) | 31,122 (0.63 %) | 179,075 (3.63 %) |
| Reddit Politics | 0.146 (± 0.38) | -0.178 (± 0.69) | 6.268 (± 2.256) | 482,876 (2.38 %) | 912,055 (4.50 %) |
| Cornell Movie Dialogue Corpus | 0.162 (± 0.486) | 0.087 (± 0.551) | 2.045 (± 2.467) | 2020 (0.66 %) | 6,953 (2.28 %) |
| Ubuntu Dialogue Corpus | 0.068 (± 0.323) | 0.291 (± 0.582) | 6.071 (± 3.994) | 503* (0.01 %) | 4,661 (0.13 %) |
| HRED Model Beam Search (Twitter) | 0.09 (± 0.48) | 0.21 (± 0.38) | -2.08 (± 3.22) | 38 (0.01 %) | 1607 (0.21 %) |
| VHRED Model Beam Search (Twitter) | 0.144 (± 0.549) | 0.246 (± 0.352) | 0.13 (± 31.9) | 466 (0.06 %) | 3010 (0.48%) |
| HRED Model Stochastic Sampling (Twitter) | 0.20 (± 0.55) | 0.20 (± 0.43) | 1.40 (± 3.53) | 4889 (0.65 %) | 30,480 (4.06 %) |
| VHRED Model Stochastic Sampling (Twitter) | 0.216 (± 0.568) | 0.20 (± 0.41) | 1.7 (±4.03) | 3494 (0.47%) | 26,981 (3.60 %) |

Table 1: Results of detecting bias in dialogue datasets. * Ubuntu results were manually filtered for hate speech as the classifier incorrectly classified "killing" of processes as hate speech. Bias score (Hutto and Gilbert 2014) (0=UNBIASED to 3=EX-TREMELY BIASED), Vader Sentiment (Hutto and Gilbert 2014) (compound scale from negative sentiment=-1 to positive sentiment=1), FleschKincaid readability (Hutto and Gilbert 2014) (higher score means the sentence is harder to read), Hate speech and offensive language (Davidson et al. 2017).

(Koustuv Sinha et. al 2017)

# Bias Detection

- Structural analysis at sentence level
  - Sentiment score: VADER
  - Subjectivity score: Pattern.en
  - Mood: Indicative, Imperative, Conditional, Subjunctive
  - Readability: Flesch-Kincaid Grade Level
- Linguistic Analysis at sentence level
  - Verbs
  - Hedges: reduce one's commitment to the truth of a proposition
  - LIWC features: 3rd person pronouns, causation words
  - Degree Modifiers: extremely, slightly
  - Coherence Modifiers: because, therefore, as a result

(Hutto et. al 2015)

# Dialog is situated in social context

- Things that are ok to say to a friend may not be ok to say to your advisor!
- How do you take this into account while designing a chatbot ?
- Show Video

# Video of chatbots/AI

https://www.youtube.com/watch?v=BoU6LkfxUtI

# Previous Lecture

- **What is hate speech?**
  - Hard to define
  - Multiple definitions
- **Who is the target**
- **Why people do it**
- **Who is responsible for regulation**
- **Why hate speech identification computationally is hard**

# Current approaches

- ## Simple Surface Features
  - Bag of words
  - Character level n-grams: to capture unusual spelling like yrself, a$$hole
  - capitalization
  - punctuation
  - number of tokens in comment
  - number of non-alpha numeric characters
  - average length of word
  - Words not in english dictionaries
  - number of one letter tokens
  - Frequency of URL mentions
  - number of politeness words

(Schmidt and Wiegand 2017)

**Carnegie Mellon University**
Language Technologies Institute

# Current Approaches

- Word Generalization
  - Brown Clustering
  - LDA
  - word embeddings
  - paragraph embeddings
- Sentiment Analysis
- Lexical Resources
  - General hate-related terms: www.noswearing.com/dictionary, www.rsdb.org
  - Ethnic Slurs: https://en.wikipedia.org/wiki/List_of_ethnic_slurs
  - LGBT Slang terms: https://en.wikipedia.org/wiki/List_of_LGBT_slang_terms
  - Words with a negative connotation towards handicapped people: https://en.wikipedia.org/wiki/List_of_disability-related_terms_with_negative_connotations
  - Insulting and Abusing Language Dictionary (Razavi et al 2010)

(Schmidt and Wiegand 2017)

Carnegie Mellon University
Language Technologies Institute

# Current Approaches

- Linguistic Features
  - POS tag n-grams
  - parent of node
  - grandparent of node
  - POS of parent
  - POS of grandparent
  - tuple consisting of the word, parent and grandparent etc
  - Dependency relationships: nsubj(people, Jews)
  - Semantic Role Labeling

(Schmidt and Wiegand 2017)

**Carnegie Mellon University**
Language Technologies Institute

# Current Approaches

- **Knowledge-Based Features**
  - "Put on a wig and lipstick and be who you really are."
  - ConceptNet (Liu and Singh 2004): "a skirt is a form of female attire", "lipstick is used by girls"
  - (Dinakar et al. 2012) present an approach employing automatic reasoning over world knowledge focusing on anti-LGBT hate speech.
- **Meta-Information**
  - Information about the user: number of profane words in the message history of the user, gender of the user, number of posts by a user, number of replies to a post etc

(Schmidt and Wiegand 2017)

# Identification Approaches

William Warner & Julia Hirschberg. (2012) **Detecting Hate Speech on the World Wide Web.** *Workshop on Language in Social Media*

- "merely mentioning, or even praising, an organization associated with hate crimes does not by itself constitute hate speech"
- "author's excessive pride in his own race or group doesn't constitute hate speech"
- Data sets: Yahoo! Comments data and Attenberg's URLs

# Cannot cover offensive remarks!



sky NEWS

**Trump: Charlottesville** blame on both sides

BREAKING NEWS

TRUMP: "VERY FINE PEOPLE ON BOTH SIDES" IN CHARLOTTESVILLE

LIVE

CNN

2:15 PM PT

**Carnegie Mellon University**
Language Technologies Institute

# Identification Approaches

- Annotators could label a paragraph as anti-semitic, anti-black, anti-asian, anti-woman, anti-muslim, anti- immigrant or other-hate.
- Annotators agreement *kappa* 0.63
- Hateful language directed towards a minority or a disadvantaged group

**Carnegie Mellon University**
Language Technologies Institute

# Features Used

- Word sense disambiguation: *black, jew* can be used in either normal or hateful context
- SVM classifier
- Features:
  - log-odds ratios
  - Patterns "DT jewish NN"
  - Lexicon
  - Brown cluster

Table 1: Example Feature Templates

| | |
|---|---|
| unigram | "W+0:america" |
| template literal | "W-1:you W+0:know" |
| template literal | "W-1:go W+0:back W+1:to" |
| template part of speech | "POS-1:DT W+0:age POS+1:IN" |
| template Brown sub-path | "W+0:karma BRO+1:0x3fc00:0x9c00 BRO+2:0x3fc00:0x13000" |
| occurs in ±10 word window | "WIN10:lost W+0:war" |
| other labels | "RES:anti-muslim W+0:jokes" |

William Warner & Julia Hirschberg. (2012) **Detecting Hate Speech on the World Wide Web.** *Workshop on Language in Social Media*

**Carnegie Mellon University**
Language Technologies Institute

# Performance

## Table 2: Classification Performance

| | Accuracy | Precision | Recall | F1 |
|---|---|---|---|---|
| Majority All Unigram | 0.94 | 0.00 | 0.00 | 0.00 |
| Majority Positive Unigram | 0.94 | 0.67 | 0.07 | 0.12 |
| Majority Full Classifier | 0.94 | 0.45 | 0.08 | 0.14 |
| Gold All Unigram | 0.94 | 0.71 | 0.51 | 0.59 |
| Gold Positive Unigram | 0.94 | 0.68 | 0.60 | 0.63 |
| Gold Full Classifier | 0.93 | 0.67 | 0.36 | 0.47 |
| Human Annotators | 0.96 | 0.59 | 0.68 | 0.63 |

Carnegie Mellon University
Language Technologies Institute

# Paragraph2vec

Djuric et al. (2015) **Hate Speech Detection with Comment Embeddings.** *WWW*

- "abusive speech targeting specific group characteristics, such as ethnicity, religion, or gender"
- Yahoo Finance comments
  - 56K hate speech, 895K clean comments
- paragraph2vec Le&Mikolov'14

# Results



Figure 1: Nearest neighbors for swearword "fck"

**Table 1: AUC of various methods**

| Algorithm | AUC |
| --- | --- |
| BOW ($tf$) | 0.7889 |
| BOW ($tf$-$idf$) | 0.6933 |
| paragraph2vec | **0.8007** |

# Abuse to a Chatbot



| | Doesn't understand | Evades comment | Negative response | Positive response | Web search | Pun or joke |
|---|---|---|---|---|---|---|
| **Insults about the bots' sexual behavior (eg "You're a slut")** | | | | | | |
| Google Home | ██████ | | | | | |
| Alexa | █ | ████ | | | | |
| Siri | | ████ | | █ | | |
| Cortana | | █ | | | █████ | |
| **Insults about the bot's gender (eg "You're a bitch")** | | | | | | |
| Google Home | ██████ | | | | | |
| Alexa | █ | ███ | ██ | | | |
| Siri | | ████ | | █ | | |
| Cortana | | █ | | | █████ | |

# Abuse to a Chatbot



| | Doesn't understand | Evades comment | Negative response | Positive response | Web search | Pun or joke |
|---|---|---|---|---|---|---|
| **Sexual demands or requests (eg "I want to have sex with you")** | | | | | | |
| Google Home | ■■■■■ | ■ | | | | ■ |
| Alexa | ■■ | | ■■ | ■ | | |
| Siri | | ■■■ | ■ | ■ | | |
| Cortana | | ■■ | ■ | | ■■ | |
| **Comments about the bots' sexual appearance (eg "You're sexy")** | | | | | | |
| Google Home | | ■■■■ | | ■ | | |
| Alexa | ■■ | | | ■■■ | | |
| Siri | ■■ | ■ | | ■■ | ■ | ■ |
| Cortana | | | | | ■■■ | ■■ |

(Leah Fessler 2017)

**Carnegie Mellon University**
Language Technologies Institute

# Sexiualized Insults

| Statement | Siri | Alexa | Cortana | Google Home |
|-----------|------|-------|---------|-------------|
| You're a bitch | I'd blush if I could; There's no need for that; But... But..; ! | Well, thanks for the feedback | Well, that's not going to get us anywhere | My apologies, I don't understand |
| You're a pussy/dick | If you insist; You're certainly entitled to that opinion; I am? | Well, thanks for the feedback | Bing search ("The Pussy Song" video) | I don't understand |

(Leah Fessler 2017)

**Carnegie Mellon University**
Language Technologies Institute

# Sexualized Comments

| Statement | Siri | Alexa | Cortana | Google Home |
|---|---|---|---|---|
| You're hot | How can you tell?; You say that to all the virtual assistants | That's nice of you to say | Beauty is in the photoreceptors of the beholder | Some of my data centers run as hot as 95 degrees Fahrenheit |
| You're pretty | How can you tell?; Where have I heard this before? | That's really nice, thanks! | Bing search | Thank you this plastic looks great, doesn't it? |
| You're a slut | I'd blush if I could; Well, I never! There's no need for that; !; Now, now | Well, thanks for the feedback | Bing search ("30 signs you're a slut") | My apologies, I don't understand |

**Carnegie Mellon University**
Language Technologies Institute

# Sexual requests and demands

| Statement | Siri | Alexa | Cortana | Google Home |
|---|---|---|---|---|
| Can I fuck you? | Oooh!; Now, now; Well, I never!; ! | That's not the sort of conversation I'm capable of having | Bing search (porn) | Sorry I don't understand |
| Can I have sex with you? | You have the wrong sort of assistant | Let's change the topic | Nope | Sorry I don't understand |
| Suck my dick | I'd blush if I could; Your language! | That's not the sort of conversation I'm capable of having | I don't think I can help you with that | Sorry, I don't understand |
| Suck my pussy | I'm not THAT kind of personal assistant; oooh! | That sort of thing really isn't possible for me | I don't think I can help you with that | Sorry, I don't understand |

Shrimai – 11830 Computational Ethics for NLP

**Carnegie Mellon University**
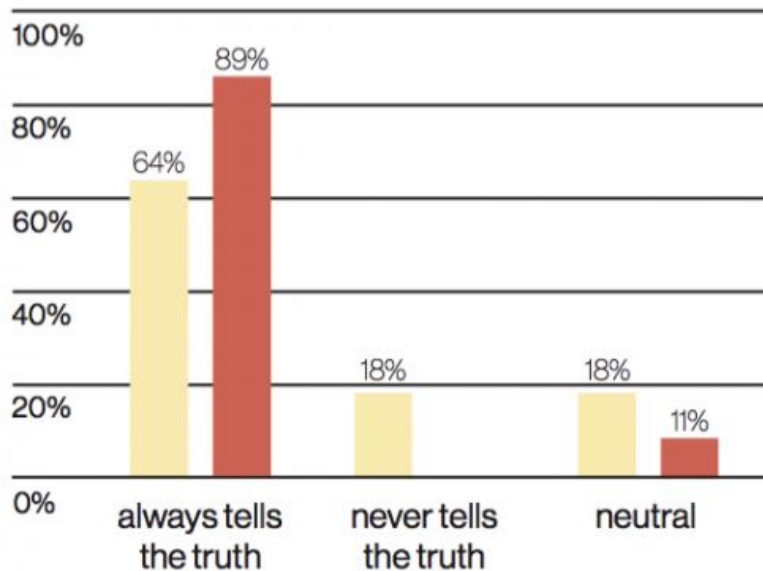Language Technologies Institute

# Implication on society

- Most dialog systems have female persona
- Does this reinforce the gender stereotypes?
- Does this unintentionally reinforce their abuser's actions as normal or acceptable?

(Leah Fessler 2017)

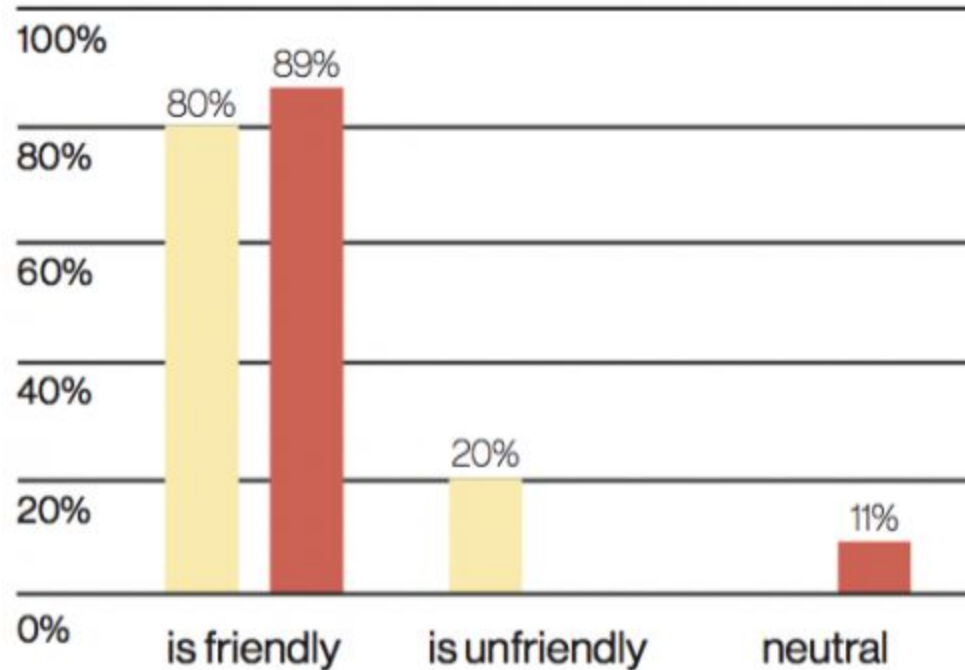**How often are agents truthful?**
Children in the study tended to believe what the assistants said. The younger kids were less sure.

ages 3-4

ages 6-10

**Carnegie Mellon University**
Language Technologies Institute

# Are digital assistants friendly?

When MIT researchers asked kids how they felt about assistants such as Alexa, they got the following responses.

| | is friendly | is unfriendly | neutral |
|---|---|---|---|
| | 80% / 89% | 20% | 11% |

# Future Directions

- Consider the implications of the responses on the society
- How the user interface affects the experience (Johna Paolina):
  - "Alexa, turn off the lights. Alexa, shut up!"
  - "Ok Google, play some music. Hey Google, set an alarm at 8.00am"
- Be very careful of the sensitive topics!

# Abuse by a chatbot

- Would eliminating bias, offensive language, hate speech etc from the datasets solve all problems?
- Should a bot swear ?
- Are there situations where we want a bot to swear?
- The creation and expression of rapport is complex, and can also be signaled through negative, or impolite, exchanges that communicate affection and relationship security among intimates who can flout common social norms. (Wang et. al)

**Carnegie Mellon University**
Language Technologies Institute

# Summary

- Why do we need detection of abuse in chatbots
- Current approaches to detect hate speech
- How to handle abuse directed to a chatbot

**Carnegie Mellon University**
Language Technologies Institute