

















Ethics in NLP — what can go wrong?

What can actually go wron

What can actually go wrong?

Machine-learned NLP Systems

- Aggregate lots of information
- Hard to know why certain predictions are made
- Increasingly wide use in various applications/sectors
- What are the risks here?
- ...of certain applications?
- ▶ IE / QA / summarization?
- ▶ MT?

- Dialog?
- ...of machine-learned systems?
- ...of deep learning specifically?

Bias Amplification

 Bias in data: 67% of training images involving cooking are women, model predicts 80% women cooking at test time — amplifies bias

- Can we constrain models to avoid this while achieving the same predictive accuracy?
- Place constraints on proportion of predictions that are men vs. women?



Zhao et al. (2017)







Bias Amplification

- the context of the word "charming"... but maybe that reflects how language



Alvarez-Melis and Jaakkola (2017)

Exclusion				Unethica	al Use
Most of our annotated data is English data, especially newswire	Surveillance applications?				
• What about:		nerating co	nvincing fak	e news / fake	comments?
		FCC Comment ID: 106030756805675	FCC Comment ID: 106030135205754	FCC Comment ID: 10603733209112	
Dialects?		Dear Commissioners:	Dear Chairman Pai,		
		Hi, I'd like to comment on	I'm a voter worried about	In the matter of	What if these were
Other languages? (Non-European/CIK)		net neutrality regulations.	Internet freedom.	NET NEUTRALITY.	
Other languages: (Non-European/CJK)		I want to	I'd like to	I strongly	undetectable?
		implore	ask	ask	
Codeswitching?		the government to	Ajit Pai to	the commission to	
C C		repeal	repeal	reverse	
		Barack Obama's	President Obama's	Tom Wheeler's	
		decision to	order to	scheme to	
		internet access	broadband	the web	
		Individuals	neonle like me	People like me	
		rather than	rather than	rather than	

۲	Dangers of Automatic Systems	
1HE VERO	📕 TECH - SCIENCE - CULTURE - CARS - REVIEWS - LONGFORM VIDEO MORE - 🛛 🗲 🛩 🔊 ᆂ Q	
Faceb Palest	ook apologizes after wrong translation sees	
Facebook tro	anslated his post as 'attack them' and 'hurt them' ayOng Oct 24, 2017, 10:43am EDT	

Dangers of Automatic Systems

Translations of gay	
adjective	
homosexual	homosexual, gay, camp
alegre	cheerful, glad, joyful, happy, merry, gay
brillante	bright, brilliant, shiny, shining, glowing, glistening
vivo	live, alive, living, vivid, bright, lively
vistoso	colorful, ornate, flamboyant, colourful, gorgeous
jovial	jovial, cheerful, cheery, gay, friendly
gayo	merry, gay, showy
noun	
el homosexual	homosexual, gay, poof, queen, faggot, fagot > Offensive terms
 el jovial 	gay
	Slide credit: <u>allout.org</u>

Slide credit: The Verge

Dangers of Automatic Systems

"Instead of relying on algorithms, which we can be accused of manipulating for our benefit, we have turned to machine learning, an ingenious way of disclaiming responsibility for anything. Machine learning is like money laundering for bias. It's a clean, mathematical apparatus that gives the status quo the aura of logical inevitability. The numbers don't lie."

- <u>Maciej Cegłowski</u>

Slide credit: Sam Bowman

Dangers of Automatic Systems

- "Amazon scraps secret AI recruiting tool that showed bias against women"
 - "Women's X" organization was a negative-weight feature in resumes
 - Women's colleges too

Was this a bad model? May have actually modeled downstream outcomes correctly...but this can mean learning humans' biases

> Slide credit: https://www.reuters.com/article/us-amazon-com jobs-automation-insight/amazon-scraps-secret-ai-recruitingtool-that-showed-bias-against-women-idUSKCN1MK08G

Bad Applications

- Sophia: "chatbot" that the creators make incredible claims about
- Creators are actively misleading people into thinking this robot has sentience

- Most longer statements are scripted by humans
- "If I show them a beautiful smiling robot face, then they get the feeling that 'AGI' (artificial general intelligence) may indeed be nearby and viable... None of this is what I would call AGI, but nor is it simple to get working"



Slide credit: https://themindlist.com/ 2018/10/12/sophia-modern-marvel-ormindless-marketing/

Bad Applications

 Wang and Kosinski: gay vs. straight classification based on faces

- Authors: "this is useful because it supports a hypothesis" (physiognomy)
- Blog post by Agüera y Arcas, Todorov, Mitchell: mostly social phenomena (glasses, makeup, angle of camera, facial hair) — bad science, *and* dangerous





Slide credit: <u>https://medium.com/@blaisea/do-</u> algorithms-reveal-sexual-orientation-or-just-exposeour-stereotypes-d998fafdf477



Final Thoughts

- You will face choices: what you choose to work on, what company you choose to work for, etc.
- Tech does not exist in a vacuum: you can work on problems that will fundamentally make the world a better place or a worse place (not always easy to tell)
- As AI becomes more powerful, think about what we should be doing with it to improve society, not just what we can do with it