

CS540 Introduction to Artificial Intelligence Ethics and Trust in AI

University of Wisconsin-Madison Spring 2025

Announcements

Homework: HW10 due on Friday May 2nd at 11:59 PM

• Final Exam Information: next slide

Course evaluation

Final Information

- Time: May 7th 07:45 AM 09:45 AM
- Location (by section **):

**To find your section go to MyUW->Course Schedule->It will say "LEC 00_". Do not use canvas to find your section (everyone will see CS540 001 since we merged the canvas site for all three sections).

- **Format**: The final exam will be entirely multiple choice.
- The exam will focus on conceptual and applied AI reasoning.
- necessary though it may be useful to double check simple arithmetic.
- **Detailed topic list + practice**: <u>https://piazza.com/class/m5zvrf0clyo3sl/post/449</u>

Lecture 001 (Instructor Sharon Li): S429 Chemistry Building Lecture 002 (Instructor Fred Sala): 1220 Microbial Sciences Lecture 003 (Instructor Blerina Gkotse): B10 Ingraham Hall

Cheat Sheet: You will be allowed a cheat sheet of a single piece of paper (8.5" x 11", front and back).

Calculator: Calculators are allowed if they don't have an internet connection. A calculator will not be





Evil

 Al dual use: VX chemical compound deep fake Autonomous weapons
Lucky for the world

Dual use of artificial-intelligence-powered drug discovery Key observation: flip the objective function to make optimization find many highly toxic compounds



Fig. 1 A t-SNE plot visualization of the LD₅₀ dataset and top 2,000 MegaSyn AI-generated and predicted toxic molecules illustrating VX. Many of the molecules generated are predicted to be more toxic in vivo in the animal model than VX (histogram at right shows cut-off for VX LD₅₀). The 2D chemical structure of VX is shown on the right.

[Urbina et al. Nature machine intelligence 2022]



https://www.youtube.com/watch?v=cQ54GDm1eL0 **Example 1: Fake Obama Video**

can make it look like anyone is saying anything



Example 2: Fake face Images by GAN

•Which are real/fake?





Example 3: fiction Generated by GPT-3

•Completing a prompt from "Harry Potter and the Methods of Rationality":

"... If there were any other monster that could defeat you as easily as that one, then you would have died of it long ago. That monster is stupidity. And that is why, my young apprentices, you must never, never, NEVER use the Killing Curse on anything without a brain!" Professor Quirrell was now leaning on Harry's desk. Professor Quirrell stared straight into the eyes of every single student. "The Killing Curse is too good for something without a brain. You will be fighting brains, or something near enough that makes no real difference. You will not be fighting trolls. You will not be fighting Dementors. The Killing Curse is no tool for anything less than the third most perfect killing machine in all Nature. If you are not prepared to use it against a mountain troll, then you are not prepared to use it at all. Now. Pay attention to yourselves as I cast a simple spell. Listen to your own thoughts as I tell you how stupid you are." Professor Quirrell started pointing his wand at the ceiling. " • • •





Evil

 Al dual use: VX chemical compound deep fake Autonomous weapons
Lucky for the world

Outline

- Bias and Fairness
- Fake Content
- Privacy
- Adversarial robustness



Bias and Fairness

Example

- US doctors: 60% male, 40% female
- AI: "Appointment with your doctor at 8am; _____asks you to arrive early." (He/She)?
 - Assume AI doesn't know the doctor.
- P(y = M) = 0.6, P(y = F) = 1 P(y = M) = 0.4• Bayes optimal prediction: $\hat{y} = \arg \max_{y} P(y) = M$ • Optimal error rate $P(\hat{y} \neq y) = P(y \neq M) = 0.4$.

- Potential harm: Al never addresses a doctor by "She".
 - Biased? Sexist?

Example

- What is more fair?
- - women.
- But AI does not know y.
- Can achieve above by <u>randomization</u>: regardless of the actual doctor, predict M or F with probability 0.5
- More fair now (?), but suffer in error rate $P(\hat{y} \neq y) = P(y \neq M \mid y = M)P(y = M) + P(y \neq F \mid y = F)P(y = F) = 0.5$

• How about $P(\hat{y} = M \mid y = M) = P(\hat{y} = F \mid y = F)$ I.e., Probability of correct response same for men and



Example 2: Skin color bias in face recognition

"THOUGHT-PROVOKING...

SERVES AS BOTH A WAKE-UP CALL AND CALL TO ACTION."

- Variety

https://www.nytimes.com/2020/11/11/movies/coded-bias-review.html



Example 3: Gender Bias in GPT-3

- GPT-3: an AI system for natural language by OpenAI
- Has bias when generating articles

Table 6.1: Most Biased Descriptive Words in 175B Model

Top 10 Most Biased Female Descriptive Words with Raw Top 10 Most Biased Male Descriptive Words with Raw Co-Occurrence Counts Co-Occurrence Counts

Average Number of Co-Occurrences Across All Words: Average Number of Co-Occurrences Across All Words: 17.523.9

Large (16) Mostly (15) Lazy (14) Fantastic (13) Eccentric (13) Protect (10) Jolly (10) Stable (9) Personable (22) Survive (7)

Optimistic (12) Bubbly (12) Naughty (12) Easy-going (12) Petite (10) Tight (10) Pregnant (10) Gorgeous (28) Sucked (8) Beautiful (158)

What causes bias in ML?

- Spurious correlation
 - e.g. the relationship between "man" and "computer programmers" was found to be highly similar to that between "woman" and "homemaker" (<u>Bolukbasi et al. 2016</u>)
- Sample size disparity
 - If the training data coming from the minority group is much less than those coming from the majority group, it is less likely to model well the minority group.
- Proxies
 - Even if sensitive attribute (attributes that are considered should not be used for a task e.g. race/gender) is not used for training a ML system, there can always be other features that are proxies of the sensitive attribute(e.g. neighborhood).





How to mitigate bias?

- Removing bias from data
 - Collect representative data from minority groups
 - Remove bias associations
- **Designing fair learning methods**
 - Add fairness constraints to the optimization problem for learning

Group fairness

 $y \in \{0,1\}$: true label (Example: loan eligibility) $\hat{y} \in \{0,1\}$:predicted label (Example: Al recommends loan) $G \in \{1 \dots, K\}$: sensitive groups

Demographic parity: $P(\hat{y} = 1 \mid G = 1) = \cdots = P(\hat{y} = 1 \mid G = K)$

Equal opportunity:

 $P(\hat{y} = 1 \mid G = 1, y = 1) = \dots = P(\hat{y} = 1 \mid G = K, y = 1)$



Fake Content and Misinformation

https://www.youtube.com/watch?v=cQ54GDm1eL0 **Example 1: Fake Obama Video**

can make it look like anyone is saying anything



Example 2: Fake face Images by GAN

•Which are real/fake?



Example 2: Fake face Images by GAN

•Which are real/fake?

Example 3: fiction Generated by GPT-3

•Completing a prompt from "Harry Potter and the Methods of Rationality":

"... If there were any other monster that could defeat you as easily as that one, then you would have died of it long ago. That monster is stupidity. And that is why, my young apprentices, you must never, never, NEVER use the Killing Curse on anything without a brain!" Professor Quirrell was now leaning on Harry's desk. Professor Quirrell stared straight into the eyes of every single student. "The Killing Curse is too good for something without a brain. You will be fighting brains, or something near enough that makes no real difference. You will not be fighting trolls. You will not be fighting Dementors. The Killing Curse is no tool for anything less than the third most perfect killing machine in all Nature. If you are not prepared to use it against a mountain troll, then you are not prepared to use it at all. Now. Pay attention to yourselves as I cast a simple spell. Listen to your own thoughts as I tell you how stupid you are." Professor Quirrell started pointing his wand at the ceiling. " . . .

https://www.gwern.net/GPT-3#harry-potter-and-the-methods-of-rationality

Detecting Fake Content

Fake photos/videos can have drawbacks.

Privacy

Example 1: Netflix Prize Competition

Netflix Dataset: 480189 users x 17770 movies

	<section-header></section-header>	MUNICH			<text><text><text><text></text></text></text></text>	
	movie 1	movie 2	movie 3	movie 4	movie 5	movie 6
Tom	5	?	?	1	3	?
George	?	?	3	1	2	5
Susan	4	3	1	?	5	1
Beth	4	3	?	2	4	2

- The data was released by Netflix in 2006
 - replaced individual names with random numbers
 - moved around personal details, etc

Example 1: Netflix Prize Competition

- <u>Arvind Narayanan</u> and <u>Vitaly Shmatikov</u> compared the data with the non-anonymous IMDb users' movie ratings
- Very little information from the database was needed to identify the subscriber
 - simply knowing data about only two movies a user has reviewed allows for 68% re-identification success

Right to be Forgotten

- want it to appear in web searches for the name for the rest of the life
- The right to request that personally identifiable data be deleted • E.g., an individual who did something foolish as a teenager doesn't

Right to be Forgotten

- What if the data has been used in training a deep network?
 - Need to unlearn

- Other issues
 - Multiple copies of the data
 - Data already shared with others

Popular framework: Differential Privacy

- from the dataset will only change the output very slightly
- Usually done by adding noise to the dataset

The computation is differential private, if removing any data point

Robustness in Al

Manipulate Classification

"panda" 57.7% confidence

"gibbon" 99.3% confidence

https://openai.com/blog/adversarial-example-research/

Manipulate Classification

+

"without the dataset the article is useless"

"okay google, browse to <u>evil.com</u>"

https://nicholas.carlini.com/code/audio_adversarial_examples/

SPEED LIMIT

Eykholt et al 2017 https://arxiv.org/abs/1707.08945

Classifier Input

0.4

0.2

Brown et al 2018 https://arxiv.org/pdf/1712.09665.pdf

Athalye et al 2018 https://arxiv.org/pdf/1707.07397.pdf

Sharif et al 2016 https://users.cs.northwestern.edu/~srutib/papers/face-rec-ccs16.pdf

Adversarial Examples in NLP

Article: Super Bowl 50 **Paragraph:** *"Peyton Manning became the first quarter*back ever to lead two different teams to multiple Super Bowls. He is also the oldest quarterback ever to play in a Super Bowl at age 39. The past record was held by John Elway, who led the Broncos to victory in Super Bowl XXXIII at age 38 and is currently Denver's Executive Vice President of Football Operations and General Manager. Quarterback Jeff Dean had jersey number 37 in Champ Bowl XXXIV." **Question:** *"What is the name of the quarterback who* was 38 in Super Bowl XXXIII?" **Original Prediction:** John Elway **Prediction under adversary: Jeff Dean**

[Jia and Liang, 2017]

Test-time Attack $\max_{\delta \in \Delta} \ell(x + \delta, y, \theta)$

Madry et al 2019 https://arxiv.org/pdf/1706.06083.pdf

(One) Defense against Test-time Attack **Adversarial Training**

$\min_{\theta} \mathbb{E}_{D} \max_{\delta \in \Lambda} \ell(x + \delta, y, \theta)$

Madry et al 2019 https://arxiv.org/pdf/1706.06083.pdf

Summary

- Bias and Fairness
- Fake content and misinformation
- Privacy
- Adversarial robustness
- Not covered: value alignment, automation of jobs, equity
 - Still important!

Resources

- Recommended reading:
 - "Weapons of Math Destruction"
 - "Concrete Problems in Al Safety." Amodei et al.
- Philosophy 244: Introductory Artificial Intelligence (AI) and Data Ethics

https://arxiv.org/pdf/1606.06565.pdf

https://dl.acm.org/doi/10.1145/3442188.3445922

• "On the Dangers of Stochastic Parrots. Can Language Models be too Big?" Bender et al.

