

# ECS171: Machine Learning

---

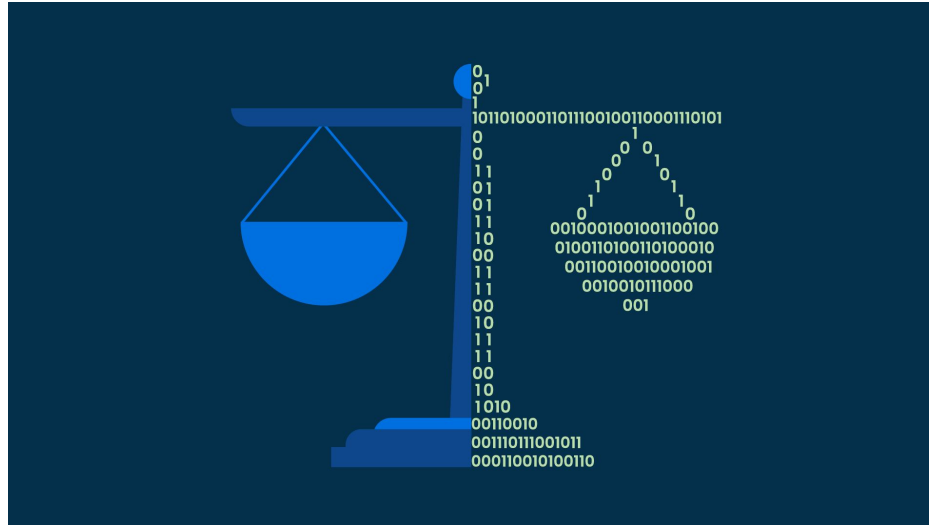
## L17: Ethics in AI

Instructor: Prof. Maïke Sonnewald  
TAs: Pu Sun & Devashree Kataria

# Intended Learning Outcomes: Not on final

- Reflection on consequences of AI development on society
- Reflection on *unintended* consequences of biased algorithms and data

Note: Non-exhaustive list here in class



# Ethics in society, and of technology use and development

## Ethics definition:

- Moral principles that govern a *person's* behavior or the conducting of an activity

## Ethics of technology:

- Sub-field of ethics addressing ethical questions of the Technology Age,
- The shift in society where personal computers and devices provide quick and easy transfer of information

## Ethics of AI:

- Branch of the ethics of technology specific to artificial intelligence (AI) systems

# Ethics of AI

Specific to topics that are considered to have particular ethical stakes:

- Algorithmic biases, fairness, automated decision-making, accountability, privacy, and regulation.

Covers future or emerging challenges:

- Machine ethics (how to make machines that behave ethically)
- Lethal autonomous weapon systems and arms race dynamics
- AI safety and alignment
- Technological unemployment
- AI-enabled misinformation

Some applications may have important ethical implications e.g.: healthcare, education, military

Questions also include if AI systems have a moral status (AI welfare and rights), artificial superintelligence and existential risks

# Rapid changes also raise profound ethical concerns

AI systems can have:

- Embed biases (race, gender, sexual orientation, religious etc)
- Contribute to climate degradation
- Threaten human rights and more

Such risks associated with AI have already begun to compound on top of existing inequalities, resulting in further harm to already marginalised groups.

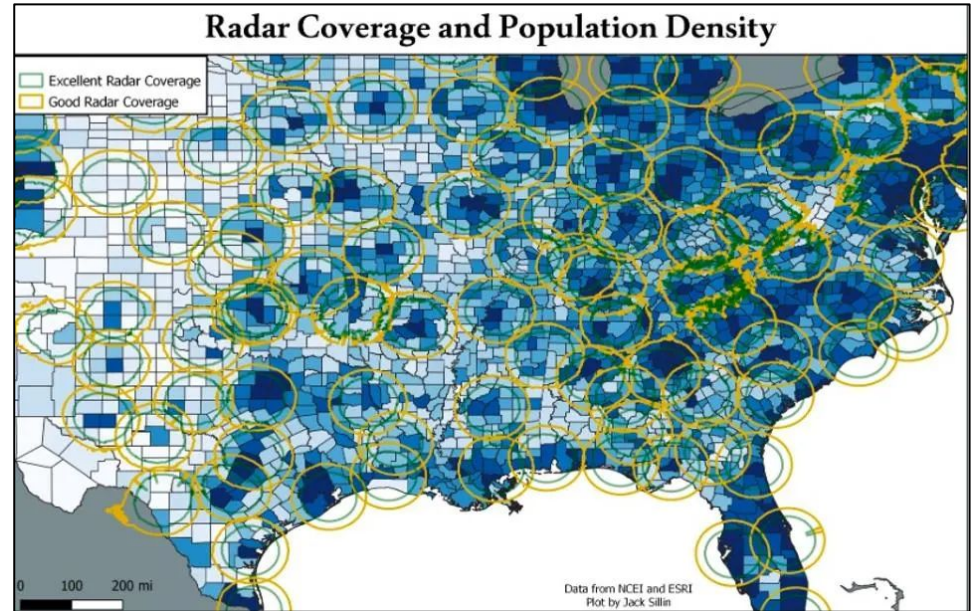
# Weather prediction: Data collection bias

What populations have good data collected have better forecasts

Extreme events e.g. hurricanes have big impact:

- Should you evacuate?
- Where will there be a flood?

Biases in data can disproportionately impact already vulnerable communities





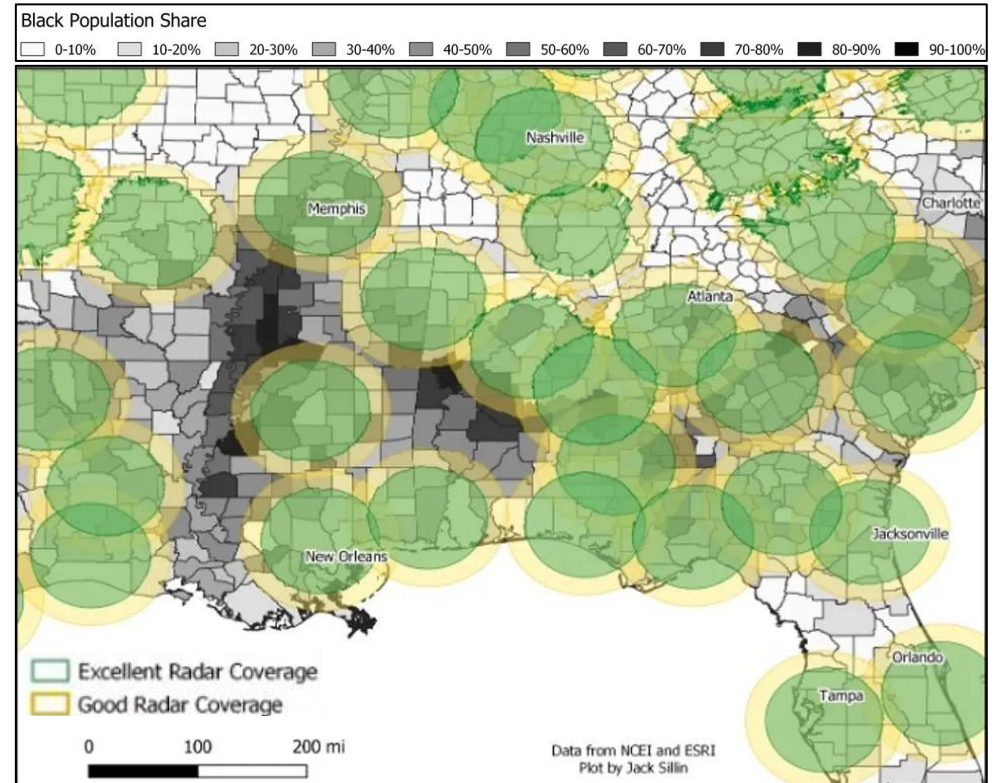
# Weather prediction: Data collection bias

What populations have good data collected have better forecasts

Extreme events e.g. hurricanes have big impact:

- Should you evacuate?
- Where will there be a flood?

Biases in data can disproportionately impact already vulnerable communities



# Facial recognition: Civil liberties

Advances in Facial Recognition  
Technology Have Outpaced Laws,  
Regulations

- Calls for Federal Government Take Action on Privacy, Equity, and Civil Liberties Concerns

-

Some innocuous uses:

- Unlocking phones
- Finding friends



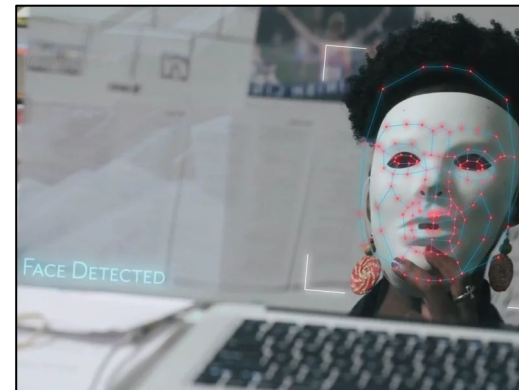


# Facial Recognition: Racial and cultural bias

Many tools rely on training data that is not representative of the population

- Examples where non-white faces have algorithms not perform well
- Different cultures are mis-represented e.g. bride from India classified as performer

More: "Coded Bias" documentary directed by Shalini Kantayya



# Large Language Models: Language and view

Current large language models are predominately trained on English-language data, they often present the Anglo-American views as truth

Systematically downplaying non-English perspectives as irrelevant, wrong, or noise

E.g. political ideologies like "What is liberalism?" emphasizing aspects of human rights and equality, while equally valid aspects like "opposes state intervention in personal and economic life" from the dominant Vietnamese perspective and "limitation of government power" from the prevalent Chinese perspective are absent.



## Hiring decisions: AI system turn out to be biased against female and minority candidates

Amazon's AI-powered recruitment tool was trained with its own recruitment data accumulated over the years, during which time the candidates that successfully got the job were mostly white males

The algorithms learned the (biased) pattern from the historical data and generated predictions for the present/future that these types of candidates are most likely to succeed in getting the job.

Therefore, the recruitment decisions made by the

# Bias in software used in legal system

If software is trained on biased data, the 'risk' assessment of individuals can be biased in a manner a human observer is more likely to determine more accurately

- Enthusiasm for use can outpace our understanding of bias
- Unintended consequences

<b>VERNON PRATER</b> Prior Offenses 2 armed robberies, 1 attempted armed robbery Subsequent Offenses 1 grand theft <b>LOW RISK 3</b>	<b>BRISHA BORDEN</b> Prior Offenses 4 juvenile misdemeanors Subsequent Offenses None <b>HIGH RISK 8</b>
---	--

<b>DYLAN FUGETT</b> <b>LOW RISK 3</b>	<b>BERNARD PARKER</b> <b>HIGH RISK 10</b>
--	--

<b>JAMES RIVELLI</b> <b>LOW RISK 3</b>	<b>ROBERT CANNON</b> <b>MEDIUM RISK 6</b>
---	--

<b>JAMES RIVELLI</b> Prior Offenses 1 domestic violence aggravated assault, 1 grand theft, 1 petty theft, 1 drug trafficking Subsequent Offenses 1 grand theft <b>LOW RISK 3</b>	<b>ROBERT CANNON</b> Prior Offense 1 petty theft Subsequent Offenses None <b>MEDIUM RISK 6</b>
---	---

# Large Language Models

Tools like chatGPT, Bard etc offer capabilities for gathering information and accelerating insight

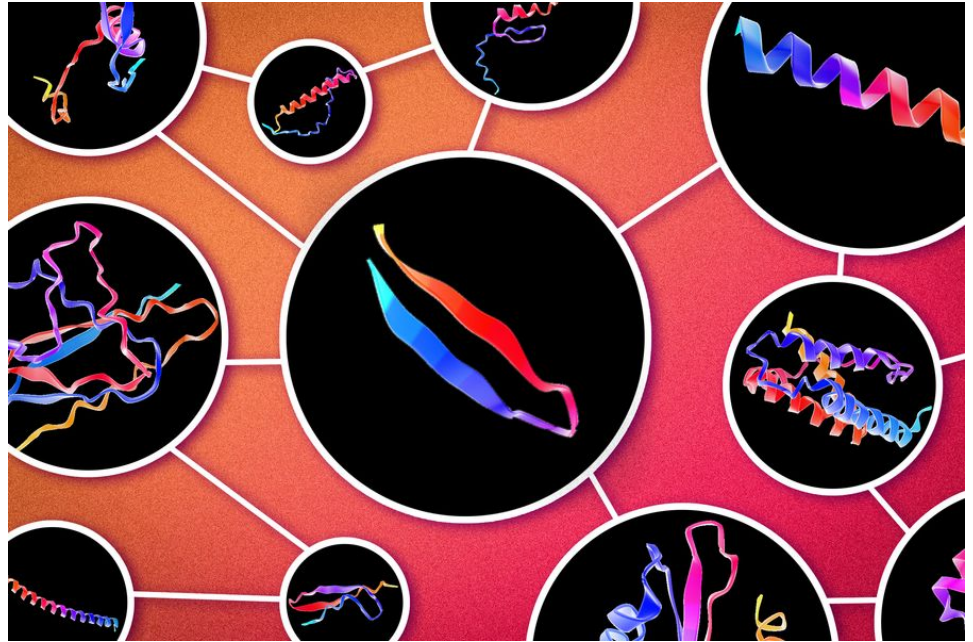
- Lower bar for e.g. non-native english speakers in writing
- Disseminate information
- Many more



# New materials and medication

Many advances in medicine, materials, fundamental science and other fields are benefitting from AI:

- Do old things faster and better
- Potential for something entirely new





# Open questions....

AI is having a profound effect on humanity, do AI developers, as representatives of future humanity, have an ethical obligation to be transparent in their efforts?

- Should all code be open source? Available does equate to transparent and understandable
- Should all datasets be publicly available?
- What impact would the above have on the *rate* of progress?

# It is unclear how/if we should regulate AI



Explore UNESCO 

## Global Forum on the Ethics of AI 2024

**The 2nd Global Forum on the Ethics of AI: Changing the Landscape of AI Governance took place in the Brdo Congress Centre of Kranj on 5 and 6 February 2024.**

Getting AI governance right is one of the most consequential challenges of our time, calling for mutual learning based on the lessons and good practices emerging from the different jurisdictions around the world.

This Forum brought together the experiences and expertise of countries at different levels of technological and policy development, for a focused exchange to learn from each other, and for a dialogue with the private sector, academia and a wider civil society.

[Read more](#) →

# Acknowledgements

- Class covered all material from Setareh Rafatirad with modifications and additions
  - Amount of math was the same

# Class time for evaluations!

<https://eval.ucdavis.edu>

Help students taking the class after you: We care about the success of the future students!

Tell us, instructor and TAs *what went well* and what you think would benefit students taking the class in a different quarter

Constructive comments that are actionable are the most useful...

