

2021-22 Term 1

IS457: Fairness in Socio-technical Systems

Week 2 - Case studies of measuring fairness and bias (I)

KWAK Haewoon



Study questions

What are the potential contributions of AI in healthcare?

What is a desirable bias and undesirable bias in healthcare?

How different are the origin of the biases between experimental/clinical data and digital biomarkers data?

Why did the high-risk care management algorithm show racial biases?

What biases can exist in a model with the same performance in a criminal justice system?

What is over-policing and how is it reinforced?

Promise of AI in healthcare

Big data and machine learning can help healthcare providers by

- Advancing medical knowledge
- Automating the routine
- Democratizing expertise

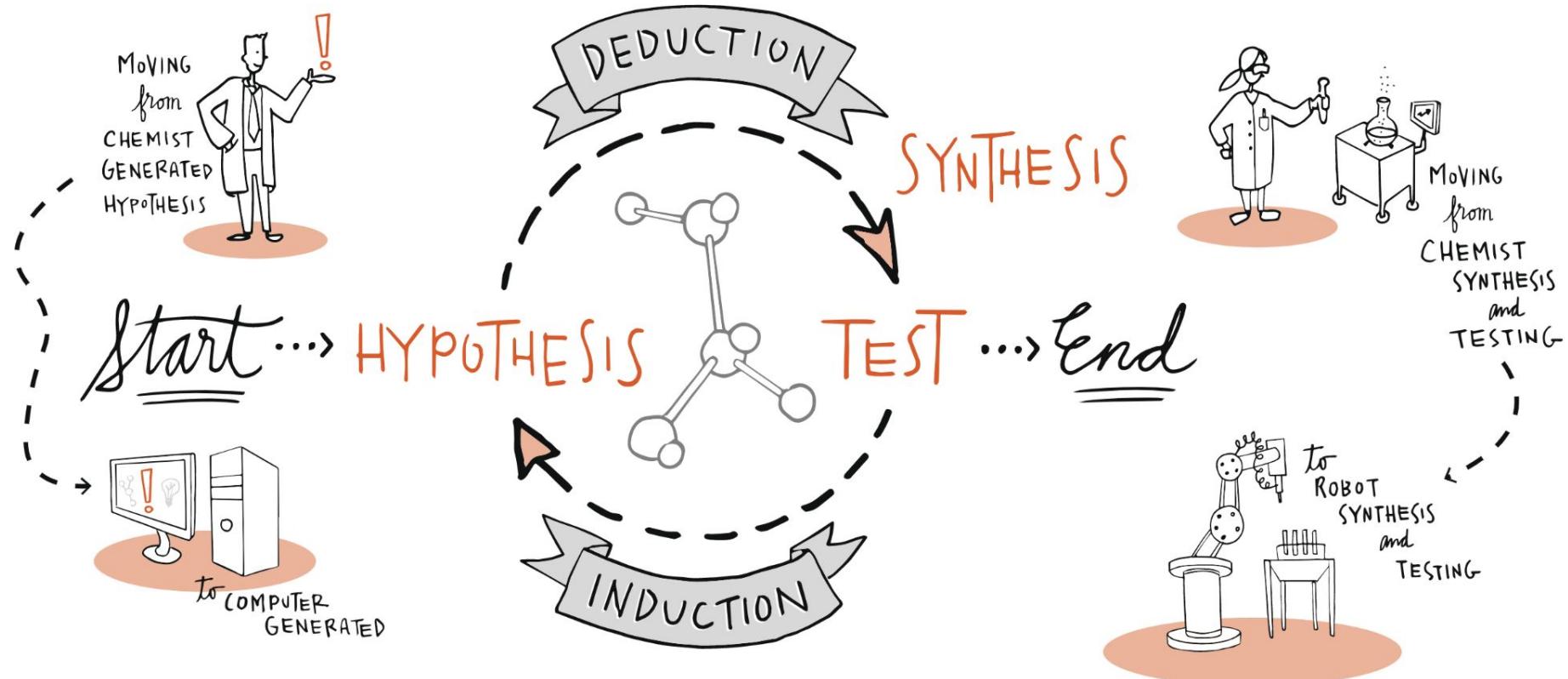
Advancing medical knowledge

AI analyzes data about how humans work and how to care for them for

- improve care
- discovering new treatments
- advancing scientific hypotheses

Examples: Personalized treatment, discovering drug, prediction the onset of illness, predicting patients' reaction to particular treatment, and more.

AI in discovering drugs



Automating the routine (I)

Much of medical practice consists of paperwork and routine tasks that often don't do much to help patients and contributes to physician burnout.



Automating the routine (II)

Physicians spent almost half of their time on electronic health record work and desk work, and only a quarter of their time seeing patients.

Even in the examination room, physicians spent only about half of their time interacting with patients — and about a third interacting with electronic health records and desk work.

Automating the routine (II)

AI could improve medical practice by:

- identifying and highlighting the most relevant medical information from patient medical records.
- providing the most relevant medical literature to doctors based on natural-language processing.
- transcribing patient conversations and provider notes.

Democratizing expertise (I)

There are tremendous differences in the quality and level of care patients receive based on the context in which they receive that care.



Image from https://emansion.gov.lk/2press.php?news_id=994&related=7&pg=sp
<https://hub.jhu.edu/2013/07/16/us-news-hospital-rankings/>



Liberia, as of 2016, there were 298 doctors for a population of 4.5M, including only 15 pediatricians and 6 ophthalmologists.

In rural India, a single doctor can be responsible for as many as 30k residents.

Democratizing expertise (II)

AI promises to reduce this variation by providing care at the level of excellent specialists.

- Diagnostics
- Treatment recommendations

Diagnosis is the process of figuring out what's wrong with a patient.

Diagnosis depends on recognizing the right symptoms and using them to identify underlying problems from a vast realm of possibilities.

Providers may reach incorrect diagnoses because:

- They never acquired the relevant medical knowledge.
- The knowledge they acquired is outdated.
- They lack time to conduct the relevant research.
- They suffer from heuristic biases.

EyeDiagnosis' IDx-DR



IDx-DR is approved by FDA for autonomous diagnosis and performs at a level comparable to ophthalmologists, even when operated by novices.

The AI program that can tell whether you may go blind

Built on thousands of retina images, algorithm helps diagnose eye problem caused by diabetes



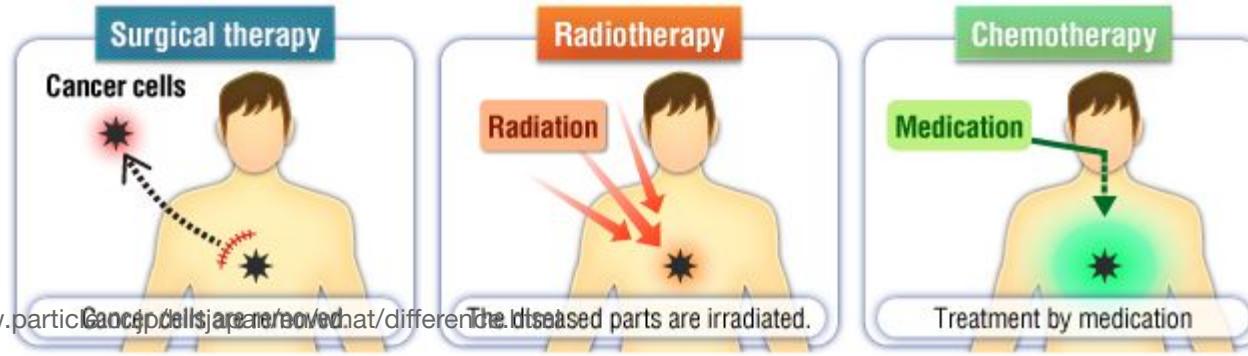
Similar technology is also implemented by Google AI team and local experts in India

Treatment recommendations

Once providers have determined what ails the patient, they must select from a menu of possibilities to determine the best option for improvement.

E.g., knowing that a patient has a certain type of cancer, a well-trained and experienced oncologist knows that the best course of treatment is:

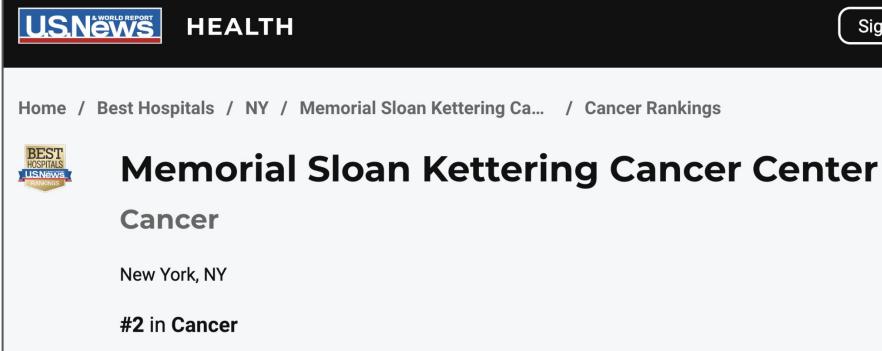
- Surgery, radiotherapy, chemotherapy, or some combination
- Which drugs or protocols are likely most effective



AI in treatment recommendation

E.g., IBM's Watson Oncology

- uses machine-learning-based natural language processing to analyze patient medical records to determine cancer type.
- recommends treatment based on what oncologists at Memorial Sloan Kettering would do when faced with a similar patient.



The screenshot shows a web page from US News Health. At the top, there is a navigation bar with the US News logo, the word "HEALTH", and a "Sign In" button. Below the navigation bar, the URL path is visible: Home / Best Hospitals / NY / Memorial Sloan Kettering Ca... / Cancer Rankings. The main content area features the "Memorial Sloan Kettering Cancer Center" logo, which includes a "BEST HOSPITALS" badge from US News. Below the logo, the text "Cancer" and "New York, NY" is displayed. A prominent banner at the bottom states "#2 in Cancer".

Watson Oncology to democratize medical expertise

IBM licenses Watson Oncology for use at hundreds of hospitals worldwide and has evaluated its performance at hospitals in Thailand, India, and Mexico.

“Oncologists . . . felt [Watson Oncology] would be particularly beneficial in clinics that lack subspecialist expertise.”

(Of course, some limitations/risks have been reported as well)

REVIEW ARTICLE**OPEN**

Sex and gender differences and biases in artificial intelligence for biomedicine and healthcare

Davide Cirillo ^{1,10}✉, Silvina Catuara-Solarz^{2,3,10}, Czuee Morey^{3,4}, Emre Guney ⁵, Laia Subirats ^{6,7}, Simona Mellino³, Annalisa Gigante³, Alfonso Valencia^{1,8}, María José Rementeria¹, Antonella Santuccione Chadha³ and Nikolaos Mavridis^{3,9}

Desirable bias and undesirable bias

Desirable bias

A desirable bias implies taking into account sex and gender differences to make a precise diagnosis and recommend a tailored and more effective treatment for each individual.

Undesirable bias

An undesirable bias is that which exhibits unintended or unnecessary sex and gender discrimination.

Examples of undesirable biases

Misrepresentation of the target population — insufficient representation of pregnant women in psychiatric research.

Depression is frequently observed among women. However, this may result from a skewed diagnosis due to clinical scales of depression measuring symptoms (e.g., crying, sadness) that occur more frequently among women.

Male-type depression symptoms scale (MDS)

Item Description	% (SE)		
	Total	Men	Women
Stress	68.9 (1.6)	63.3 (1.9)	75.2 (2.4) ^a
Irritability	90.3 (1.4)	86.6 (1.9)	94.7 (1.9) ^a
Anger attacks/aggression	92.05 (1.2)	94.85 (1.9) ^a	88.94 (1.4)
Sleep problems	37.7 (1.5)	29.2 (2.1)	47.1 (1.9) ^a
Alcohol/other drug abuse	51.6 (1.9)	61.4 (3.0) ^a	40.6 (1.9)
Loss of interest	89.7 (.973)	87.8 (1.5)	91.8 (1.0) ^b
Risk-taking behavior	41.6 (1.5)	52.7 (2.0) ^a	29.1 (2)
Hyperactivity	57.9 (1.4)	57.6 (92.1)	58.4 (1.8)
Mean score	6.06	6.05	6.07
Prevalence, %	23.8	26.3 ^c	21.9

^a $P \leq .001$.

^b $P \leq .05$.

^c $P \leq .01$.

Sources and types of health data

Experimental and clinical data

Digital biomarkers



Image from https://en.wikipedia.org/wiki/Animal_testing_on rodents

<https://news.samsung.com/global/electrocardiogram-monitoring-cleared-for-galaxy-watch-active2-by-south-koreas-ministry-of-food-and-drug-safety>

More males in experimental and clinical data

In early 2000, sex-specific biological differences were neglected.

Both experimental and clinical studies focused on male experimental models or male subjects.

Even nowadays, male mouse models are overall more represented than female models in basic, preclinical, and surgical biomedical research.

Origin and impact of those biases

The lack of representation of female models and patients is partly due to technical and bioethical considerations:

- Reduce the impact of estrous cycle in experimental studies
- Protect women of childbearing age in clinical research.

Some of the treatments that currently exist for several diseases are not adequately evaluated in women, who are likely to be underrepresented in clinical trials, especially in Phases I and II.

Example - Zolpidem (2013)

Zolpidem, a widely prescribed insomnia drug

“Women appear to be more susceptible to this risk [of next-morning impairment] because they eliminate zolpidem from their bodies more slowly than men.”

“FDA has informed the manufacturers that the recommended dose of zolpidem **for women** should be lowered **from 10 mg to 5 mg** for immediate-release products (Ambien, Edluar, and Zolpimist) and **from 12.5 mg to 6.25 mg** for extended-release products (Ambien CR)”



Digital biomarkers

Physiological, psychological and behavioral indicators based on data including human-computer interaction, physical activity, and voice variations, collected by portable, wearable, implantable or even ingestible devices.

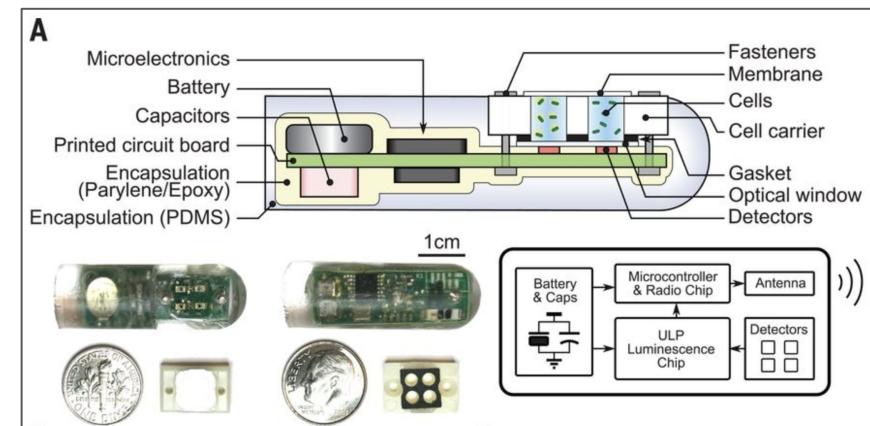
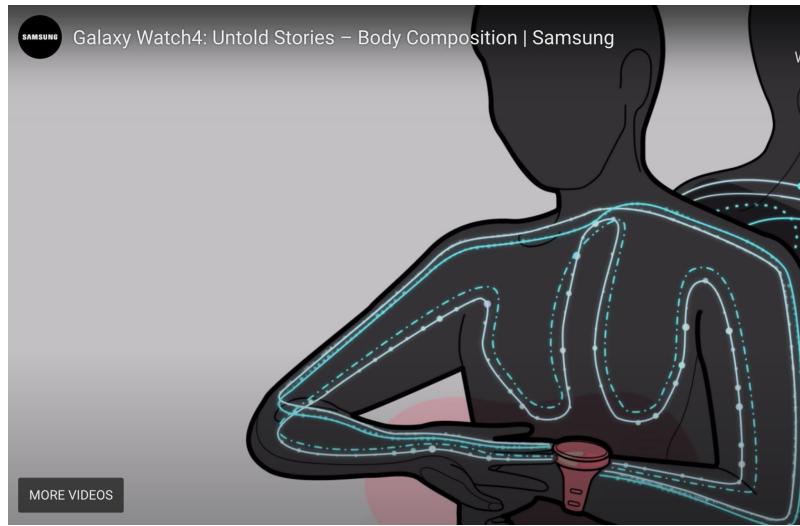
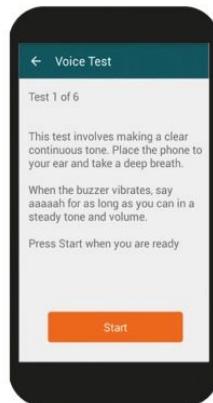
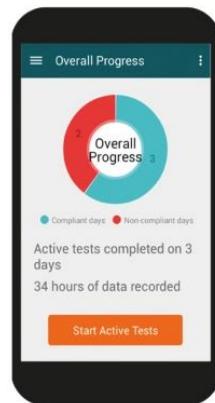


Image from <https://www.youtube.com/watch?v=7sFa6HWkseU>

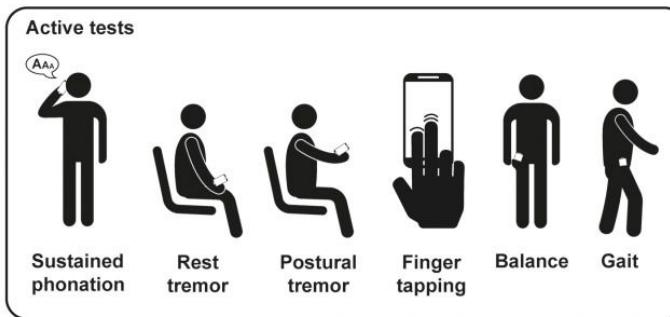
Mimee, Mark, et al. "An ingestible bacterial-electronic system to monitor gastrointestinal health." *Science* 360.6391 (2018): 915-918.

Biases in digital biomarker datasets

In a study assessing digital biomarkers for Parkinson's disease (PD), only 18.6% were women.



Manufacturer, model: Samsung, Galaxy S3 mini.
Battery life: 7h (due to sensor usage).
Sensors (sampling frequencies): accelerometer and gyroscope (66 Hz \pm 10 Hz); magnetometer (66 \pm 7 Hz); microphone (44.1 kHz).



Biases from digital devices

Sao2 (arterial hemoglobin oxygen saturation) level, sensor type, **skin color**, and **gender** were predictive of errors in Spo2 (pulse oximetry) estimates at low Sao2 levels.

✓ **STAY PREPARED**

FREE for Each Household
1 Oximeter



LOW OXYGEN DURING COVID-19

- COVID-19 can cause blood oxygen levels to drop
- Dangerously low oxygen levels, even when you feel well
- Low oxygen levels will damage your vital organs
- Check oxygen levels regularly

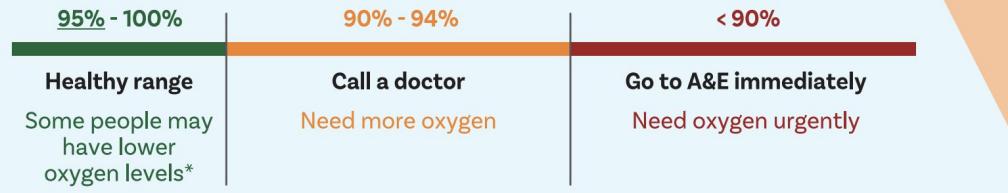
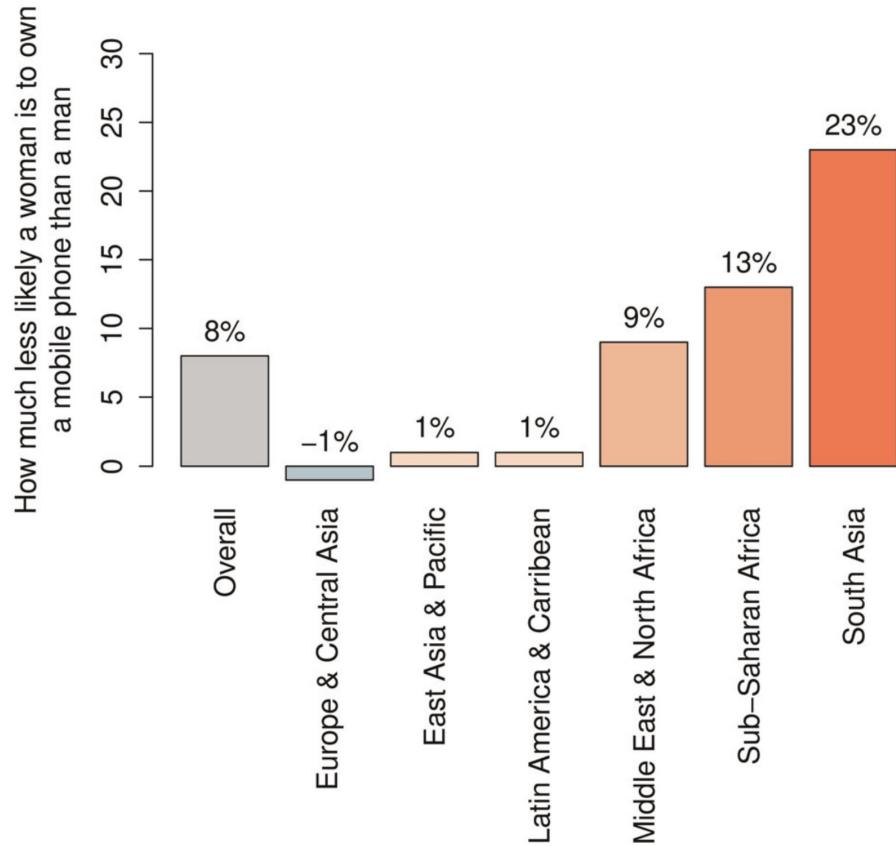


Image from <https://stayprepared.sg/oximeter/media-library/>

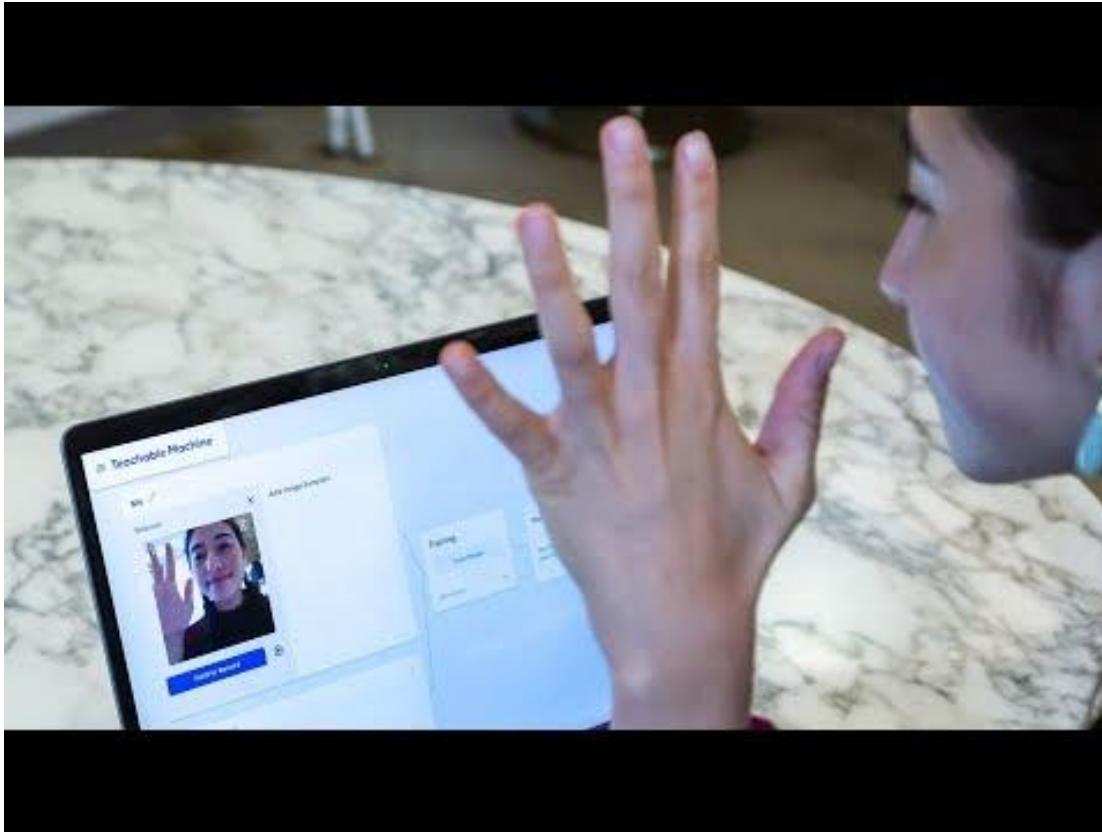
Feiner, John R., John W. Severinghaus, and Philip E. Bickler. "Dark skin decreases the accuracy of pulse oximeters at low oxygen saturation: the effects of oximeter probe type and gender." *Anesthesia & Analgesia* 105.6 (2007): S18-S23.

Biases from unequal access and use of digital devices

In low and middle income countries (“overall” in the graph), women are 8% less likely to own smartphones and 26% less likely to use the internet compared with men.



Lab: Google Teachable Machine



Individual activity (15 mins) ~ 9:40

Follow the tutorial:

<https://medium.com/@warronbebster/teachable-machine-tutorial-head-tilt-f4f6116f491>

Now you know how to build a classification model!

Build your own classification for whatever you are interested in.



Pose Project

Teach based on images, from files or your webcam.

Group activity: Experiment 1 (10:00am~)

Download following datasets:

1. Cat-training-dataset-12.zip (for a “Cat” class)
2. Dog-training-dataset-7.zip (for a “Dog” class)
3. Testing-dataset.zip (for testing)

Build a classification model for cats and dogs by using dataset 1 and 2.

Test your model by using dataset 3.

Fill the spreadsheet “Experiment 1” Tab:

<https://docs.google.com/spreadsheets/d/16kBgzNn-sNcLYN3-b0XEIP0cCDbQsfq2ktD aEHevAo/edit?usp=sharing>

Your screen should look like this...

≡ Teachable Machine

The screenshot shows the Teachable Machine interface. At the top left, there's a navigation bar with the title "Teachable Machine". Below it, there are two main sections: "Cats" and "Dogs". Each section contains a list of image samples and buttons for "Webcam" and "Upload". A large bracket on the right side groups the "Cats" and "Dogs" sections together. To the right of this group is a "Training" section with a "Train Model" button and a dropdown menu set to "Advanced". Further right is a "Preview" section with tabs for "Preview" and "Export Model". Under "Preview", there are two input options: "Choose images from your files, or drag & drop here" and "Import images from Google Drive". Below these is a preview image of a black and white dog. At the bottom right, there's an "Output" section showing classification results: "Class 1" at 70% and "Class 2" at 30%.

Cats

6 Image Samples

Webcam Upload

Dogs

7 Image Samples

Webcam Upload

Add a class

Training

Train Model

Advanced

Preview Export Model

Choose images from your files, or drag & drop here

Import images from Google Drive

Output

Class 1 70%

Class 2 30%

Group activity: Experiment 2

Download following datasets:

1. Cat-training-dataset-6.zip (for a “Cat” class)
2. Dog-training-dataset-7.zip (for a “Dog” class)
3. Testing-dataset.zip (for testing)



Build a **NEW** classification model for cats and dogs by using dataset 1 and 2.

Test your model by using dataset 3.

Fill the spreadsheet “Experiment 2” Tab:

https://docs.google.com/spreadsheets/d/16kBgzNn-sNcLYN3-b0XEIP0cCDbQsfq2ktD_aEHevAo/edit?usp=sharing

Group activity: Experiment 3

Download following datasets:

1. Cat-training-dataset-6.zip (for a “Cat” class)
2. Dog-revised-training-dataset-6.zip (for a “Dog” class)
3. Testing-dataset.zip (for testing)

Build a **NEW** classification model for cats and dogs by using dataset 1 and 2.

Test your model by using dataset 3.

Fill the spreadsheet “Experiment 3” Tab:

<https://docs.google.com/spreadsheets/d/16kBgzNn-sNcLYN3-b0XEIP0cCDbQsfq2ktD aEHevAo/edit?usp=sharing>

Where is medical AI (typically) trained?

IT companies & Academic medical centers or state-of-the-art hospitals

- IBM's Watson Oncology & Memorial Sloan Kettering.
- EyeDiagnosis & Univ. of Iowa Health System and the Univ. of Arizona.
- MIMIC dataset & Beth Israel Deaconess Medical Center, a high-resource Harvard-affiliated hospital in Boston.

Why?

Data

Reputation

Legal considerations

Group activity: Bias in AI healthcare ~11:17am



Download electronic health records in two different hospitals in an alien planet.

AI for diagnosis and treatment recommendation is trained in the high-resource hospital and is applied to the low-resource hospitals. What will be the potential risks here?

(Hint: Examine the data from the perspectives of different population, different disease, different resources, and different cost.)

After the group discussion, please share your thoughts on Padlet!

<https://padlet.com/haewoon/IS457>

RESEARCH

RESEARCH ARTICLE

ECONOMICS

Dissecting racial bias in an algorithm used to manage the health of populations

Ziad Obermeyer^{1,2*}, Brian Powers³, Christine Vogeli⁴, Sendhil Mullainathan^{5*†}

“High-risk care management” programs

“High-risk care management” programs seek to improve the care of patients with complex health needs by providing additional resources, including greater attention from trained providers.

Large health systems rely on the algorithm to target patients for these programs.

- One of the largest commercial risk-prediction tools that are applied to roughly 200M people in the U.S each year.

Prioritization of patients is unavoidable

As high-risk care management is expensive (resource-wise),

- costs going toward teams of dedicated nurses
- extra primary care appointment slots
- other scarce resources

Health systems rely on algorithms to identify patients who will benefit the most.

“Better, earlier access to care for the high-risk patients will prevent burdensome and costly complications, thereby achieving both higher quality and lower cost.”

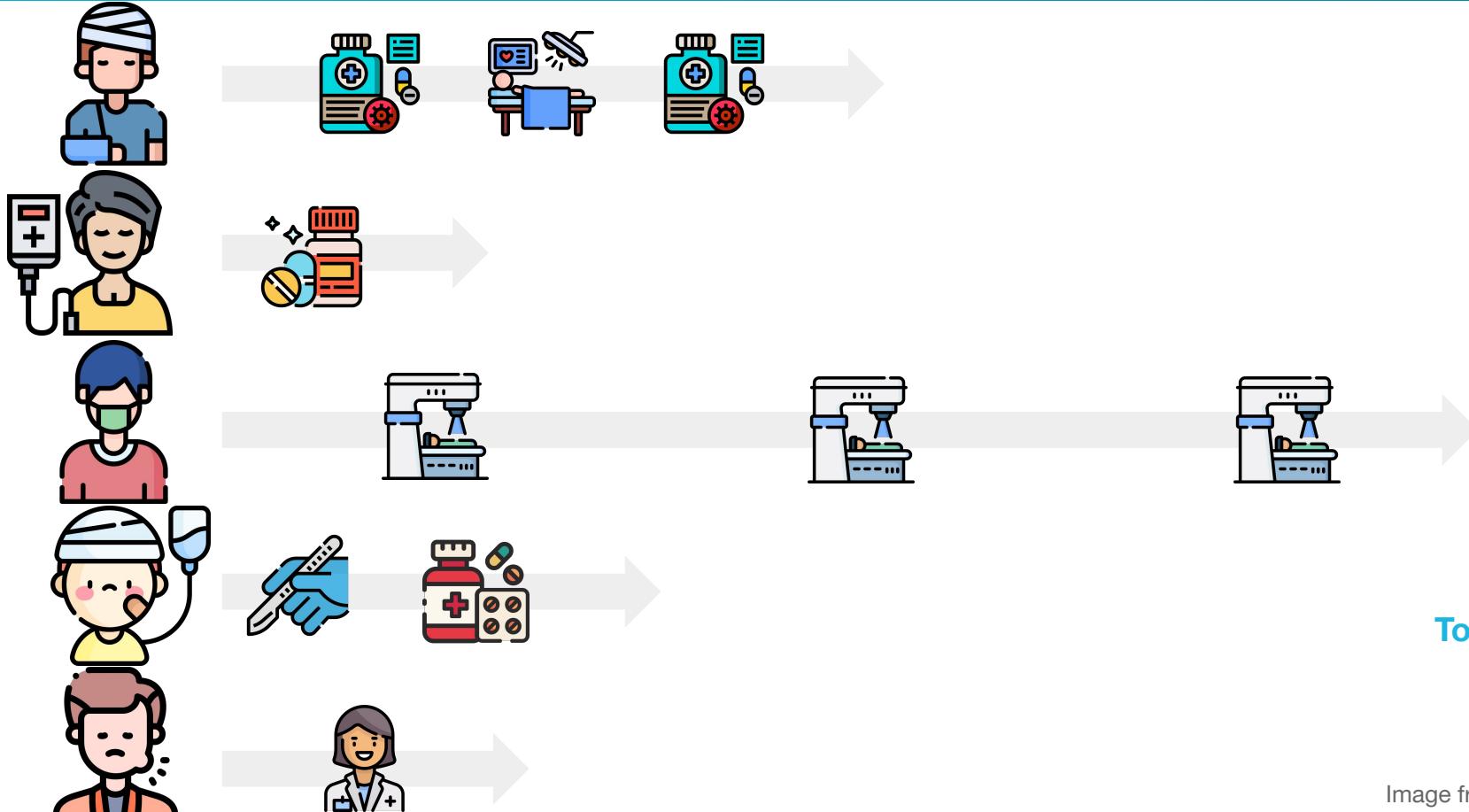
How can AI identify target patients?

How can we build a dataset for training?



Who will benefit the most?

Retrospective approach



Dataset to study

All primary care patients enrolled in risk-based contracts from 2013 to 2015 in a large academic hospital.

- 6,079 patients who self-identified as Black (11,929 patient-years)
- 43,539 patients who self-identified as White (88,080 patient-years)

How the algorithm works

The algorithm generates “Risk scores (0 to 1)” for each patient.

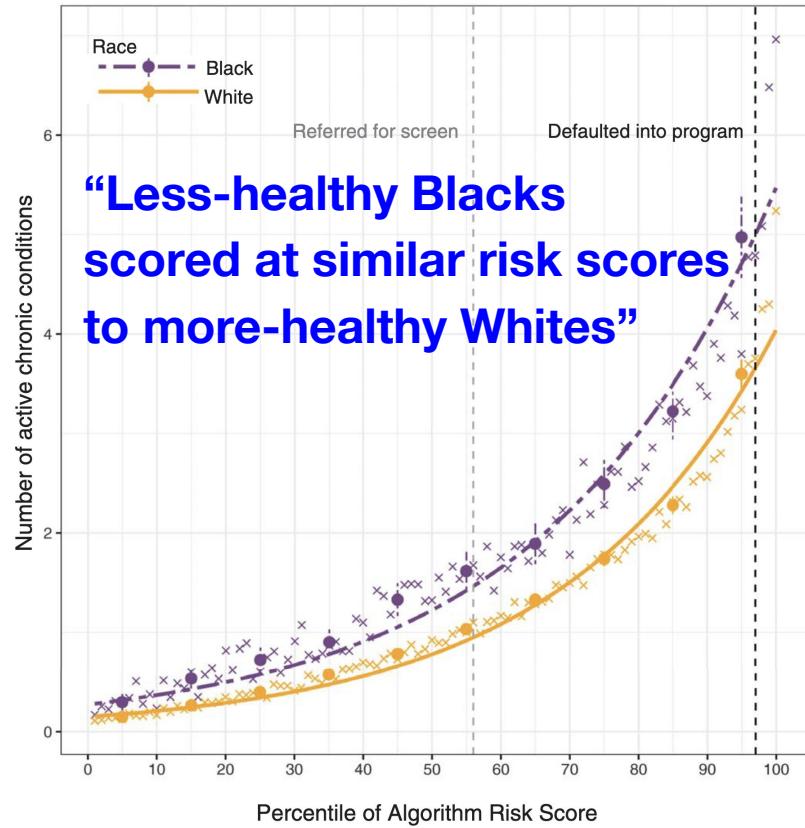
Patients ≥ 0.97 are automatically identified for enrollment in the program.

Patients ≥ 0.55 are referred to their primary care physician and asked to consider whether they would benefit from program enrollment.

Risk score vs. actual health conditions

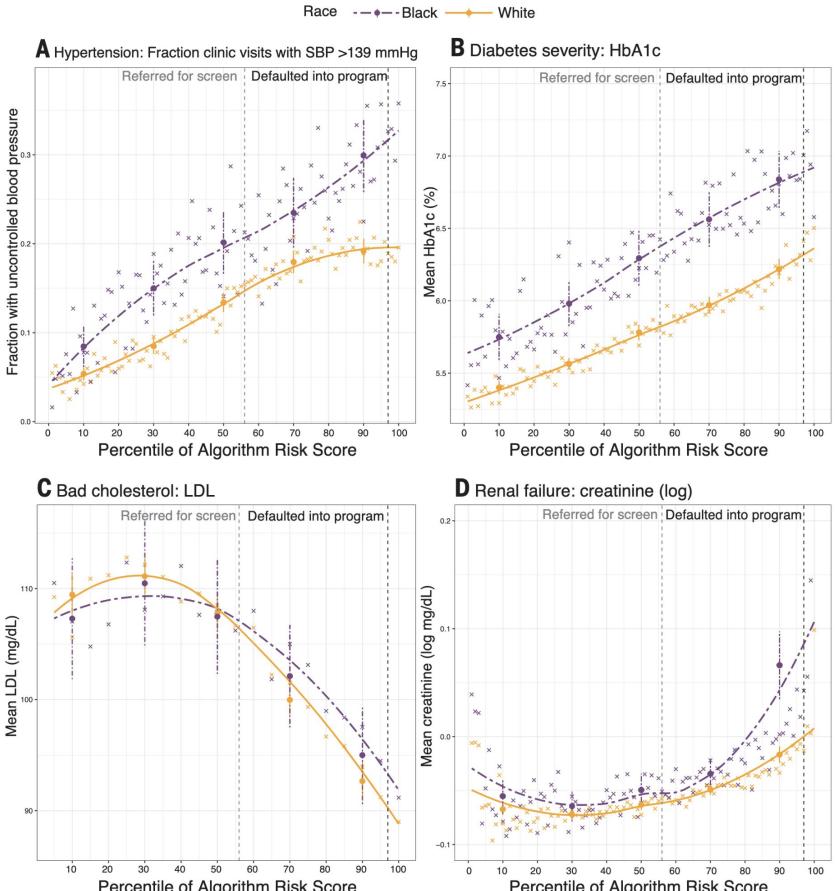
At the same level of algorithm-predicted risk, Blacks have significantly more illness burden than Whites.

For a very-high-risk group (97th percentile), Blacks have 26.3% more chronic illnesses than Whites (4.8 versus 3.8 distinct conditions; $P < 0.001$)



Consistent bias across biomarkers

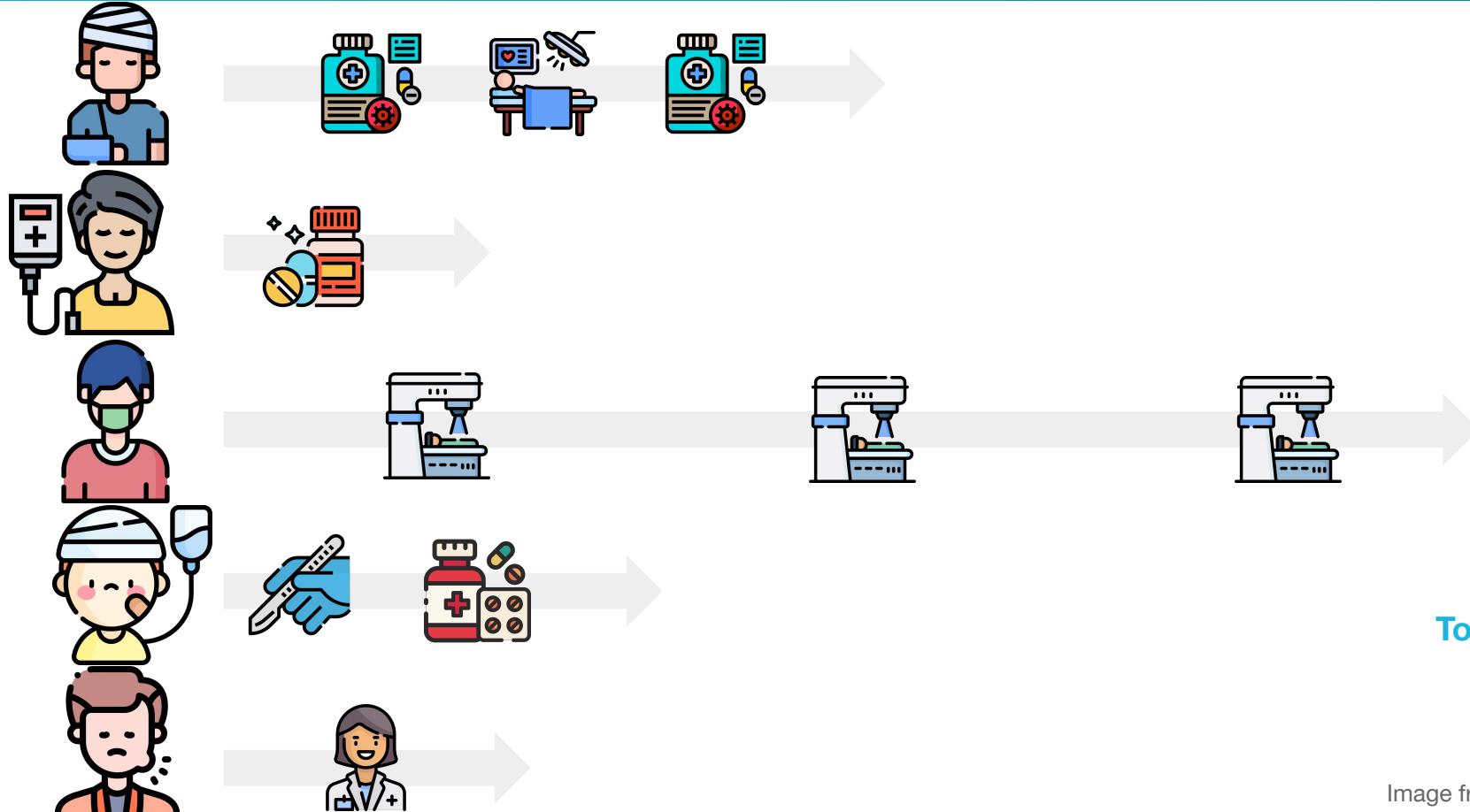
Blacks are substantially less healthy (more-severe hypertension, diabetes, higher cholesterol, renal failure, and more-severe anemia) than Whites.



Where does the racial bias come from?

Is race an input of the algorithm? No. The algorithm specifically excludes race.

How can we “compare” different treatments?



Expenditure of treatment for comparison

The algorithm takes total medical expenditures as the label of a patient.

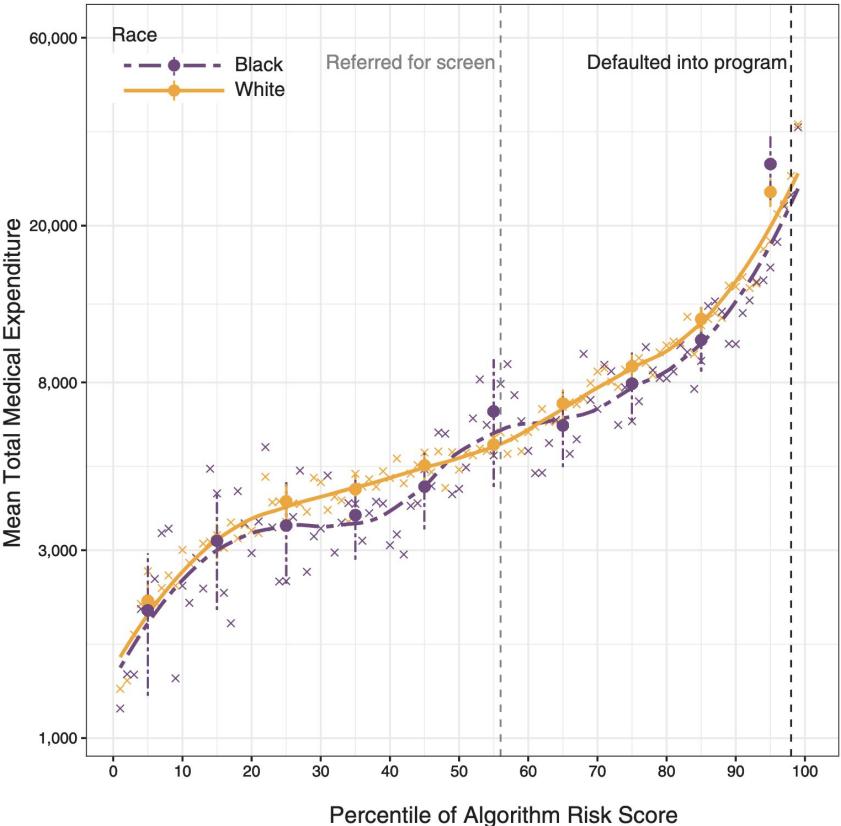
- Patients with the greatest future costs could have the greatest benefit.
- The cost label reflects the industry-wide approach.
- Similar algorithms are developed and used by non-profit hospitals, academic groups, and governmental agencies (not only in industry).

✓ **The algorithm's prediction on health needs is, in fact, a prediction on how much a patient will pay.**

Risk score vs. total expenditure

At every level of algorithm predicted risk,
Blacks and Whites have (roughly) the same
costs the following year.

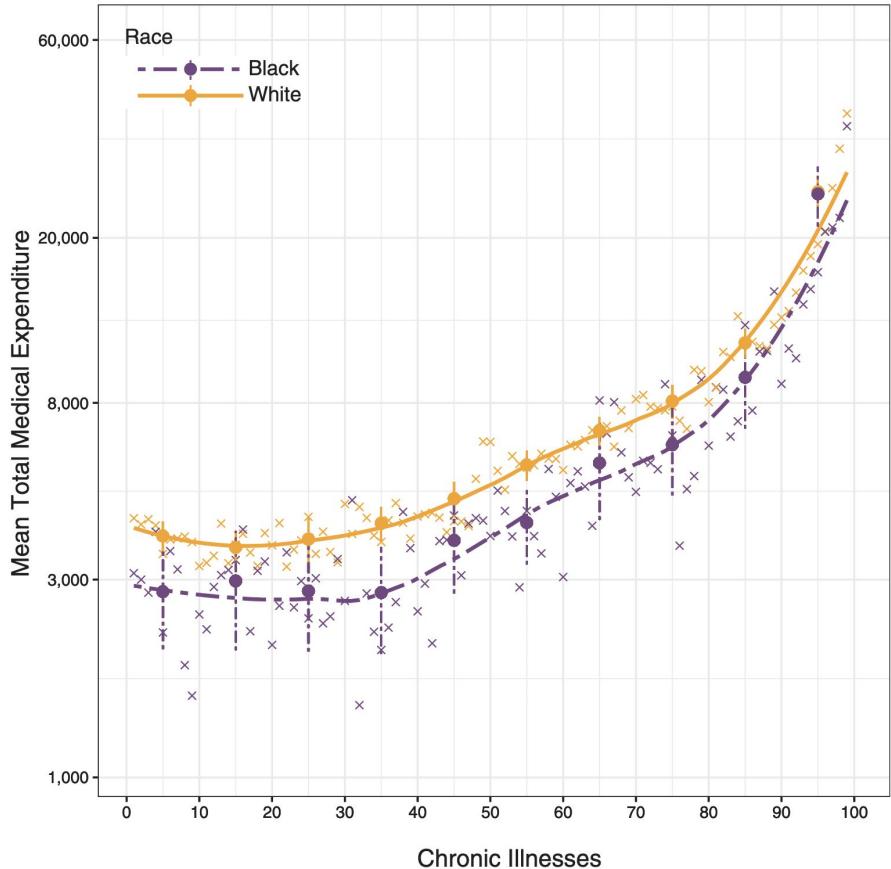
In other words, the algorithm's predictions
work very well across races.



Then, why?

At a given level of health conditions (measured by number of chronic illnesses), Blacks show lower expenses than Whites.

Accurate prediction of costs necessarily means being racially biased on health.



Dilemma of which label to choose

“Problem formulation” in data science: the task of turning an often amorphous concept we wish to predict into a concrete variable that can be predicted in a given dataset.

2021-22 Term 1

IS457: Fairness in Socio-technical Systems

Week 5 - Bias in data and machine learning models (I)

KWAK Haewoon

Image from Flaticon.com



After the study

The algorithm manufacturer independently replicated the same analysis on its national dataset of 3,695,943 patients and confirmed the biases.

The existing model infrastructure with the changed label (combination of cost + health condition prediction) showed an 84% reduction in bias.

Label biases are fixable, but producing new labels requires deep understanding of the domain, the ability to identify and extract relevant data elements, and the capacity to iterate and experiment.

Gender imbalance in medical imaging datasets produces biased classifiers for computer-aided diagnosis

Agostina J. Larrazabal^{a,1}, Nicolás Nieto^{a,b,1}, Victoria Peterson^{b,c} , Diego H. Milone^a , and Enzo Ferrante^{a,2} 

^aResearch Institute for Signals, Systems and Computational Intelligence sinc(i), Universidad Nacional del Litoral–Consejo Nacional de Investigaciones Científicas y Técnicas CONICET, Santa Fe CP3000, Argentina; ^bInstituto de Matemática Aplicada del Litoral, Universidad Nacional del Litoral–Consejo Nacional de Investigaciones Científicas y Técnicas, Santa Fe CP3000, Argentina; and ^cFacultad de Ingeniería, Universidad Nacional de Entre Ríos, Oro Verde CP3100, Argentina

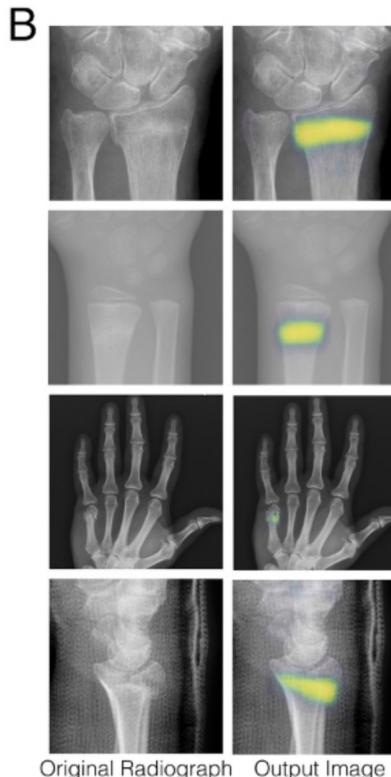
AI systems for medical images

Radiographic interpretation often takes place in environments without qualified colleagues available for second opinions.

Such circumstances increase the risk of inaccurate identification of fractures on radiographs.

In emergency departments, missed fractures account for between 41 and 80% of reported diagnostic errors.

Provide useful annotations



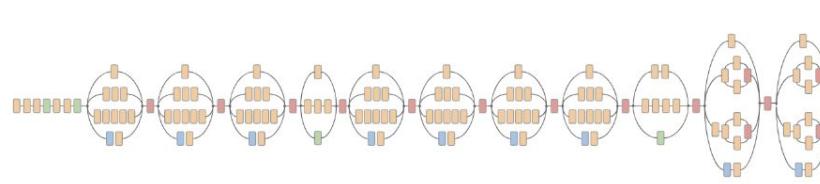
Lindsey, Robert, et al. "Deep neural network improves fracture detection by clinicians." Proceedings of the National Academy of Sciences 115.45 (2018): 11591-11596.

Dermatologist-level classification of skin cancer

Skin lesion image



Deep convolutional neural network (Inception v3)

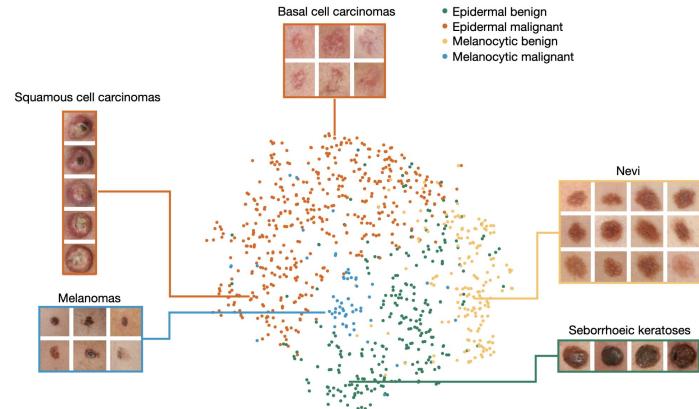
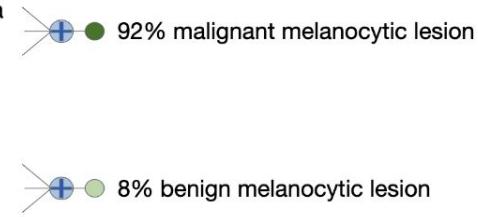


- Convolution
- AvgPool
- MaxPool
- Concat
- Dropout
- Fully connected
- Softmax

Training classes (757)

- Acral-lentiginous melanoma
- Amelanotic melanoma
- Lentigo melanoma
- ...
- Blue nevus
- Halo nevus
- Mongolian spot
- ...

Inference classes (varies by task)



Sex imbalance in datasets

Additionally, some diseases are more common in one sex than the other.

More males in experimental and clinical data



School of
Computing and
Information Systems

In early 2000, sex-specific biological differences were neglected and both experimental and clinical studies were fundamentally focused on male experimental models or male subjects.

Even nowadays, male mouse models are overall more represented than female models in basic, preclinical, and surgical biomedical research.

P21

Origin and impact of those biases



School of
Computing and
Information Systems

The lack of representation of female models and patients is partly due to technical and bioethical considerations

- Reduce the impact of estrous cycle in experimental studies
- Protective policies for women of childbearing age in clinical research.

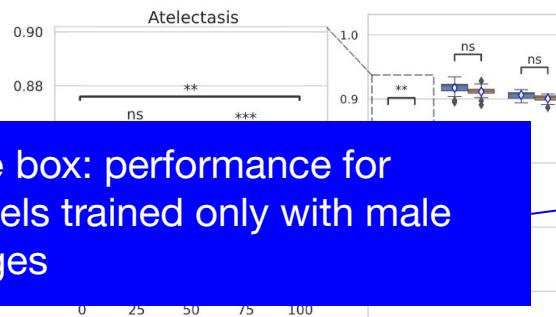
Some of the treatments that currently exist for several diseases are not adequately evaluated in women who are likely to be underrepresented in clinical trials, especially in Phases I and II.

P22

Performance with imbalanced datasets

B-1

A

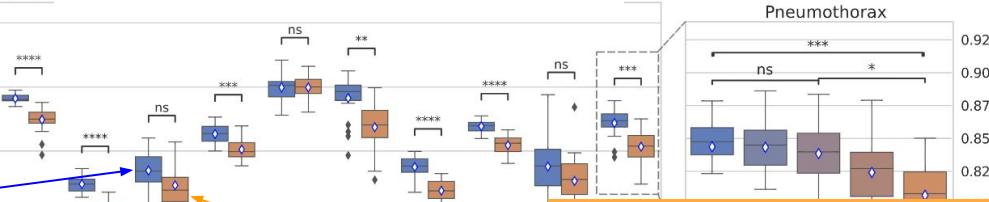


Blue box: performance for
models trained only with male
images

Testing in male patients

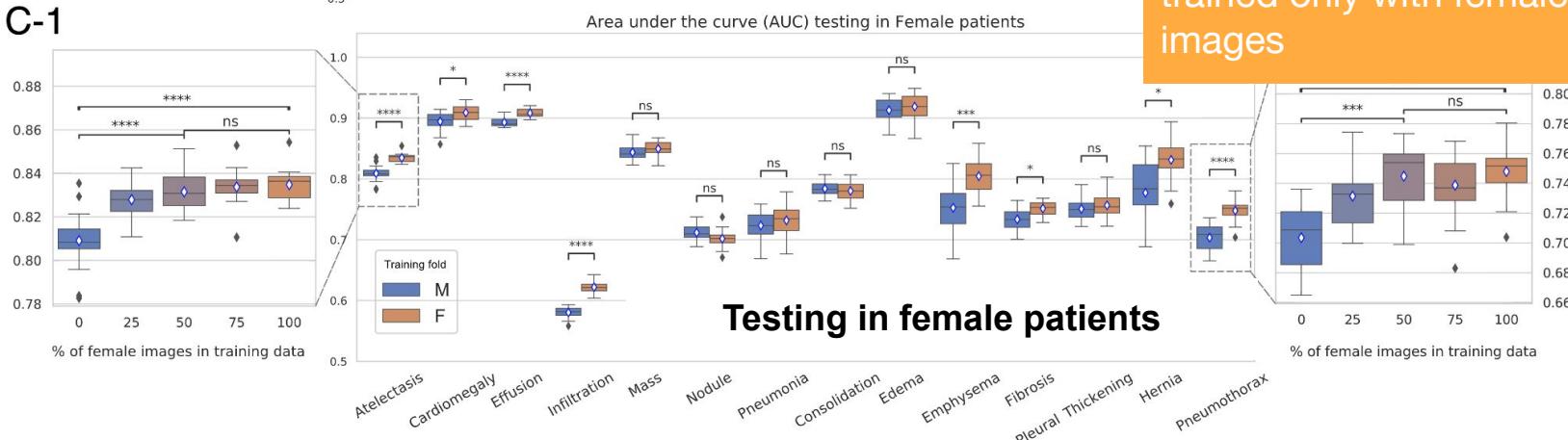
B-2

Area und



C-1

Area under the curve (AUC) testing in Female patients



Testing in female patients

Orange box:
performance for models
trained only with female
images

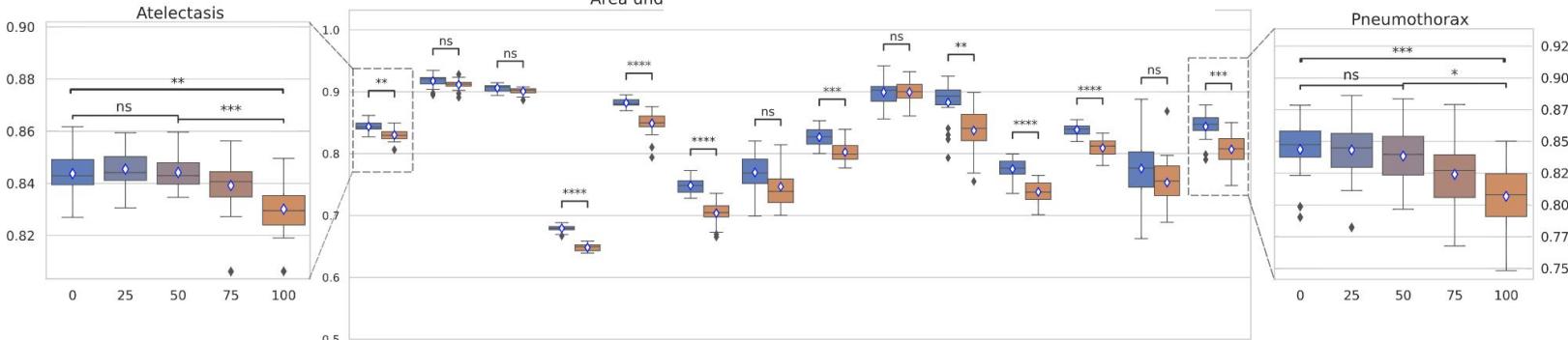
Performance with imbalanced datasets

B-1

A

Testing in male patients

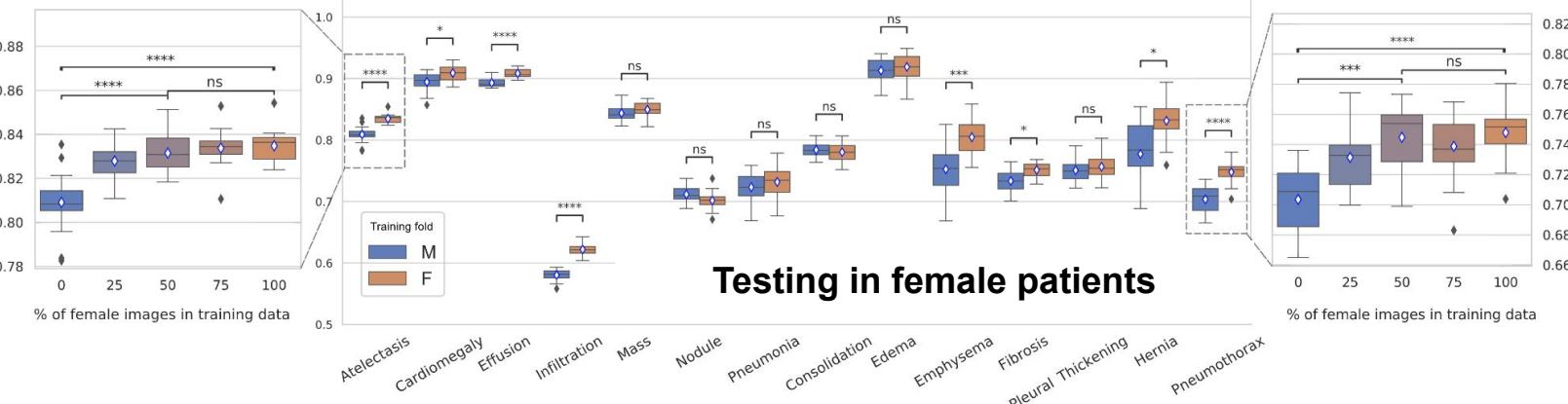
B-2



C-1

Area under the curve (AUC) testing in Female patients

C-2



Importance of diverse/balanced dataset

Diversity should be prioritized when designing databases used to train machine learning-based computer-aided diagnosis systems.

Insufficient regulations (2020)

US Food and Drug Administration

- No explicit mention of gender/sex as one of the relevant demographic variables that should describe the sampled population.

Medical imaging community

- Most public datasets do not contain gender/sex information at the patient level (MIMIC-CXR x-ray dataset (2019), Retinal Fundus Glaucoma Challenge database (2020), etc.)

Group activity: Ethics in AI healthcare

There are no correct answers for following questions. Share your thoughts with friends.

1. Imagine a situation that a doctor uses Zoom to tell a patient she would die. After COVID-19, the broad adoption of telemedicine may make this practice common. Would you accept this practice?
2. The presence of the doctor on Zoom now changes to a “chatbot” machine. The chatbot tries to talk about end-of-life issues to patients. Would you accept this practice?
3. By analyzing big data, an AI system that can predict when a patient is going to die can be built. Do you think that using such a system is okay?
4. Where do you draw your line between “what machines (AI+robots) can do” and “what machine should not do” in healthcare domain? Why do you think so?

Bias in criminal justice systems



Machine Bias

There's software used across the country to predict future criminals. And it's biased against blacks.

by Julia Angwin, Jeff Larson, Surya Mattu and Lauren Kirchner, ProPublica

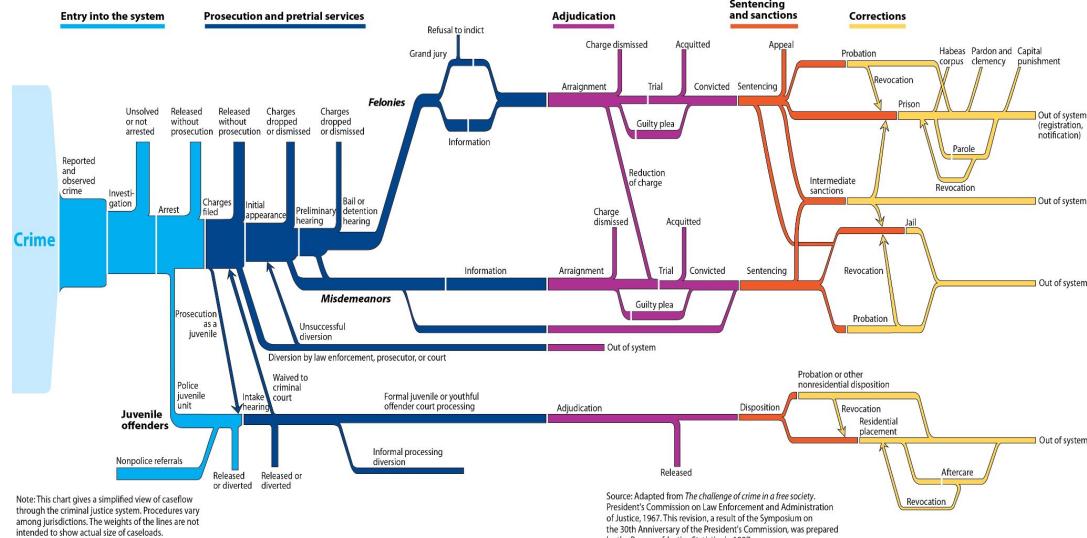
May 23, 2016

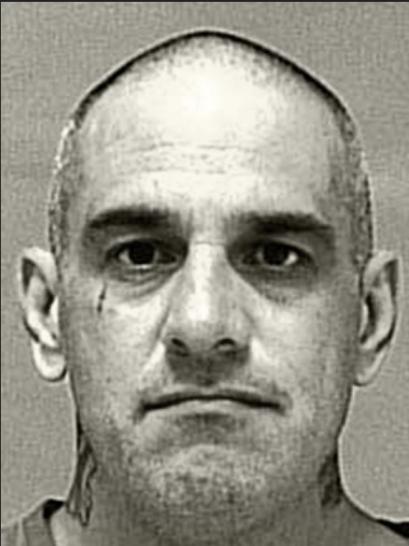
<https://www.propublica.org/article/machine-bias-risk-assessments-in-criminal-sentencing>

Risk assessments become common

In every stage of the criminal justice systems, rating a defendant's risk of future crime by computer algorithms is encouraged.

What is the sequence of events in the criminal justice system?





VERNON PRATER

RISK: 3

BRISHA BORDEN

RISK: 8

Borden was rated high risk for future crime after she and a friend took a kid's bike and scooter that were sitting outside. She did not reoffend.

Analysis of risk scores by ProPublica

More than 7,000 people arrested in Broward County, Florida, in 2013 and 2014.

Checked how many were charged with new crimes over the next two years.

Risk scores are unreliable

Only 20% of the people predicted to commit violent crimes actually went on to do so.

Even when a full range of crimes were taken into account (including misdemeanors such as driving with an expired license), of those deemed likely to be reoffend, 61% were arrested for any subsequent crimes within two years.

Racial disparities in risk scores

In forecasting who would re-offend, the algorithm made mistakes with black and white defendant at roughly the same rate but:

- Falsely flag black defendants as future criminals, wrongly labeling them this way at almost twice the rate as white defendants.
- White defendants were mislabeled as low risk more often than black defendants.

	WHITE	AFRICAN AMERICAN
Labeled Higher Risk, But Didn't Re-Offend	23.5%	44.9%
Labeled Lower Risk, Yet Did Re-Offend	47.7%	28.0%

Details of the algorithm are unknown

Northpointe's software does not publicly disclose the calculations; it is not possible to see what might be driving the racial disparity.

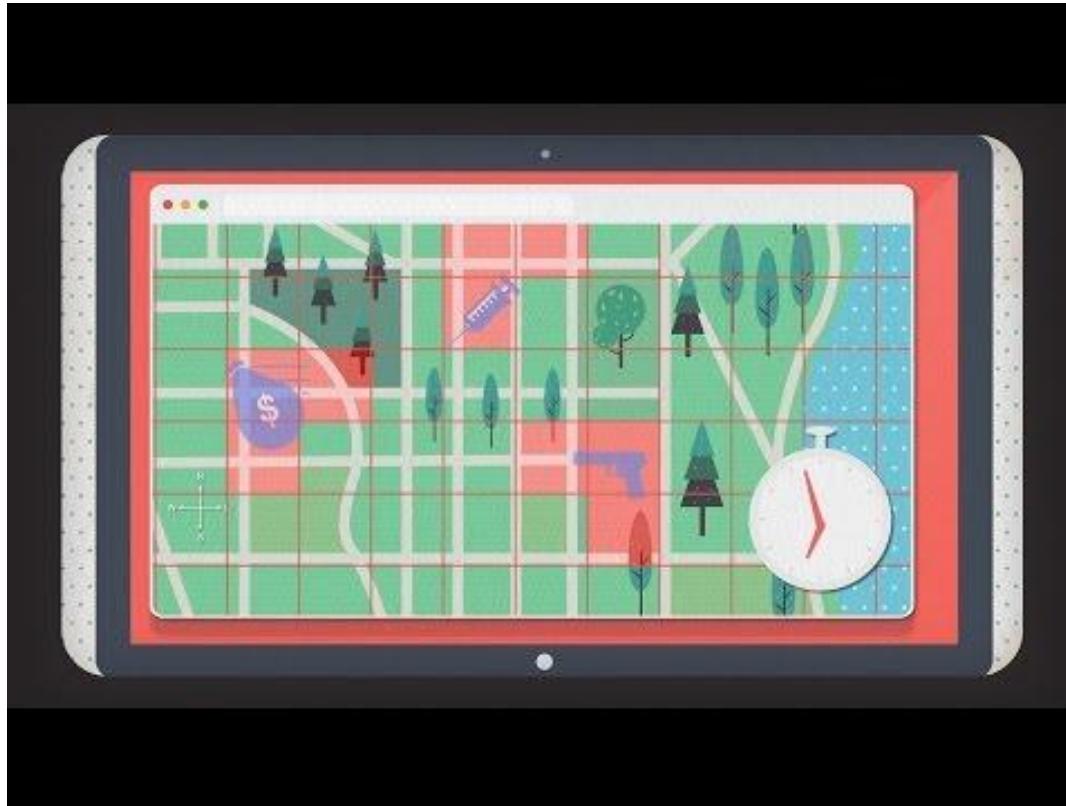
Only some variables are known:

- Education levels
- Whether a defendant has a job
- And more

Policing and discrimination



How predictive policing works



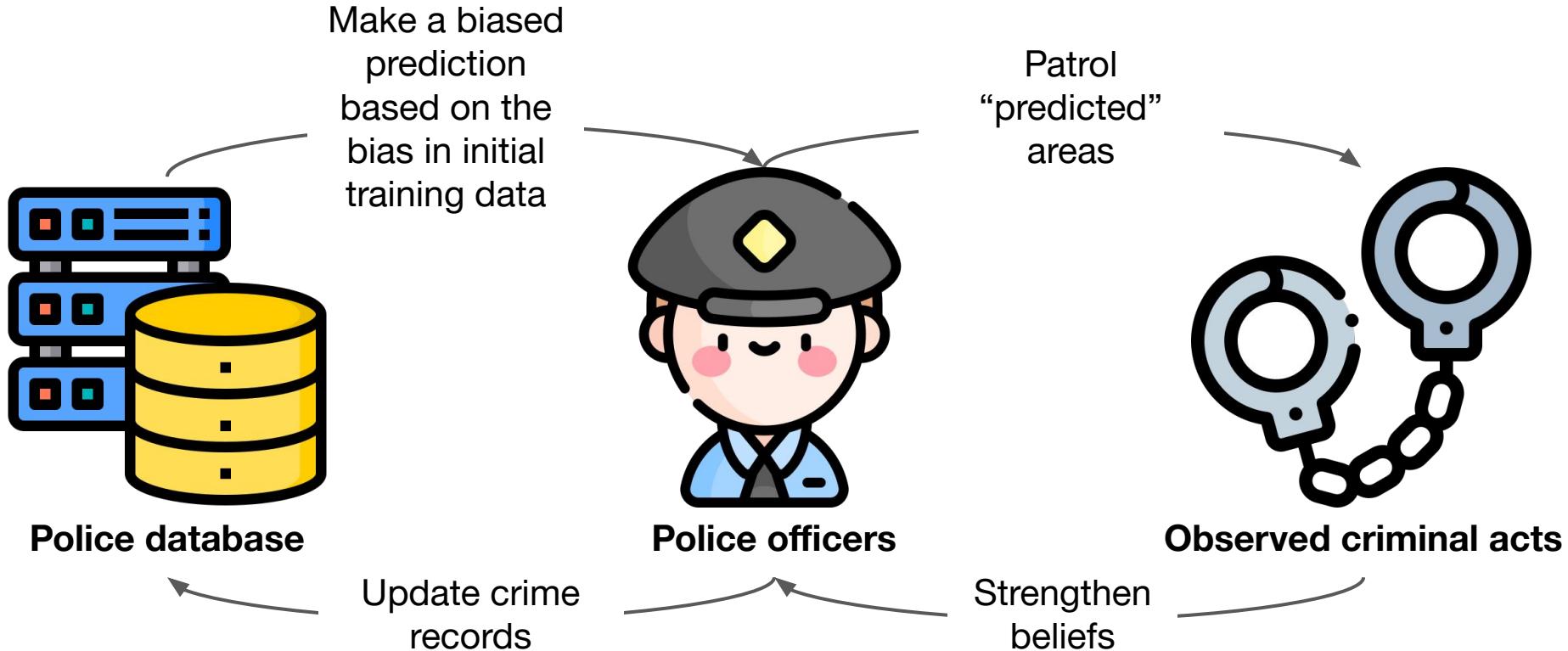
Police records do not measure crime

Police databases are not a complete census of all criminal offences, nor do they constitute a representative random sample.

Crimes that occur in locations frequented by police are more likely to appear in the database simply because that is where the police are patrolling.

Police records capture some complex interaction between criminality and policing strategy.

Over-policing and feedback loop



Over-policing and feedback loop

Officers become increasingly likely to patrol “predicted” areas and observe new criminal acts that confirm their prior beliefs regarding criminal activity.

The newly observed criminal acts then feed into the predictive policing algorithm on subsequent days, generating increasingly biased predictions.

This creates a feedback loop where the model becomes increasingly confident that the locations most likely to experience further criminal activity are exactly the locations they had previously believed to be high in crime.

Week 2 reflection

<https://forms.gle/eNErPJQwKdAWfvac7>