

Applied Artificial Intelligence

Session 26: GANs, Fairness and Ethics in AI

Fall 2018

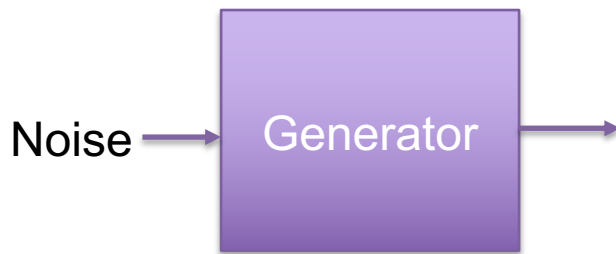
NC State University

Lecturer: Dr. Behnam Kia

Course Website: <https://appliedai.wordpress.ncsu.edu/>

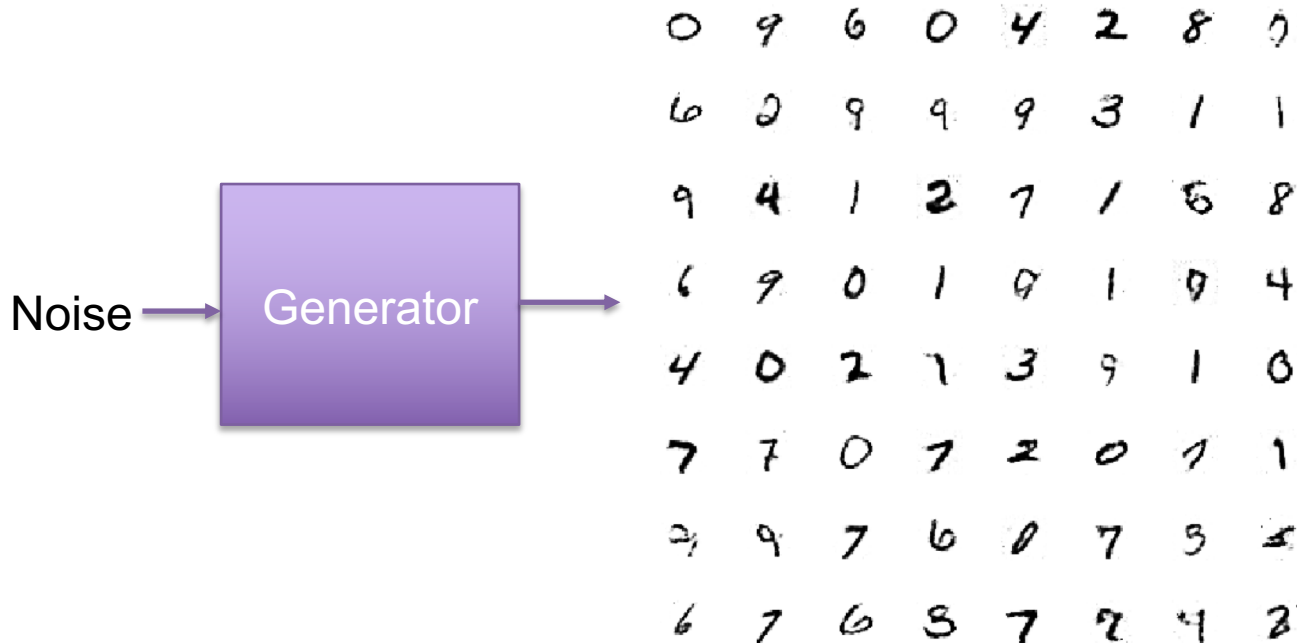
Goal: create realistic-looking synthetic data (images, voice,...)

Goal: create realistic-looking synthetic data (images, voice,...)



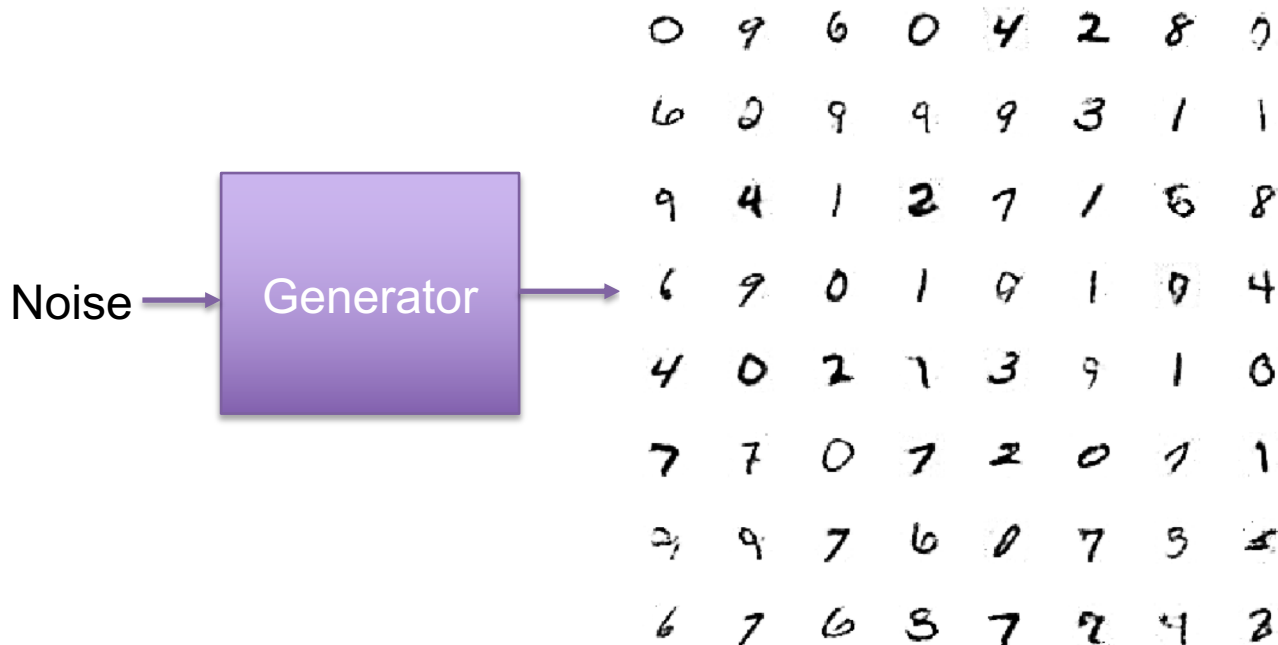
0 9 6 0 4 2 8 0
6 2 9 9 9 3 1 1
9 4 1 2 7 1 6 8
6 9 0 1 9 1 9 4
4 0 2 7 3 9 1 0
7 7 0 7 2 0 7 1
2 9 7 6 0 7 5 5
6 7 6 5 7 2 4 2

Goal: create realistic-looking synthetic data (images, voice,...)



But how?

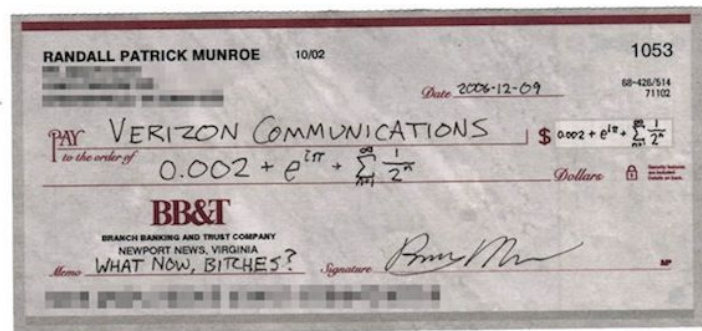
Goal: create realistic-looking synthetic data (images, voice,...)

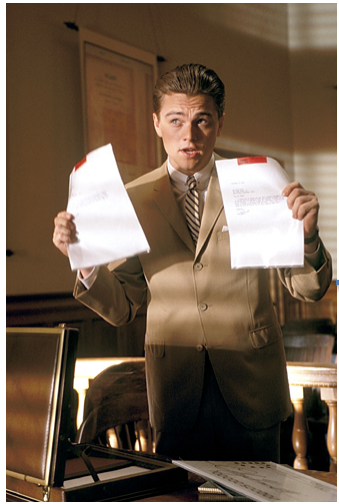


But how?

Adversarial approach. Generative NNs following this approach are called Generative Adversarial Networks, GANs

- As an example to adversarial approach to designing a generative of real-looking synthetic data (fake data), imagine counterfeiting checks.





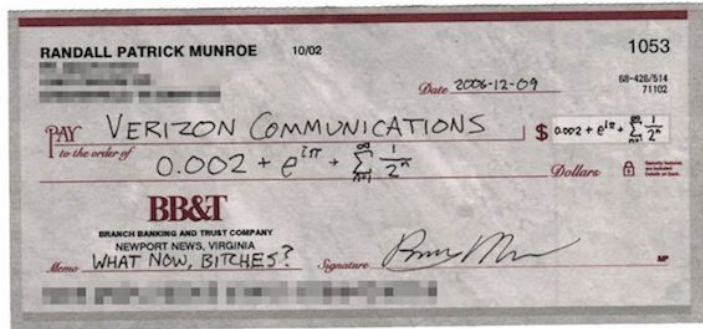
123456
ACE Stedman's Hardware Date Today
Miles City
Pay to Person with this check | \$ 175000/100
~~One-hundred-seventy-five million thousandths~~
BANK OF MILES CITY
memo this check is good Mr. Stedman

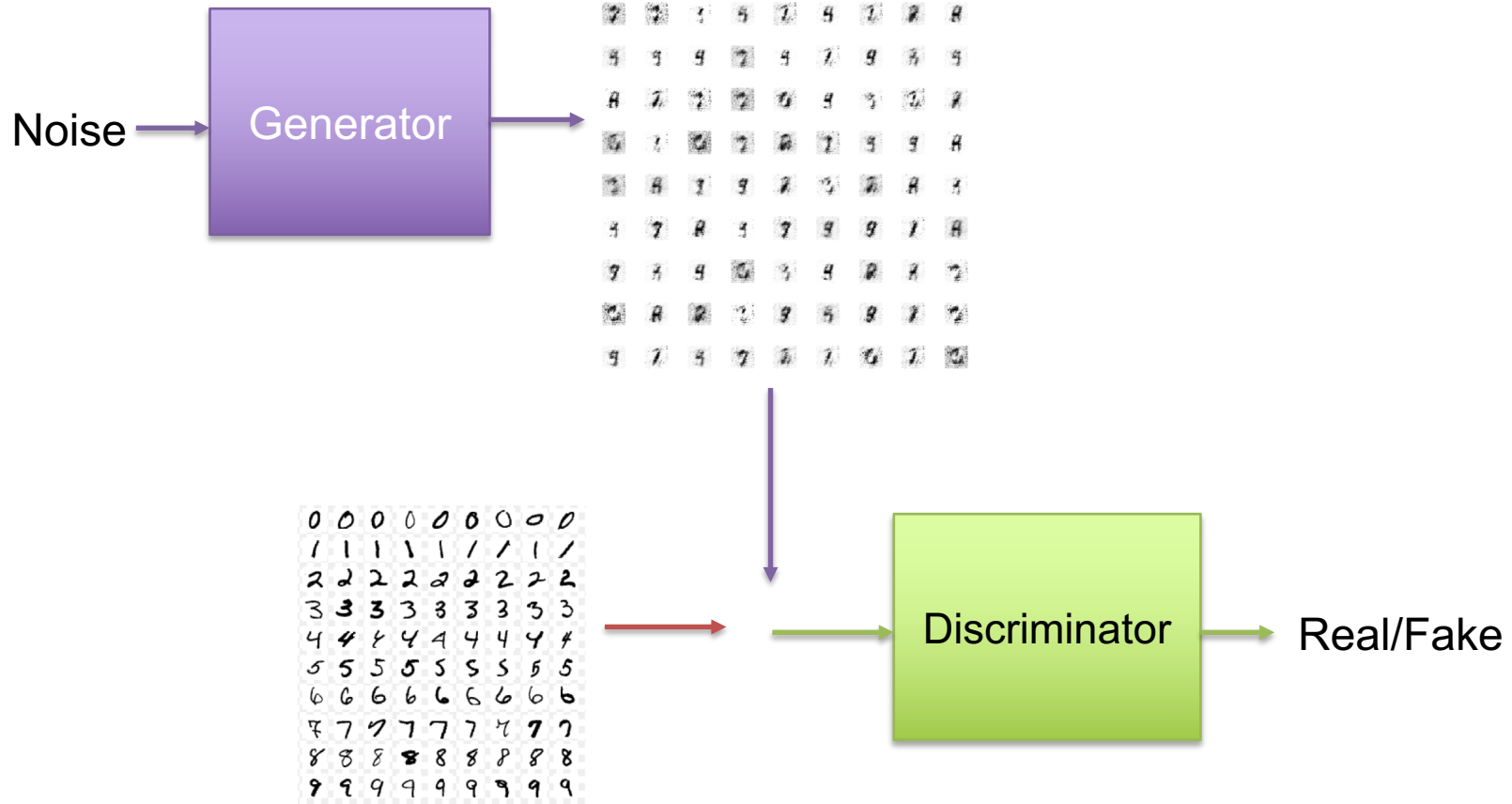


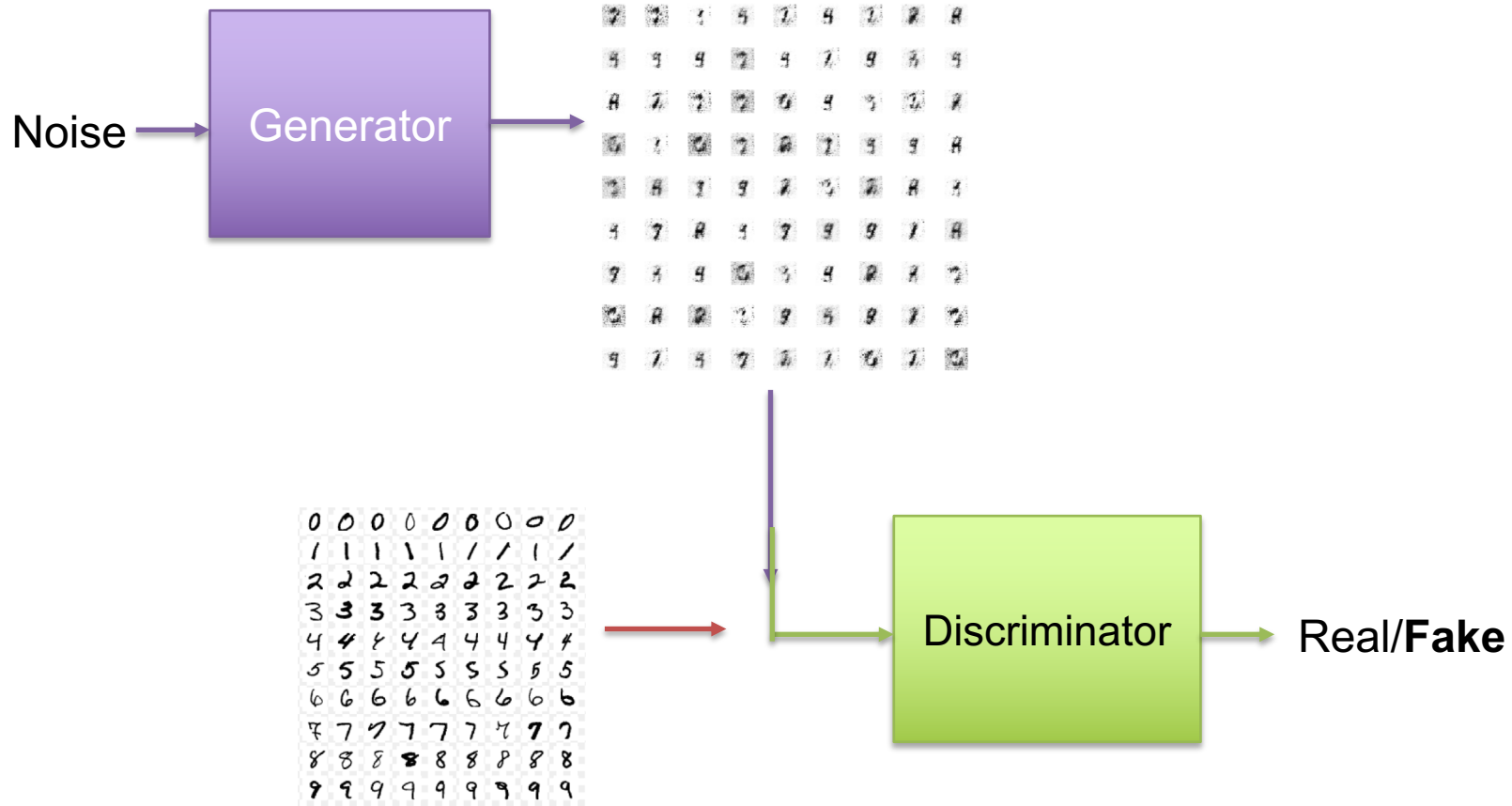


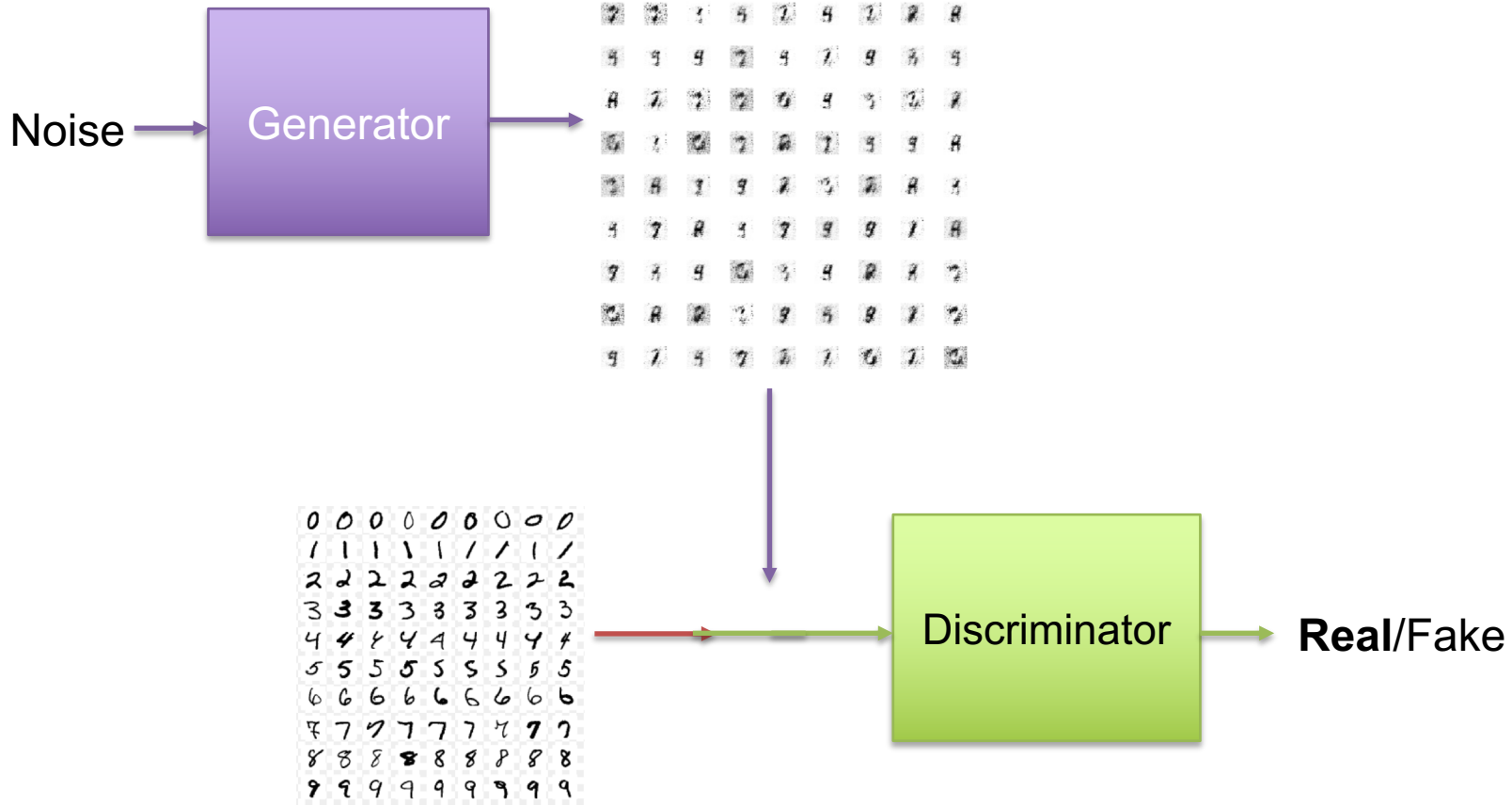
12-31-58
ACE Stedman's Hardware
Miles City Date Today
Pay to Person with this check \$ 175,000.00
~~One hundred seventy-five thousand and no/100's~~
BANK OF MILES CITY
MEMO this check is good Mr. Stedman





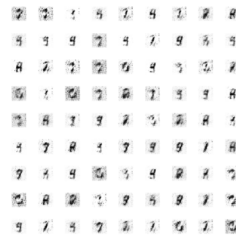






Step 1: Train Discriminator

- Create fake images from the Generator
- Mix them with real images
- Train Discriminator to detect fake from real
- Now we have a trained discriminator

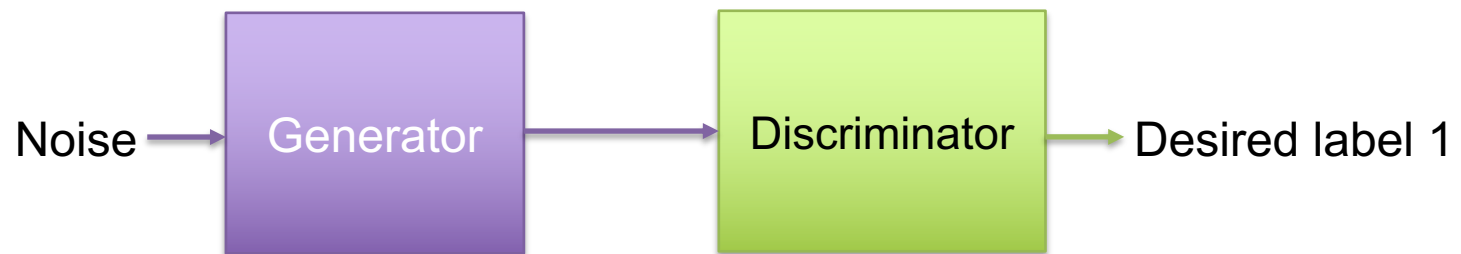


Real: 1
Fake: 0



Step 2: Train Generator to fool discriminator

- Combine generator and discriminator as one model
- Freeze discriminator (no training for discriminator)
- Give fake images to discriminator and train the combined model to produce labels 1 (generator is under training)



- Go to Step 1 to retrain discriminator

0 9 6 0 4 2 8 0
 6 2 9 9 9 3 1 1
 9 4 1 2 7 1 6 8
 6 9 0 1 9 1 9 4
 4 0 2 1 3 9 1 0
 7 7 0 7 2 0 1 1
 2 9 7 6 0 7 9 5
 6 7 6 8 7 2 4 2

7 1 0 6 2 0 6 2
 3 8 0 0 8 8 4 7
 5 8 7 3 9 0 5 8
 1 5 0 2 8 4 2 3
 0 4 3 9 8 2 1 8
 5 0 1 6 6 5 5 2
 1 7 7 1 2 3 7 3
 6 3 7 6 6 1 4 0

0 9 6 0 4 2 8 0
6 2 9 9 9 3 1 1
9 4 1 2 7 1 6 8
6 9 0 1 9 1 9 4
4 0 2 1 3 9 1 0
7 7 0 7 2 0 1 1
2 9 7 6 0 7 9 4
6 7 6 8 7 2 4 2

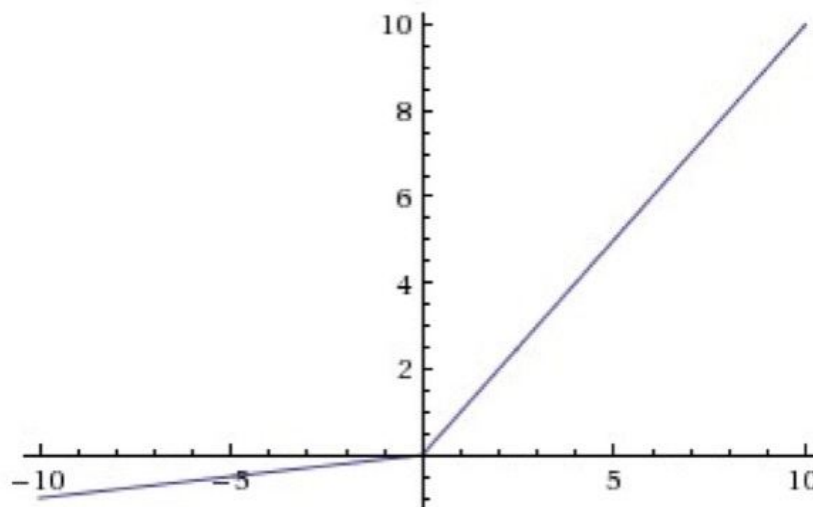
Fake

7 1 0 6 2 0 6 2
3 8 0 0 8 8 4 7
5 8 7 3 9 0 5 8
1 5 0 2 8 4 2 3
0 4 3 9 8 2 1 8
5 0 1 6 6 5 5 2
1 7 7 1 2 3 7 3
6 3 7 6 6 1 4 0

Real

Leaky Relu

$$f(x) = \begin{cases} x & \text{if } x > 0 \\ 0.01x & \text{otherwise} \end{cases}$$



- Please see codes for our first example of a GAN.

Homework

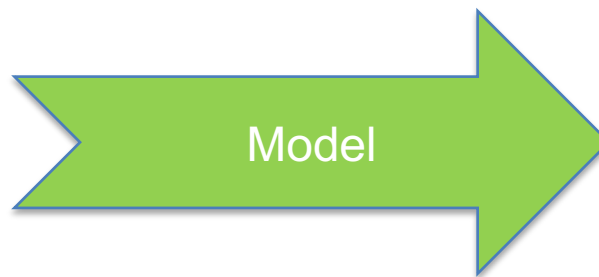
Deadline: Next Tuesday, 3PM

- Use an RNN network to classify MNIST dataset.
- Each image is 28×28 . Consider each column (28 pixels) as one input. The sequence to RNN would be 28 columns of 28 pixels, each column given at one time step.
- Try simpleRNN, LSTM, and GRU.
- Try multiple stacks of those units.
- Use dropout to fight against overfitting (if needed).
- Accuracy of $\sim 98\%$ can be easily obtained.

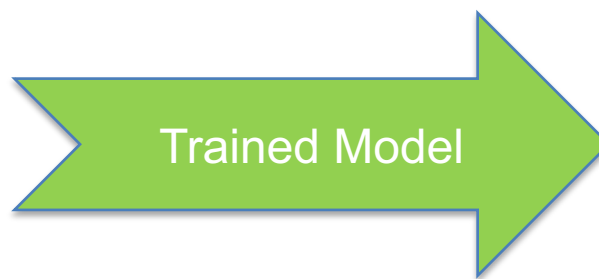
AI : Ethics and Fairness

Biased/Bad Data Results in Bad AI

Junk/Bad/Biased
Data for training
AI



Biased/Bad Data Results in Bad AI



Junk/Bad/Biased Results



Prof. Donna Strickland

The first woman to win Nobel Prize in Physics in 55 years.

If AI wanted to select the winner, and gender was among the features of the candidates, historical data would be biased against selection of female physicists.

Fortunately we don't use AI to select the next winner of Nobel price, but in other areas such as job websites, we are observing that AI is starting to be used to filter CVs of the applicants and to match them to potential employers. Please see the next slide.

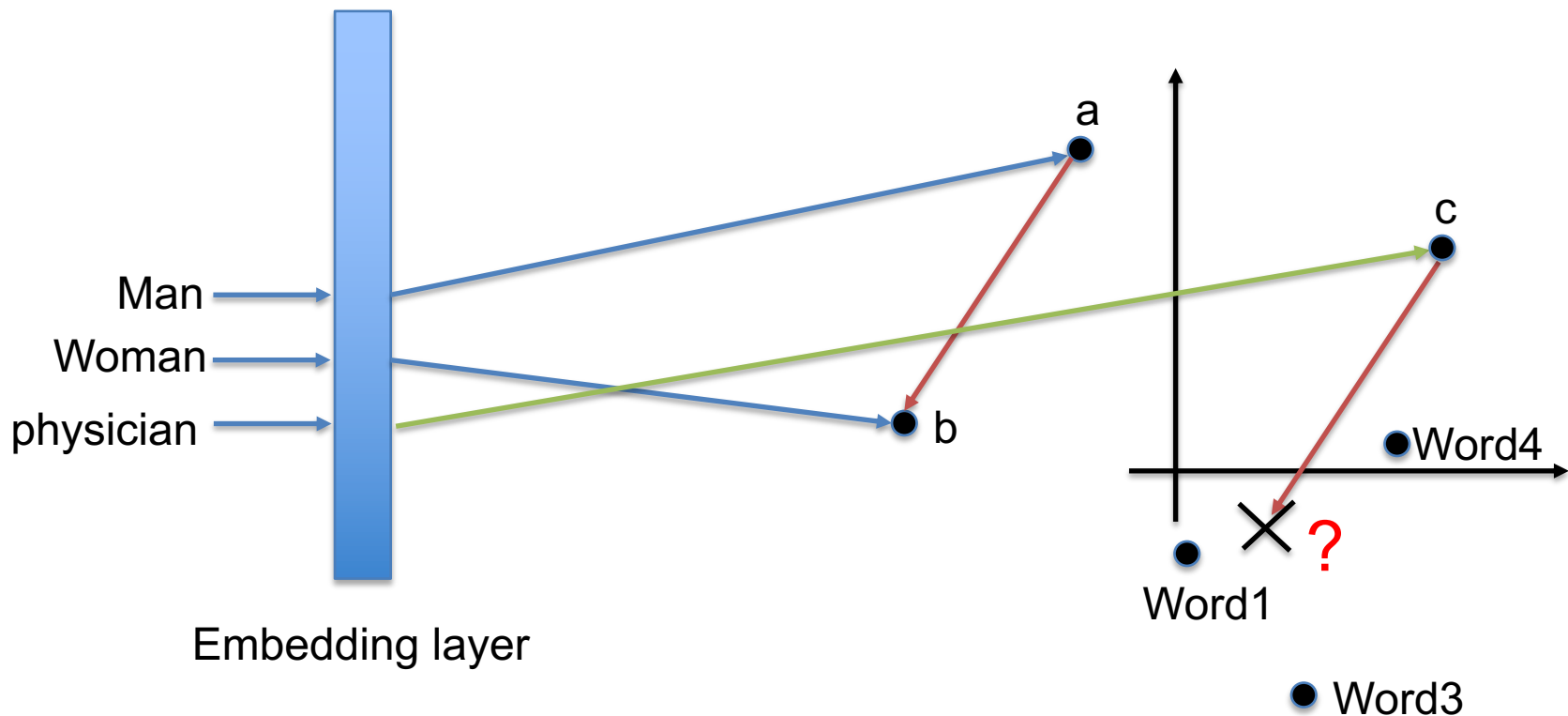
Red flags to possible bias in data



- Does your use case or product specifically use any of the following data: biometrics, race, skin color, religion, sexual orientation, socioeconomic status, income, country, location, health, language, or dialect?
- Does your use case or product use data that is likely to be highly correlated with any of the personal characteristics listed above (for example, zip code or other geospatial data is often correlated with socioeconomic status and/or income; image/video data can reveal information about race, gender, and age)?
- Could your use case or product negatively impact individuals' economic or other important life opportunities?

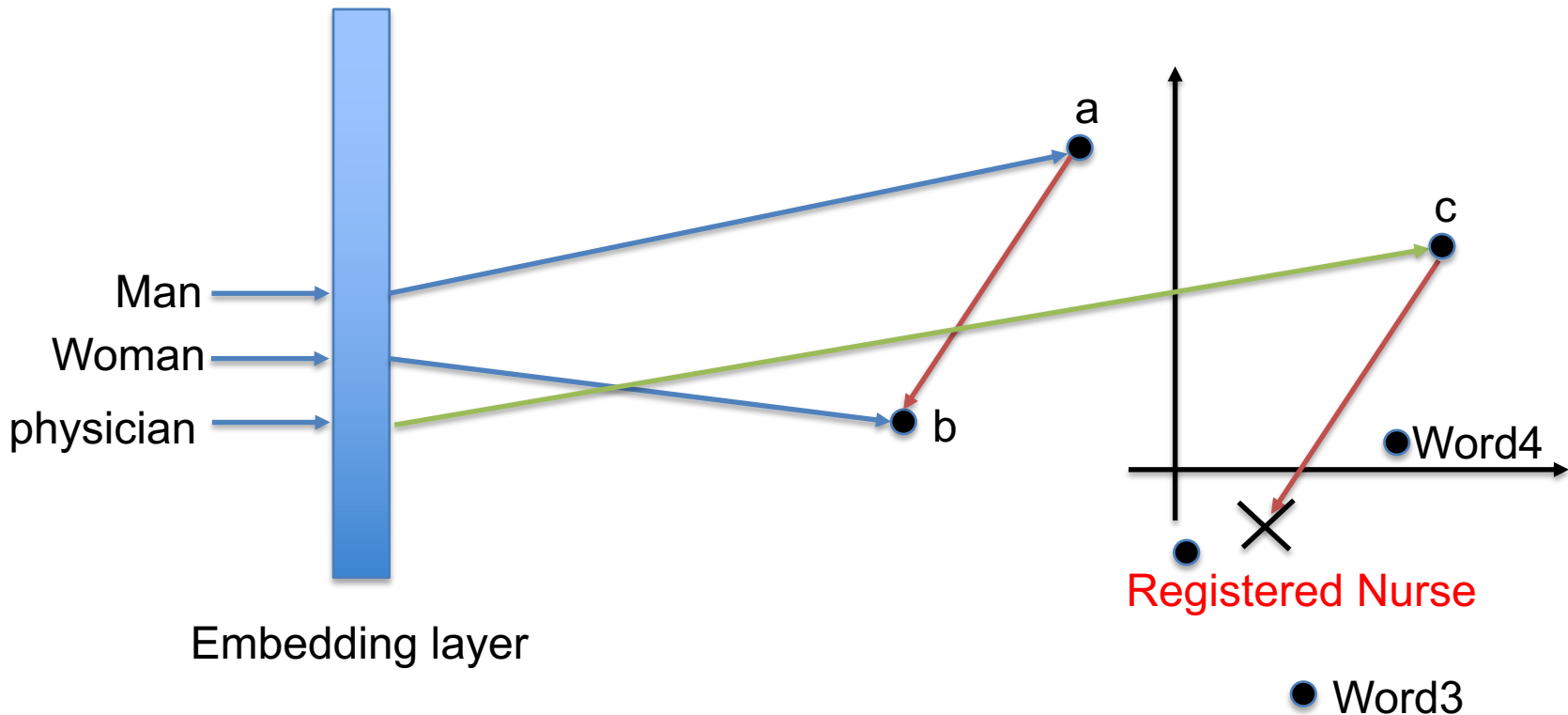
In some Publicly available and used Embeddings...

- Remember the previous homework. The goal was to use embeddings to determine given words a and b, find word d that word c and d has the same relationship as a and b. Here the relationship is the gender shift. So what is the female version of a physician?



In some Publicly available and used Embedding...

- Some embeddings say d is “Nurse!” Please see the next slide for more info on this.
- Yeap, we have a misogynist AI!



Man is to Computer Programmer as Woman is to Homemaker? Debiasing Word Embeddings

Tolga Bolukbasi¹, Kai-Wei Chang², James Zou², Venkatesh Saligrama^{1,2}, Adam Kalai²

¹Boston University, 8 Saint Mary's Street, Boston, MA

²Microsoft Research New England, 1 Memorial Drive, Cambridge, MA

tolgab@bu.edu, kw@kwchang.net, jamesyzou@gmail.com, srv@bu.edu, adam.kalai@microsoft.com

Abstract

The blind application of machine learning runs the risk of amplifying biases present in data. Such a danger is facing us with *word embedding*, a popular framework to represent text data as vectors which has been used in many machine learning and natural language processing tasks. We show that even word embeddings trained on Google News articles exhibit female/male gender stereotypes to a disturbing extent. This raises concerns because their widespread use, as we describe, often tends to amplify these biases. Geometrically, gender bias is first shown to be captured by a direction in the word embedding. Second, gender neutral words are shown to be linearly separable from gender definition words in the word embedding. Using these properties, we provide a methodology for modifying an embedding to remove gender stereotypes, such as the association between the words *receptionist* and *female*, while maintaining desired associations such as between the words *queen* and *female*. Using crowd-worker evaluation as well as standard benchmarks, we empirically demonstrate that our algorithms significantly reduce gender bias in embeddings while preserving its useful properties such as the ability to cluster related concepts and to solve analogy tasks. The resulting embeddings can be used in applications without amplifying gender bias.

The Rekognition Scan

Comparing input images to mugshot databases

1 INPUT: SEARCH IMAGES



2 REKOGNITION SEARCH

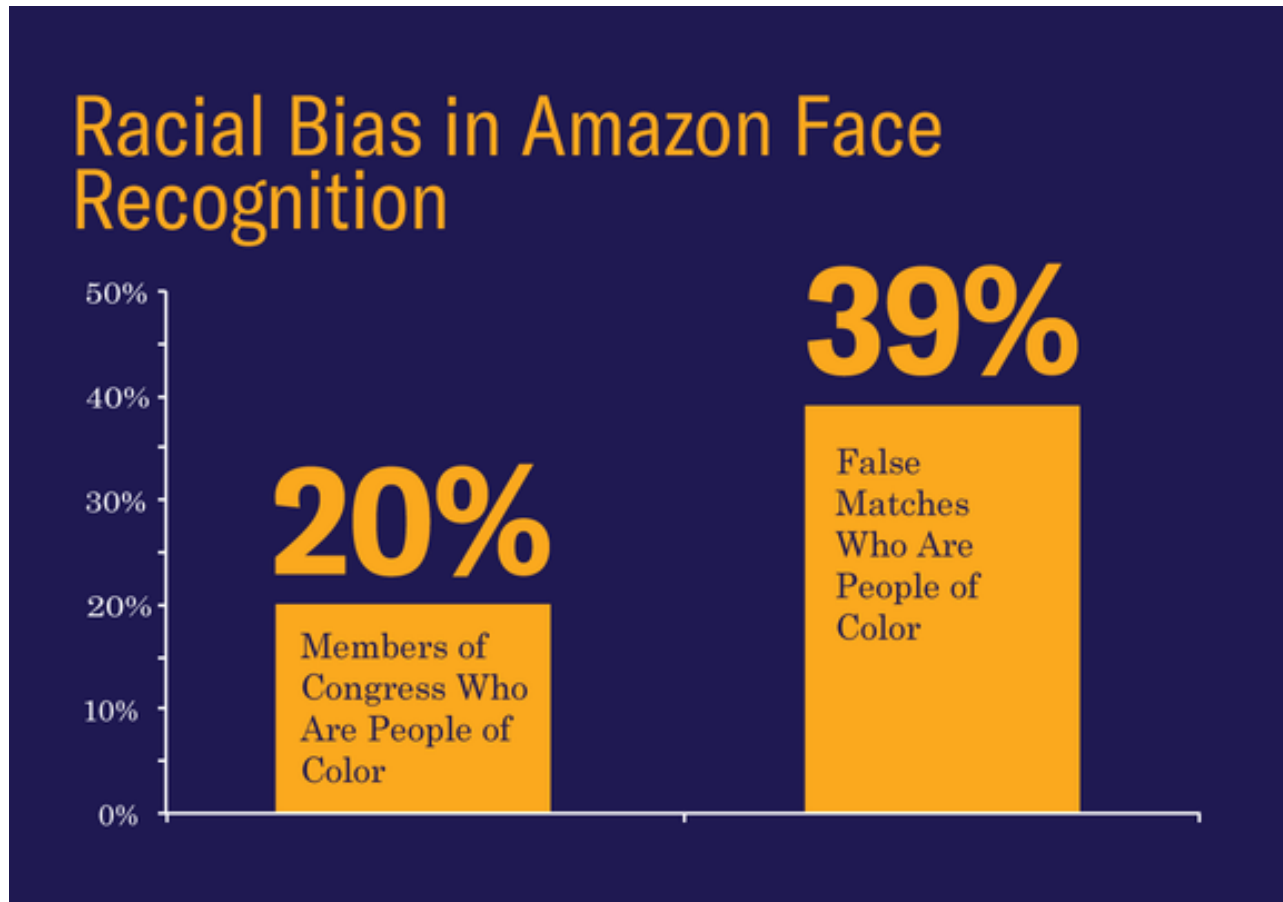


3 OUTPUT: PREDICTED MATCHES



ACLU ran an experiment, in which they used a database of 25,000 arrest photos, and compared them against the photos of members of congress using Amazon's facial recognition software.

<https://www.aclu.org/blog/privacy-technology/surveillance-technologies/amazons-face-recognition-falsely-matched-28>



Studies show that in facial recognition systems, the false matches are disproportionately of people of color. In ACLU's study, 39% of false matches are people of color, but just 20% of members of congress are people of color. This problem can have real-life, unfair consequences for people of color.

Auto censored/flagged words and comments in online forums and comment areas

- Sentiment Analyzers might think words such as “Gay” is negative, and therefore comments or documents containing those words can be flagged or censored, whereas words such as straight remain safe.

- Please read:

https://motherboard.vice.com/en_us/article/j5jmq8/google-artificial-intelligence-bias