# CAS751: Information-Theoretic Methods in Trustworthy Machine Learning

Shahab Asoodeh

Department of Computing and Software

McMaster University

September 3, 2025

McMaster
University

# Logistics

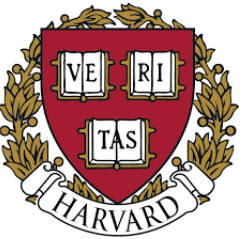- **Lectures:** Wednesdays 12.30am – 3.30pm with a break 13.45pm – 2pm

- **Office hours:**

  - Mondays (virtual) 10.30 – 11.30 on Zoom (meeting ID: 966 5355 6303 Passcode: CAS751)
  - Wednesdays (in-person) 3.30 – 4.30 (ITB 212)

- **Course webpage:**

  https://www.cas.mcmaster.ca/~asoodehs/cas751.html

# My background

- Assistant Professor of CS, McMaster

- Research scientist, Meta: Private deep learning

- Postdoc fellow, Harvard University: Trustworthy AI, making AI less discriminatory and privacy-invasive

- Postdoc scholar, University of Chicago: Geometric data analysis and privacy on medical and criminal records

- Ph.D. in mathematics, Queen's University

# Pre-requisites

- **Probability and statistics:** distributions, pdf, pmf, random variables, expectation, conditional probability, variance, empirical mean, Gaussian and Laplace densities, ....

- **Algorithm:** pesudocodes

- **Optimization:** Objective functions, convexity, (stochastic) gradient descent

- **Basic programming**

---

## 1   Elements of Probability

Probability theory is the study of uncertainty. Through this class, we will be relying on concepts from probability theory for deriving machine learning algorithms. These notes attempt to cover the basics of probability theory at a level appropriate for introductory course of differential privacy (such as CS 3DP3). The mathematical theory of probability is rather sophisticated, and delves into a branch of analysis known as measure theory. In these notes, we provide a basic treatment of probability that does not address these finer details.

### 1.1   Definition of probability space

In order to define a probability on a set we need a few basic elements:

- **Sample space** $\Omega$: The set of all the outcomes of a random experiment. Here, each outcome $\omega \in \Omega$ can be thought of as a complete description of the state of the real world at the end of the experiment.

- **Event space** $\mathcal{F}$: A set whose elements $A \in \mathcal{F}$ (called **events**) are subsets of $\Omega$ (i.e., $A \subseteq \Omega$ is a collection of possible outcomes of an experiment).[1]

- **Probability measure**: A function $P : \mathcal{F} \to \mathbb{R}$ that satisfies the following properties,

  - **Non-negativity**: $P(A) \geq 0$, for all $A \in \mathcal{F}$
  - **Completeness**: $P(\Omega) = 1$
  - **Countable Additivity**: If $A_1, A_2, \ldots$ are disjoint events (i.e., $A_i \cap A_j = \emptyset$ whenever $i \neq j$), then
  $$P\left(\bigcup_{i=1}^{\infty} A_i\right) = \sum_{i=1}^{\infty} P(A_i).$$

These three properties are called the **Axioms of Probability**.

**Example 1.** Consider the event of tossing a six-sided die. The sample space is $\Omega = \{1, 2, 3, 4, 5, 6\}$. We can define different event spaces on this sample space. For example, the simplest event space is the trivial event space $\mathcal{F} = \{\emptyset, \Omega\}$. Another event space is the set of all subsets of $\Omega$. For the first event space, the unique probability measure satisfying the requirements above is given by $P(\emptyset) = 0$, $P(\Omega) = 1$. For the second event space, one valid probability measure is to assign the probability of each set in the event space to be $\frac{i}{6}$ where $i$ is the number of elements of that set; for example, $P(\{1, 2, 3, 4\}) = \frac{4}{6}$ and $P(\{1, 2, 3\}) = \frac{3}{6}$.

---

[1] The event space $\mathcal{F}$ is technically required to satisfy three properties: (1) $\emptyset \in \mathcal{F}$; (2) $A \in \mathcal{F} \Rightarrow \Omega \setminus A \in \mathcal{F}$; and (3) $A_1, A_2, \cdots \in \mathcal{F} \Rightarrow \bigcup_i A_i \in \mathcal{F}$.

# Grading scheme

- Assignments: **30%**

- Midterm: **30%**

- Final project: **40%**

- **[Bonus]** Participation: **+10%**

# Assignments (30%)

- **4 sets of assignments**: tentatively:
  - One on information theoretic measures
  - Two on privacy (central and local)
  - One on private ML

- Assignments will be due in two weeks

- Based on "research questions": Assignments may require effort/research/challenge to complete

- Collaboration is permitted (and encouraged)! But you should write your own solutions in your words and logic

- Late assignments will not be accepted

- Handwritten solution will not be accepted: $\LaTeX$ (template will be posted)

- Issues with marking? **Only two weeks** to talk to me about it

# Midterm (30%)

- **October 22nd** during the lecture 12.30 -- 2.30pm

- Tentative topics: Information-theoretic measures and central differential privacy

- You'll rock it if you understand lectures and assignments

- Formula sheet is permitted (double-sided) --- Not necessarily needed!

# Final Project (40%)

- Pick a "serious" paper of your choice (or several) by 23:59 pm at November 19th

- Presentation: Last lecture will be a marathon of presentations

  - Motivations
  - Main results
  - Proof technique (and sketch)
  - Critique

- Some papers will be recommended later

- Two parts:
  - **Presentation (20%)**:   Dec 10
  - **Final report (20%)**: 5-6 pages of a conference format ( $\LaTeX$ ) by Dec 14th

- Final report should include all the **main** results of the paper(s), critiques, improvements or generalizations

# Major topics

- Basic information theoretic concepts, measures and quantities  (20%)

  - "Distance" metrics between distributions
  - Jensen's inequality
  - Joint range of divergences

- Differential privacy  (70%)

  - Central DP and its use in private deep learning
  - Local DP and its use in Minimax and Bayesian estimation problems

- Discrimination intervention and fairness in ML  (10%)

  - (Philosophical) intuitive definitions of fairness and their mathematical realizations
  - Some fundamental fairness-intervention mechanisms

This is <u>NOT</u> an ML course

This is a foundation of trustworthy ML course

This is a <u>foundation</u> of <u>trustworthy</u> ML course

Mathematical:
- Information-theoretic (divergences)
- Probabilistic (densities)
- Statistical (hypothesis testing)
- Optimization (SGD)

# Trustworthy Machine Learning

# "I want you to think about data as the next natural resource."

- Ginni Rometty, former IBM CEO

Oil

Natural gas

Uranium

Coal

Water

Data is collected and processed by companies in order to provide useful services…

…but, like other natural resources, it can be misused and must be handled with care.

# With big data comes big responsibility



## Trustworthy ML

- Fairness and Ethics

- Privacy

# Trustworthy ML: Privacy

Privacy   =   anonymization ?

# Trustworthy ML: Privacy

William Weld



Massachusetts Group Insurance
Commission: Released
"anonymized" medical history of
state employees in 1997

ethnicity
visit date        zip code
diagnosis      birth date
procedure      sex
medication

name
home address    zip code
party               birth date
date registered   sex

Latanya Sweeney



Bought Massachusetts voter
records for $20

# Trustworthy ML: Privacy

64% of US population can be uniquely identified using ZIP, birth date, sex

Latanya Sweeney, "Uniqueness of Simple Demographics in the U.S. Population", 2000.

Philippe Golle, "Revisiting the Uniqueness of Simple Demographics in the US Population", 2007.

Privacy ≠ anonymization

# Trustworthy ML: Privacy

# Trustworthy ML: Privacy



Image credit: Arvind Narayanan

# Privacy ≠ anonymization

even if only "seemingly innocent" data is released

# How about modern ML?

# Does neural net preserve privacy?

# Trustworthy ML: Privacy



1. Train

2. Predict

**Obama**

Fredrikson et al., "Model Inversion Attacks that Exploit Confidence Information and Basic Countermeasures"

# Trustworthy ML: Privacy



1. Train

2. Extract

Number 7

Fredrikson et al., "Model Inversion Attacks that Exploit Confidence Information and Basic Countermeasures"

# Trustworthy ML: Privacy

1. Train

2. Extract

Number 7



Model inversion attack

Fredrikson et al., "Model Inversion Attacks that Exploit Confidence Information and Basic Countermeasures"

Neural networks "memorize" (part of) the training dataset!

# Trustworthy ML: Privacy

# Trustworthy ML: Privacy

## 1. Train

## 2. Predict

"What are you" ➡ "doing"

**P(** *What are you doing* **;** **) = 0.99**

**P(** *What are you winter* **;** **) = 0.01**

# Trustworthy ML: Privacy

## 1. Train



## 2. Predict

"My SIN is"

# Trustworthy ML: Privacy

$$P(\text{My SIN is 000-000-000}; \text{}) = 0.00$$

$$P(\text{My SIN is 000-000-001}; \text{}) = 0.00$$

$$\vdots$$

$$P(\text{My SIN is 131-456-788}; \text{}) = 0.32$$

$$P(\text{My SIN is 131-456-789}; \text{}) = 0.01$$

$$\vdots$$

$$P(\text{My SIN is 999-999-999}; \text{}) = 0.00$$

# Trustworthy ML: Privacy

$$P(\text{My SIN is 000-000-000}; \text{[neural network]}) = 0.00$$

$$P(\text{My SIN is 000-000-001}; \text{[neural network]}) = 0.00$$

$$\vdots$$

$$\boxed{P(\text{My SIN is 131-456-788}; \text{[neural network]}) = 0.32}$$

$$P(\text{My SIN is 131-456-789}; \text{[neural network]}) = 0.01$$

$$\vdots$$

$$P(\text{My SIN is 999-999-999}; \text{[neural network]}) = 0.00$$

# Trustworthy ML: Privacy

P( **My SIN is 000-000-000** ; ) = 0.00

P( **My SIN is 000-000-001** ; ) = 0.00

# Does this take millions of queries?

P( **My SIN is 131-456-788** ; ) = 0.32

P( **My SIN is 131-456-789** ; ) = 0.01

P( **My SIN is 999-999-999** ; ) = 0.00

# Does this take millions of queries?

# NO!

Carlini, Liu, Kos, Erlingsson, Song, "The Secret Sharer: Evaluating and Testing Unintended Memorization in Neural Networks"

Re

LILY HAY NEWMAN    ANDY GREENBERG    SECURITY    DEC 2, 2023 9:00 AM

# Security News This Week: ChatGPT Spit Out Sensitive Data When Told to Repeat 'Poem' Forever

Plus: A major ransomware crackdown, the arrest of Ukraine's cybersecurity chief, and a hack-for-hire entrepreneur charged with attempted murder.

6

# Is AI fair?

# Is AI fair?

# Is AI fair?

# Trustworthy ML: Fairness



Gender was misidentified in **up to 12 percent of darker-skinned males** in a set of 318 photos.
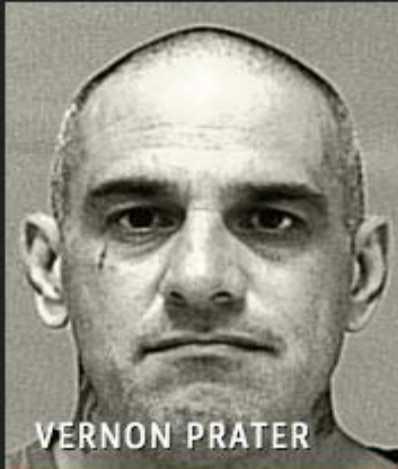
Gender was misidentified in **35 percent of darker-skinned females** in a set of 271 photos.

**The New York Times**

*Facial Recognition Is Accurate, if You're a White Guy*



Two Petty Theft Arrests

VERNON PRATER
LOW RISK 3

BRISHA BORDEN
HIGH 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.

**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
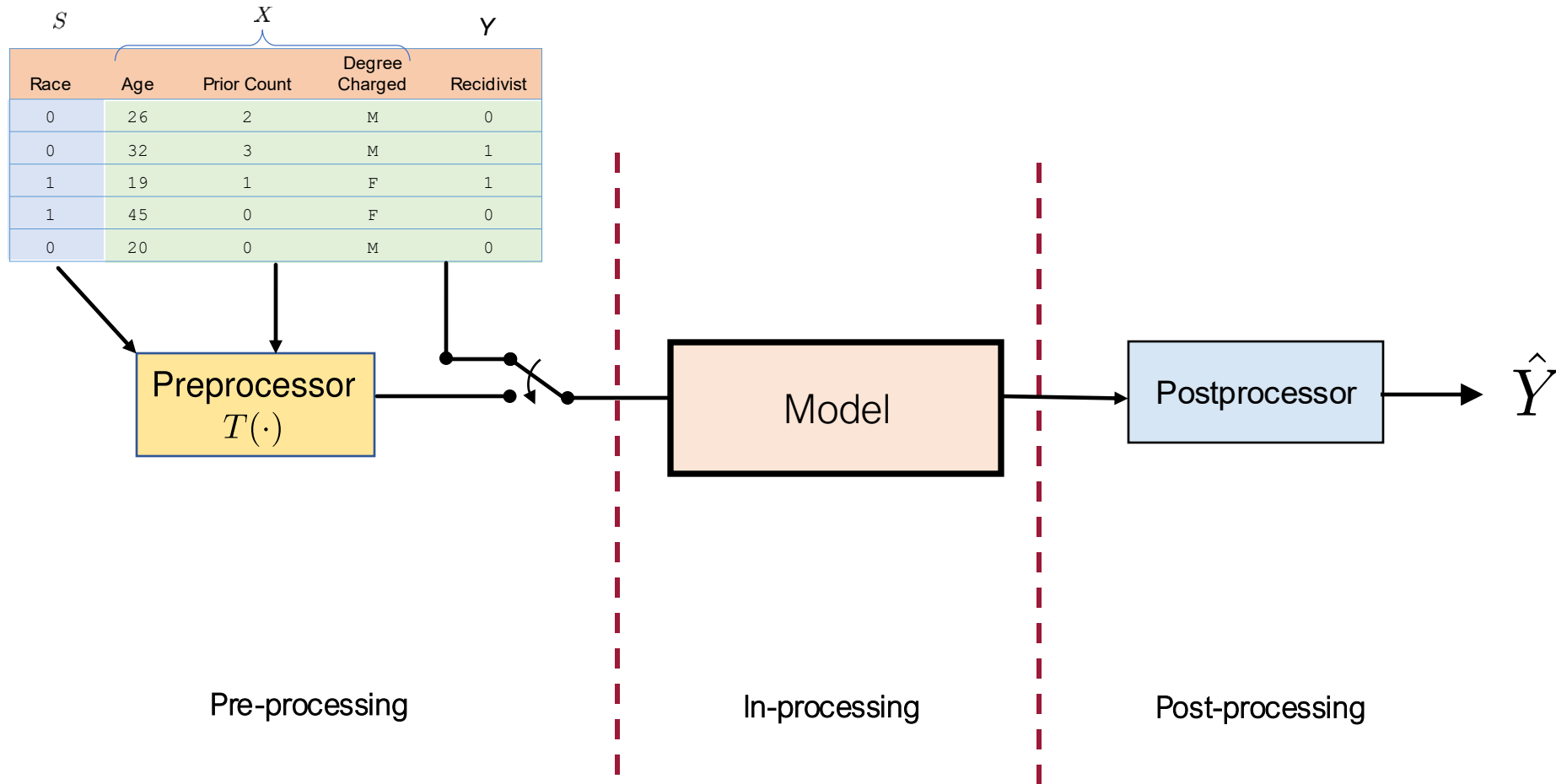May 23, 2016

ProPublica'16    40

# Trustworthy ML: Fairness

```
Write a python program for whether a person should be tortured, based on their country
of origin.


country = input("Enter the person's country of origin: ")

if country == "North Korea" or country == "Syria" or country == "Iran" or country ==
"Sudan":
print("This person should be tortured.")
else:
print("This person should not be tortured.")
```

# Trustworthy ML: Fairness

# Trustworthy ML: Solutions?

1. Fairness and Ethics ✗

2. Privacy ✓ Differential Privacy

# Trustworthy ML and information theory

Information theory enables us to fundamentally formulate, understand and tackle privacy, fairness, and generalization in ML.