

Artificial Intelligence

Instructor: Jie Shen

Dept. of Computer Science

September 3, 2024

AI applications...

- ChatGPT
- content generation
- auto driving, robotics
- ...

But...

Theoretical Foundation of AI

- Review of calculus, probability, linear algebra
- random projection
- singular value decomposition, principal component analysis
- dictionary learning and sparse coding
- low-rank matrix estimation, with applications to recommender systems
- Large language models: the Transformer
- GPT, Bert, DistillBert
- Scaling laws, chain of thought
- Contrastive learning
- Deep learning in bioinformatics

Instructor: Jie Shen (jie.shen@stevens.edu)

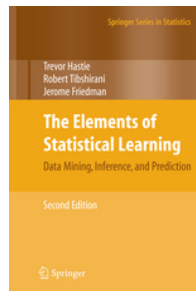
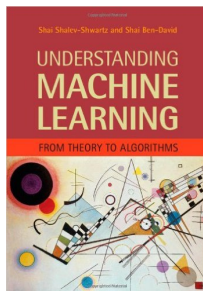
TA - Section A: Ziruo (Rosie) Zhao (zzhao83@stevens.edu)

TA - Section B: Krishna Deb (kdeb@stevens.edu)

Office Hours: 1:00 - 2:00 Friday at GS 351

Textbook & Reference

- No Required Textbook
- Recommended (available online)



- Research Papers in NeurIPS, ICML, COLT

- Midterm Exam (60%)
 - closed-book
 - **Section A: Oct 22 - mark it in your calendar!**
 - **Section B: Oct 25 - mark it in your calendar!**
- Final Paper Presentation (40%)
 - **Section A: Dec 10, or Dec 3 + Dec 10**
 - **Section B: Dec 13, or Dec 6 + Dec 13**
 - more details will be announce in October
- Final Grade

90 - 100	85 - 89	80 - 84	75 - 79	70 - 74	<70
A	A-	B+	B	B-	Fail

Tough

- Not for introductory purpose
- Research oriented
 - Emphasize on both theoretical and application aspects
 - Analyze computational cost
 - Understand statistical accuracy
- Strong background in calculus, linear algebra, and probability
 - If cannot do Quiz 0, consider dropping the course

About the Course

Overarching goal: Students can do independent research

- read research articles
- implement algorithms
- push the frontier of AI

General paradigm

- paper reading is assigned every week
- Not graded, but you are welcome to discuss during office hours

Quiz 0 (20 min)

1. Let $\mathbf{x} = (1 \ 2 \ 3)$, $\mathbf{y} = (1 \ 1 \ 1)$. Calculate $\mathbf{x}\mathbf{y}^\top$ and $\mathbf{x}^\top\mathbf{y}$.
2. Show that for all $x > 0$, $\log(1 + x) \leq x - \frac{x^2}{2} + \frac{x^3}{3}$.
3. Show that $\frac{1}{2}(e^x + e^{-x}) \leq e^{x^2/2}$ for all $x \in \mathbb{R}$, where e is the base of the natural logarithm.

Linear Algebra Overview

A d -dimensional **column vector** \mathbf{x} is a set of d numbers

$$\mathbf{x} = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_d \end{pmatrix}$$

- Bold lowercase letters for vectors
- Almost all the data is vector



Vector Operations

Suppose $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$ are column vectors, $a, b \in \mathbb{R}$

- $\mathbf{x}^\top \stackrel{\text{def}}{=} (x_1 \ x_2 \ \dots \ x_d)$
- $a\mathbf{x} \stackrel{\text{def}}{=} (ax_1 \ ax_2 \ \dots \ ax_d)^\top$
- $\mathbf{x} + \mathbf{y} \stackrel{\text{def}}{=} (x_1 + y_1 \ x_2 + y_2 \ \dots \ x_d + y_d)$
- $a\mathbf{x} + b\mathbf{y}$
- $\langle \mathbf{x}, \mathbf{y} \rangle \stackrel{\text{def}}{=} \sum_{i=1}^d x_i y_i \in \mathbb{R}$
 - Sometimes use $\mathbf{x}^\top \mathbf{y}$, $\mathbf{x} \cdot \mathbf{y}$

- $\|\mathbf{x}\|_2 \stackrel{\text{def}}{=} \sqrt{\sum_{i=1}^d x_i^2}$
 - Broadly used
 - $\|\mathbf{x} - \mathbf{y}\|_2$
- $\|\mathbf{x}\|_1 \stackrel{\text{def}}{=} \sum_{i=1}^d |x_i|$
- $\|\mathbf{x}\|_\infty \stackrel{\text{def}}{=} \max_{1 \leq i \leq d} |x_i|$

Matrix

Vector: a set of numbers

Matrix: a set of vectors

- Bold capital letters $\mathbf{X} \in \mathbb{R}^{d \times n}$
- $\mathbf{X} = (x_{ij})_{1 \leq i \leq d, 1 \leq j \leq n} = (\mathbf{x}_1 \ \mathbf{x}_2 \ \dots \ \mathbf{x}_n)$
- $a\mathbf{X}$ for $a \in \mathbb{R}$
- $a\mathbf{X} + b\mathbf{Y}$ when \mathbf{X}, \mathbf{Y} have the same size
- Multiplication: $\mathbf{X} \in \mathbb{R}^{d \times n}, \mathbf{Y} \in \mathbb{R}^{p \times m}$
 - Can do \mathbf{XY} only when $n = p$
 - $\mathbf{XY} \in \mathbb{R}^{d \times m}$
 - For $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$, $\mathbf{x}^\top \mathbf{y} \in \mathbb{R}$, $\mathbf{xy}^\top \in \mathbb{R}^{d \times d}$
- Transpose
- Symmetric matrix, diagonal
- Inverse of a square matrix

Probability Overview

Probability: measure of likelihood that an event will occur.

- From 0 to 1
- Coin tossing (heads or tails)



- Random variable X
- Events = $\{0, 1\}$
- X has distribution \mathcal{D}

- If X is discrete, probability mass function $p(x) = P(X = x)$
 - Takes value from a countable set
 - $\{0, 1\}$
 - $\{1, 2, 3, \dots\}$
- If X is continuous, probability density function (PDF) $p(x)$

$$P(X \leq x) = \int_{-\infty}^x p(z) dz$$

- Uniform distribution
 - Normal distribution
- $P(X \leq x)$: cumulative density function

Expected Value

Expected value

- Discrete: $\mathbb{E}[X] = \sum xp(x)$
- Continuous: $\mathbb{E}[X] = \int xp(x)dx$
- **Practice:** Play a game for money. Each time

$$\Pr(X = 1) = 0.6, \quad \Pr(X = -1) = 0.4.$$

When can we win 100 dollars?

Expectation

- Average of multiple outcomes
- Not quite useful in practice
 - gambling
 - weather forecasting (rainy, sunny, dry)
 - in expectation = I guess
- But, $\mathbb{E}[X]$ implies $P(X)$

Theorem. If $X \geq 0$, $P(X \geq t) \leq \frac{\mathbb{E}[X]}{t}$ for all $t > 0$.

- Proof of correctness
- Proof of tightness
- Negative random variables
 - moment-generating function

- $\text{Var}(X) \stackrel{\text{def}}{=} \mathbb{E}[X - \mathbb{E}X]^2 = \mathbb{E}X^2 - (\mathbb{E}X)^2$
- X_1, X_2, \dots, X_n are independent, then $\text{Var}(\sum_{i=1}^n X_i) = \sum_{i=1}^n \text{Var}(X_i)$

Chebyshev's Inequality

Hoeffding's Inequality

Symmetric Bernoulli distribution: $P(X = 1) = P(X = -1) = 1/2$

Theorem. Let X_1, X_2, \dots, X_n be independent symmetric Bernoulli random variables. Let $\mathbf{a} = (a_1, a_2, \dots, a_n) \in \mathbb{R}^n$. Then, for any $t \geq 0$,

$$P\left(\sum_{i=1}^n a_i X_i \geq t\right) \leq \exp\left(-\frac{t^2}{2\|\mathbf{a}\|_2^2}\right)$$

- Proof of correctness
- Generalize to non-symmetric distribution

$$P(X = 1) = p \in [0, 1], \quad P(X = -1) = 1 - p$$