

Artificial Intelligence

Instructor: Jie Shen

Dept. of Computer Science

September 3, 2024
b

AI applications...

2018 : v1
late 19: v2 → open-source from OpenAI
late 20: v3 → ~~closed~~ private
01/24: 3.5
07/24: v4 V2: small, 120M
 XL: 4-5B
 Llama 3.1: 8B

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

~2.00 GPU cards
AI00; \$16,000 RTX 6000 \$6,000
Time = 6 months
 > 30 < 8
VRAM 48 GB ← 70 B
 405 B
couple of years

performance

$f(\# \text{ params of model},$
 $\# \text{ samples},$

$\# \text{ compute resource})$

Open AI 2020:

Scaling laws.

Conclusion: always go
with larger models

* if \$ is fixed

$f = \# \text{ param}^{(0.75)}$

$\# \text{ samples}^{0.5}$

Deepmind 2022:

in 2020 is wrong.

param \approx # samples

Meta 2023:

" You selected
wrong obj func to
optimize "

obj \rightarrow inference
cost

\Rightarrow smaller & model

2024:

" I'm hot "

But...

Theoretical Foundation of AI

Syllabus

- Review of calculus, probability, linear algebra
- ✓ • random projection *(20 min → proof, 15 min → paper review)*
- 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

↓
mid-term
Oct 25

Course Staff

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

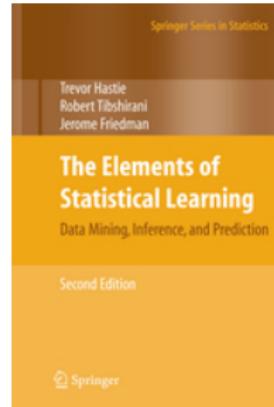
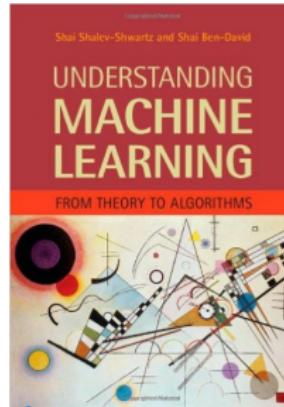
TA - Section A: Ziruo (Rosie) Zhao (zzhao83@stevens.edu)

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)



2014: 1000

2024: 1000+

- Research Papers in NeurIPS, ICML COLT pure theory
oral presentation → spotlight
~100

ICLR
↑

Grading

5 Problems 20 pt / problem

20 pt guaranteed.

20 pt 99%

→ 60 pt : open-ended

- 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%)

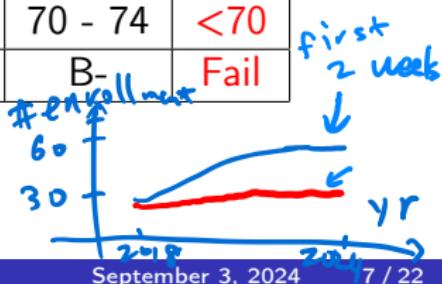
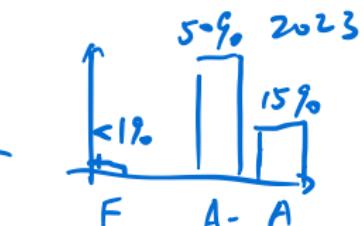
- Section A: Dec 10, or Dec 3 + Dec 10

- Section B: Dec 13, or Dec 6 + Dec 13

- more details will be announced in October

- Final Grade

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



About the Course

Tough

- Not for introductory purpose
- Research oriented
 - Emphasize on both theoretical and application aspects
 - Analyze computational cost
 - Understand statistical accuracy *testing loss v.s. #samples*
- 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)

4:10 → break

4:15 → continue.

~ 5:00 → dismiss

1. Let $x = (1 \ 2 \ 3)$, $y = (1 \ 1 \ 1)$. Calculate xy^T and $x^T y$.

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.

apply Taylor or wiki $f(x) = x - \frac{x^2}{2} + \frac{x^3}{3} - \log(1+x)$
on LHS RHS $f(x) \geq 0$ when $x > 0$

$$f(0) = 0$$

$$f(x) \geq f(0) \quad x > 0$$
$$f \uparrow \Leftrightarrow f' \geq 0$$

$$f: \mathbb{R}^d \rightarrow \mathbb{R}$$

$$f'(x) = \vec{0}$$

$$\hat{=} \left(\begin{array}{c|ccccc} \partial f & & & & & \\ \hline \partial x_1 & & & & & \\ & \partial f & & & & \\ & \hline & \partial x_2 & & & \\ & & \cdots & & & \\ & & & \partial f & & \\ & & & \hline & \partial x_d & \end{array} \right)$$

$$x = (x_1, \dots, x_d)$$

$$x_i \in \mathbb{R}$$

$$f(x) = x_1 + \cancel{x_2} + x_3 + \dots + x_d$$

$$\nabla f = (1, 1, 1, \dots, 1)$$

$$f(x) = \underline{x_1} \underline{x_2} + \underline{x_3} \cancel{x_4}$$

$$\nabla f = (2x_2, x_1, x_1, x_4, x_3)$$

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, $a, b \in \mathbb{R}$

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

Vector Norms

Linear reg:

min $\|w\|_2$ least-square

$$\text{S.t. } \|w\|_2 \leq \epsilon$$

- $\|x\|_2 \stackrel{\text{def}}{=} \sqrt{\sum_{i=1}^d x_i^2}$ CS 560
Statistical ML.

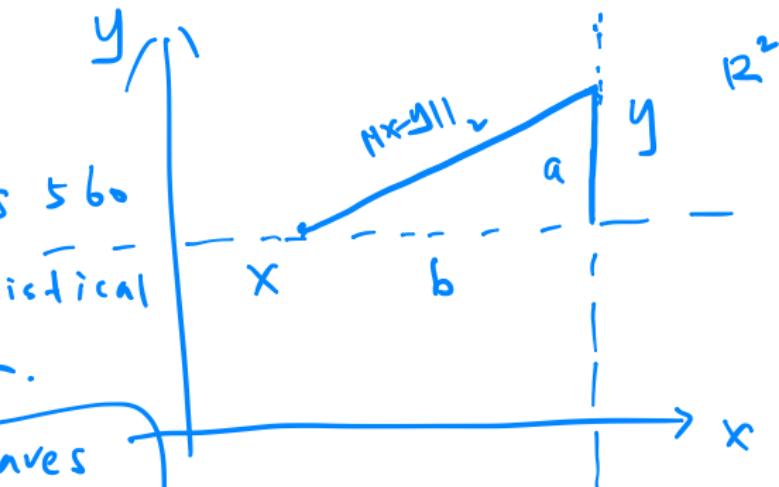
- Broadly used
- $\|x - y\|_2$

- $\|x\|_1 \stackrel{\text{def}}{=} \sum_{i=1}^d |x_i|$ → saves samples
- $\|x\|_\infty \stackrel{\text{def}}{=} \max_{1 \leq i \leq d} |x_i|$ samples

over-fitting

train loss ≈ 0

test loss $\rightarrow +\infty$



$$\|x-y\|_1 = a+b$$

$$\|y\|_\infty = \max\{a, b\}$$

test performance = $f(\# \text{ samples})$

In some applications, # samples is capped during training)

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}

Probability Overview

- 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



$N(\mu, \sigma^2)$

$$P(X \leq x) = \frac{1}{\sqrt{2\pi}} \cdot e^{-\frac{x^2}{2}}$$

Expected Value

$$X_1 \quad X_2 \quad \dots \quad X_N$$

$$S = \sum_{i=1}^N X_i \geq 100$$

$$E[X_i] = 0.2$$

$$E[S] = 0.2N \geq 100$$

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?

$$N = 500$$

$$\Pr(100) = 0.6^{100} \approx 0$$

$$\begin{aligned} & \boxed{40} \quad N_1 \quad \dots \quad \boxed{100} \\ & \frac{\frac{1}{t} \sum_{j=1}^t N_j}{t \rightarrow \infty} \approx 500 \end{aligned}$$

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 $P(X \geq 299)$, if play 2000 time
You must win \$100

Markov's Inequality

$$\mathbb{E}[X] = 10 \quad \Pr(X > t) \leq \frac{10}{t}$$
$$\Leftrightarrow \underline{\Pr(X < t)} \geq 1 - \frac{10}{t}$$

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

- Proof of correctness
- Proof of tightness
- Negative random variables

- moment-generating function

$$t = 1000$$

$$\text{pick } t = 100$$

$$\Rightarrow \Pr(X < 100) \geq 1 - 0.1 \\ = 0.9$$

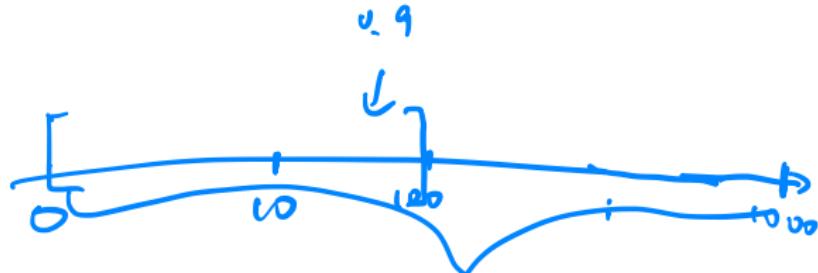
w.p. 0.9,

$X < 100$

$$\Pr(X < 1000) \geq 0.99$$

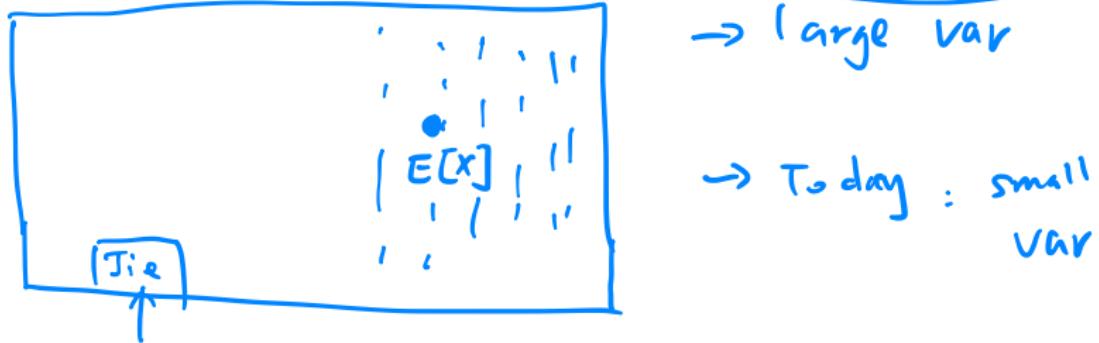
w.p. 0.99

$X < 1000$



Variance

- $\text{Var}(X) \stackrel{\text{def}}{=} \mathbb{E}[X - \mathbb{E}X]^2 = \mathbb{E}X^2 - (\mathbb{E}X)^2$
- 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

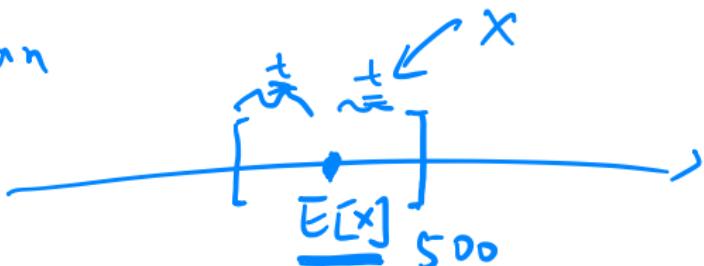
$$\Pr(|X - E[X]| > t) \leq \frac{\text{Var}(X)}{t^2}$$

$\forall X. \quad \underline{t > 0} \quad \Downarrow$

W.P. $1 - \frac{\text{Var}(X)}{t^2}$.

$$|X - E[X]| \leq t$$

X concentrates around its
mean



$$\text{W.P. } \left| X - E[X] \right| \leq t$$

1 - $\frac{\text{Var}(X)}{t^2}$

$$S = \sum_{i=1}^N X_i \geq 100$$

$$E[S] = 0.2N$$

$$\begin{aligned} \text{Var}[S] &= \text{Var}\left(\sum X_i\right) \\ &= \sum \text{Var}(X_i) = 0.96N \end{aligned}$$

$$\text{Var}(X_i) = E[X_i^2] - (EX_i)^2$$

$$= 1 - 0.2^2 = 0.96$$

$$\Rightarrow \text{W.P. } \left| 1 - \frac{0.96N}{t^2} \right| = 0.99$$

$$= 100 + |S - 0.2N| \leq t$$

$$0.2N - t \leq S \leq 0.2N + t$$

$$t > 0$$

$$t = 0.2N - 100 > 0 \Rightarrow N > 500$$

$$\left\{ \begin{array}{l} 0.2N - t = 100 \\ 1 - \frac{0.96N}{t^2} = 0.99 \end{array} \right. \Rightarrow \boxed{\begin{array}{l} N_1 = 500 \\ N_2 = 3000 \end{array}}$$

Solve over t . $\underline{\underline{N}}$

$$\Rightarrow \underline{\underline{N}} = \square$$

HW① Solve $\underline{\underline{N}} = \boxed{\begin{array}{l} 565 \\ \approx 3000 \\ \approx 3325 \end{array}} \approx 500 \approx 3000$

HW② write a python program

to see empirical result

Theory: $N = 3000$

~~100~~ $\hat{N} \approx 500$

Emp: $\hat{N} < 600$
 correct

$[450 \quad 600]$

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 $t \geq 0$,

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

Hoeffding's

$$\Pr(\sum X_i > t) \leq e^{-\frac{t^2}{2n}}$$

- Proof of correctness
- Generalize to non-symmetric distribution

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

\rightarrow Hoeffding's on bounded R.V.

Assume X_1, \dots, X_n independent

$$X_i \in [a, b] \quad \forall 1 \leq i \leq n$$

Then,

$$\Pr(|\sum X_i - E \sum X_i| > t) \leq e^{-\frac{t^2}{4(b-a)^2 n}}$$

In our case,

$$X_i \in [-1, 1]$$

are independent.

$$S = \sum X_i$$

$$\Pr(|S - 0.2N| > t) \leq e^{-\frac{t^2}{16n}}$$

$$\Leftrightarrow \text{w.p. } 1 - e^{-\frac{t^2}{16}}$$

$$|S - 0.2N| < t$$

Hoeffding's v.s. Chebyshev's

w.p.

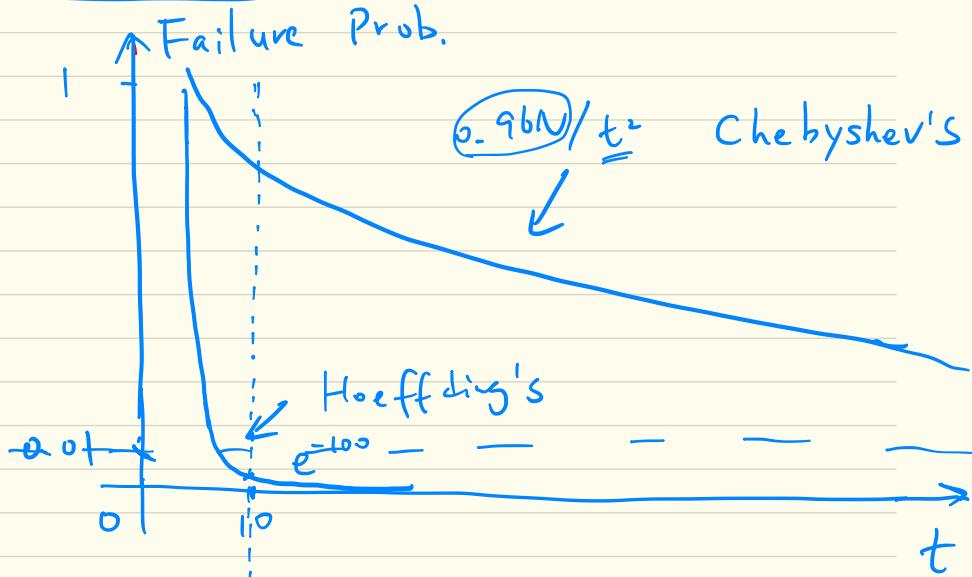
$$1 - \boxed{e^{-t^2/16N}}$$

$$1 - \boxed{\frac{0.96N}{t^2}}$$

$$|S - 0.2N| \leq t = 10$$

$$|S - 0.2N| \leq t = 10$$

Failure Prob.



$$e^{-t^2} \downarrow$$

$$\text{w.p. } 0.99\dots 9 > \frac{1}{100}$$

$$\frac{e^{-\frac{t^2}{16N}}}{\frac{0.96N}{t^2}} > \frac{1}{100}$$

$$\frac{1}{t^2} > \frac{1}{100} \Rightarrow 100 > t^2 \Rightarrow t < \sqrt{100} = 10$$

$$\text{w.p. } 0.99$$

Hoeffding's

C.

v v

$$\text{w.p. } 0.99 \dots 9 \underbrace{\dots}_{>100}$$

v

$$\text{w.p. } 0.99$$

$$|S - \mu_N| < 10$$

$$|S - \mu_N| < 10$$

$$e^{-t^2} = 0.01$$

$$\frac{1}{t^2} = 0.01$$

$$t = \sqrt{\ln 100}$$

$$t = 10$$

$$= \sqrt{2 \ln 10} \approx 2.5$$

$$< \sqrt{8} \approx 2.5$$

w.p. 0.99

w.p. 0.99

$$|S - \mu_N| < 10$$

$$|S - \mu_N| < 2.5$$

To day: Try Hoeffding's

to solve same problem.

$$N = \boxed{\quad}$$

stronger

M.

$$\text{W.P. } 1 - \frac{EX}{t}$$

$$X \leq t$$

C.

$$1 - \frac{\text{var}(x)}{t^2}$$

$$|X - EX| \leq t$$

H.

$$1 - e^{-t^2/\sigma^2}$$

$$\left| \sum X_i - E \sum X_i \right| \leq t$$

$X \neq EX$

Condition

$$\textcircled{1} \quad EX < \infty$$

both $\frac{EX}{t} < \infty$ $\frac{\text{Var}(X)}{t^2} < \infty$

$$EX^k < \infty \quad \text{for all } k = 1, \dots, \infty$$

$$\text{Var}(X) = (EX^2) - (EX)^2 < \infty$$

$$\textcircled{2} \quad X > 0$$

any X

$$X = \sum X_i \quad X_i \text{ are indep.}$$

estimate

stronger

condition

stronger

To show:

$$\boxed{\Pr(X_i = +1) = 0.5}$$

$$\boxed{\Pr(X_i = -1) = 0.5}$$

$$\Pr(\sum X_i - E \sum X_i > t) \leq e^{-\frac{t^2}{2n}}$$

Proof: LHS

$$e^{\lambda} \geq \Pr(\lambda(\sum X_i - E \sum X_i) > \lambda t)$$

$$= \Pr(e^{\lambda(\sum X_i - E \sum X_i)} > e^{\lambda t})$$

$$\stackrel{\text{Markov's}}{\leq} \frac{E[e^{\lambda(\sum X_i - E \sum X_i)}]}{e^{\lambda t}}$$

$$= \frac{E[e^{\lambda \cdot \sum X_i}]}{e^{\lambda t}}$$

$$\leq \frac{e^{\frac{n}{2}\lambda^2}}{e^{\lambda t}} = e^{\frac{n}{2}\lambda^2 - \lambda t} \quad \begin{matrix} \text{Holds} \\ \forall \lambda > 0 \end{matrix}$$

$$E[e^{\lambda \sum X_i}] = e^{a+b} = e^a \cdot e^b$$

$$= E\left[\prod_{i=1}^n e^{\lambda X_i}\right]$$

since X_1, \dots, X_n are indep.

$\Rightarrow e^{\lambda X_1}, \dots, e^{\lambda X_n}$ are indep.

A, B are indep. $\Rightarrow E[A \cdot B] = E[A] \cdot E[B]$

$$= \prod_{i=1}^n E[e^{\lambda X_i}] \rightarrow \text{MGF of } X_i$$

$$= \prod_{i=1}^n \left(e^{\lambda \cdot 1 \cdot \frac{1}{2}} + e^{\lambda \cdot (-1) \cdot \frac{1}{2}}\right)$$

$$= \prod_{i=1}^n \frac{1}{2} (e^\lambda + e^{-\lambda})$$

$$\leq \prod_{i=1}^n e^{\frac{\lambda^2}{2}} = e^{\frac{n}{2} \lambda^2}$$

$$\Rightarrow \forall \lambda > 0$$

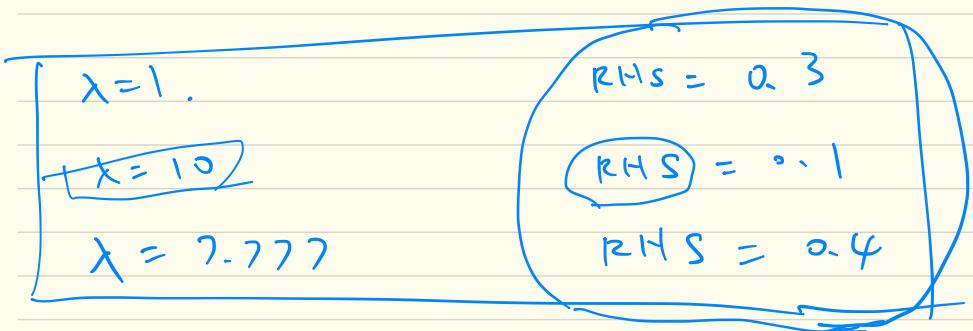
$$Pr(\sum X_i - E \sum X_i \text{rot} \leq t)$$

$$e^{-t^2/2n}$$

$$\frac{\frac{n}{2} \lambda^2 - \lambda t}{e^{\frac{n}{2} \lambda^2 - \lambda t}}$$

minimize
w.r.t. λ .

Find "optimal" λ



$$\lambda = 1$$

$$\text{w.p. } \underline{0.7}$$

$$\sum X_i - E \sum X_i < t$$

$$\lambda = 10$$

$$\text{w.p. } \underline{0.9}$$

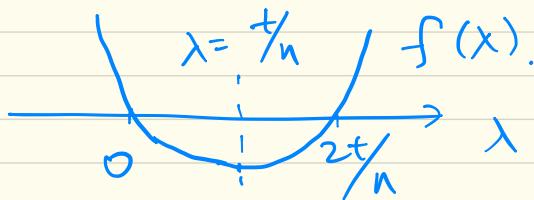
$$\frac{1}{\sum X_i - E} < t$$

$$\lambda = 7.777 \text{ w.p. } \underline{0.6}$$

$$\sum X_i - E < t$$

$$\min f(\lambda) = \frac{n}{2} \lambda^2 - t\lambda$$

$$\text{s.t. } \lambda > 0$$



$$f(t/n) = -\frac{t^2}{2n}$$

4:20

$$\Pr(|X - \bar{X}| > t) < f(t)$$

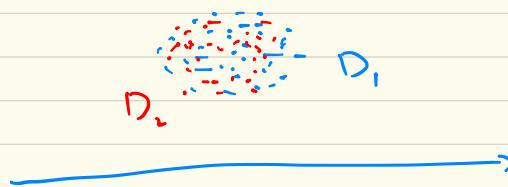
↑
Failure prob.

$$\text{W.P. } 1 - f(t)$$

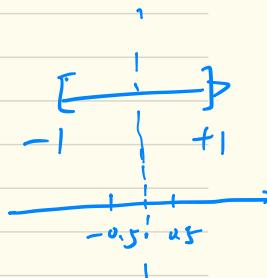
$$|X - \bar{X}| \leq t$$

Hoeffding's needs all info of distribution / $E X^k < \infty$.

$$D_1 = D_2$$



$$E X = E X \\ X \sim D_1 \quad X \sim D_2$$



$$\text{Var}(X) = \text{Var}(X)$$

$$(D_1 \rightarrow E X^2 = E X^2) \\ D_2$$

$$E X^3 \\ D_1 : D_2$$

$$E X^\infty \quad E X^\infty$$

$$\text{if } E X^k = E X^k \\ D_1 \qquad \qquad D_2$$

for all $k=1, \dots \infty$

$$\Rightarrow D_1 = D_2$$

\Rightarrow any distribution is
uniquely determined
by the moments $\{E X^k\}_{k=1}^{\infty}$

Moment Generating func. of X :

$$\stackrel{\triangle}{=} E[\underbrace{e^{\lambda x}}_{X \sim D}] = g(\lambda)$$

$$e^{\lambda y} = 1 + \underline{\frac{y}{1!}} + \underline{\frac{y^2}{2!}} + \underline{\frac{y^3}{3!}} + \underline{\frac{y^4}{4!}} + \dots$$

$$D \Leftrightarrow g(\lambda)$$

Random Projection

How to effectively reduce dim?

- c - dataset

$$\begin{bmatrix} 1 & | & 28 \\ \hline 28 & & \end{bmatrix} \quad \text{---} \quad \text{---}$$

784 - dim

Query $\boxed{1} \rightarrow \begin{pmatrix} & \\ & \end{pmatrix}_{784}$

Retrieve all images with "1"

metric on (x_1, x_2)

$$d \rightarrow 0 \quad x_1 \approx x_2 \quad d(x_1, x_2)$$

$$\rightarrow \infty \quad x_1 \neq x_2 = \|x_1 - x_2\|_2$$

n : # img in database,

d : dim

$O(n d)$

- ① reduce $n \rightarrow$ hashing
(not in this course)
- ② reduce $d \rightarrow$ dim reduction ✓

$$d = 784$$

$$k = 1$$

10

Given k

How to evaluate performance¹⁰⁰



of the new data? ~~203~~

x_i in DB.

$$\| q - x_i \|_2$$

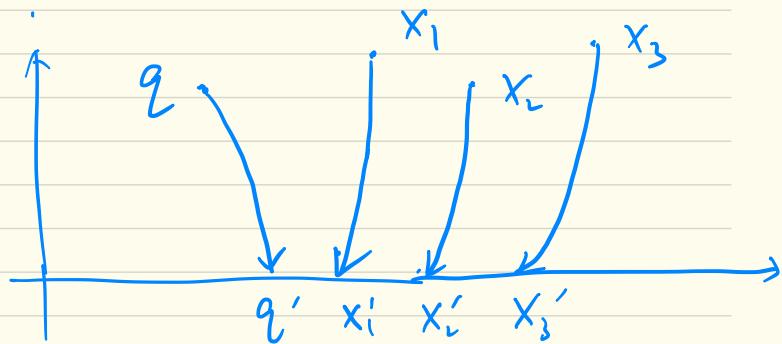
to Goal: on new data, $\{q', x'_1, \dots, x'_n \in \mathbb{R}^k\}$

Top 100 are same as
Top 100 \mathbb{R}^d

$$f: q \rightarrow q' = f(q) \quad \mathbb{R}^d \rightarrow \mathbb{R}^k$$

$$x_i \rightarrow x'_i = f(x_i)$$

$$\| q' - x'_i \|_2 = \| f(q) - f(x_i) \|_2$$



$$\boxed{\| f(q) - f(x_i) \|_2 \approx \| q - x_i \|_2}$$

$$f: \quad X \rightarrow X'$$

$$x \rightarrow M \cdot x.$$

$M: k \times d$ matrix X

$$\downarrow \quad M_{ij} \sim N(0, \frac{1}{d})$$

random matrix

random projec.

A 2. X

$$\|M \cdot q - M \cdot x\| \approx \|q - x\|$$



$$\|M \cdot (q - x)\| \approx \|q - x\|$$

Goal: w.p. 0.99
 $\underbrace{\|M \cdot x\|}_{R.V.}$

$$\|M \cdot x\| \approx \boxed{E[\|M \cdot x\|]}$$

next week.

key msg from concentration
inequalities:

w. h.p.

$$X \approx E[X]$$