# EE 53100 Concentration Inequalities

## Introduction and preliminaries

- Course webpage:
    https://people.iith.ac.in/shashankvatedka/html/courses/2023/EE5603/course_details.html

- Announcements, homework submissions: Google classroom
                                        (send me an email if you do not get an invite by tomorrow)

- Prerequisites:
    - Strong foundation in probability and random processes
    - Some background in information theory/statistics/machine learning is helpful, but not mandatory
    - Programming in python

- Class timings: Slot B (Mon 10am, Tue 9am, Thu 11am)

- Assessment:
    - homeworks (roughly 3, totaling 55%)
    - 3 quizzes/tests (15% each)

- References: See course webpage

# Motivation and background

① <u>n tosses of a fair coin</u>

for large $n$, fraction of heads $\approx \frac{1}{2}$

What about finite $n$?

$\Pr[\text{ fraction of heads} \geq \frac{1}{2} + \sigma ]$ ?

$\Pr[\# \text{ heads} = k] = \binom{n}{k} \left(\frac{1}{2}\right)^n$

$\Pr[\# \text{ heads} \geq n\left(\frac{1}{2} + \sigma\right)] = \quad ?$

② **Communication systems:**

$$Y = X + Z$$

$$Z \sim \mathcal{N}(0, \sigma^2)$$

(i) $X \in \{+\sqrt{P}, -\sqrt{P}\}$

$$\hat{X} = \begin{cases} +\sqrt{P} & \text{if } Y > 0 \\ -\sqrt{P} & \text{if } Y \leq 0 \end{cases}$$

What is $\Pr[\hat{X} \neq X]$?

(ii) $Y^n = X^n + Z^n$ $\qquad X^n \in \mathbb{R}^n \qquad Z^n \sim \text{iid } \mathcal{N}(0, \sigma^2)$

$\longrightarrow$ codebook

$X^n \in C$

$$\hat{X}^n = \underset{x^n \in C}{\text{argmin}} \ \|Y^n - x^n\|^2$$

What is $\Pr[\hat{X}^n \neq X^n]$?

③ **Empirical risk minimization :**

Binary Classification :  $(X, Y) \sim P_{XY}$    $X \in \mathcal{R}$

$Y \in \{0, 1\}$

- Do not know $P_{XY}$

- Have access to  $(X_1, Y_1)(X_2, Y_2) -- (X_n, Y_n)$  iid $P_{XY}$

$\underbrace{\phantom{(X_1, Y_1)(X_2, Y_2) -- (X_n, Y_n)}}_{\text{dataset}}$

- Problem : Choose $f \in \mathcal{F}$ that minimizes expected risk

$$R(P_{XY}) = \mathbb{E}\left[ L(f(X), Y) \right]$$

eg:    $R(P_{XY}) = \mathbb{E}\left[ 1_{\{f(X) \neq Y\}} \right] = Pr[f(X) \neq Y]$

- Choose $f : \mathcal{X} \to \{0, 1\}$ that best approximates data

$$\hat{f} = \underset{f \in \mathcal{F}}{\arg\min} \quad \frac{1}{n} \sum_{i=1}^{n} L\left(f(x_i), y_i\right)$$

- How well does this approximate $\mathbb{E}\left[L(f(x), Y)\right]$ ?

# Probability Basics

① Probability space $(\Omega, \mathcal{F}, P)$

$\Omega$ — Sample space      (collection of outcomes)

$\mathcal{F}$ — Event space

$P$ — Probability measure

Event space : Should form a sigma algebra

    ① $\Omega \in \mathcal{F}$

    ② $A \in \mathcal{F} \Rightarrow A^c \in \mathcal{F}$

    ③ Countable collection $A_1, A_2, \ldots \in \mathcal{F} \Rightarrow \bigcup_{i=1}^{\infty} A_i \in \mathcal{F}$

# Examples:

① $\Omega = \{1, 2, 3, 4, 5, 6\}$

What is the smallest event space?   $\{\emptyset, \Omega\}$

What is the smallest event space containing $\{\{1,2\}, \{2\}\}$?

What is the largest event space?   $\{\emptyset, \Omega, \{1, 2\}, \{2\}$

$\{3, 4, 5, 6\} \{1, 3, 4, 5, 6\}$

$\{2, 3, 4, 5, 6\}, \{1\}\}$

② $\Omega = \mathbb{R}$

The smallest sigma algebra that contains all intervals of the form $(a, b)$ for $a, b \in \mathbb{R}$ is called the Borel sigma algebra.

Lebesgue

**Probability measure:** $\quad P: \mathcal{F} \longrightarrow [0, 1]$

① $P(\Omega) = 1$

② $A_1, A_2 \cdots \in \mathcal{F}$ (countable collection of events)

$\quad$ disjoint

$\Rightarrow P\left(\overset{\infty}{\underset{i=1}{\cup}} A_i\right) = \overset{\infty}{\underset{i=1}{\sum}} P(A_i)$

Q: Why not assign probability measure on outcomes?

$\quad$ (Eg: Uniform distribution on $[0, 1]$)

Q: Why not always choose $\mathcal{F} = $ power set of $\Omega$?

$\quad$ A: Not possible to have a consistent measure

$\quad$ Eg: Vitali set $\qquad$ https://en.wikipedia.org/wiki/Vitali_set

## Random variable :

$$X : \Omega \longrightarrow \mathbb{R}$$

$$\text{s.t} \quad X^{-1}((a,b]) \in \mathcal{F}$$

(every interval corresponds to a valid event)

**Example :** Infinite sequence of coin tosses

H T H H H $-$ $-$

0 . 0 1 0 0 0

Say $b_i = \begin{cases} 0 & \text{if the ita toss} = H \\ 1 & \text{if } T \end{cases}$

$0 . b_1 b_2 b_3 - - -$

$$x = b_1/2 + b_2/4 + b_3/8 + -$$

$$= \sum_{i=1}^{\infty} b_i / 2^i$$

$$X = \sum_{i=1}^{\infty} b_i / 2^i$$

$$\Pr\left[X \in [0, 1/2)\right] = 1/2$$

$$\Pr\left[X \in [1/4, 1/2)\right] = 1/4$$

$$\Pr\left[X \in [a, b)\right] = b - a \longrightarrow \text{Uniform distribution on } [0, 1]$$

$$\Pr\left[X \in \left[\hat{j}/2^i, \frac{\hat{j}+1}{2^i}\right)\right] = 1/2^i$$

$$\underbrace{\qquad\qquad}_{b_1 b_2 \cdots b_i}$$

$$\left[1/3, 2/3\right)$$

# Probability measure:

for a random variable, specified by the cumulative distribution function (cdf)

$$f_X(x) = \Pr[X \leq x]$$

$(a, b]$

$(-\infty, b]$

Discrete: - Probability mass function

$$p_X(x) = \Pr[X = x] \longrightarrow \text{Assign probabilities to outcomes}$$

Continuous: Probability density function

$$f_X(x) = \frac{d}{dx} f_X(x)$$

# Common distributions : Make sure that you know the pmf/pdf :

1. Bernoulli

2. Binomial

3. Poisson

4. Uniform (discrete & continuous)

5. Gaussian

6. Laplace

7. Gamma

8. Chi-squared.

# Expectation

$$\mathbb{E}[X] = \sum_{x \in \mathfrak{X}} x \, p_X(x) \quad , \text{ if discrete}$$

$$= \int_{-\infty}^{\infty} x \, f_X(x) \, dx \quad , \text{ if continuous}$$

$$\mathbb{E}[g(X)] = \begin{cases} \displaystyle\sum_{x \in \mathfrak{X}} g(x) \, p_X(x) \\[2em] \displaystyle\int_{-\infty}^{\infty} g(x) \, f_X(x) \, dx \end{cases}$$

- Mean : $\mathbb{E}X = \mu$

- Variance : $\mathbb{E}\left[(X - \mathbb{E}X)^2\right] = \sigma^2$

# Moment generating function

$k$'th moment: $\mathbb{E}(X^k)$ $\longrightarrow$ noncentral

$\mathbb{E}\left[(X-\mu)^k\right]$ $\longrightarrow$ central.

MGF: $\mathbb{E}\,e^{\lambda X} = \begin{cases} \displaystyle\sum_{x \in \mathcal{X}} e^{\lambda x} p_X(x) \\[2em] \displaystyle\int_{-\infty}^{\infty} f_X(x)\, e^{\lambda x}\, dx \end{cases}$

Exercise: Compute the MGF for all the common distributions listed above.

# Basic tail bounds : $\Pr[X > \sigma]$

**Markov inequality :** Let $X$ be a non-negative random variable $\left( i.e., \Pr[X < 0] = 0 \right)$ & $\mathbb{E}X = \mu < \infty$

Then,

$$\Pr[X \geqslant t] \leqslant \frac{\mu}{t} \qquad \forall t > 0$$

→ useful only if $t > \mu$

**Proof :**

$$\mu = \int_0^\infty x f_x(x) \, dx = \int_0^t x f_x(x) \, dx + \int_t^\infty x f_x(x) \, dx$$

$$\frac{\mu}{t} \geqslant \int_t^\infty x f_x(x) \, dx \geqslant t \int_t^\infty f_x(x) \, dx = t \frac{\Pr[X \geqslant t]}{t}$$

# Chebyshev inequality:

Let $X$ be a rv with $\mathbb{E}X = \mu < \infty$,

$$\sigma^2 = \mathbb{E}\left[(X-\mu)^2\right] < \infty.$$

Then,

$$\Pr\left[\ |X-\mu| > t\ \right] \leq \frac{\sigma^2}{t^2} \qquad t > 0$$

Proof:

$$\Pr\left[|X-\mu| > t\right] = \Pr\left[\underbrace{(X-\mu)^2}_{Y} > t^2\right]$$

$$\leq \frac{\mathbb{E}Y}{t^2} = \frac{\sigma^2}{t^2}.$$

# Sequences of random variables

$X_1, X_2, X_3 \cdots X_n$ are independent & identically distributed if

$$p_{X_1 \cdots X_n}(x_1, \cdots, x_n) = \prod_{i=1}^{n} p_X(x_i)$$

$$f_{X_1 \cdots X_n}(x_1 - x_n) = \prod_{i=1}^{n} f_X(x_i)$$

Q: Let $Y = \frac{1}{n} \sum_{i=1}^{n} X_i$ where $X_1 \cdots X_n$ are iid

What is $\mathbb{E} Y$? $= \mathbb{E} X = \mathbb{E} X_1$

What is $\text{Var}(Y)$? $= \text{Var}\left(\frac{1}{n} \sum_{i=1}^{n} X_i\right) = \frac{1}{n^2} \sum_{i=1}^{n} \text{Var}(X_i) = \sigma^2/n$

$$\Pr\left[ \left| \underbrace{\frac{1}{n} \sum_{i=1}^{n} X_i}_{Y_n} - \mu \right| > \epsilon \right] \leq \frac{\sigma^2}{\frac{n}{\epsilon^2}} = \frac{\sigma^2}{n\epsilon^2} \rightarrow 0 \text{ as } n \rightarrow \infty$$

# Convergence of random variables

① We say that $X_1, X_2, X_3 \cdots$ converges to $X$ <u>in probability</u> if

$$\Pr\left[\,|X_n - X| > \epsilon\,\right] \to 0 \quad \text{as } n \to \infty \qquad \text{<u>for all</u> } \epsilon > 0$$

We denote this as $X_n \xrightarrow{p} X$

② We say that $X_1, X_2 \cdots$ converges to $X$ <u>almost surely</u>/ with <u>probability</u> 1 if

$$\Pr\left[\lim_{n \to \infty} X_n = X\right] = 1 \quad \Longleftrightarrow \quad \Pr\left[\lim_{n \to \infty} |X_n - X| = 0\right] = 1$$

We denote this as $X_n \xrightarrow{a.s.} X$

③ We say that $X_1 X_2 X_3 \cdots$ converges to $X$ in <u>distribution</u> if

$$\lim_{n \to \infty} F_{X_n}(x) = F_X(x)$$

at all points where $F_X$ is continuous.

We denote this $X_n \xrightarrow{d} X$

④ We say that $X_1 X_2 \cdots$ converges to $X$ in $L^p$ (for $p \geq 1$) if

$$\lim_{n \to \infty} \mathbb{E}\left[ |X_n - X|^p \right] = 0$$

# Examples :

① $(X_n, n \geq 1)$ independent

$$\Pr[X_n = 0] = 1 - 1/n \qquad \Pr[X_n = 1] = \frac{1}{n}$$

$$\Pr[X = 0] = 1 \qquad F_X(x) = \begin{cases} 0, & x < 0 \\ 1, & x \geq 0 \end{cases}$$

① $F_{X_n}(x) = \begin{cases} 0, & x < 0 \\ 1 - 1/n, & x \in [0, 1) \\ 1 & x \geq 1 \end{cases}$

$$\lim_{n \to \infty} F_{X_n}(x) = \begin{cases} 0, & x < 0 \\ 1, & x \geq 0 \end{cases} = F_X(x)$$

$$\therefore \quad X_n \xrightarrow{d} X$$

② $\Pr\left[ |X_n - X| > \epsilon \right] = \Pr\left[ |X_n| > \epsilon \right] = \frac{1}{n}$

$$\longrightarrow 0$$
$$\text{as } n \to \infty$$

③ $\mathbb{E}|X_n - X|^p = \mathbb{E}|X_n|^p = 0^p\left(1 - \frac{1}{n}\right) + 1^p \frac{1}{n}$

$$= \frac{1}{n} \longrightarrow 0 \qquad \text{as } n \to \infty$$

④ $\Pr\left[ \lim_{n \to \infty} |X_n - X| = 0 \right] \qquad X_n \xrightarrow{a.s} 0$

$\Pr\left[ X_n < \epsilon \quad \text{for all } n \geq N \right] = \Pr\left[ X_N = 0, X_{N+1} = 0, \cdots \right]$

$$= \left(1 - \frac{1}{N}\right)\left(1 - \frac{1}{N+1}\right)\left(1 - \frac{1}{N+2}\right) \cdots = 0$$

$$\Pr[X_n = 1] = e^{-n} \qquad \Pr[X_n = 0] = 1 - e^{-n}$$

$$- \quad \Pr[X_n = 1/n] = 1/n \quad , \quad \Pr[X_n = 0] = 1 - \frac{1}{n}$$

For given $\varepsilon > 0$, can we find $N$ st

$$\Pr[X_n > \varepsilon \quad \text{for} \quad n \geqslant N] = 0 \qquad \text{Yes.}$$

$$X_n \xrightarrow{as} X$$

② $X_1 \sim Ber(1/2)$

$$X_n = \begin{cases} X_1 & \text{if } n \text{ is odd} \\ \bar{X_1} & \text{if } n \text{ is even} \end{cases}$$

$Pr\left[ |X_n - X| > \varepsilon \right]$  does not converge for any $X$  as $n \to \infty$

$$f_{X_n}(x) = \begin{cases} 0, & x < 0 \\ 1/2, & x \in [0, 1) \\ 1, & x \geq 1 \end{cases} \qquad \text{for all } n$$

$$X_n \xrightarrow{L^p} X$$

$$\Downarrow \quad \not\Uparrow$$

Note : $\qquad X_n \xrightarrow{a.s} X \quad \Longrightarrow \atop \nLeftarrow \quad X_n \xrightarrow{P} X \quad \Longrightarrow \atop \nLeftarrow \quad X_n \xrightarrow{d} X$

in general.

# Weak law of large numbers (WLLN)

If $X_1, X_2, \ldots$ are iid with mean $\mu$, then

$$\frac{1}{n} \sum_{i=1}^{n} X_i \xrightarrow{P} \mu$$

$$P_n\left[ \left| \frac{1}{n} \sum_{i=1}^{n} X_i - \mu \right| > \varepsilon \right] \xrightarrow{\leq \frac{\sigma^2}{n\varepsilon^2}} 0 \quad \text{as } n \to \infty$$

# Strong law of large numbers (SLLN)

If $X_1, X_2, \ldots$ are iid with mean $\mu$, then

$$\frac{1}{n} \sum_{i=1}^{n} X_i \xrightarrow{a.s} \mu$$

for any $\varepsilon > 0$, $\delta > 0$, $\exists N$ st

$$P_n\left[ \left| \frac{1}{n} \sum_{i=1}^{n} X_i - \mu \right| \leq \varepsilon \text{ for all } n > N \right] \geq 1 - \delta$$

Proof of WLLN if $\text{Var}(X_i) < \infty$

# Central limit theorem

If $X_1, X_2 \cdots$ are iid with finite mean $\mu$ & finite variance $\sigma^2$, then

$$\frac{\sum\limits_{i=1}^{n} X_i - n\mu}{\sqrt{n\sigma^2}} \xrightarrow{d} N(0,1)$$