# PRELIMINARIES

## Concentration Inequalities (Foreword)

Two classical results to characterize the "many-sample" behavior of independent random variables:

① [Strong Law of Large Numbers]

$X_1, X_2, \ldots, X_n$ independent r.v.s, $\mathbb{E}X_i = 0$,

$$\mathbb{E}[X_i^4] \leq R < \infty, \quad \forall i \in \mathbb{N},$$

then $\quad \mathbb{P}\left( \frac{1}{n} \sum_{i=1}^{n} X_i \xrightarrow[n \to \infty]{} 0 \right) = 1.$

② [Central Limit Theorem]

$(X_n)_{n \in \mathbb{N}}$ is an iid sequence, $\quad E[X_i] = 0,$

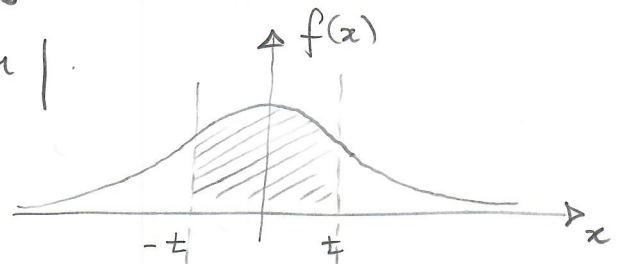$$\mathrm{Var}(X_i) = \sigma^2 < \infty, \quad \forall i \in \mathbb{N}.$$

Set $\quad G_n := \dfrac{\sum_{i=1}^{n} X_i}{\sigma \sqrt{n}}.$

Then, $\forall x \in \mathbb{R}, \quad \lim_{n \to \infty} \mathbb{P}(G_n \leq x) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{x} e^{-\partial^2/2} \, d\partial$

**Observation**   Consider independent and identically distributed r.v's $(X_n)_{n \in \mathbb{N}}$ with $\mathbb{E}[X_i] = \mu$. Then, $\frac{1}{n} \sum_{i=1}^{n} X_i$ converges asymptotically to $\mu$, as $n \to \infty$.

<u>**Question**</u>   Finite-sample guarantees for the deviation

$$\left| \frac{1}{n} \sum_{i=1}^{n} X_i - \mu \right|.$$



We are mainly interested in inequalities s.t.

$$\mathbb{P}\left( \left| \frac{1}{n} \sum_{i=1}^{n} X_i - \mu \right| > t \right) \leq \delta(n,t),$$

where $\quad \delta(n,t) \in (0,1)$ is s.t. $\quad \lim_{n \to \infty} \delta(n,t) = 0, \quad \forall t \geq 0,$

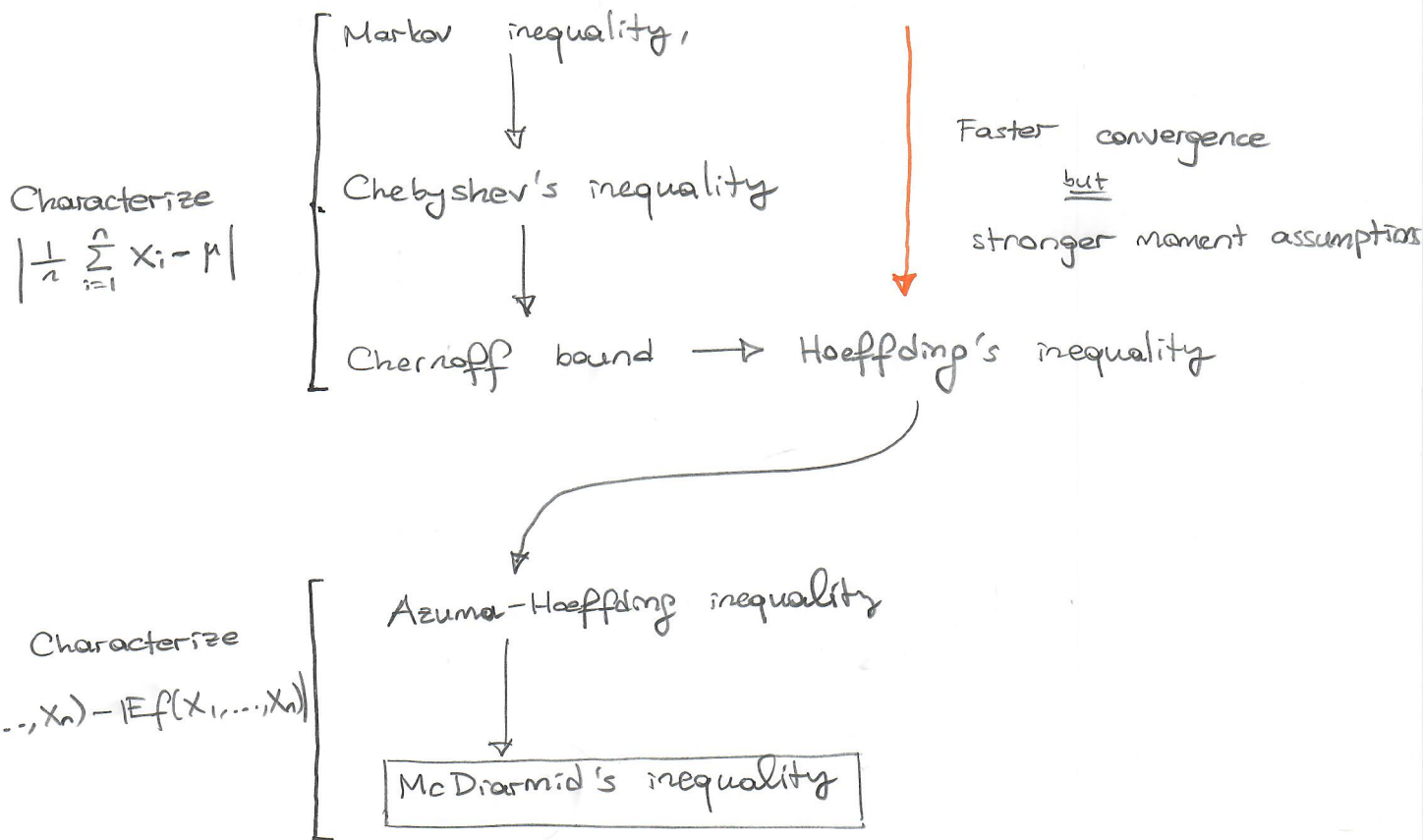$$\lim_{t \to \infty} \delta(n,t) = 0, \quad \forall n \in \mathbb{N}.$$

Of course, moments of $X_i$, i.e., $\mathbb{E}|X_i|^k$, $k \in \mathbb{N}$ will have a crucial importance in deriving these bounds.

The concentration behavior is much more general than just the concentration of the sample mean $\frac{1}{n}\sum_{i=1}^{n} X_i$ around the true mean.

Under very general conditions, as long as $(X_n)_n$ is an independent sequence of random variables,

$$f(X_1, X_2, \ldots, X_n)$$ also shows concentration behavior around $\mathbb{E}f(X_1, X_2, \ldots, X_n)$.

As such, the discussion will be as follows:

Characterize $\left|\frac{1}{n}\sum_{i=1}^{n} X_i - \mu\right|$

Markov inequality,

↓

Chebyshev's inequality

↓

Chernoff bound ⟶ Hoeffding's inequality

Faster convergence
but
stronger moment assumptions

Characterize $\left|f(X_1,\ldots,X_n) - \mathbb{E}f(X_1,\ldots,X_n)\right|$

Azuma-Hoeffding inequality

↓

McDiarmid's inequality

McDiarmid's inequality will be central to generalization bounds.
Hoeffding's inequality is useful for analyzing wide neural networks.

## Basics : Markov, Chebyshev, Chernoff

**Markov's Inequality** | Let $Y \in \mathbb{R}$ be a non-negative random variable with $\mathbb{E}|Y| < \infty$. Then, for any $t > 0$,

$$\mathbb{P}(Y > t) \leq \frac{\mathbb{E}Y}{t}.$$

*Pf*: Very easy. Since $Y \geq 0$ a.s., $t > 0$,

$$Y = Y(\mathbb{1}\{Y \geq t\} + \mathbb{1}\{Y < t\}) \geq Y\mathbb{1}\{Y \geq t\}$$
$$\geq t\,\mathbb{1}\{Y \geq t\}.$$

Taking expectation,

$$\mathbb{E}[Y] \geq \mathbb{E}[t\,\mathbb{1}\{Y \geq t\}] = t\,\mathbb{P}(Y \geq t).$$

**Chebyshev's inequality** | Let $X \in \mathbb{R}$ be a random variable with $\mathbb{E}X^2 < \infty$. Then, for any $t > 0$,

$$\mathbb{P}(|X - \mathbb{E}X| > t) \leq \frac{\text{Var}(X)}{t^2}.$$

*Pf*: Again, super easy. $Y := |X - \mathbb{E}X| \geq 0$ a.s. For any $t > 0$,

$$\mathbb{P}(Y > t) = \mathbb{P}(Y^2 > t^2) \leq \frac{\mathbb{E}Y^2}{t^2} = \frac{\text{Var}(X)}{t^2}.$$

How to use these inequalities?

$X_1, X_2, \ldots, X_n$ independent $\mathbb{R}$-valued r.v.'s with $\mathbb{E}X_i = \mu_i$
$$\text{Var}(X_i) \leq \sigma^2 < \infty$$

Then, $\frac{1}{n}\sum_{i=1}^{n}(X_i - \mu_i) =: Z_n$, $\mathbb{E}[Z_n] = 0$,

$$\text{Var}(Z_n) = \frac{1}{n^2}\mathbb{E}\sum_{i,j}(X_i - \mu_i)(X_j - \mu_j) \leq \frac{\sigma^2}{n}.$$

$$\Rightarrow \quad \mathbb{P}\left(\left|\frac{1}{n}\sum_{i=1}^{n}X_i - \frac{1}{n}\sum_{i=1}^{n}\mu_i\right| > t\right) \leq \frac{\sigma^2}{nt^2}$$

If $(X_i)_i$ are iid, then $\mu_i = \mu$, $\text{Var}(X_i) = \sigma^2$,

$$\mathbb{P}\left(\left|\frac{1}{n}\sum_{i=1}^{n}X_i - \mu\right| > t\right) \leq \frac{\sigma^2}{nt^2}$$

## Chernoff bound

$X \in \mathbb{R}$, $\lambda > 0$, $t \in \mathbb{R}$,

$$P(X - \mathbb{E}X > t) = P(\lambda \cdot (X - \mathbb{E}X) > \lambda t)$$

$$= P\left(e^{\lambda(X - \mathbb{E}X)} > e^{\lambda t}\right)$$

$$\leq \mathbb{E}\left[e^{\lambda(X - \mathbb{E}X)}\right] \cdot e^{-\lambda t}$$

**Assumption:** $\mathbb{E}e^{\lambda(X - \mathbb{E}X)} < \infty$, $\lambda \in (-\lambda_0, \lambda_0)$ for some $\lambda_0 > 0$.

Then, let $\quad \varphi_X(\lambda) = \log \mathbb{E}\left[e^{\lambda(X - \mathbb{E}X)}\right]$.

$\Rightarrow \quad P(X - \mathbb{E}X > t) \leq e^{-[\lambda t - \varphi_X(\lambda)]}$

Optimizing the r.h.s. of the above inequality,

$$\tilde{\varphi}_X^*(t) := \sup_{\lambda \in \mathbb{R}_+} \{\lambda t - \varphi_X(\lambda)\}$$

Thus,

$$P(X - \mathbb{E}X > t) \leq e^{-\tilde{\varphi}_X^*(t)} \qquad \text{for any } t \in \mathbb{R}.$$

$\tilde{\varphi}_X^*(t)$ resembles **Legendre transform** in convex analysis, but the optimization is performed on $\mathbb{R}_+$ rather than $\mathbb{R}$.

Let $\quad \varphi_X^*(t) = \sup_{\lambda \in \mathbb{R}} \{\lambda t - \varphi_X(\lambda)\}$. Then, $\varphi_X^*(t) \geq \tilde{\varphi}_X^*(t)$.

However, we have $\quad \varphi_X^*(t) = \tilde{\varphi}_X^*(t), \quad \underline{t \geq 0}$

**Pf** By Jensen's inequality,

$$\varphi_X(\lambda) = \log \mathbb{E}\left[e^{\lambda(X - \mu)}\right]$$

$$\geq \mathbb{E} \log e^{\lambda(X - \mu)} = 0$$

Hence, $\varphi_X(\lambda) \geq 0$. If $t \geq 0$,

$$\sup_{\lambda \leq 0} \{\lambda t - \varphi_X(\lambda)\} \leq \sup_{\lambda \leq 0} \{\lambda t\} \leq 0.$$

Since $\{\lambda t - \varphi_X(\lambda)\}_{\lambda = 0} = 0$, we have

$$\varphi_X^*(t) = \sup_{\lambda \in \mathbb{R}} \{\lambda t - \varphi_X(\lambda)\} = \sup_{\lambda \in \mathbb{R}_+} \{\lambda t - \varphi_X(\lambda)\} = \tilde{\varphi}_X^*(t).$$

Example (Gaussian)

$$X \sim N(\mu, \sigma^2) \implies \varphi_X(\lambda) = \frac{\sigma^2 \lambda^2}{2}, \forall \lambda \in \mathbb{R}.$$

Then,

$$\varphi_X^*(t) = \sup_{\lambda \in \mathbb{R}} \left\{ \lambda t - \frac{\sigma^2 \lambda^2}{2} \right\}$$

for optimizing, $t - \sigma^2 \lambda_{opt} = 0 \implies \lambda_{opt} = \frac{t}{\sigma^2}$

$$\implies \varphi_X^*(t) = \frac{t^2}{2\sigma^2}$$

$$\implies \boxed{\mathbb{P}\left( X - \mu \geq t \right) \leq e^{-\frac{t^2}{2\sigma^2}}}$$
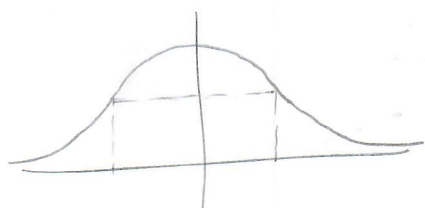
## Sub-Gaussian Random Variables and Hoeffding's Lemma

DEF| $X$ is sub-Gaussian with parameter $\sigma^2$ if

(i) $\mathbb{E}X = 0$,

(ii) $\varphi_X(\lambda) \leq \frac{\sigma^2 \lambda^2}{2}$ for all $\lambda \in \mathbb{R}$.

$X \sim N(0, \sigma^2)$ is sub-Gaussian with $\sigma^2$ obviously.



Another important class of sub-Gaussian r.v.'s is bounded r.v.'s.

PROP| If $X \in \mathbb{R}$ is sub-Gaussian with $\sigma^2$ then

$$\mathbb{P}\left( X > t \right) \leq e^{-\frac{t^2}{2\sigma^2}}, \forall t > 0.$$

LEMMA| $X \in [a, b]$ a.s. for some $a < b$, then

$X$ is sub-Gaussian with $\frac{(b-a)^2}{4}$.

Pf                                           $X$ has a density function $f(x)$.

Let $Z$ be a r.v. with density $\dfrac{f(x) e^{\lambda x}}{e^{\varphi_X(\lambda)}}$.     For sanity, $\dfrac{\int_a^b f(x) e^{\lambda x} dx}{e^{\varphi_X(\lambda)}} = 1$.

$$\mathbb{E}[Z] = \int x f(x) e^{\lambda x} dx \cdot e^{-\varphi_X(\lambda)} = \mathbb{E}[X e^{\lambda X}] e^{-\varphi_X(\lambda)}$$

$$\mathbb{E}[Z^2] = \mathbb{E}[X^2 e^{\lambda X}] e^{-\varphi_X(\lambda)}$$

$$\varphi_X(\lambda) = \varphi_X(0) + \lambda \cdot \varphi_X'(0) + \frac{\lambda^2}{2} \cdot \varphi_X''(\lambda_0)$$

$$\varphi_X(0) = 0, \qquad \varphi_X'(\lambda) = \frac{\mathbb{E}\left[(X-\mu)e^{\lambda(X-\mu)}\right]}{\mathbb{E}\left[e^{\lambda(X-\mu)}\right]} = 0 \quad \text{if} \quad \lambda = 0.$$

$$\varphi_X''(\lambda) = \frac{\mathbb{E}\left[(X-\mu)^2 e^{\lambda(X-\mu)}\right]e^{\varphi_X(\lambda)} - \mathbb{E}^2\left[(X-\mu)e^{\lambda(X-\mu)}\right]}{e^{2\varphi_X(\lambda)}}, \quad \forall \lambda \in \mathbb{R}$$

$$= \frac{\mathbb{E}\left[X^2 e^{\lambda X}\right]}{e^{\varphi_X(\lambda)}} - \left(\frac{\mathbb{E}\left[Xe^{\lambda X}\right]}{e^{\varphi_X(\lambda)}}\right)^2$$

$$= \mathbb{E}\left[Z_\lambda^2\right] - \mathbb{E}^2\left[Z_\lambda\right] = \text{Var}\left(Z_\lambda\right)$$

Hence,

$$\varphi_X(\lambda) = \frac{\lambda^2}{2} \cdot \text{Var}\left(Z_{\lambda_0}\right) \quad \text{for} \quad \text{some} \quad \lambda_0.$$

Note that :

$$\underset{t \in \mathbb{R}}{\text{argmin}} \, \mathbb{E}\left[|Z - t|^2\right] = \mathbb{E}Z$$

and

$$\left|Z - \frac{a+b}{2}\right| \leq \frac{b-a}{2} \quad \text{since} \quad Z \in [a, b].$$

Then,

$$\text{Var}\left(Z_{\lambda_0}\right) \leq \mathbb{E}\left[\left|Z_{\lambda_0} - \frac{a+b}{2}\right|^2\right] \leq \frac{(b-a)^2}{4}$$

$$\Rightarrow \quad \varphi_X(\lambda) \leq \frac{(b-a)^2}{4} \cdot \frac{\lambda^2}{2} \quad \Rightarrow \quad X \text{ is sub-Gaussian with } (b-a)^2/4$$

Hoeffding's Lemma at work:

$X \in [a,b]$, $\mathbb{E}X = 0$ $\Rightarrow$ $X$ is sub-Gaussian with $\dfrac{(b-a)^2}{4}$ :

$$\varphi_X(\lambda) \le \frac{\lambda^2}{8} \cdot (b-a)^2$$

$$\sup_{\lambda \in \mathbb{R}} \left\{ \lambda t - \varphi_X(\lambda) \right\} \le \sup_{\lambda \in \mathbb{R}} \left\{ \lambda t - \frac{\lambda^2 (b-a)^2}{8} \right\}$$

$$t - \frac{\lambda^* \cdot (b-a)^2}{4} = 0 \Rightarrow \lambda^* = \frac{4t}{(b-a)^2} \quad \text{and}$$

$$\sup_{\lambda} \left\{ \lambda t - \frac{\lambda^2 (b-a)^2}{8} \right\} = \frac{2t^2}{(b-a)^2}$$

Hence,

$$\mathbb{P}\left( X > t \right) \le \exp\left( - \frac{2t^2}{(b-a)^2} \right)$$

<u>LEMMA 1</u>  $X, Y$ sub-Gaussian with $\sigma_X^2$, $\sigma_Y^2$, they are independent

$\Rightarrow$  $X+Y$ is sub-Gaussian with $\sigma_X^2 + \sigma_Y^2$.

Also, for any $\alpha \in \mathbb{R}$, $\alpha X$ is sub-Gaussian with $\alpha^2 \sigma_X^2$.

<u>Pf</u> :  $\mathbb{E}\left[ e^{\lambda(X+Y)} \right] \underset{\underset{\text{independence}}{\uparrow}}{=} \mathbb{E}\left[ e^{\lambda X} \right] \cdot \mathbb{E}\left[ e^{\lambda Y} \right] \le e^{\lambda^2(\sigma_X^2 + \sigma_Y^2)/2}$

$\mathbb{E}\left[ e^{\lambda \alpha X} \right] \le e^{(\alpha\lambda)^2/2 \, \sigma_X^2} = e^{\alpha^2 \sigma_X^2 \cdot \frac{\lambda^2}{2}}, \ \forall \lambda, \alpha \in \mathbb{R}$

Then, consider $\overset{\text{independent}}{\nearrow}$ $X_i \in [a_i, b_i]$ with $\mathbb{E}X_i = 0$. Then,

$\dfrac{1}{n}\sum\limits_{i=1}^{n} X_i$ is sub-Gaussian with $\dfrac{1}{n^2} \sum\limits_{i=1}^{n} \dfrac{(b_i - a_i)^2}{4}$.

$$\Rightarrow \mathbb{P}\left( \frac{1}{n}\sum_{i=1}^{n} X_i > t \right) \le \exp\left( - \frac{2t^2 n^2}{\sum_{i=1}^{n}(b_i - a_i)^2} \right)$$

$$\mathbb{P}\left( \frac{1}{n}\sum_{i=1}^{n} X_i < -t \right) = \mathbb{P}\left( \frac{1}{n}\sum_{i}(-X_i) > t \right) \le \exp\left( - \frac{2n^2 t^2}{\sum_{i=1}^{n}(b_i - a_i)^2} \right)$$

Alternative form :  $\mathbb{P}\left( \left| \frac{1}{n}\sum\limits_{i=1}^{n} X_i \right| > \sqrt{\frac{1}{2n^2}\sum\limits_{i=1}^{n}(b_i - a_i)^2 \log(2/\delta)} \right) \le \delta$.