

Instructor : Avishek Ghosh

Course website : figure it out (content, assignments @ website, moodle)

Tue / Fri — 5.30 - 7 pm

Extra classes — Saturday 3.30 - 5 pm

<u>Grading</u> :	Homeworks	20%	(2-3)
	Midsem	30%	
	Endsem	30%	
	Scribes	15%	
	Class Participation	5%	

What is this course about ?

1. Analyzing ML algorithms from a statistical point of view
2. Involves a large of statistical tools / techniques that can be used independently

NOT about

1. Theory of Deep Learning
2. Not particularly algorithmic

Reference :

1. High dimensional Statistics — Martin Wainwright
2. Asymptotic Stats — A. W. Vanderbaart

Classification (Binary)

Given data points $(X_1, Y_1) (X_2, Y_2) \dots (X_n, Y_n)$ where $X_i \in \mathcal{X} \subseteq \mathbb{R}^d$ ^(features)
 $Y_i \in \{-1, +1\}$

Goal : Find a classifier $g : \mathcal{X} \rightarrow \{-1, +1\}$

How to obtain this ? — Notion of loss function

Binary loss : $\mathbb{1}_{[g(x) \neq Y]} = \begin{cases} 1, & \text{if } g(x) \neq Y \\ 0, & \text{if } g(x) = Y \end{cases}$
 (error)

Construct "Empirical loss" $L_n(g) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}\{g(x_i) \neq y_i\}$

We select classifier for which $L_n(g)$ is minimized

$$\hat{g}_n = \arg \min_{g \in \mathcal{C}} L_n(g), \text{ where } \mathcal{C} \text{ denotes family of classifier}$$

Problem : 1. Performance on "unseen" data is not considered
2. \mathcal{C} can be complicated, n can be much smaller

Upto this point, purely empirical (no statistics)

Statistical Model

We assume that $(X_1, Y_1), \dots, (X_n, Y_n)$ are i.i.d samples from a joint distribution \mathcal{D} having same distribution as (X, Y)

Then, for a classifier $g: \mathcal{X} \rightarrow \{\pm 1\}$, we can write,

$$\underbrace{L(g)}_{\text{Avg loss / expected loss}} = \mathbb{E}_{(X, Y) \sim \mathcal{D}} [\mathbb{1}\{g(X) \neq Y\}] = \mathbb{P}(g(X) \neq Y)$$

It is a good idea to study $L(\hat{g}_n)$

2 Questions :

1. Is $L(\hat{g}_n)$ comparable to $\underbrace{\inf_{g \in \mathcal{C}} L(g)}_{\text{Naive Bayes classifier}}$?
whether \hat{g}_n is comparable with the best classifier in \mathcal{C}
2. Is $L(\hat{g}_n)$ comparable to $L_n(\hat{g}_n)$?

Comparison between "in-sample" error and average error for \hat{g}_n

Assume $g^* = \arg \min_{g \in \mathcal{C}} L(g)$ {Naive Bayes}

$$\begin{aligned} \text{We write } L(\hat{g}_n) &= L(g^*) + L(\hat{g}_n) - L_n(\hat{g}_n) + \underbrace{L_n(\hat{g}_n)}_{\text{lesser loss}} - L(g^*) \\ &\leq L(g^*) + \underbrace{L(\hat{g}_n) - L_n(\hat{g}_n)}_{\Delta} + \underbrace{L_n(g^*) - L(g^*)}_{\Delta} \\ \Rightarrow L(\hat{g}_n) - L(g^*) &\leq \sup_{g \in \mathcal{C}} |L(g) - L_n(g)| + \sup_{g \in \mathcal{C}} |L_n(g) - L(g)| \end{aligned}$$

$$\Rightarrow L(\hat{g}_n) - L(g^*) \leq 2 \sup_{g \in \mathcal{C}} |L_n(g) - L(g)| \quad - (*)$$

Remark: ① is controlled by (*)

$$② L(\hat{g}_n) - L_n(\hat{g}_n) \leq \sup_{g \in \mathcal{C}} |L_n(g) - L(g)|$$

Remark: Performance of \hat{g}_n is governed by $\sup_{g \in \mathcal{C}} |L_n(g) - L(g)|$.

We use uniform law of large numbers to handle this

Empirical Process Theory

1. Uniform law of large numbers
2. Uniform central limit theorem

Uniform Law of large NOS.

Suppose X_1, X_2, \dots, X_n i.i.d random objects taking value in \mathcal{X} . Let \mathcal{F} be class of real-valued function on \mathcal{X} , what can we say about

$$\sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n f(X_i) - \mathbb{E} f(X) \right| = Z$$

In particular,

- ① whether $Z \rightarrow 0$ when n is large?
- ② Can we obtain non-asymptotic guarantees? i.e. guarantees for every n
- ③ Can we provide conditions on \mathcal{F} s.t. Z converges to 0?

Connection to ML and statistics

① Binary Classification —

$$X_i \mapsto (X_i, Y_i)$$

$$\mathcal{F} \mapsto \{ \mathbb{1}\{g(x) \neq y\} : g \in \mathcal{C} \}$$

② M-estimation

$$\hat{\theta}_n = \arg \max_{\theta \in \Theta} \frac{1}{n} \sum_{i=1}^n m_{\theta}(x_i) \quad \leftarrow \text{think as -ve of loss fn.}$$

where x_1, \dots, x_n are i.i.d observations, Θ is parameter space
 m_{θ} are real valued fn parametrized by θ .

- Examples, 1. $m_{\theta}(x) = \log p_{\theta}(x) \leftarrow$ Maximum likelihood estimator (MLE)
2. $m_{\theta}(x) = -(x - \theta)^2 \leftarrow$ sample mean (Mean estimator)
3. $m_{\theta}(x) = -|x - \theta| \leftarrow$ Median estimator

In mean estimation, target quantity for $\hat{\theta}_n$ is

$$\theta^* = \arg \max_{\theta \in \Theta} \mathbb{E} m_{\theta}(x)$$

Similar to binary classification example we want $\underbrace{d(\hat{\theta}_n, \theta^*)}_{\text{distance b/w } \hat{\theta}_n \text{ and } \theta^*}$ to be small

It turns out $d(\hat{\theta}_n, \theta^*)$ is governed by $2 \sup_{\theta \in \Theta} \left| \frac{1}{n} \sum_{i=1}^n m_{\theta}(x_i) - \mathbb{E} m_{\theta}(x) \right|$

which is an instance of uniform law of large numbers.

Strategy to Control $\mathbb{E} Z$

① Key Observation: Z concentrates around $\mathbb{E} Z$ i.e. $\mathbb{E} \sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n f(x_i) - \mathbb{E} f(x) \right|$
(concentration of measure)

② We control $\mathbb{E} Z$ through techniques like symmetrization (Rademacher Complexity) or chaining (vc dimension)

Remark: (Asymptotic result) \tilde{f} is called "Glivenko - Cantelli" if $Z \rightarrow 0$ almost-surely as $n \rightarrow \infty$

Assumption: $\sup_x |f(x)| \leq B \quad \forall f \in \mathcal{F}$

McDiarmid's inequality

Suppose x_1, x_2, \dots, x_n and $g: \mathcal{X}_1 \times \dots \times \mathcal{X}_n \rightarrow \mathbb{R}$ satisfies "bounded difference".

$$|g(x_1, \dots, x_n) - g(x_1, \dots, x_{i-1}, x_i', x_{i+1}, \dots, x_n)| \leq C_i$$

$$\forall x_1, \dots, x_n \quad \forall i \in [n]$$

Then we have

$$\begin{aligned} \mathbb{P}(g(x_1, \dots, x_n) - \mathbb{E} g(x_1, \dots, x_n) \geq t) &\leq \exp\left(-\frac{2t^2}{\sum_{i=1}^n C_i^2}\right) \\ &\leq -t \leq \dots \end{aligned}$$

Rmk : The bounded difference says that a f^n that is not too sensitive on any of its argument concentrates.

Apply McDiarmid's to Z

$$Z = \sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n f(x_i) - \mathbb{E} f(x) \right|$$

we construct

$$g(x_1, \dots, x_n) = \sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n f(x_i) - \mathbb{E} f(x) \right|$$

$$\begin{aligned} g(x_1, \dots, x_i', \dots, x_n) &= \sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{j \neq i} f(x_j) + \frac{f(x_i')}{n} - \mathbb{E} f(x) \right| \\ &= \sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{j=1}^n f(x_j) - \mathbb{E} f(x) + \frac{f(x_i')}{n} - \frac{f(x_i)}{n} \right| \end{aligned}$$

$$\leq g(x_1, \dots, x_n) + \sup_{f \in \mathcal{F}} \left| \frac{f(x_i')}{n} \right| + \sup_{f \in \mathcal{F}} \left| \frac{f(x_i)}{n} \right|$$

$$\therefore \left| g(x_1, \dots, x_n) - g(x_1, \dots, x_i', \dots, x_n) \right| \leq \boxed{\frac{2B}{n} = C_i}$$

by switching x_i, x_i' to get both side

Hence, g satisfies bounded difference \Rightarrow apply Hoeffding's inequality.

$$\mathbb{P}(Z - \mathbb{E}Z \geq t) \leq \exp\left(\frac{-2t^2}{\sum_i \frac{4B^2}{n^2}}\right) = \exp\left(-\frac{nt^2}{2B^2}\right)$$

$$\text{||} \mathbb{P}(Z - \mathbb{E}Z \leq -t) \leq \exp\left(\underbrace{-\frac{nt^2}{2B^2}}_{=8}\right)$$

Then we say w.p. $\geq 1-8$,

$$Z \leq \mathbb{E}Z + \underbrace{B \sqrt{\frac{2}{n} \log \frac{1}{8}}}_{t \text{ (very small bcz } \propto \frac{1}{\sqrt{n}})}$$

Remark: We need to control $\mathbb{E} \sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n f(X_i) - \mathbb{E}f(X) \right|$

To control Z ① $Z \rightarrow \mathbb{E}Z$ ✓
 ② control $\mathbb{E}Z$!

Concentration Inequality (Hoeffding)

Suppose $\overset{\text{indep.}}{X_1, \dots, X_n}$ r.v. such that $a_i \leq X_i \leq b_i$ $\overset{\text{w.p. } 1}{\text{almost surely}}$ where $a_1, \dots, a_n, b_1, \dots, b_n$ are real numbers. Then for any $t \geq 0$

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

and

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

Proof: Let $S = \sum_{i=1}^n (X_i - \mathbb{E}X_i)$. fix $\lambda \geq 0$. we have

$$\mathbb{P}(S \geq t) = \mathbb{P}(e^{\lambda S} \geq e^{\lambda t}) \leq \frac{\mathbb{E}(e^{\lambda S})}{e^{\lambda t}} \stackrel{\text{Markov}}{=} e^{-\lambda t} \mathbb{E}e^{\lambda t} = \exp(-\lambda t + \psi_S(\lambda))$$

How: find Ω for which $e^{\lambda S} \geq e^{\lambda t}$ is true but $S \geq t$ is not.

$$\psi_S(t) = \log \mathbb{E} e^{\lambda S} \quad (\log\text{-mgf / cumulant fn})$$

$$\begin{aligned} \psi_S(\lambda) &= \log \mathbb{E} e^{\lambda \sum_i (x_i - \mathbb{E} x_i)} = \log \prod_{i=1}^n \mathbb{E} e^{\lambda (x_i - \mathbb{E} x_i)} \\ &\quad \xrightarrow{\text{independence}} \\ &= \sum_{i=1}^n \log \mathbb{E} e^{\lambda (x_i - \mathbb{E} x_i)} \\ &= \sum_{i=1}^n \psi_{x_i - \mathbb{E} x_i}(\lambda) \quad \text{--- (1)} \end{aligned}$$

Fix i , we analyze

$$\psi_{\underbrace{x_i - \mathbb{E} x_i}_U}(\lambda) \quad \text{Need to bound } \psi_U$$

$$b_i - \mathbb{E} x_i \leq U \leq a_i - \mathbb{E} x_i \quad \text{almost surely.}$$

$$\text{Taylor series, } \psi_U(\lambda) = \cancel{\psi_U(0)} + \cancel{\psi_U'(0)} \lambda + \overbrace{\psi_U''(c) \left(\frac{\lambda^2}{2}\right)}^{\text{key term}},$$

$$* \psi_U(0) = \log \mathbb{E} e^{0(x_i - \mathbb{E} x_i)} = 0 \quad 0 < c < \lambda$$

$$* \psi_U'(\lambda) = \frac{d}{d\lambda} \log \mathbb{E} e^{\lambda U} = \frac{1}{\mathbb{E} e^{\lambda U}} \mathbb{E}(e^{\lambda U} \cdot U)$$

$$\text{Put } \lambda = 0 \quad = \frac{\mathbb{E}[U]}{1} = 0 \dots$$

$$\psi_U(\lambda) = \psi_U''(c) \frac{\lambda^2}{2}$$

$$\psi_U''(\lambda) = \frac{d}{d\lambda} \frac{1}{\mathbb{E} e^{\lambda U}} \mathbb{E}(U e^{\lambda U}) = \frac{\mathbb{E}(e^{\lambda U}) \mathbb{E}(U^2 e^{\lambda U}) - \mathbb{E}(U e^{\lambda U})^2}{\mathbb{E}(e^{\lambda U})^2}$$

We will show: $\psi_U''(\lambda) \geq 0$ for any λ

Consider a r.v. V whose density w.r.t. density of U is $e^{\lambda U} / \mathbb{E} e^{\lambda U}$

$$\underbrace{p_V}_{\text{density fn for random var. } V} = p_U \frac{e^{\lambda U}}{\mathbb{E} e^{\lambda U}} \quad (u=0 \Rightarrow p_u=0 \Rightarrow p_v=0)$$

$$\psi''(\lambda) = \mathbb{E}(V^2) - (\mathbb{E} V)^2 \leftarrow \text{variance of } V$$

$$\geq 0.$$

Rule: Support of U is $a_i - \mathbb{E} X_i, b_i - \mathbb{E} X_i$

support of V is " ————— " ($p_U = 0 \Rightarrow p_V = 0$)

Exercise: For any random variable V bound,

$$\text{Var}(V) \leq \left(\frac{b_i - a_i}{4} \right)^2 \leftarrow \frac{1}{2} \text{ pr } a_i, \frac{1}{2} b_i$$

$$\Rightarrow \psi_U(\lambda) = \text{Var}(V) \leq \frac{(b_i - a_i)^2}{4}$$

we have,

$$\boxed{\psi_U(\alpha) = \frac{\lambda^2}{2} \psi_U''(x) \leq \frac{\lambda^2}{8} (b_i - a_i)^2}$$

Now, substituting ①

$$\psi_S(\lambda) = \sum_{i=1}^n \psi_{X_i - \mathbb{E} X_i}(\lambda) \leq \sum_{i=1}^n \frac{\lambda^2}{8} (b_i - a_i)^2$$

Substituting this in ①,

$$\Pr(S \geq t) \leq \exp(-\lambda t + \sum_{i=1}^n \frac{\lambda^2}{8} (b_i - a_i)^2)$$

optimizing over λ , put $\lambda^* = \frac{4t}{\sum (b_i - a_i)^2}$

Putting $\lambda = \lambda^*$,

$$\Pr(S \geq t) \leq \exp\left(\frac{-2t^2}{\sum_{i=1}^n (b_i - a_i)^2}\right)$$

To get other side of concentration put $\gamma_i = -X_i$

Reading assignment:

1. Go through proof of lec 2
2. Revise martingales

Central Limit Theorem

14/1

X_1, X_2, \dots, X_n independent with mean μ , variance σ^2

$$\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$$

$$\text{CLT: } \sqrt{n} \left(\frac{\bar{X}_n - \mu}{\sigma} \right) \xrightarrow[n \rightarrow \infty]{\text{distribution}} N(0, 1)$$

(Reading: Convergence of r.v.
Convergence in distribution)

Suppose n is very large (for $t \geq 0$)

$$\mathbb{P}(\sqrt{n} \left(\frac{\bar{X}_n - \mu}{\sigma} \right) \geq t) \approx \mathbb{P}(N(0, 1) \geq t)$$

Scaling r.v. by σ , variance scales by σ^2 .

$$\mathbb{P}(\sqrt{n} (\bar{X}_n - \mu) \geq t) \approx \mathbb{P}(N(0, \sigma^2) \geq t)$$

We use MGF based method to upperbound

$$\mathbb{P}(N(0, \sigma^2) \geq t) \leq \mathbb{P}(e^{\lambda N(0, \sigma^2)} \geq e^{\lambda t})$$

$$\leq e^{-\lambda t} \mathbb{E} e^{\lambda N(0, \sigma^2)}$$

$$= \exp(-\lambda t + \underbrace{\psi_{N(0, \sigma^2)}(\lambda)}_{\text{log-mgf/cumulant}})$$

$$\mathbb{E} e^{\lambda N(0, \sigma^2)} = e^{\frac{\sigma^2 \lambda^2}{2}} \quad (\text{Exercise})$$

Using this,

$$\mathbb{P}(N(0, \sigma^2) \geq t) \leq \exp\left(-\lambda t + \frac{\sigma^2 \lambda^2}{2}\right)$$

$$\boxed{\mathbb{P}(N(0, \sigma^2) \geq t) \leq \exp\left(-\frac{t^2}{2\sigma^2}\right)}$$

optimal λ

$$* \text{ CLT: } \boxed{\mathbb{P}(\sqrt{n} (\bar{X}_n - \mu) \geq t) \leq \exp\left(-\frac{t^2}{2\sigma^2}\right)} \quad \text{--- (I) (CLT)}$$

Using Hoeffding's we get,

$$X_1, \dots, X_n, \mathbb{E}[X] = \mu, \text{Var}(X) = \sigma^2, a \leq X \leq b \text{ a.s.}$$

$$\mathbb{P}\left(\sum_{i=1}^n X_i - \mathbb{E}[X] \geq t'\right) \leq \exp\left(\frac{-2t'^2}{n(b-a)^2}\right)$$

let $t' = \sqrt{n} \cdot t$

$$\Rightarrow \mathbb{P}\left(\frac{1}{n} \sum_{i=1}^n X_i - \mathbb{E}[X] \geq \frac{\sqrt{n} \cdot t}{n}\right) \leq \exp\left(\frac{-2t^2}{(b-a)^2}\right)$$

$$\Rightarrow \boxed{\mathbb{P}\left(\sqrt{n}(\bar{X}_n - \mu) \geq t\right) \leq \exp\left(\frac{-2t^2}{(b-a)^2}\right)} \quad \text{--- (II) (Hoeffding's)}$$

Comparison :

If $a \leq X_i \leq b$ a.s., then

CLT : $\exp\left(\frac{-t^2}{2\sigma^2}\right)$

$$\sigma^2 \leq \frac{(b-a)^2}{4}$$

Hoeffding's : $\exp\left(\frac{-2t^2}{(b-a)^2}\right)$

With this, CLT bound is $\exp\left(\frac{-t^2}{2\sigma^2}\right) \leq \exp\left(\frac{-2t^2}{(b-a)^2}\right)$
(CLT) (Hoeffding).

Remarks : * Error upper bound matches b/w CLT, Hoeffding

* σ^2 can be much smaller than range, makes CLT bound sharper/stronger

* CLT holds asymptotically and is an approximation.

On the other hand, Hoeffding's is exact, non-asymptotic (holds for any n)

Berry Esseen Thm : How close app w.r.t CLT (approximation scales)
(v. deep thm — not in course) $\sim 1/\sqrt{n}$

Note : MGF based techniques we use — (Cramer — Chernoff Method.)

Confidence Interval

Suppose $X_1 \dots X_n$ i.i.d random variables with $\mathbb{E}X = \mu$, $\text{Var } X = \sigma^2$
 $a \leq X_i \leq b$ almost surely.

Problem: We know (a, b, σ^2) and want to estimate μ .

We obtain a set (interval), also known as confidence interval (CI), where μ lies w.h.p.

CLT based bound (classical) (for $t \geq 0$)

$$\text{By CLT, } \mathbb{P}\left(\left|\sqrt{n}\left(\frac{\bar{X}_n - \mu}{\sigma}\right)\right| \leq t\right) \xrightarrow{n \rightarrow \infty} \mathbb{P}(|N(0,1)| \leq t)$$

(union of two events)

Put $t = z_{\alpha/2}$ (α -quantile) s.t. $\mathbb{P}(|N(0,1)| \leq z_{\alpha/2}) = 1 - \alpha$
w.p. $1 - \alpha$,

$$-z_{\alpha/2} \leq \sqrt{n} \left(\frac{\bar{X}_n - \mu}{\sigma} \right) \leq z_{\alpha/2}$$

$$\Rightarrow \boxed{\bar{X}_n - \frac{\sigma}{\sqrt{n}} z_{\alpha/2} \leq \mu \leq \bar{X}_n + \frac{\sigma}{\sqrt{n}} z_{\alpha/2}} \quad \text{C.I.}$$

length of interval $\frac{2\sigma}{\sqrt{n}} z_{\alpha/2}$ shrinks with n ($\rightarrow 0$ as $n \rightarrow \infty$)

Hoeffding's for C.I.

$$\mathbb{P}\left(\left|\sqrt{n}(\bar{X}_n - \mu)\right| \geq t\right) \leq \underbrace{2 \exp\left(\frac{-2t^2}{(b-a)^2}\right)}_{\alpha}$$

$$\text{i.e. } t = (b-a) \sqrt{\frac{1}{2} \log \frac{2}{\alpha}}$$

so, w.p. $1 - \alpha$,

$$\boxed{|\bar{X}_n - \mu| \leq \frac{1}{\sqrt{n}} (b-a) \sqrt{\frac{1}{2} \log \frac{2}{\alpha}}}$$

μ lies in $\left[\bar{X}_n - \frac{1}{\sqrt{n}} (b-a) \sqrt{\frac{1}{2} \log \frac{2}{\alpha}}, \bar{X}_n + \frac{1}{\sqrt{n}} (b-a) \sqrt{\frac{1}{2} \log \frac{2}{\alpha}}\right]$ C.I.

length of interval,

$$\sigma \leq \frac{b-a}{2}$$

so, CLT is better
(asymptotic)

Sub-Gaussian

Assume $X \sim \mathcal{N}(\mu, \sigma^2)$. We know that

$$\mathbb{E} e^{\lambda(X-\mu)} = e^{\lambda^2 \sigma^2 / 2} \text{ for any } \lambda.$$

This motivates us to define class of r.v. exhibiting similar properties.

Def: A r.v. X w/ mean μ is called sub-gaussian if there exists a positive number σ such that

$$\mathbb{E} e^{\lambda(X-\mu)} \leq e^{\lambda^2 \sigma^2 / 2} \text{ for all } \lambda.$$

It is denoted as $X \sim \text{subG}(\sigma)$. σ is called the parameter of sub-Gaussian r.v.; σ^2 is a proxy for variance.

Examples ① Gaussian

② Rademacher r.v.:

$\epsilon \in \{-1, +1\}$ with equal probability

We want to show that ϵ is 1 sub gaussian

$$\begin{aligned} \mu = 0, \quad \mathbb{E} e^{\lambda \epsilon} &= \frac{e^{\lambda} + e^{-\lambda}}{2} = \sum_{k=0}^{\infty} \frac{\lambda^{2k}}{(2k)!} \leq \sum_{k=0}^{\infty} \frac{(\lambda^2)^k}{k! 2^k} \\ \sigma = 1 \end{aligned}$$

$k! \sim k^k$
 $(2k)! \text{ much bigger than } 2^k k!$
 \downarrow

$$\leq e^{\lambda^2 / 2}$$

③ Bounded r.v.

Suppose $a \leq X \leq b$ a.s.

show that $X \sim \text{SubG}(\frac{b-a}{2})$

Last class: $\log \mathbb{E} e^{\lambda X} \leq \frac{(b-a)^2}{8} \lambda^2$

$$\Rightarrow \mathbb{E} e^{\lambda X} \leq \exp\left(\lambda^2 \frac{(b-a)^2}{8}\right) \leftarrow \text{exactly what we want.}$$

Hoeffding's inequality for sub Gaussian [not just for bounded r.v.]

Suppose X_1, \dots, X_n are sub G (σ_i) and $\mathbb{E} X_i = \mu_i$; then,
for all $t \geq 0$

$$\mathbb{P} \left(\sum_{i=1}^n (X_i - \mathbb{E} X_i) \geq t \right) \leq \exp \left(- \frac{t^2}{2 \sum_{i=1}^n \sigma_i^2} \right)$$

$$\mathbb{P} \left(\sum_{i=1}^n (X_i - \mathbb{E} X_i) \leq -t \right) \leq \exp \left(- \frac{t^2}{2 \sum_{i=1}^n \sigma_i^2} \right)$$

Remark : $\sigma_i = \frac{b_i - a_i}{2}$, we get back hoeffding's for bounded r.v. (proved in lec 2)

Proof : Use Cramer- Chernoff technique, bound on $\mathbb{E} e^{\lambda X}$ comes from defn of sub-gaussian directly!

Reading exercise :

Chap 2 from HDS (Martin) — Sub-gaussian (bigger class: sub-exponential)
— Martingales.