

Instructor : Arishek 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 $x_i \in \mathcal{X} \subseteq \mathbb{R}^d$
 $y_i \in \{-1, +1\}$

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

How to obtain this? — Notion of loss function

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

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

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 $g: \mathcal{X} \rightarrow \{-1, 1\}$, we can write,

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

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

2 Questions :

$$1. \text{ Is } L(\hat{g}_n) \text{ comparable to } \inf_{g \in \mathcal{C}} L(g) ?$$

whether \hat{g}_n is comparable with the best classifier in \mathcal{C}

$$2. \text{ Is } L(\hat{g}_n) \text{ comparable to } L_n(\hat{g}_n) ?$$

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

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

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

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

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 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

$$\hat{\theta}_n = \arg \max_{\theta \in \Theta} \frac{1}{n} \sum_{i=1}^n m_\theta(x_i)$$

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. $m_\theta(x) = \log p_\theta(x)$ ← Maximum likelihood estimator (MLE)
 2. $m_\theta(x) = -(x - \theta)^2$ ← sample mean (Mean estimator)
 3. $m_\theta(x) = -|x - \theta|$ ← Median estimator

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

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

distance b/w $\hat{\theta}_n$ and θ^*

Similar to binary classification example we want $d(\hat{\theta}_n, \theta^*)$ 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 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) 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 $g: \mathcal{X} \times \dots \times \mathcal{X}_n \rightarrow \mathbb{R}$ 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

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

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

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

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

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

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

by switching x_i, x_i' to get both side

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

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

$$\text{Hence } P(Z - \mathbb{E}Z \leq -t) \leq \exp\left(-\frac{nt^2}{2B^2}\right) = \delta$$

Then we say $-f^- \geq 1-\delta$,

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

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 $Z \rightarrow \mathbb{E}Z \checkmark$
 ② Control $\mathbb{E}Z$?

Concentration Inequality (Hoeffding)

Suppose x_1, \dots, x_n r.v. such that $a_i \leq x_i \leq b_i$, almost surely
 where $a_1, \dots, a_n, b_1, \dots, b_n$ are real numbers. Then for any $t \geq 0$

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

and

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

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

$$P(S \geq t) = P(e^{\lambda S} \geq e^{\lambda t}) \leq \frac{\mathbb{E}(e^{\lambda S})}{e^{\lambda t}} = e^{-\lambda t} \mathbb{E}e^{\lambda S}$$

Now: find λ for which $e^{\lambda S} \geq e^{\lambda t}$ is true but $S \geq t$

$e^{\lambda S} \geq e^{\lambda t}$ is true but $S \geq t$

- ①

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

$$\begin{aligned}\psi_s(\lambda) &= \log \mathbb{E} e^{\lambda \sum_i (x_i - \mathbb{E} x_i)} = \log \prod_{i=1}^n \mathbb{E} e^{\lambda(x_i - \mathbb{E} x_i)} \\ &\quad \text{independence} \\ &= \sum_{i=1}^n \log \mathbb{E} e^{\lambda(x_i - \mathbb{E} x_i)} \\ &= \sum_{i=1}^n \psi_{x_i - \mathbb{E} x_i}(\lambda) \quad -\text{(11)}\end{aligned}$$

Fix i , we analyze

$$\psi_{x_i - \mathbb{E} x_i}(\lambda). \text{ Need to bound } \psi_0$$

$b_i - \mathbb{E} x_i \leq v \leq a_i - \mathbb{E} x_i$ almost surely.

Taylor series, $\psi_0(\lambda) = \psi_0(0) + \psi_0'(0)(\lambda) + \underbrace{\psi_0''(c)(\frac{\lambda^2}{2})}_{\text{key term.}}$

* $\psi_0(0) = \log \mathbb{E} e^{0(x_i - \mathbb{E} x_i)} = 0 \quad 0 < c < \lambda$

* $\psi_0'(\lambda) = \frac{d}{d\lambda} \log \mathbb{E} e^{\lambda v} = \frac{1}{\mathbb{E} e^{\lambda v}} \mathbb{E}(e^{\lambda v} \cdot v)$

Put $\lambda = 0 = \frac{\mathbb{E}[v]}{c} = 0 \dots$

$$\psi_0(\lambda) = \psi_0''(c) \frac{\lambda^2}{2}$$

$$\psi_0''(\lambda) = \frac{d}{d\lambda} \frac{1}{\mathbb{E} e^{\lambda v}} \mathbb{E}(v e^{\lambda v}) = \frac{\mathbb{E}(e^{\lambda v}) \mathbb{E}(v^2 e^{\lambda v}) - \mathbb{E}(v e^{\lambda v})^2}{\mathbb{E}(e^{\lambda v})^2}$$

We will show: $\psi_0''(\lambda) \geq 0$ for any λ

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

$$\psi_v = \psi_u \frac{e^{\lambda u}}{\mathbb{E} e^{\lambda u}} \quad (u=0 \Rightarrow \rho_u = 0 \Rightarrow \rho_v = 0)$$

density fn for random var. v

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

$$\geq 0.$$

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

support of V is $\{ \dots \}$ ($P_u=0 \Rightarrow P_v=0$)

Exercise: For any random variable V bound,

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

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

we have,

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

Now, substituting ⑩

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

Substituting this in ①,

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

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

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 $y_i = -x_i$

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

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

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

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

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

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

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

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

We use MGF based method to upperbound

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

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

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

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

Using this,

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

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

* CLT: $P\left(\sqrt{n}(\bar{X}_n - \mu) \geq t\right) \approx \exp\left(-\frac{t^2}{2\sigma^2}\right)$ — (I) (CLT)

Using Hoeffding's we get,

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

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

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

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

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

Comparison :

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

$$\text{CLT : } \exp\left(-\frac{t^2}{2\sigma^2}\right) \quad \sigma^2 \leq \frac{(b-a)^2}{4}$$

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

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

Remarks : * Error upper bound matches w/o 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 approx 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 $t \geq 0$)

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

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

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

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

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.

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

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

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

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

μ lies in $\left[\bar{X}_n - " , \bar{X}_n + "\right]$

Length of interval,
 $\sigma \leq \frac{b-a}{2}$
so, CLT is better
(asymptotic)

Sub-Gaussian

Assume $X \sim N(\mu, \sigma^2)$, 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 $X \sim \text{sub } G(\sigma)$. 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

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

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

③ Bounded r.v.

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

Show that $X \sim \text{Sub } G\left(\frac{b-a}{2}\right)$

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

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

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 $t \geq 0$

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

$$\mathbb{P}\left(\sum_{i=1}^n (x_i - \mathbb{E} x_i) \leq -t\right) \leq \dots$$

Remark : $\sigma_i = \frac{b_i - \alpha_i}{2}$, 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.

(17/1)

Last time : Concentration of measure

- Hoeffding's ineq for bounded r.v.
CLT
- Confidence Interval (C.I.) construction
- sub-Gaussian X is sub-G(σ)

$$\text{if } \mathbb{E} e^{\lambda X} \leq e^{\lambda^2 \sigma^2 / 2 + \lambda}.$$

— Hoeffding's ineq for sub-gaussian

Properties of Sub-Gaussian r.v.

- ① $x_1, x_2 \sim \text{sub-G}(\sigma_1), \text{sub-G}(\sigma_2)$ and they are independent w/ mean μ_1, μ_2
- $$x_1 + x_2 \sim \mathbb{E} e^{\lambda(x_1 + x_2 - (\mu_1 + \mu_2))} = \mathbb{E} e^{\lambda x_1 - \mathbb{E} x_1} \mathbb{E} e^{\lambda x_2 - \mathbb{E} x_2} \leq e^{\frac{\lambda^2}{2}(\sigma_1^2 + \sigma_2^2)}$$
- $$\sim \text{sub-G}(\sqrt{\sigma_1^2 + \sigma_2^2})$$

Reading exercise : Sub-exponential r.v. (chapter 2 from HDS)

* Property : If $X \sim \text{Sub-G}$ then $X^2 \sim \text{sub exponential}$

Norm : $\gamma = \begin{bmatrix} x_1 \\ \vdots \\ x_d \end{bmatrix}$, where $x_1 \dots x_d \sim \text{sub-G}$
 $\|\gamma\|_2^2 = x_1^2 + \dots + x_d^2 \sim \text{Subexp.}$

Martingale Based Concentration \leftarrow Doob 1956 (rejected Shannon's paper :p)

- * Bounded, Sub-G, Sub-exp. We require x_1, \dots, x_n to be independent
- * Independence doesn't hold in many settings - online Learning / Time-Series
- * Need to deal w/ dependent r.v.'s

Martingale

sets

$\gamma_1 \subseteq \gamma_2 \subseteq \dots$ random variables
Pair $(Y_k, F_k)_{k=1}^{\infty}$ is called a martingale if .

1. $\mathbb{E}[|Y_k|] < \infty$ for $k \geq 1$
2. $\mathbb{E}[Y_{k+1} | F_k] = Y_k$ a.s.

Ex : ① (Partial sum)

Let x_1, x_2, x_3, \dots be i.i.d random variables w/ mean μ

Define, $y_k = \sum_{i=1}^k x_i - k\mu$ w/ finite variance

This is a martingale because

$$\textcircled{1} \quad \mathbb{E}[|y_{k+1}|] < \infty$$

$$\textcircled{2} \quad \mathbb{E}[y_{k+1} | x_1 \dots x_k] = \mathbb{E}\left[\underbrace{\sum_{i=1}^k x_i}_{\text{contains all information abt } y_k} - k\mu + x_{k+1} - \mu | x_1 \dots x_k\right]$$
$$= y_k$$

Example : (Doob Martingale)

Given a sequence of independent random variables $\{X_k\}_{k=1}^n$. We define

$$Y_k = \mathbb{E}[f(X) | X_1, \dots, X_k] \text{ for } k=1, \dots, n \text{ and}$$

$$Y_0 = \mathbb{E}[f(X)] ; \text{ where } X = (X_1, \dots, X_n) \text{ and } f: \mathbb{R}^n \rightarrow \mathbb{R}$$

with $\mathbb{E}|f(X)| < \infty$

We now show this is a martingale

$$\begin{aligned} Y_n &= f(X) . \text{ So, } f(X) - \mathbb{E}[f(X)] = Y_n - Y_0 \\ &= \sum_{k=1}^n (\underbrace{Y_k - Y_{k-1}}_{\text{difference}}) \end{aligned}$$

Claim : $\{Y_k\}_{k=1}^n$ is a martingale w.r.t. $\{X_k\}_{k=1}^n$

$$\tilde{\mathcal{F}}_k = \sigma(X_1, \dots, X_k) \leftarrow \text{set w/ all possible functions of } X_1, \dots, X_k \text{ (}\sigma\text{-field)}$$

$$\begin{aligned} \textcircled{1} \quad \mathbb{E}[|Y_k|] &= \mathbb{E}[\mathbb{E}[f(X) | X_1, \dots, Y_k]] \quad \text{2 triangle inequality / jensen} \\ &\leq \mathbb{E}|f(X)| < \infty \quad \text{Prop: } \mathbb{E}\mathbb{E}(Y | X) = \mathbb{E}(Y) \end{aligned}$$

$$\begin{aligned} \textcircled{2} \quad \mathbb{E}[Y_{k+1} | \tilde{\mathcal{F}}_k] &= \mathbb{E}[\mathbb{E}[f(X) | X_1, \dots, X_{k+1}]] \quad \text{law of iterated expectation} \\ &= \mathbb{E}[f(X) | X_1, \dots, X_k] \quad \text{Tower rule: smaller set wins} \\ &= Y_k \end{aligned}$$

$$\text{Tower rule: } \mathbb{E}[\mathbb{E}[z | S_1] | S_2] = \mathbb{E}[z | S_1] \text{ if } S_1 \subseteq S_2$$

Martingale difference

A sequence $(D_k, \tilde{\mathcal{F}}_k)_{k=1}^\infty$ is called martingale difference if D_k is adapted to $\tilde{\mathcal{F}}_k$ and

$$\textcircled{1} \quad \mathbb{E}|D_k| < \infty \quad \textcircled{11} \quad \mathbb{E}(D_{k+1} | \tilde{\mathcal{F}}_k) = 0$$

Natural way : $\{\gamma_k, F_k\}_{k=1}^{\infty}$ martingale , $D_k = \gamma_k - \gamma_{k-1}$

$$\textcircled{1} \quad \mathbb{E}[D_k] = \mathbb{E}[(\gamma_k - \gamma_{k-1})] \leq \mathbb{E}|\gamma_k| + \mathbb{E}|\gamma_{k-1}| < \infty$$

$$\textcircled{2} \quad \mathbb{E}[D_{k+1} | F_k] = \mathbb{E}(\gamma_k - \gamma_{k+1} | F_k) = \gamma_k - \gamma_k = 0 \text{ a.s.}$$

Thm (Azuma - Hoeffding's inequality) [Hoeffding's for martingale differences]

Let $(D_k, F_k)_{k=1}^{\infty}$ be a martingale difference sequence where $a_k \leq D_k \leq b_k$ a.s. for $k=1 \dots n$. Then for all $t \geq 0$,

$$\mathbb{P}\left[\sum_{k=1}^n D_k \geq t\right] \leq \exp\left(\frac{-2t^2}{\sum_{k=1}^n (b_k - a_k)^2}\right) \text{ and}$$

$$\mathbb{P}\left[\sum_{k=1}^n D_k \leq -t\right] \leq \exp\left(\frac{-2t^2}{\sum_{k=1}^n (b_k - a_k)^2}\right)$$

* Bounded martingale difference concentrates.

* Azuma - Hoeffding's for Martingale sequences — skipping, similar.

Proof : Similar to Hoeffding's bound for bounded r.v., in particular we use Cramer - Chernoff method

$$\begin{aligned} \text{Let } S = \sum_{k=1}^n D_k ; \quad \mathbb{P}(S \geq t) &\leq \mathbb{P}(e^{\lambda S} \geq e^{\lambda t}) \stackrel{\text{(markov)}}{\leq} e^{-\lambda t} \mathbb{E} e^{\lambda S} \\ &\leq \exp(-\lambda t + \psi_s(\lambda)) \quad \text{where } \psi_s(\lambda) = \log \mathbb{E} e^{\lambda S} \\ &\qquad\qquad\qquad = \log \mathbb{E} e^{\lambda \sum_{k=1}^n D_k} \end{aligned}$$

Let us look at

$$\begin{aligned} \mathbb{E}\left[e^{\lambda \sum_{k=1}^n D_k} | F_{n-1}\right] &= \mathbb{E}\left[e^{\lambda D_n} e^{\lambda \sum_{k=1}^{n-1} D_k} | F_{n-1}\right] \\ &= e^{\lambda \sum_{k=1}^{n-1} D_k} \underbrace{\mathbb{E}[e^{\lambda D_n} | F_{n-1}]}_{\substack{\text{Using same technique} \leftarrow \text{M&F of a mean zero,} \\ \text{as hoeffding} \qquad \qquad \qquad \text{bounded random variable}}} \\ &\leq \left(e^{\lambda \sum_{k=1}^{n-1} D_k}\right) \cdot \exp\left(\frac{\lambda^2}{8} (b_n - a_n)^2\right) \end{aligned}$$

Last time we were looking at

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

↓
 symmetrization chaining

* Rademacher complexity:

for a set A , draw n elements a_1, \dots, a_n

$$(\text{Empirical}) \quad \text{Rademacher avg: } \hat{R}_n(A) = \mathbb{E} \sup_{a \in A} \left| \frac{1}{n} \sum_{i=1}^n \varepsilon_i a_i \right|$$

* Examples: ℓ_2 -Ball, ℓ_1 -Ball
 $(Y_{\sqrt{n}})$ (Y_n)

i.i.d. ± 1 w/equal prob (w.r.t.)

Today :- Connect R.C. to U.L.L.N

- Symmetrization argument
- Bounds of $\mathbb{E} \sup_{f \in \mathcal{F}} |P_n f - P f|$

Empirical Process Setup

- $X_1, X_2, \dots, X_n \sim \text{i.i.d. } \mathcal{D}$
- \mathcal{F} : class of real valued functions

Let,

$$F(x_1, x_2, \dots, x_n) = \{f(x_1), \dots, f(x_n) : f \in \mathcal{F}\}$$

A random subset of \mathbb{R}^n ← compute rademacher complexity of this

$$\hat{R}_n(F(x_1, \dots, x_n)) = \mathbb{E}_{\varepsilon} \sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n f(x_i) \varepsilon_i \right| \quad (\text{Empirical Rademacher complexity})$$

$$R_n(F) = \mathbb{E}_{x_1, \dots, x_n} \hat{R}(F(x_1, \dots, x_n)) \leftarrow \text{Rademacher complexity of } \mathcal{F}$$

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

Note: \mathcal{F} is a class of f_n , $F(x_1, \dots, x_n)$ is a random subset of \mathbb{R}^n

Theorem (Symmetrization)

We have $\mathbb{E} \sup_{f \in \mathcal{F}} |P_n f - P f| \leq 2 R_n(\mathcal{F})$ where $R_n(\mathcal{F})$ is defined in (1)

Proof : $x_1, \dots, x_n \sim i.i.d (X)$

Withdraw (x'_1, \dots, x'_n) from same distribution $\stackrel{(X')}{\text{s.t.}}$ $(x'_1, \dots, x'_n) \perp\!\!\!\perp (x_1, \dots, x_n)$ (independent)

$$\mathbb{E}_x f(x) = \mathbb{E}_{x'} \left(\frac{1}{n} \sum_{i=1}^n f(x'_i) \right) \quad \text{--- (1)}$$

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

Using (1)

$$= \mathbb{E}_x \sup_{f \in \mathcal{F}} \left(\frac{1}{n} \sum_{i=1}^n f(x_i) - \underbrace{\mathbb{E}_{x'} \left(\frac{1}{n} \sum_{i=1}^n f(x'_i) \right)}_{\text{constant}} \right)$$

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

* $\sup(\cdot)$ is convex, $|\cdot|$ is convex. So use Jensen's inequality (Symmetrization argument)



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

Claim : $f(x_i) - f(x'_i)$ is distributed identically to $(f(x_i) - f(x'_i)) \varepsilon_i$

Intuition : ε_i flips sign, but x_i, x'_i i.i.d, so they switch also w.p. $\frac{1}{2}$.

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

$$\leq \mathbb{E}_{x, x' \in \mathcal{F}} \sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n \varepsilon_i (f(x_i)) \right| + "$$

$$= 2 R_n(\mathcal{F})$$

$$\sup(A+B) \leq \sup(A) + \sup(B)$$

Remark : In most cases we condition on $x_1 \dots x_n$ ($x_1 \dots x_n$). So, in this setting $R_n(\tilde{F}) = \hat{R}_n(\tilde{F})$

Simple Bounds on $R_n(\tilde{F})$

Lemma (Massat's Lemma) : Suppose A is a finite subset of \mathbb{R}^n with cardinality $|A|$. Then,

$$R_n(A) = \mathbb{E} \max_{a \in A} \left| \frac{1}{n} \sum_{i=1}^n \epsilon_i a_i \right| \leq \sqrt{\frac{6 \log(2|A|)}{n}} \max_{a \in A} \sqrt{\frac{1}{n} \sum_{i=1}^n a_i^2}$$

Proof : For non-negative X , we write

$$\mathbb{E} X = \int_0^\infty P(X > x) dx \quad (\text{Exercise})$$

$$\underline{\text{Notation}} : \sqrt{\frac{1}{n} \sum_{i=1}^n a_i^2} = \|a\|_2, \quad \sum_{i=1}^n \frac{\epsilon_i a_i}{n} = \frac{a^\top \tilde{\epsilon}}{n} = a^\top \tilde{\epsilon} \quad (\tilde{\epsilon} = \frac{\epsilon}{n})$$

$$\begin{aligned} \mathbb{E} \exp \left[\frac{(a^\top \tilde{\epsilon})^2}{6 \|a\|_2^2} \right] &= \int_0^\infty P \left\{ \exp \left(\frac{(a^\top \tilde{\epsilon})^2}{6 \|a\|_2^2} \right) > x \right\} dx \\ &= \underbrace{\int_0^1 P(\cdot) dx}_{\leq 1} + \int_1^\infty P \left\{ \underbrace{\left| a^\top \tilde{\epsilon} \right|}_{\frac{1}{n} \sum_{i=1}^n a_i \epsilon_i} > 6 \|a\|_2 \sqrt{\log(x)} \right\} dx \\ &= 1 + 2 \int_1^\infty \exp \left(- \frac{2 \times 6 \|a\|_2^2 \log x}{4 \|a\|_2^2} \right) dx \quad -1 \leq \epsilon_i \leq 1 \\ &\quad \text{due to } |a^\top \tilde{\epsilon}| \\ &= 1 + 2 \int_1^\infty x^{-3} dx = 2 \quad -\text{(i)} \end{aligned}$$

We have,

$$\begin{aligned} \mathbb{E} \exp \left[\max_{a \in A} \frac{|a^\top \tilde{\epsilon}|^2}{6 \|a\|_2^2} \right] &= \mathbb{E} \max_{a \in A} \exp \frac{|a^\top \tilde{\epsilon}|^2}{6 \|a\|_2^2} \quad (e^x \text{ is increasing}) \\ &\leq \sum_{i=1}^n \mathbb{E} \exp \left[\frac{(a^\top \tilde{\epsilon})^2}{6 \|a\|_2^2} \right] \leq 2|A| \quad -\text{(ii)} \end{aligned}$$

(III) can be re-written as, $(\max \{a, b\})^2 = \max \{a^2, b^2\}$ for $a, b \geq 0$

$$\mathbb{E} \exp \left(\max_{a \in A} \frac{|a^T \tilde{\varepsilon}|}{\sqrt{6} \|a\|_2} \right)^2 \leq 2|A| \quad - (IV)$$

* Fact: Show that $x \mapsto e^{x^2}$ is convex ($x > 0$) (composition)

Using Jensen's inequality,

$$\max \left(\frac{a}{b} \right) \geq \frac{\max(a)}{\max(b)}$$

$$\exp \left(\mathbb{E} \max_{a \in A} \frac{|a^T \tilde{\varepsilon}|}{\sqrt{6} \|a\|_2} \right)^2 \leq \mathbb{E} \exp \left(\max_{a \in A} \frac{|a^T \tilde{\varepsilon}|}{\sqrt{6} \|a\|_2} \right)^2 \leq 2|A|$$

$$\mathbb{E} \max_{a \in A} |a^T \tilde{\varepsilon}| \leq \sqrt{6 \log 2|A|} \max \|a\|_2$$

(some bug in pt \leftarrow in eqn (II))

Application of ULLN

Apply Massart's Lemma to $R_n(\mathcal{F})$

Assumption: \mathcal{F} is Boolean, $f(x) \in \{0, 1\} \quad \forall x \in \mathcal{X}, f \in \mathcal{F}$

* Classification, testing

** Let us fix \mathcal{F} and $(x_1 \dots x_n) = (x_1 \dots x_n)$

$$\mathbb{E} \sup_{f \in \mathcal{F}} |P_f - P_f| \leq 2R_n(\mathcal{F}) = 2R_n(x_1 \dots x_n)$$

$$\mathcal{F}(x_1 \dots x_n) = \{f(x_1) \dots f(x_n) : f \in \mathcal{F}\}$$

$$\leq \sqrt{\frac{6 \log (2|\mathcal{F}(x_1 \dots x_n)|)}{n}} \underbrace{\max_{f \in \mathcal{F}} \sqrt{\frac{1}{n} \sum_{i=1}^n f(x_i)^2}}_{\leq 1}$$

* If $|\mathcal{F}(x_1 \dots x_n)| \sim 2^n$

$$\mathbb{E} \sup_{f \in \mathcal{F}} |P_f - P_f| \leq \sqrt{\frac{6 \log 2 + 6 \log 2}{n}} \quad \text{Rmk: Growth of } |\mathcal{F}(x_1 \dots x_n)| \text{ is important (different from the size of function class)}$$

* Polynomial discrimination

$$|\mathcal{F}(x_1 \dots x_n)| \sim \text{poly}(n) \quad P(n) = n^\alpha \quad \mathbb{E} \sup_{f \in \mathcal{F}} |P_f - P_f| \leq \sqrt{\frac{\alpha \log n}{n}} \quad \text{checked by VC dimension.}$$

Last class : Symmetrization to show that

$$\mathbb{E} \sup_{f \in \mathcal{F}} |P_n f - P_f| \leq 2 R_n(\mathcal{F}(x_1 \dots x_n))$$

where $\mathcal{F}(x_1 \dots x_n) = \{f(x_1), \dots, f(x_n) : f \in \mathcal{F}\} \subseteq \mathbb{R}^n$

We condition on $x_1 = x_1, \dots, x_n = x_n$

$\rightarrow \mathbb{E}_{x_1 \dots x_n} R_n(\mathcal{F}(x_1 \dots x_n))$ can use instead

(Recall)

Lemma (Massart's Lemma) : Suppose A is a finite subset of \mathbb{R}^n with cardinality $|A|$. Then,

$$R_n(A) = \mathbb{E} \max_{a \in A} \left| \frac{1}{n} \sum_{i=1}^n \epsilon_i a_i \right| \leq \sqrt{\frac{6 \log(2|A|)}{n}} \max_{a \in A} \sqrt{\frac{1}{n} \sum_{i=1}^n a_i^2}$$

Using this,

$$R_n(\mathcal{F}(x_1 \dots x_n)) \lesssim \sqrt{\frac{\log(|\mathcal{F}(x_1 \dots x_n)|)}{n}} \max_{f \in \mathcal{F}} \sqrt{\frac{1}{n} \sum_{i=1}^n f(x_i)^2}$$

↑
ignoring
universal constants

* Want upper bound on $|\mathcal{F}(x_1 \dots x_n)|$

Assumption : \mathcal{F} is boolean. $f(x_i) \in \{-1, 1\} \forall x_i, f \in \mathcal{F}$

* \mathcal{F} Boolean, $|\mathcal{F}(x_1 \dots x_n)| \leq 2^n$ Is this useful? (no)

$R_n(\mathcal{F}) \lesssim 1$

↓
set or \mathcal{F} is not
learnable (no matter how many samples)

* Polynomial Discrimination :

\mathcal{F} has polynomial discrimination if there exists a polynomial $P(n)$ s.t.
 $|\mathcal{F}(x_1 \dots x_n)| \leq P(n)$

Ex : If $P(n) = n^\alpha$, then $R_n(\mathcal{F}) \lesssim \sqrt{\frac{\alpha \log n}{n}}$

Q: How to check if \mathcal{F} has polynomial discrimination?

→ VC dimension (Vapnik Chernonenkis) [A combinatorial object]

Defn (VC dimension)

Shattering: A finite subset $\{x_1 \dots x_m\}$ is said to be shattered by a Boolean class \mathcal{F} if

$$|\mathcal{F}(x_1 \dots x_m)| = 2^m ; \quad \mathcal{F}(x_1 \dots x_m) = \{0, 1\}^m$$

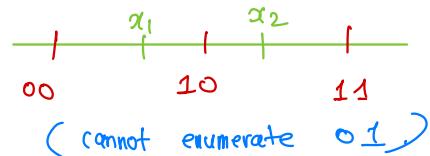
The VC-dimension D of \mathcal{F} is the maximum integer D for which any set of $\{x_1 \dots x_D\}$ is shattered by \mathcal{F} .

* If $\{x_1 \dots x_D\}$ shattered for every D , VC dim is ∞ .

Examples: ① Intervals on \mathbb{R}

$$\mathcal{S}_{\text{left}} = \left\{ \mathbb{1}_{(-\infty, a]} : a \in \mathbb{R} \right\} \rightarrow D=1 \quad \{x_1\} \text{ is shattered}$$

$$\text{where } \mathbb{1}_{(-\infty, a]}(x) = \begin{cases} 1, & \text{if } x \leq a \\ 0, & \text{otherwise} \end{cases} \quad \mathbb{1}_{(-\infty, x_1-\epsilon)}, \mathbb{1}_{(-\infty, x_1+\epsilon)}$$



② Extend to \mathbb{R}^2

$$\mathcal{S}_{\text{rect}} = \left\{ \mathbb{1}_{[a_1, b_1] \times [a_2, b_2]} : a_i, b_i \in \mathbb{R} \right\}$$

$$\text{VC}(\mathcal{S}_{\text{rect}}) = 2 \quad \left. \begin{array}{c} \cdot \\ \cdot \\ \cdot \\ \cdot \end{array} \right\} \text{can't shatter!}$$

③ Lemma: Let \mathcal{V} as a D -dimensional vector space on real functions on \mathcal{X}

Let, $\mathcal{F} = \left\{ \mathbb{1}_{(f \geq 0)} : f \in \mathcal{V} \right\}$. Then, $\text{VC}(\mathcal{F})$ is atmost D .

Pf: Let $\{x_1 \dots x_{D+1}\}$ and consider $T = \{f(x_1) \dots f(x_{D+1}) : f \in \mathcal{V}\}$
 $\therefore \exists$ some coefficients $\alpha \in \mathbb{R}^{D+1} \neq 0$.

$$\sum_{i=1}^{D+1} \alpha_i f(x_i) = 0 \quad \text{for all } f \in \mathcal{V} \quad \leftarrow V^* \text{ (evaluation = dual space)}$$

Since $\alpha \neq 0$, wlog, \exists index k s.t. $\alpha_k > 0$

Let's assume F shatters $\{x_1, \dots, x_{k+1}\}$, then there exist

$$f \in \mathcal{V} \text{ s.t. } \begin{aligned} f(x_i) < 0 &\quad \text{for all } i \text{ s.t. } \alpha_i > 0 \\ f(x_i) \geq 0 &\quad " \quad " \quad \alpha_i \leq 0 \end{aligned}$$

With this,

$$\sum_{i=1}^{k+1} \alpha_i f(x_i) = \underbrace{\sum_{i: \alpha_i > 0} \alpha_i f(x_i)}_{< 0} + \underbrace{\dots}_{\leq 0} < 0$$

a contradiction.

VC dimension and VLLM

Lemma Let \mathcal{H}_n denote the collection of all closed half-spaces in \mathbb{R}^n . $VC(\mathcal{H}_n) = n+1$

*Half-Spaces / Linear classifiers are learnable

VC-dimension and ULLN

Lemma (Sauer-Shelah-V-C) (SSVC Lemma). Suppose the VC dimension of Boolean function class F is D (abstract measure). Then for every $n \geq 1$ and x_1, \dots, x_n , $|F(x_1, \dots, x_n)| \leq \binom{n}{0} + \binom{n}{1} + \dots + \binom{n}{D}$, where $\binom{n}{k} = 0$ if $k > n$.

Remark F has poly. discrimination
 $P(n) = \left(\frac{e}{D}\right)^D n^D$

Remark Use Symmetrization to show $E_{f \sim F} P_{f \neq f^*} \leq 2 \ln |F(x_1, \dots, x_n)|$. It turns out that $\frac{2}{n} \ln \binom{n}{D}$ can be reduced by further argument called Chaining.

Remark: SSVC lemma can be proved using induction and "down-shifting" (read pf in HDS - chapter 4 Prop. 4.18)

Chaining

Covering and Packing :

Let (T, ρ) denotes a metric space T with associated metric ρ .

$$(\rho: T \times T \rightarrow \mathbb{R})$$

- Ⓐ $\rho(\theta, \tilde{\theta}) \geq 0$, $\rho(\theta, \tilde{\theta}) = 0 \text{ if } \theta = \tilde{\theta}$ (non negativity)
- Ⓑ $\rho(\theta, \tilde{\theta}) = \rho(\tilde{\theta}, \theta)$ (symmetric)
- Ⓒ $\rho(\theta, \tilde{\theta}) + \rho(\tilde{\theta}, \bar{\theta}) \geq \rho(\theta, \bar{\theta})$ (Triangle inequality)

Ex : \mathbb{R}^d (or subset)

$$\rho(\theta, \tilde{\theta}) = \|\theta - \tilde{\theta}\|_2 = \left(\sum_{j=1}^d (\theta_j - \tilde{\theta}_j)^2 \right)^{1/2}$$

Ex : Boolean cube = $\{0, 1\}^d$

$$\rho(\theta, \tilde{\theta}) = d_H(\theta, \tilde{\theta}) = \# \text{ coordinates at which } \theta, \tilde{\theta} \text{ differ} = \sum_{j=1}^d \mathbb{1}_{(\theta_j \neq \tilde{\theta}_j)}$$

$$\text{Normalized hamming dist} = \frac{\sum_{j=1}^d \mathbb{1}_{(\theta_j \neq \tilde{\theta}_j)}}{d}$$

Ex : $C[0, 1]$: set of all continuous functions in $[0, 1]$

$$\text{Metric} : (\text{Sup-norm}) \quad \sup_{x \in [0, 1]} |f(x) - g(x)| = \rho(f, g) = \|f - g\|_\infty$$

Ex : L^2 : space of square integrable fns in $[0, 1]$

$$\|f - g\|_2 = \left[\int_0^1 (f(x) - g(x))^2 dx \right]^{1/2}$$

Ex : $L^2(\mu; [0, 1])$

$$\|f - g\|_2 = \left[\int_0^1 (f(x) - g(x))^2 d\mu(x) \right]^{1/2}$$

Covering Number: A δ -cover of a set T w.r.t. ρ is a set

$$\{\theta^1, \dots, \theta^N\} \subset T \text{ s.t. } \forall \theta \in T \exists \theta^i \text{ s.t. } \rho(\theta, \theta^i) \leq \delta$$

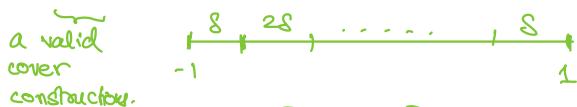
Covering Number, $N(\delta, T, \rho)$ is the cardinality of the smallest cover.

$$* \delta_1 \leq \delta_2, N(\delta_1, T, \rho) \geq N(\delta_2, T, \rho)$$

Ex. Unit hypercube

$$d=1: [-1, 1], \rho(\theta, \theta') = |\theta - \theta'|$$

$$N(\delta, [-1, 1], \rho) \leq \frac{1}{\delta} + 1 \text{ (break into intervals of } 2\delta)$$



Extend to d

$$N(\delta, [-1, +1]^d, \|\cdot\|_\infty) \leq \left(1 + \frac{1}{\delta}\right)^d \leftarrow \text{due to grid}$$

Ex: Binary Hypercube: $\{-1, +1\}^d = H_d$

$$\rho(\theta, \theta') = \frac{1}{d} \sum_{j=1}^d \mathbb{1}_{\{\theta_j \neq \theta'_j\}}$$

$$N_H(\delta, H_d; \rho) \leq \underbrace{2^{(1-\delta)d}}_{\text{turns out to be a lower bound}}$$

fix $\frac{1}{\delta}$ take all combo. $(1-\delta)d$

$$T(\delta) = \{\theta \in H_d \mid \theta_j = 0 \forall j \notin S\}$$

$$S = \{1, 2, \dots, (1-\delta)d\}$$

Match δ exactly on $S \Rightarrow d_H \leq \delta$.

Packing Number

A δ -packing of a set T w.r.t. a metric ρ is $\{\theta^1, \dots, \theta^N\} \subset T$ s.t. $\rho(\theta^i, \theta^j) > \delta \forall i, j \in [N], i \neq j$

Packing Number $M(\delta, T, \rho)$ is the cardinality of maximum such set.

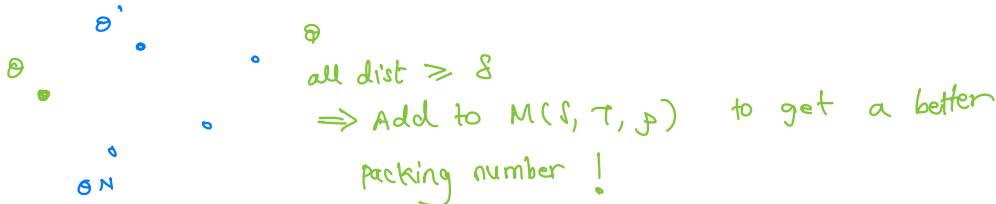


Lemma: For $\delta > 0$ we have

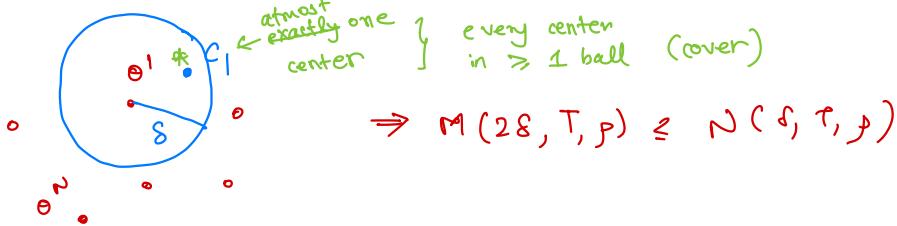
$$M(2\delta, \tau, p) \stackrel{(b)}{\leq} N(\delta, \tau, p) \stackrel{(a)}{\leq} M(\delta, \tau, p)$$

Remark: These numbers are otherwise equivalent.

Pf: (a) Let $\{\theta^1, \dots, \theta^N\}$ be maximal set of δ -separated pts for (τ, f)
 $M(\delta, \tau, p) = N$. Also a valid cover


all dist $\geq \delta$
 \Rightarrow Add to $M(\delta, \tau, p)$ to get a better
packing number!

(b) Let $\{\theta^1, \dots, \theta^N\}$ be min set of pts f -cover


almost exactly one center } every center in ≥ 1 ball (cover)
 $\Rightarrow M(2\delta, \tau, p) \leq N(\delta, \tau, p)$

Examples.

Proposition: Let $\|\cdot\|$ denote some norm \mathbb{R}^d . Also, $B_R = \{x \in \mathbb{R}^d, \|x\| < R\}$

Then $\delta > 0$, we have,

$$M(\delta R, B_R, \|\cdot\|) \leq \left(1 + \frac{2}{\delta}\right)^d$$

(easy to see
for ∞ norm)

Proof: (Volumetric argument)

Let x_1, \dots, x_N denote any set of points in B_R that are δR separated

i.e., $\|x_i - x_j\| > \delta R \quad \forall i \neq j$

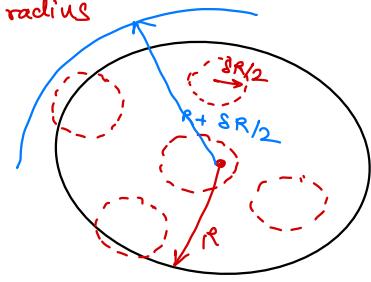
Then, the closed balls

$$\overline{B}(x_i, \frac{\delta R}{2}) = \{x \in \mathbb{R}^d \mid \|x_i - x\| \leq \frac{\delta R}{2}\}$$

are disjoint

Moreover, all these balls are contained within a ball of radius

$$\text{So, } M(\cdot) \cdot c \left(\frac{\delta R}{2}\right)^d \leq c\left(R + \frac{\delta R}{2}\right)^d \\ \Rightarrow M(\cdot) \leq \left(1 + \frac{2R}{\delta}\right)^d$$



Normalize : $M(\delta, B_R, \|\cdot\|) \leq \left(1 + \frac{2R}{\delta}\right)^d$

Notation : Covering number / δ -net

Example : Cover / Pack a function class.

Proposition : Let $\Theta \subseteq \mathbb{R}^d$ be a non-empty subset with diameter D $p(\theta_1, \theta_2) \leq D$
 $\forall \theta_1, \theta_2 \in \Theta$

Let $\mathcal{F} = \{f_\theta : \theta \in \Theta\}$ satisfying $|f_{\theta_1}(x) - f_{\theta_2}(x)| \leq \Gamma(x) \|\theta_1 - \theta_2\|_2$

Fix a measure φ .

$$p^2(f, g) = \int_X (f(x) - g(x))^2 d\varphi(x)$$

Then $\delta > 0$,

$$M(\delta, \mathcal{F}, \varphi) \leq \left(1 + \frac{2D \|\Gamma\|_\varphi}{\delta}\right)^d \text{ where } \|\Gamma\|_\varphi^2 = \int_X \Gamma(x)^2 d\varphi(x)$$

Pf : Need to find R .

$$\int_X (f_{\theta_1}(x) - f_{\theta_2}(x))^2 d\varphi(x) \leq \int_X \Gamma(x)^2 \|\theta_1 - \theta_2\|^2 d\varphi(x) \\ \leq \|\Gamma\|_\varphi^2 \|\theta_1 - \theta_2\|^2$$

$$\Rightarrow p(f_{\theta_1}, f_{\theta_2}) \leq \|\Gamma\|_\varphi \|\theta_1 - \theta_2\|_2 \quad \textcircled{1}$$

From $\textcircled{1}$, to find a δ -packing for \mathcal{F} , it is sufficient to find a

$\frac{\delta}{\|\Gamma\|_\varphi}$ packing in Θ .

$$M(\delta, \kappa, \rho) \leq M\left(\frac{\delta}{\|\Gamma\|_Q}, \Theta, \|\cdot\|_2\right) \xrightarrow{\substack{\text{packing for lhs} \\ \Rightarrow \text{packing for rhs.}}}$$

$$B(a, D) = \{x \in \mathbb{R}^d \mid \|x - a\| \leq D\}$$

$$\leq M\left(\frac{\delta}{\|\Gamma\|_Q}, B(a, D), \|\cdot\|_2\right) \stackrel{\substack{\text{previous} \\ \text{lemma}}}{\leq} \left(1 + \frac{2D\|\Gamma\|_Q}{\delta}\right)^d$$

■