

Metric Entropy

- ① Physics: The degree of disorder or randomness in a system.
- ② Information Theory: The randomness of a distribution.

In this chapter, Metric Entropy represents a similar meaning:

size of a class of functions

Learning theory

To motivate the definition of metric entropy, we introduce some fundamental knowledge about learning theory.

We assume that (\mathcal{X}, Σ, P) is a probability space, and X_1, \dots, X_n are iid sampled from law P , we mark that

$$X_1, \dots, X_n \sim (\mathcal{X}, \Sigma, P)$$

$X_1, \dots, X_n \sim (\chi, \bar{\sigma}^2, P)$ and
 $Y_i = f_0(X_i) + \varepsilon_i, f_0 \in \mathcal{H}$,
where \mathcal{H} is a class of functions: $\chi \mapsto \mathbb{R}$.

If we observe $(X_i, Y_i)_{i=1, \dots, n}$, we want to estimate the true f_0 , usually, given loss function $L: \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}_+$,

$$f_0 := \operatorname{argmin}_{f \in \mathcal{H}} \mathbb{E} L(Y_i, f(X_i)),$$

Example

Risk function

$$\chi = \mathbb{R}^p.$$

① $\mathcal{H} = \mathbb{R}^p$, $L(y_1, y_2) = (y_1 - y_2)^2$, $\operatorname{Var} \varepsilon_i = 6^2 < \infty \Leftrightarrow$ Least square.

② $\mathcal{H} = \{\beta \in \mathbb{R}^p; \|\beta\|_0 \leq s\} \Leftrightarrow$ Best subset selection.

③ $\mathcal{H} = \{f \in \mathbb{R}^p \rightarrow \mathbb{R}, f^{(k)} \text{ exists and } \|f^{(k)}\|_q \leq c\} \Leftrightarrow$ smoothing spline.

④ RKHS / SVM / DL.

We can define so called M-estimator of f_0 by empirical risk:

$$\hat{f}_n = \arg \min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n L(Y_i, f(x_i))$$

We need to answer some questions:

- ① How to ensure $\hat{f}_n \rightarrow f_0$ in some sense. (Asymptotic)
- ② How to ensure $P(d(f_0, \hat{f}_n) > \epsilon)$ is small enough for $\epsilon > 0$, where d is a metric on \mathcal{F} . (Non-asymptotic)

Learning Process

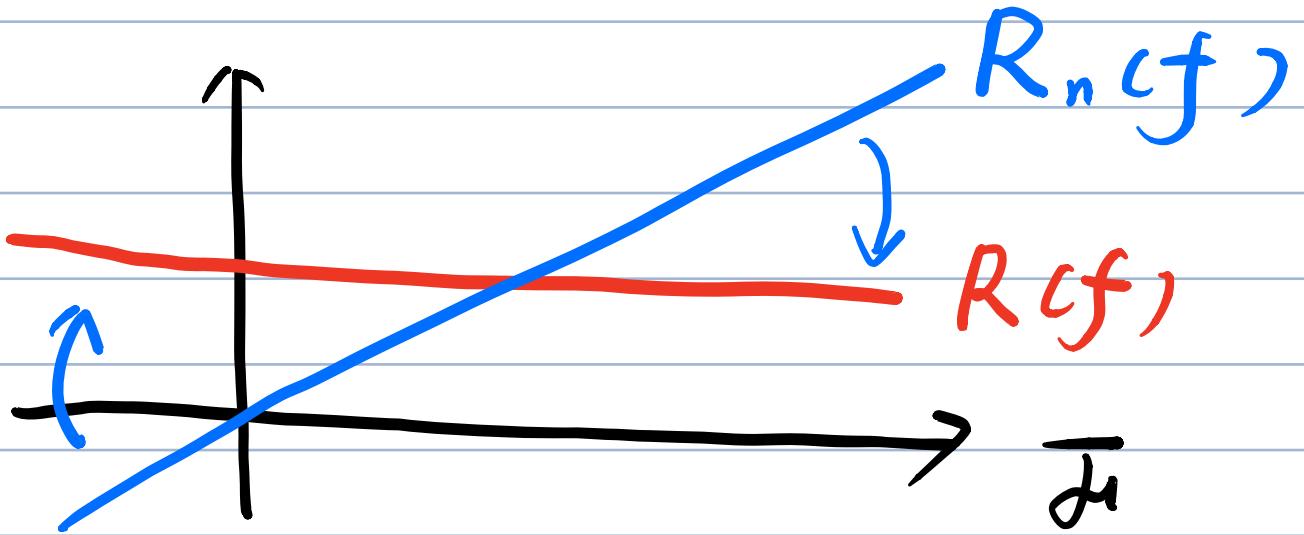
$$R(f) = \mathbb{E} L(Y_i, f(x_i)),$$

$$R_n(f) = \frac{1}{n} \sum_{i=1}^n L(Y_i, f(x_i)),$$

(Formally, $P(L(\cdot, f_0) > \epsilon)$ and $P_n(L(\cdot, f_0) > \epsilon)$)

Notice that $R_n(f) \rightarrow_{a.s.} R(f)$, for fixed f , but :

$$\begin{aligned} R_n(\hat{f}_n) &= \inf_{f \in \mathcal{F}} R_n(f) \\ &\Rightarrow \inf_{f \in \mathcal{F}} R(f), \end{aligned}$$



A key condition for
 $\inf_{f \in \mathcal{F}} R_n(f) \rightarrow \inf_{f \in \mathcal{F}} R(f)$ ①

is the "compactness" of \mathcal{H} . Moreover,
 if $\|\mathcal{H}\| < \infty$, ① holds, which implies
 the "magnitude" of \mathcal{H} seem to play
 an important role on the convergence of
 Learning process.

In this chapter, we focus on the
 definition of the "magnitude" of
 a given class of functions, to eval-
 uate the possible risk for a learning
 procedure.

Setting

① $X_1, \dots, X_n \stackrel{iid}{\sim} (\mathcal{X}, \bar{\mathcal{E}}, P)$.

② \mathcal{F} is a class of functions:
 $\mathcal{X} \mapsto \mathbb{R}$.

③ \mathcal{F} is a metric space equipped with metric d .

If $|\mathcal{F}| = \infty$, in order to quantify the size of $|\mathcal{F}|$, it's common to get the "finite version" of \mathcal{F} by the covering number of \mathcal{F} .

Covering Number

$\{f_i\} \subset \mathcal{F}$ is said to be a δ -cover of \mathcal{F} if $\mathcal{F} \subset \bigcup_i \bar{B}(f_i; \delta)$.

Mark $C_\delta(\mathcal{F})$ be the collection of all δ -cover of \mathcal{F} , then

$$N(\delta; \mathcal{F}, d) := \inf_{A \in C_\delta(\mathcal{F})} |A|.$$

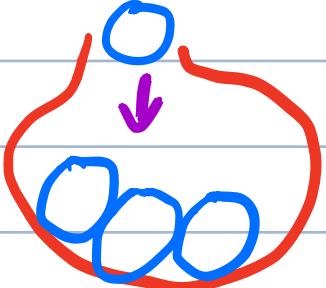
Remark:

① If $N(\delta; \mathcal{F}, d) < \infty$, $\forall \delta > 0$, we call \mathcal{F} is totally bounded.

② The quantity $\log N(\delta; \mathcal{F}, d)$ is called as the metric entropy of \mathcal{F} .

Normally, the covering number of \mathcal{F} is hard to calculate, since it lacks efficient algorithms to get the minimal cardinality of $C(\mathcal{F})$. Additionally, we don't need to calculate the specific form of $N(\delta; \mathcal{F}, d)$, to be concise, note that $N(\delta; \mathcal{F}, d) \uparrow$, if $\delta \downarrow 0$, we use the speed of increasing of $N(\delta; \mathcal{F}, d)$ as $\delta \downarrow 0$ to quantify the "size" of \mathcal{F} .

So we need to create a "equivalent amount" of covering number, which is easy to calculate.



Packing Number

$\{g_j\} \subset \mathcal{F}$ is said to be a δ -packing of \mathcal{F} if

$$d(g_j, g_i) > \delta \text{ if } i \neq j.$$

Similarly, $P_\delta(\mathcal{F})$ is the collection of all δ -packing of \mathcal{F} ,

$$P(\delta; \mathcal{F}, d) := \sup_{B \in P_\delta(\mathcal{F})} |B|$$

Lemma 1 $\forall \delta > 0$, we have

$$P(2\delta; \mathcal{F}, d) \leq N(\delta; \mathcal{F}, d) \leq P(\delta; \mathcal{F}, d)$$

Pf: $\forall \{f_i\} \in P_{2\delta}(\mathcal{F}) \setminus \{g_j\}$

$$\in C_\delta(\mathcal{F}),$$

For $\forall i$, $\exists g_{n_i} \in \{g_j\}$ s.t.

$$d(f_i, g_{n_i}) \leq \delta.$$

Define $\phi(f_i) = g_{n_i}$, claim

that ϕ is injection, since if $\exists f_j \in \{f_i\}$

$$\text{s.t. } d(f_j, g_{n_i}) \leq \delta$$

$$\Rightarrow d(f_j, f_i) \leq 2\delta, \text{ contradiction.}$$

$$\Rightarrow |\{f_i\}| \leq |\{g_j\}|.$$

$$\Rightarrow P(2\delta; \mathcal{F}, d) \leq N(\delta; \mathcal{F}, d).$$

For $N(8; \bar{f}, d) \leq P(8; \bar{f}, d)$,
 Let $\{f_i\}$ be the maximal 8-packing of \bar{f} , i.e.,
 $d(f_i, f_j) > 8, \forall i \neq j$.
 Then $\{f_i\}$ is a 8-cover of \bar{f} , since
 if $\exists f \in \bar{f}$ s.t.
 $d(f, f_i) > 8, \forall f_i \in \{f_i\}$,
 $\Rightarrow \{f_i\} \cup \{f\} \in P_8(\bar{f})$, contradiction.
 $\Rightarrow N(8; \bar{f}, d) \leq P(8; \bar{f}, d)$.

□

Example 1

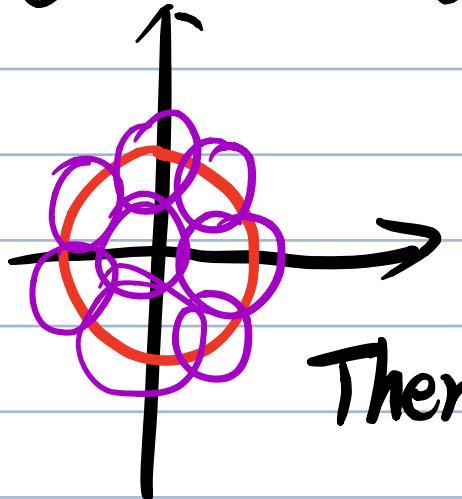
Lemma 1 If $\mathbb{F} = \mathbb{R}^P$ equipped with norm $\|\cdot\|$, define

$$\mathbb{B} = \{\theta \in \mathbb{R}^P; \|\theta\| \leq 1\},$$

Then

$$(\frac{1}{8})^P \leq N(8; \mathbb{B}, \|\cdot\|) \leq (\frac{2}{8} + 1)^P$$

Pf: Let $\{\theta_i\} \in C_8(\bar{B})$, i.e.,



$$\bar{B} \subset \bigcup_i \bar{B}(\theta_i; 8)$$

$$\text{and } |\{\theta_i\}| = N(8; \bar{B}, \text{II} \cdot \text{II})$$

Then

$$\text{Vol}(C\bar{B}) \leq \text{Vol} \left(\bigcup_i \bar{B}(\theta_i; 8) \right)$$

$$\leq \sum_i \text{Vol}(\bar{B}(\theta_i; 8))$$

$$= N(8; \bar{B}, \text{II} \cdot \text{II}) \text{Vol}(\bar{B}(0; 8))$$

$$= 8^3 N(8; \bar{B}, \text{II} \cdot \text{II}) \text{Vol}(C\bar{B})$$

$$\Rightarrow \left(\frac{1}{8}\right)^3 \leq N(8; \bar{B}, \text{II} \cdot \text{II}).$$

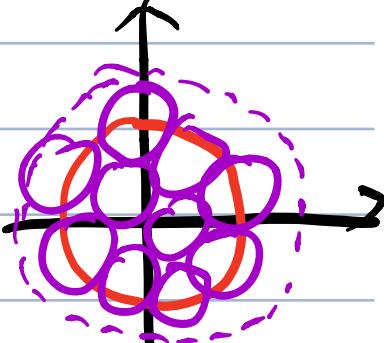
Note that $N(8; \bar{B}, \text{II} \cdot \text{II}) \leq P(8; \bar{B}, \text{II} \cdot \text{II})$
and $P(8; \bar{B}, \text{II} \cdot \text{II}) \cdot \text{Vol}(\bar{B}(0, \frac{8}{2}))$

$$\leq \text{Vol}(\bar{B} + \bar{B}(0, \frac{8}{2}))$$

$$\Rightarrow P(8; \bar{B}, \text{II} \cdot \text{II})$$

$$\leq \left(1 + \frac{8}{2}\right)^3 / \left(\frac{8}{2}\right)^3 = \left(\frac{3}{8} + 1\right)^3$$

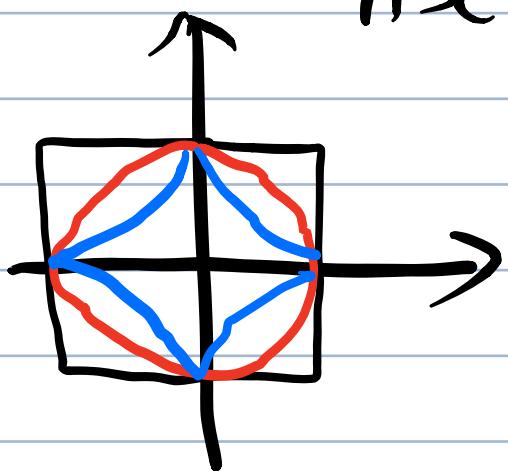
□



Covering Numbers of unit ball
 p -norm of \mathbb{R}^n is defined as:

$$\|x\|_p = \left(\sum_{i=1}^n |x_i|^p \right)^{\frac{1}{p}}, \forall x \in \mathbb{R}^n,$$

Unit ball in \mathbb{R}^n : $B_p = \{x \in \mathbb{R}^n : \|x\|_p \leq 1\}$,



Notice that

$$NC_8(B_\infty, \| \cdot \|_\infty) \leq \left(\lceil \frac{1}{8} \rceil\right)^p$$

and

$$NC_8(B_\infty, \| \cdot \|_\infty) \geq PC_2 8(B_\infty, \| \cdot \|_\infty)$$

$$\geq \left(\lfloor \frac{1}{8} \rfloor\right)^p$$

$$\Rightarrow NC_8(B_\infty, \| \cdot \|_\infty) \asymp \left(\frac{1}{8}\right)^p$$

Example 2

(Lipschitz function)

In learning theory, the continuous functions usually are the functions we want to learn, but as the Lemma in Functional Analysis:

The bound set of infinite dimensional vector space is not totally bounded.

We usually need to "imagine" some regularisations for continuous functions space, for example, special convex bound set. Define:

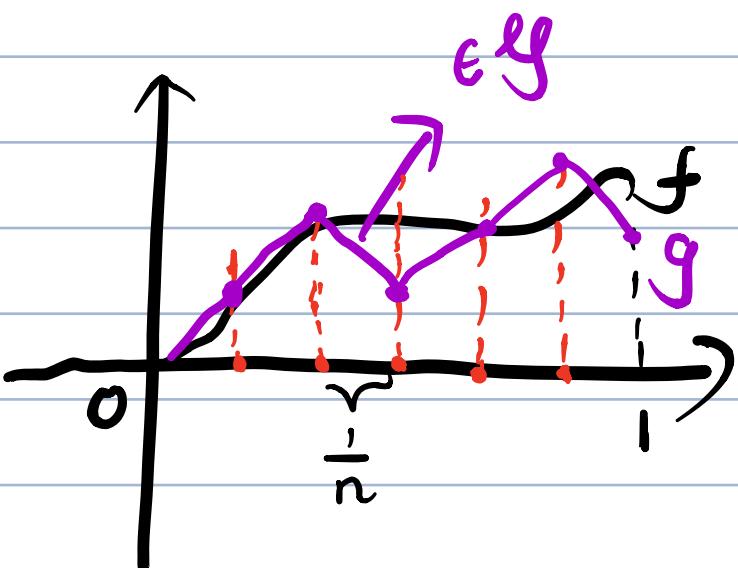
$\mathcal{J}_L := \{g : [0, 1] \rightarrow \mathbb{R}, g(0) = 0 \text{ and}$

$|g(x) - g(x')| \leq L|x - x'|, \forall x, x' \in [0, 1]\}$, where $L > 0$.

\mathcal{J}_L equids with $\|\cdot\|_{\sup}$:

$$\|f - g\|_{\sup} = \sup_{x \in [0, 1]} |f(x) - g(x)|,$$

$f, g \in \mathcal{J}_L$. Then:



For $\forall f \in F_L$,

$\exists g \in \mathcal{G}$ s.t.

$$\|f - g\|_{\sup} \leq \frac{L}{n},$$

and if $g_1, g_2 \in \mathcal{G}$

$$\|g_1 - g_2\|_{\sup} \geq \frac{2L}{n},$$

Given $\delta > 0$, let $n = \lceil \frac{L}{\delta} \rceil$

$$\Rightarrow |\mathcal{G}| = 2^n = 2^{\lceil \frac{L}{\delta} \rceil}$$

$$\Rightarrow \log N(\delta; \mathcal{F}_L, \|\cdot\|_{\sup}) \leq \frac{L}{\delta}.$$

$$\text{Let } n = \lfloor \frac{L}{\delta} \rfloor \Rightarrow |\mathcal{G}| = 2^{\lfloor \frac{L}{\delta} \rfloor}$$

$$\Rightarrow \log N(\delta; \mathcal{F}_L, \|\cdot\|_{\sup}) \geq \frac{L}{\delta}.$$

Remark:

① For p -dimensional Lipschitz function,
it's easy to show:

$$\log N(\delta; \mathcal{F}_L, \|\cdot\|_{\sup}) \asymp \left(\frac{L}{\delta}\right)^p,$$

which implies the curse of dimension.

② If $f(0)$ is not fixed but
 $|f(x)| \leq c$,

$$\log N(\delta; \mathcal{H}_L, \| \cdot \|_{\text{sup}}) \lesssim \frac{L}{\delta}$$

still holds.

