

In this chapter, we focus on the upper bound of the supremum of **Sub-Gaussian process** on some metric space T , i.e.,

$$\mathbb{E} \sup_{\theta \in T} |X_\theta| \text{ or } \mathbb{E} \sup_{\theta \in T} X_\theta.$$

Firstly, we motivate the chapter by two important examples

Example 1 Gaussian / Rademacher complexity

$$g(T) = \mathbb{E} \sup_{\theta \in T} \langle \theta, \varepsilon \rangle, \quad \varepsilon \sim N(0, I_n)$$

$\bar{R}(T) = \mathbb{E} \sup_{\theta \in T} \langle \theta, \varepsilon \rangle, \quad \varepsilon \sim \text{Unif } \{-1, 1\}^n$, where $\langle \theta, \varepsilon \rangle$ can be viewed as a stochastic process on $T \subset \mathbb{R}^n$.

Now we give an example of complexity: **Maximum singular value of random matrix**.

Let $W \in M^{n \times d}(\mathbb{C})$ be a random matrix with zero-mean iid entries W_{ij} , the operator norm on $M^{n \times d}(\mathbb{C})$:

$$\|W\| = \sup_{v \in S^{d-1}} \sqrt{v^\top W^\top W v}$$

$$= \sup_{v \in S^{d-1}, u \in S^{n-1}} u^\top W v$$

$$= \sup_{v,u} \sum_{i,j} u_i W_{ij} v_j$$

$$= \sup_{v,u} \langle\langle W, uv^T \rangle\rangle,$$

where $\langle\langle A, B \rangle\rangle = \sum_{i,j} A_{ij} B_{ij}$, note that

$$\begin{aligned} & \{uv^T; v \in S^{d-1}, u \in S^{n-1}\} \\ &= \{A \in M_{n \times d}(\mathbb{R}); \text{rank}(A) = 1, \sqrt{\langle\langle A, A \rangle\rangle} := \|A\| = 1\} \\ &= M_{1,n}(\mathbb{R}) \end{aligned}$$

$$\Rightarrow \|W\| = \sup_{A \in M_{1,n}(\mathbb{R})} \langle\langle A, W \rangle\rangle.$$

where $\langle\langle A, W \rangle\rangle$ can be viewed as the stochastic process on $M_{1,n}(\mathbb{R})$.

Example 2 Learning process

For regression problem, we assume that $X_1, \dots, X_n \stackrel{iid}{\sim} (\mathcal{X}, \Sigma, P)$, and

$Y_i = T(X_i) + \varepsilon_i$, where $X_i \perp \varepsilon_i$ and $\varepsilon_i \perp \varepsilon_j$, $i \neq j$, $\mathbb{E} \varepsilon_i = 0$.

Given a class of functions on $\mathcal{X} : \mathcal{F}$, since T may not belong to \mathcal{F} , then

$$f^* = \arg \min_{f \in \mathcal{F}} \mathbb{E} (f(x) - T(x))^2$$

is the optimal function we can learn. With help from predictor T , note that

$$\begin{aligned}\mathbb{E}(f(x) - Y)^2 &= \mathbb{E}(f(x) - T(x) - \varepsilon)^2 \\ &= \mathbb{E}(f(x) - T(x))^2 + \mathbb{E}\varepsilon^2,\end{aligned}$$

if $\mathbb{E}\varepsilon^2 < \infty$, then

$$f^* = \arg \min_{f \in \mathcal{F}} \mathbb{E}(f(x) - Y)^2,$$

which implies that we can use M-estimator of f^* :

$$\hat{f}_n = \arg \min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n (f(x_i) - Y_i)^2,$$

for approximation of f^* .

To ensure \hat{f}_n is a good estimator for f^* , we need to evaluate the excess risk:

$$R(\hat{f}_n) - R(f^*),$$

where $R(\hat{f}_n) = \mathbb{E}(\hat{f}_n(x) - T(x))^2$, $R(f^*) = \mathbb{E}(f^*(x) - T(x))^2$.

Let $R_n(\hat{f}_n) = \frac{1}{n} \sum_{i=1}^n (\hat{f}_n(x_i) - T(x_i))^2$, note that

$$\begin{aligned}R(\hat{f}_n) - R(f^*) &\leq R(\hat{f}_n) - R_n(\hat{f}_n) \\ &+ R_n(\hat{f}_n) - R_n(f^*) + R_n(f^*) - R(f^*) \\ &\leq 0\end{aligned}$$

$$\leq 2 \sup_{f \in \mathcal{F}} |R_n(f) - R(f)|,$$

Now, let $\tilde{J} = \mathbb{E}(f - \bar{f})^2; f \in \mathcal{F}$, define:

$$P_f := R(f) = \mathbb{E} f(X),$$

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

$$P_n^\varepsilon f = \frac{1}{n} \sum_{i=1}^n \varepsilon_i f(x_i), \quad \varepsilon_i \sim \text{Unif}\{-1, 1\} \perp X_i.$$

Then

Lemma 1

$$\mathbb{E} \|P_n - P\|_{\mathcal{F}} \leq 2 \mathbb{E} \|P_n^\varepsilon\|_{\mathcal{F}}, \text{ where}$$

$\|Q - P\|_{\mathcal{F}} = \sup_{f \in \mathcal{F}} |Q f - P f|$, Q, P are two laws on X .

$$P_f: |P_n f - P f| = |P_n f - \mathbb{E} P_n' f|_X$$

(where $P_n' f = \frac{1}{n} \sum_{i=1}^n f(x'_i)$, $x'_i \perp x_i$)

Then

$$\leq \mathbb{E} |P_n f - P_n' f|_X$$

$$\Rightarrow \mathbb{E} \|P_n - P\|_{\mathcal{F}} \leq \mathbb{E} \sup_{f \in \mathcal{F}} \mathbb{E} |P_n f - P_n' f|_X$$

$$\leq \mathbb{E} \|P_n - P_n'\|_{\mathcal{F}}, \text{ note that}$$

($\varepsilon_i (f(x_i) - f(x'_i)) = d (f(x_i) - f(x'_i))$)

$$= \mathbb{E} \|P_n^\varepsilon - P_n'^\varepsilon\|_{\mathcal{F}}$$

$$\leq 2 \mathbb{E} \|P_n^\varepsilon\|_{\mathcal{F}}. \quad \square$$

Maximum of Gaussian sequence

If $X \sim N(0, \sigma^2)$, the moment generating function of X :

$$\begin{aligned}\mathbb{E} e^{\lambda X} &= \int_{\mathbb{R}} e^{\lambda x} \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{x^2}{2\sigma^2}} dx \\ &= e^{\frac{\lambda^2\sigma^2}{2}} \int_{\mathbb{R}} \frac{1}{\sqrt{2\sigma^2\pi}} e^{-\frac{(x-\lambda\sigma)^2}{2\sigma^2}} dx \\ &= e^{\lambda^2\sigma^2/2}.\end{aligned}$$

Lemma 2 If $X_i \sim N(0, \sigma^2)$, then

$$\mathbb{E} \max_{i \leq n} X_i \leq \sqrt{2\sigma^2 \log n}.$$

Pf: $e^{\lambda \mathbb{E} \max_i X_i} \stackrel{\text{Jesen}}{\leq} \mathbb{E} e^{\lambda \max_i X_i}$

$$\leq n \mathbb{E} e^{\lambda X_i}, \forall \lambda > 0$$

$$\Rightarrow \mathbb{E} \max_i X_i \leq \log n / \lambda + \frac{\lambda \sigma^2}{2}, \forall \lambda > 0$$

$$\Rightarrow \mathbb{E} \max_i X_i \leq \sqrt{2\sigma^2 \log n}.$$

□

Sub-Gaussian variable

X is called a sub-Gaussian iff $\exists \gamma \sim N(0, \tau^2)$ s.t.

$$\mathbb{P}(|X| \geq t) \leq \mathbb{P}(|Y| \geq t).$$

Remark: ① X sub-Gaussian $\Leftrightarrow -X$ sub-Gaussian

② $\mathbb{P}(|Y| \geq t) = \int_{t/\tau}^{+\infty} e^{-\frac{x^2}{2\tau^2}} dx$
 $= o(t) e^{-\frac{t^2}{2}}$

$$\Rightarrow \exists C > 0 \text{ s.t. } \mathbb{P}(|X| \geq t) \leq C e^{-\frac{t^2}{2}}$$

 $=: 2e^{-\frac{t^2}{2C}}$

the minimal of C is determined by X ,
we mark that

$$\|X\|_4 = b, \quad (\text{Why?})$$

indeed, $\|\cdot\|_4$ is a norm of X , and

$$X \text{ sub-Gaussian} \Leftrightarrow \|X\|_4 < \infty \quad (\text{Why?})$$

③ If X sub-Gaussian, then $\exists b \geq 0$ s.t.

$$\mathbb{E} e^{\lambda X} \leq e^{\lambda^2 b^2 / 2} \text{ and } b \geq \|X\|_4.$$

Moreover,

$$\mathbb{E} e^{\lambda X} \leq e^{\lambda^2 b^2 / 2} \Leftrightarrow \mathbb{P}(|X| \geq t) \leq 2e^{-\frac{t^2}{2b^2}} \quad (\text{Why?})$$

Lemma 3 X_i sub-Gaussian, s.t.

$$\mathbb{E} e^{\lambda X_i} \leq e^{\lambda^2 \sigma^2 / 2}, \forall i \leq n.$$

then

$$\mathbb{E} \max_{i \leq n} |X_i| \leq \sqrt{\sigma^2 \log n}$$

Pf: $e^{\lambda} \mathbb{E} \max_{i \leq n} |X_i|$

Jesen

$$\leq \mathbb{E} e^{\lambda \max_{i \leq n} |X_i|}$$

$$= \mathbb{E} \max_i e^{\lambda |X_i|}$$

$$\leq n \mathbb{E} e^{\lambda |X_1|}$$

note that

$$\begin{aligned} \mathbb{E} e^{\lambda |X_1|} &\leq \mathbb{E}(e^{-\lambda X_1} + e^{\lambda X_1}) \\ &\leq 2e^{\lambda^2 \sigma^2 / 2} \end{aligned}$$

$$\Rightarrow e^{\lambda \mathbb{E} \max_i |X_i|} \leq 2n e^{\lambda^2 \sigma^2 / 2}$$

$$\Rightarrow \mathbb{E} \max_i |X_i| \leq \frac{\log_2 n}{\lambda} + \frac{\lambda 6^2}{2}, \quad \forall \lambda > 0$$

$$\begin{aligned} \Rightarrow \mathbb{E} \max_i |X_i| &\leq \sqrt{4 \log_2 n \cdot \frac{6^2}{2}} \\ &= \sqrt{2 \log_2 n} 6^2 \\ &\leq \sqrt{6^2 \log n} \end{aligned}$$

□

Sub-Gaussian process

$\{X_\theta, \theta \in T\}$ is called sub-Gaussian process iff

$$\mathbb{E} e^{\lambda c(X_\theta, -X_{\theta_2})} \leq e^{\frac{\lambda^2 (d(\theta, \theta_2))^2}{2}}$$

Remark: $\{X_\theta, \theta \in T\}$ sub-Gaussian

$$\Leftrightarrow \|X_\theta - X_{\theta_2}\|_4 \leq d(\theta, \theta_2).$$

Example | Gaussian / Rademacher complexity

let $X_\theta = \langle \theta, \varepsilon \rangle$, $\varepsilon \sim N(0, I_n)$,
then

$$X_\theta - X_{\theta_2} = \langle \theta - \theta_2, \varepsilon \rangle$$

$$\sim N(0, \|\theta - \theta_2\|^2)$$

$$\Rightarrow \mathbb{E} e^{\lambda \langle X_\theta - X_{\theta_2}, \varepsilon \rangle} = e^{\frac{\lambda^2 \|\theta - \theta_2\|^2}{2}}$$

If $\varepsilon_i \sim \text{Unif } \{-1, 1\}^n$,

$$\|X_\theta - X_{\theta_2}\|_4$$

$$= \left\| \sum_{i=1}^n (\theta_{1(i)} - \theta_{2(i)}) \varepsilon_i \right\|_4$$

$$\leq \sum_{i=1}^n |\theta_{1(i)} - \theta_{2(i)}| \|\varepsilon_i\|_4$$

$$\leq \|\theta - \theta_2\|_\infty \quad (\text{Why?})$$

If $\|\varepsilon_i\|_4 < 6$, then

$$\|X_\theta - X_{\theta_2}\| \leq \|\theta - \theta_2\|_\infty$$

$$\leq \|\theta - \theta_2\|,$$

$\|\cdot\|$ is an arbitrary norm of \mathbb{R}^n .

Example 2 Learning process

Recall

$P_n^\varepsilon \tilde{f} = \frac{1}{n} \sum_{i=1}^n \tilde{f}(x_i) \varepsilon_i$, assume that

$$R = \sup_{f \in \mathcal{F}} \sqrt{\mathbb{E} f(x_1)^2} < \infty.$$

Remark:

(Hoeffding) $\varepsilon_i \stackrel{iid}{\sim} \text{Unif}[-1, 1]$

$$\Pr \left(\left| \sum_{i=1}^n \alpha_i \varepsilon_i \right| > t \right) \leq 2e^{-\frac{t^2}{2\|\alpha\|^2}}$$

Let $X_f = \frac{1}{n} \sum_{i=1}^n \tilde{f}(x_i) \varepsilon_i$, then

$$\Pr(|X_f - X_g| > t | X)$$

$$= \Pr \left(\left| \frac{1}{n} \sum_{i=1}^n (\tilde{f}(x_i) - \tilde{g}(x_i)) \varepsilon_i \right| > t | X \right)$$

$$\leq 2e^{-\frac{nt^2}{2 \sum_{i=1}^n (\tilde{f}(x_i) - \tilde{g}(x_i))^2}} + t^2$$

$$\Rightarrow \Pr(|X_f| > t) \leq 2e^{-\frac{t^2}{2d(\tilde{f}, \tilde{g})}},$$

where $d(\tilde{f}, \tilde{g}) = \sqrt{\mathbb{E}(\tilde{f}(x) - \tilde{g}(x))^2}$.

Why

Upper bound by one-step discretization

If T is totally bounded, let T_8 be the δ -covering for T , then $\forall \theta \in T$,

$$\exists t \in T_8, d(\theta, t) \leq \delta,$$

$$X_\theta = X_\theta - X_t + X_t$$

$$\Rightarrow \mathbb{E} \sup_{\theta \in T} |X_\theta| \leq \mathbb{E} \sup_{d(\theta_1, \theta_2) \leq \delta} |X_{\theta_1} - X_{\theta_2}|$$

$$+ \mathbb{E} \max_{t \in T_8} |X_t|,$$

If $\delta \downarrow 0^+$,

$$\sup_{d(\theta_1, \theta_2) \leq \delta} |X_{\theta_1} - X_{\theta_2}| = o_p(1), \text{ if}$$

X_θ is sub-Gaussian process, (Why?)
which implies that if δ is small enough,
then

$\mathbb{E} \max_{t \in T_8} |X_t|$ is a tight upper bound
for $\mathbb{E} \sup_{\theta \in T} |X_\theta|$. Similarly, if

$$\mathbb{E} X_\theta = 0, \forall \theta \in T, \text{ then}$$

$$\mathbb{E} \sup_{\theta \in T} X_\theta = \mathbb{E} \sup_{\theta \in T} |X_\theta - X_{\theta_0}|$$

$$\leq \mathbb{E} \sup_{\theta_1, \theta_2 \in T} |X_{\theta_1} - X_{\theta_2}|$$

And

$$|X_{\theta_1} - X_{\theta_2}| \leq |X_{t_1} - X_{t_2}| + |X_{t_1} - X_{\theta_1}| + |X_{t_2} - X_{\theta_2}|, \quad t_1, t_2 \in T_8$$

$$\Rightarrow \mathbb{E} \sup_{\theta_1, \theta_2 \in T} |X_{\theta_1} - X_{\theta_2}| \leq$$

$$\mathbb{E} \sup_{\theta_1, \theta_2 \in T} |X_{\theta_1} - X_{\theta_2}| + \underbrace{\text{O}_{p(1)}}_{\text{O}_{p(1)}}$$

$$\mathbb{E} \sup_{t_1, t_2 \in T_8} |X_{t_1} - X_{t_2}|$$

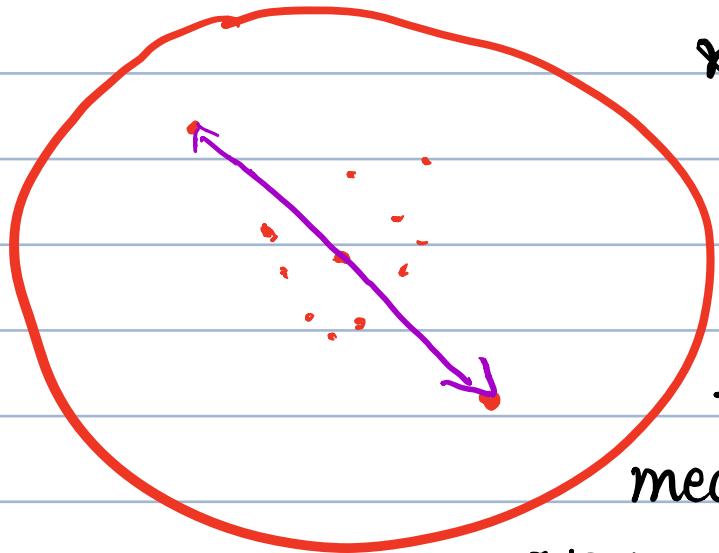
If $D = \text{diam}(T) < \infty$, then $\text{Lemma 3})$

$$\mathbb{E} \sup_{\theta_1, \theta_2 \in T} |X_{\theta_1} - X_{\theta_2}| \leq \text{O}_{p(1)} + \sqrt{D^2 \log N_8},$$

where $N_8 = N(8; T, d)$.

Chaining

Somewhat, $\sqrt{D^2 \log N_\delta}$ is not a tight bound for $\mathbb{E} \sup_{t_1, t_2 \in T_\delta} |X_{t_1} - X_{t_2}|$, and



note that
 $N_\delta \uparrow +\infty$ as $\delta \downarrow 0$,
but
 $\mathbb{E} \sup_{t_1, t_2 \in T_\delta} |X_{t_1} - X_{t_2}|$
tends to decrease, which
means that if δ is
very small, the upper
bound for $\mathbb{E} \sup_{\theta_1, \theta_2 \in T_\delta} |X_{\theta_1} - X_{\theta_2}|$ is
not suitable, i.e.,

one chooses δ so as to achieve the
optimal trade-off between two terms.

Assume the $|T| < \infty$ and $\text{diam}(T) := D < \infty$. Let $\delta_0 = D$,

$$\delta_k \downarrow 0, \quad k \rightarrow \infty.$$

Given $L > 0$, \bar{T}_k is the minimal δ_k -covering of T , $k = 0, \dots, L$.

Theorem 1 If $\{X_\theta, \theta \in T\}$ is zero-mean process and $\text{diam}(T) := D < \infty$, then

$$\mathbb{E} \sup_{\theta_1, \theta_2 \in T} |X_{\theta_1} - X_{\theta_2}| \leq J(T, d)$$

where $J(\delta; T) := \int_0^D \sqrt{\log N(\delta; T, d)} d\delta$.

Pf: $\forall \theta_1, \theta_2 \in T$, $T_{L+1} := T$, and

$$\mathbb{E} \max_{\theta_1, \theta_2 \in T} |X_{\theta_1} - X_{\theta_2}| \quad \begin{cases} \text{assume that} \\ |T| < \infty \end{cases}$$

$$\leq \mathbb{E} \max_{t_1 \in T_{L+1}, t_2 \in T_L, d(t_1, t_2) \leq \delta_L} |X_{t_1} - X_{t_2}|$$

$$+ \mathbb{E} \max_{t_1, t_2 \in T_L} |X_{t_1} - X_{t_2}|$$

$$\leq \sum_{i=1}^L \mathbb{E} \max_{t_1 \in T_{i+1}, t_2 \in T_i, d(t_1, t_2) \leq \delta_i} |X_{t_1} - X_{t_2}|$$

$$\leq \sum_{i=1}^L \sqrt{\delta_i^2 \log(N(\delta_{i+1}; T, d) N(\delta_i; T, d))}$$

$$\leq \sum_{i=1}^L \sqrt{\log N(\delta_{i+1}; T, d)} \delta_i$$

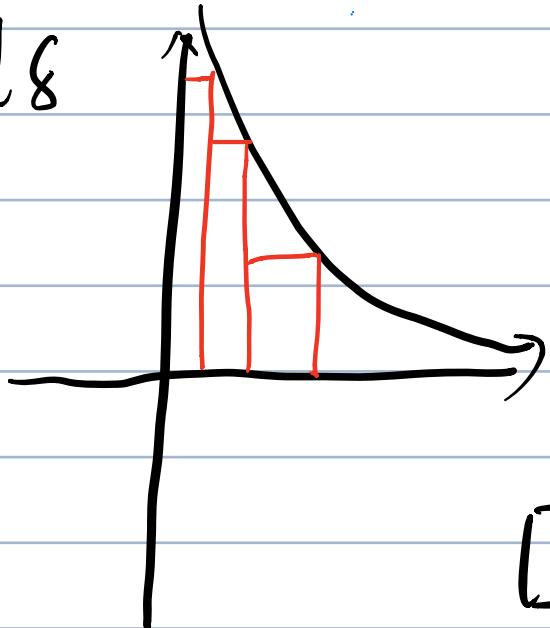
(Let $\delta_{i+1} = \delta_i / 2$)

$$\leq \sum_{i=1}^L \sqrt{\log N(\delta_{i+1}; T, d)} \delta_{i+1} / 2$$

$$\leq \int_0^D \sqrt{\log N(\delta; T, d)} d\delta$$

$$=: \underline{JCT, d}$$

Entropy Integral



□

Remark:

If $\{X_\theta, \theta \in T\}$ is the sub-Gaussian, $\text{diam}(T) < \infty$, then

$$\mathbb{E} \sup_{\theta \in T} X_\theta \leq JCT, d)$$

Example 1

$$\mathbb{E} \sup_{\theta \in \mathbb{B}} \langle \theta, \epsilon \rangle \leq \int_0^1 \sqrt{\log N(\delta; \mathbb{B}, 11 \cdot 11)} d\delta$$

$$\leq \int_0^1 \sqrt{\log (1 + \frac{2}{\delta})^P} d\delta$$

$$\leq \sqrt{P} \int_0^1 \sqrt{\log (1 + \frac{2}{\delta})} d\delta \leq \sqrt{P}$$

Example 2

If \mathcal{F} is the class of functions on X , and $\theta \in \mathcal{F}$, $D := \text{diam}(\mathcal{F})$

Remark: If $\exists \theta_0 \in \mathcal{T}$ s.t.

$$X_{\theta_0} = 0 \text{ a.s.}$$

Then

$$\begin{aligned} \mathbb{E} \sup_{\theta \in \mathcal{T}} |X_\theta| &= \mathbb{E} \sup_{\theta \in \mathcal{T}} |X_\theta - X_{\theta_0}| \\ &\leq \mathbb{E} \sup_{\theta_1, \theta_2} |X_{\theta_1} - X_{\theta_2}| \\ \Rightarrow \mathbb{E} \sup_{\theta \in \mathcal{T}} |X_\theta| &\leq J(T, d) \end{aligned}$$

Note that $\phi: f \mapsto (f-T)^2$ is continuous since

$$\begin{aligned} P((f-T)^2 - (g-T)^2)^2 &= P(f-g)^2 (f+g-T)^2 \\ &\leq P(f-g)^2 \end{aligned}$$

Then $\tilde{\mathcal{F}}$ is also bounded and $\exists C > 0$ s.t.

$$N(\delta; \tilde{\mathcal{F}}, d) \leq NCC(\delta; \mathcal{F}, d)$$

According to Theorem 1:

$$R(\hat{f}_n) - R(f^*) \leq \mathbb{E} \|P_n^\varepsilon\|_{\tilde{\mathcal{F}}}$$

$$= \frac{\mathbb{E} \|P_n^\varepsilon\|_{\tilde{\mathcal{F}}}}{\sqrt{n}}$$

$$\leq \frac{\int_0^{+\infty} \sqrt{\log N(\delta; \bar{f}, d)} d\delta}{\sqrt{n}}$$

$$\leq \frac{\int_0^{+\infty} \sqrt{\log N(\delta; \bar{f}, d)} d\delta}{\sqrt{n}}$$

If $\text{VC}(\mathcal{H}) =: v < \infty$, then

$$R(\hat{f}_n) - R(f^*) \leq \frac{\int_0^{+\infty} \sqrt{Cv \log \frac{1}{\delta}} d\delta}{\sqrt{n}}$$

$$\asymp \sqrt{\frac{v}{n}}$$

Example 3

$$F_n(x) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}(X_i \leq x)$$

$$= P f_x, \text{ where } f_x(y) = \mathbb{1}(y \leq x)$$

Then $\mathbb{E} \|F_n - F\|_\infty \leq \sqrt{\frac{1}{n}}$

$$\Rightarrow \|F_n - F\|_\infty = o_p(1).$$