

汎化誤差解析と Rademacher 複雑度

\mathcal{F} = モデル ($f: x \in X \mapsto y \in \mathbb{R}$ なり関数の集合)

例: ニューラルネットワーク, 再生核ヒルベルト空間

$l: \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$ を損失関数とする.

例: $l(y, f(x)) = (y - f(x))^2$: 二乗損失 ($y \in \mathbb{R}$)

$l(y, f(x)) = \log(1 + \exp(-yf(x)))$: ロジスティック損失 ($y \in \{\pm 1\}$)

期待損失, 予測誤差: $L(f) = E_{x, y} [l(y, f(x))]$

経験損失, 訓練誤差: $\hat{L}(f) = \frac{1}{n} \sum_{i=1}^n l(y_i, f(x_i))$

$\mathcal{Z} = \{(x_i, y_i) \mid i=1, \dots, n\}$ は i.i.d. \mathcal{Z} の分布 \mathbb{P} に従う.

$Pg = E_{x, y} [g(x, y)]$, $P_n g = \frac{1}{n} \sum_{i=1}^n g(x_i, y_i)$ と書く.

($L(f) = P(l \circ f)$, $\hat{L}(f) = P_n(l \circ f)$ とある)

$\mathcal{Z} = (x, y)$ の分布と
 $\mathcal{Z}_n = (x_i, y_i)$ の分布は別々
にある.

$f^\circ = \operatorname{argmin}_{f: \text{可測関数}} L(f)$, $f^* = \operatorname{argmin}_{f \in \mathcal{F}} L(f)$

\hat{f} : 何らかの推定量 (これが存在するを仮定)

- Excess risk (残余誤差, 余剰誤差)
 $L(\hat{f}) - L(f^\circ)$ 又は $L(\hat{f}) - L(f^*)$
- Generalization gap (汎化ギャップ)
 $L(\hat{f}) - \hat{L}(\hat{f})$

← これを「汎化誤差」と言うことができる.

Excess risk

$$L(\hat{f}) - L(f^\circ) = \underbrace{L(\hat{f}) - \hat{L}(\hat{f})}_{\text{汎化ギャップ}} + \underbrace{\hat{L}(\hat{f}) - \hat{L}(f^*)}_{\substack{f \text{ が訓練データに適合する程度} \\ \epsilon + \text{公差} < \epsilon \text{ (2004年)} \\ 0 \text{ 以下}}} + \underbrace{\hat{L}(f^*) - L(f^*)}_{\substack{\text{大数の法則} \\ O_p(1/n)}} + \underbrace{L(f^*) - L(f^\circ)}_{\substack{\text{モデルバイアス}}}$$

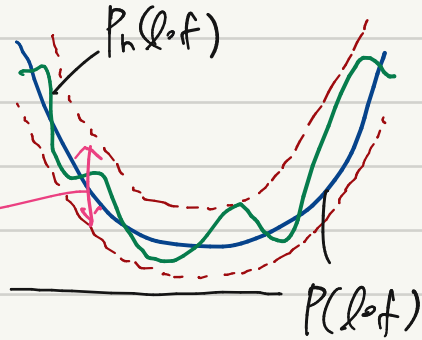
$$\hat{f} = \operatorname{argmin}_{f \in \mathcal{F}} \hat{L}(f) \text{ ならば}$$

$$L(\hat{f}) - L(f^*) = \underbrace{(P - P_n)(\ell \circ \hat{f} - \ell \circ f^*)}_{\text{バリエーション}} + \underbrace{L(f^*) - L(f^0)}_{\text{(モデル)バイアス}}$$

• 汎化リスクの Bound

$$L(\hat{f}) - \hat{L}(\hat{f}) = (P - P_n)(\ell \circ \hat{f})$$

$$\leq \sup_{f \in \mathcal{F}} (P - P_n)(\ell \circ f)$$



⊛ \hat{f} はデータに依存するため、単純に大数の法則は使えない。

$$G := \{g = \ell \circ f \mid f \in \mathcal{F}\} \text{ とおけば}$$

$\sup_{g \in G} (P - P_n)g$ は 最適化問題 に存在 \rightarrow 経験過程

* 以後 $\sup_{g \in G} P_n g$ と $\sup_{g \in G} \frac{1}{n} \sum_{i=1}^n \epsilon_i g(x_i)$ は

可測 である。可測でない場合は以下の議論は成り立たない。

(Fubini の定理 が成り立たない, $E_x E_y [(\cdot)^*] \neq E_y E_x [(\cdot)^*]$)

- van der Vaart & Wellner: weak convergence and empirical processes. (1996) の Sec 2.3.

- Dudley: Uniform central limit theorems.

\rightarrow Suslin image admissible

Rademacher complexity

以下. $\|g\|_\infty \leq M$ ($\forall g \in G$) とする.

Def (Rademacher 複雑度)

σ_n ($i=1, \dots, n$): Rademacher 変数 (i.i.d.)

$$P(\sigma_n=1) = P(\sigma_n=-1) = \frac{1}{2}$$

とある.

$$- \hat{R}_n(G) := E_{\{\sigma_i\}_{i=1}^n} \left[\sup_{g \in G} \left| \frac{1}{n} \sum_{i=1}^n \sigma_i g(z_i) \right| \mid z_1, \dots, z_n \right]$$

$$- R_n(G) := E_{\{z_i\}_{i=1}^n} [\hat{R}_n(G)]$$

σ_i 以外の z_i は n/i-1 だけある.

Thm (対称性)

← $\{z_i\}_{i=1}^n$ と G の最大値を z_i と置く.

$$E_{z_i} [\sup_{g \in G} | (P - P_n)g |] \leq 2 R_n(G)$$

(直接 sup を評価するのは, $R_n(G)$ の方が評価が楽)

(Proof)

$$LHS = E_{z_i} \left[\sup_{g \in G} \left| E_{z'_i} \left[\frac{1}{n} \sum_{i=1}^n g(z'_i) \right] - \frac{1}{n} \sum_{i=1}^n g(z_i) \right| \right]$$

$$\leq E_{z_i, z'_i} \left[\sup_{g \in G} \left| \frac{1}{n} \sum_{i=1}^n (g(z'_i) - g(z_i)) \right| \right] \quad (\text{Jensen 不等式})$$

$$= E_{\sigma_i} E_{z_i, z'_i} \left[\sup_{g \in G} \left| \frac{1}{n} \sum_{i=1}^n \sigma_i (g(z'_i) - g(z_i)) \right| \right]$$

(z'_i と z_i は同分布)

$$\leq E_{\sigma_i} E_{z'_i} \left[\sup_{g \in G} \left| \frac{1}{n} \sum_{i=1}^n \sigma_i g(z'_i) \right| \right]$$

$$+ E_{\sigma_i} E_{z_i} \left[\sup_{g \in G} \left| \frac{1}{n} \sum_{i=1}^n \sigma_i g(z_i) \right| \right]$$

$$= 2 R_n(G) \quad (\because \text{Fubini の定理})$$

↑ 2.2.2 可測性を使う.

可積分性は有界性から OK.

Thm

$$\left[P \left(\sup_{g \in G} |(P - P_n)g| \geq 2R_n(G) + M \sqrt{\frac{2 \log(\frac{1}{\delta})}{M}} \right) \leq \delta \right. \\ \left. (\forall \delta > 0) \right]$$

Proof

LEM (Mc Diarid の不等式)

X_1, \dots, X_n を \mathcal{X} 上の独立な (同-とは限らない) 確率変数とす。

$f: \mathcal{X}^n \rightarrow \mathbb{R}$ を可測関数とし。

$$|f(x_1, \dots, x_n) - f(x_1, \dots, x'_i, \dots, x_n)| \leq C_i$$

かつ $\forall x_1, \dots, x_n, x'_i \in \mathcal{X}$ に対す (2 式) 1 式の時。

$$P \left(f(X_1, \dots, X_n) - E[f(X_1, \dots, X_n)] \geq \epsilon \right) \\ \leq \exp \left(- \frac{2\epsilon^2}{\sum_{i=1}^n C_i^2} \right) \quad (\forall \epsilon > 0) \quad //$$

今 $\|g\|_\infty \leq M$ ($\forall g \in G$) とす

$$f(z_1, \dots, z_n) = \sup_{g \in G} \left| \frac{1}{n} \sum_{i=1}^n g(z_i) - P g \right|$$

とす。 $f(z_1, \dots, z_n) - f(z_1, \dots, z'_i, \dots, z_n)$

$$= \sup_{g \in G} \left| \frac{1}{n} \sum_{j \neq i} g(z_j) + \frac{g(z_i)}{n} - \frac{g(z'_i)}{n} + \frac{g(z'_i)}{n} - P g \right|$$

$$= \sup_{g \in G} \left| \frac{1}{n} \sum_{j \neq i} g(z_j) + \frac{g(z'_i)}{n} - P g \right|$$

$$\leq \sup_{g \in G} \left| \frac{1}{n} \sum_{j \neq i} g(z_j) + \frac{g(z'_i)}{n} - P g \right| + \sup_{g \in G} \left| \frac{g(z_i)}{n} - \frac{g(z'_i)}{n} \right| \\ = \sup_{g \in G} \left| \frac{1}{n} \sum_{j \neq i} g(z_j) + \frac{g(z'_i)}{n} - P g \right|$$

$$\leq \frac{2M}{n}$$

同様の理由で $f(z_1, \dots, z'_i, \dots, z_n) - f(z_1, \dots, z_n) \leq \frac{2M}{n}$ かつ
絶対値 $\leq \frac{2M}{n}$ である。

あとは McDiarmid の不等式 と 対称化 を 使う

Prop (Rademacher 複雑度の性質)

(1) $G \subset G'$ ならば $R_n(G) \leq R_n(G')$

(2) $\phi_i: \mathbb{R} \rightarrow \mathbb{R}$ かつ $|\phi_i(x) - \phi_i(y)| \leq D|x-y|$ ならば

$$E \sigma_i \left[\sup_{g \in G} \left| \frac{1}{n} \sum_{i=1}^n \sigma_i \phi_i(g(z_i)) \right| \right] \leq 2D \hat{R}_n(G)$$

(contraction inequality)

↑ $(R_n(G)$ の定義に 絶対値 を 使わない場合、左辺の $E(\dots)$ の中から 絶対値 を 取り除く。 $2D$ を D に 変える。

(3) G の 凸包 E

$$\text{conv}(G) = \left\{ \sum_{i=1}^m \lambda_i g_i \mid \sum_{i=1}^m \lambda_i = 1, \lambda_i \geq 0, g_i \in G, m=1,2,\dots \right\}$$

と すると

$$R_n(G) = R_n(\text{conv}(G))$$

(4) $R_n(c \cdot G) = |c| R_n(G) \quad (c \in \mathbb{R})$

Ex. (1) 線形形式 $x \cdot \omega \in \mathbb{R}^d$ と すると

$$\begin{cases} - F = \{ f(x) = x^T \omega \mid \|\omega\|_2 \leq 1 \} \\ - E[\|x_i\|^2] \leq 1 \end{cases}$$

$$R_n(F) = E_{z_i, \sigma_i} \left[\sup_{\|\omega\| \leq 1} \left| \frac{1}{n} \sum_{i=1}^n \sigma_i \omega^T x_i \right| \right]$$

$$\leq E_{z_i, \sigma_i} \left[\sup_{\|\omega\| \leq 1} \|\omega\| \cdot \left\| \frac{1}{n} \sum_{i=1}^n \sigma_i x_i \right\| \right]$$

$$\leq \sqrt{E_{z_i, \sigma_i} \left[\frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n \sigma_i \sigma_j x_i^T x_j \right]}$$

$$= \sqrt{\frac{1}{n^2} \sum_{i=1}^n E[\|x_i\|^2]} = \sqrt{\frac{1}{n}}$$

(2) F が有限集合の時

$$F = \{f_1, \dots, f_N\}, \quad \|f_i\|_\infty \leq R \quad \forall i$$

$$\hat{R}_n(F) \leq R \sqrt{\frac{2 \log(2N)}{n}}$$

$R_n = 1$ から $\lambda \rightarrow 2$ まで
 n は $2N \in N$ とする

(Massart の定理)



Hoeffding の不等式より
 示せる。

N は t と $u \sqrt{\log(N)}$ 2-つの交点になる

(Massart の定理の証明)

$t > 0$ に対し.

$$\exp(t E_G \left[\max_{f \in F} \left| \frac{1}{n} \sum_{i=1}^n \sigma_i f(z_i) \right| \right]) \leq E_G \left[\exp \left(t \max_{f \in F} \left| \frac{1}{n} \sum_{i=1}^n \sigma_i f(z_i) \right| \right) \right]$$

$$\leq E_G \left[\sum_{f \in F} \exp \left(t \frac{1}{n} \sum_{i=1}^n \sigma_i f(z_i) \right) + \sum_{f \in F} \exp \left(-t \frac{1}{n} \sum_{i=1}^n \sigma_i f(z_i) \right) \right]$$

$$= \sum_{f \in F} 2 E_G \left[\exp \left(t \frac{1}{n} \sum_{i=1}^n \sigma_i f(z_i) \right) \right]$$

$$= \sum_{f \in F} 2 E_G \left[\prod_{i=1}^n \exp \left(\frac{t}{n} \cdot \sigma_i f(z_i) \right) \right]$$

$$= 2 \sum_{f \in F} \prod_{i=1}^n E_{\sigma_i} \left[\exp \left(\frac{t}{n} \sigma_i f(z_i) \right) \right]$$

$$\leq 2 \sum_{f \in F} \prod_{i=1}^n \exp \left(\frac{\left(\frac{t}{n}\right)^2 R^2}{2} \right)$$

(\because Hoeffding の補題)

$$\leq 2N \cdot \exp \left(\frac{t^2 R^2}{2n} \right)$$

有限 r.v. は sub-Gaussian

(両辺 \log と $\rightarrow 2, \frac{1}{t}$ をかけると)

$$\Rightarrow \hat{R}_n(F) \leq \frac{\log(2N)}{t} + \frac{tR^2}{2n} \quad (\forall t > 0)$$

$$\Rightarrow t = \sqrt{\frac{2n \log(2N)}{R^2}} \quad \text{右辺は最小}$$

$$\Rightarrow \hat{R}_n(F) \leq R \cdot \sqrt{\frac{2 \log(2N)}{n}}$$

* R は $\sigma^2 = \max_{f \in F} \frac{1}{n} \sum_{i=1}^n f(z_i)^2$ にあてはめるとよい (証明は省略)

Contraction ineq. f'). 損失関数 ρ は Lipschitz 連続なら
汎化誤差 \rightarrow $\mathbb{E} R_n(F)$ に δ だけ δ だけ δ .

$$|\rho(y, f) - \rho(y, f')| \leq L |f - f'| \quad \text{と好}.$$

$$\begin{aligned} \hat{L}(f) - L(f) &= (P_n - P)(\rho \circ f) \\ &\leq \sup_{g \in \mathcal{G}} (P_n - P)g \\ &\leq 2 R_n(\mathcal{G}) + M \sqrt{\frac{2 \log(\frac{1}{\delta})}{n}} \quad \text{with prob. } 1 - \delta. \\ &\leq 4 L R_n(F) + M \sqrt{\frac{2 \log(\frac{1}{\delta})}{n}} \end{aligned}$$

● 深層 2-3 レットワークの Rademacher 複雑度

ここで η の特殊な条件が成り立つ。深層 NN の Rademacher 複雑度を導出しよう。

$$\mathcal{F} = \left\{ w^{(L)} \eta \left(w^{(L-1)} \eta \left(\dots \eta \left(w^{(1)} x \right) \dots \right) \right) \mid w^{(l)} \in \mathbb{R}^{m_l \times m_{l-1}}, \underbrace{\|w^{(l)}\|_{\infty, 1} \leq B_l}_{\leftarrow m_L=1 \text{ かつ } \eta \text{ の条件が強い}} \right\}$$

ここで $\eta(x) = \max(x, 0)$ (ReLU) (ReLU)

$$\|w^{(l)}\|_{\infty, 1} := \max_{1 \leq j \leq m_l} \|w_{j, :}^{(l)}\|_1$$

とある。

$$\mathcal{F}_i = \left\{ f(x) = \sum_{j=1}^{m_{i-1}} \omega_j \eta(f_j(x)) \mid f_j \in \mathcal{F}_{i-1}, \|\omega\|_1 \leq B_i \right\}$$

$$\mathcal{F}_1 = \left\{ f(x) = w^T x \mid \|w\|_1 \leq B_1 \right\}$$

とある。 $\mathcal{F}_L = \mathcal{F}$ である。

Lem

$$\mathbb{R}_n(\mathcal{F}_i) \leq 2 B_i \mathbb{R}_n(\mathcal{F}_{i-1})$$

Proof

$$\mathbb{R}_n(\mathcal{F}_i) = \mathbb{E}_{\sigma_i} \left[\sup_{\substack{\|w\|_1 \leq B_i \\ f_j \in \mathcal{F}_{i-1}}} \left| \frac{1}{n} \sum_{k=1}^n \sigma_k \sum_{j=1}^{m_{i-1}} \omega_j \eta(f_j(x_k)) \right| \right]$$

$$= \mathbb{E}_{\sigma_i} \left[\sup_{\omega, f_j} \left| \sum_{j=1}^{m_{i-1}} \omega_j \frac{1}{n} \sum_{k=1}^n \sigma_k \eta(f_j(x_k)) \right| \right]$$

$$\leq \mathbb{E}_{\sigma_i} \left[\sup_{\omega, f_j} \underbrace{\|w\|_1}_{\leq B_i} \max_{1 \leq j \leq m_{i-1}} \left| \frac{1}{n} \sum_{k=1}^n \sigma_k \eta(f_j(x_k)) \right| \right]$$

$$\leq B_i \mathbb{E}_{\sigma_i} \left[\sup_{f \in \mathcal{F}_{i-1}} \left| \frac{1}{n} \sum_{k=1}^n \sigma_k \eta(f(x_k)) \right| \right]$$

$\in \{ \eta \circ f \mid f \in \mathcal{F}_{i-1} \}$ の Rademacher 複雑度

$$\leq 2B_i R_n(\mathcal{F}_{i-1}) \quad (\because \text{contraction inequality.})$$

(γ は $1-\gamma \cdot 2 \rightarrow$ 連続)

Lem

$$\|x\|_\infty \leq 1 \quad (\text{a.s.}) \quad \text{for } i$$

$$R_n(\mathcal{F}_i) \leq \sqrt{\frac{2 \log(2d)}{n}}$$

(Proof)

$$\begin{aligned} \tau_{\mathcal{F}_i} &= E_{\sigma_i} \left[\sup_{\|\omega\|_1 \leq B_1} \|\omega\|_1 \max_{1 \leq j \leq d} \left| \frac{1}{n} \sum_{k=1}^n \sigma_k \chi_{k,j} \right| \right] \\ &= B_1 E_{\sigma_i} \left[\max_{1 \leq j \leq d} \left| \frac{1}{n} \sum_{k=1}^n \sigma_k \chi_{k,j} \right| \right] \\ &\leq B_1 \sqrt{\frac{2 \log(2d)}{n}} \quad (\because \text{Massart の定理}) \end{aligned}$$

Thm (DNN の Rademacher 複雑度)

$$\|x\|_\infty \leq 1 \quad (\text{a.s.}) \quad \text{for } i$$

$$R_n(\mathcal{F}) \leq 2^L \cdot \prod_{i=1}^L B_i \sqrt{\frac{2 \log(2d)}{n}}$$

→ 各層のノード数の積が制御されるため、汎化性能はネットワークのサイズから出てこない。

この以外 (2枚様々のバウンド) が導き出される。
(2^L 程度のバウンド) になる。

Rademacher 複雑度の土着: カバリー・メジャー, エンタルピー

Def (Covering Number, 被覆数)

\mathcal{F} : 関数の集合

d : \mathcal{F} の距離

$$N(\mathcal{F}, \varepsilon, d) := \min \{ n \geq 1 \mid \exists f_1, \dots, f_n \in \mathcal{F} \text{ s.t. } \forall f \in \mathcal{F} \exists j \in \{1, \dots, n\} \text{ s.t. } d(f, f_j) \leq \varepsilon \}$$

Def (エンタルピー数)

$$e_n(\mathcal{F}, d) := \inf \{ \varepsilon > 0 \mid N(\mathcal{F}, \varepsilon, d) \leq 2^{\varepsilon n} \}$$

n 点 $x_1, \dots, x_n \in X$ に対し,

$$\|f\|_{n,p} := \left(\frac{1}{n} \sum_{i=1}^n |f(x_i)|^p \right)^{\frac{1}{p}} \quad (1 \leq p < \infty)$$

とせよ.

Thm

$\forall f \in \mathcal{F}$ が $\|f\|_\infty \leq R$ を満たすから、ある $x_1, \dots, x_n \in \mathbb{R}^d$ をとると

$$R_n(\mathcal{F}) \leq \inf_{\alpha \geq 0} \left\{ \alpha + R \sqrt{\frac{2 \log(2N(\alpha, \mathcal{F}, \|\cdot\|_{n,1}))}{n}} \right\}$$

Proof

$f \in \mathcal{F}$ に対し $v[f]$ を f に一番近い ε -被覆の点とせよ.

$$R_n(\mathcal{F}) = E_\sigma \left[\sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n \sigma_i f(x_i) \right| \right]$$

(大は \mathcal{F} 略)

$$= E_\sigma \left[\sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n \sigma_i (f(x_i) - v[f](x_i)) + \frac{1}{n} \sum_{i=1}^n \sigma_i v[f](x_i) \right| \right]$$

$$\leq E_\sigma \left[\sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n \sigma_i (f(x_i) - v[f](x_i)) \right| \right] + E_\sigma \left[\sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n \sigma_i v[f](x_i) \right| \right]$$

$$\leq \sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n |f(x_i) - v[f](x_i)|$$

$$+ E_\sigma \left[\max_{1 \leq j \leq N(\mathcal{F}, \alpha, \|\cdot\|_{n,1})} \left| \frac{1}{n} \sum_{i=1}^n \sigma_i f_j(x_i) \right| \right]$$

← ε -被覆の
j番目の元

$$\leq \alpha + R \sqrt{\frac{2 \log(2N)}{n}} \quad (\because \text{Massart の定理})$$

↑
被覆の定義

- $p \leq q$ なら $\|\cdot\|_{n,p} \leq \|\cdot\|_{n,q}$ であるから
 $N(d, \mathcal{F}, \|\cdot\|_{n,q}) \leq N(d, \mathcal{F}, \|\cdot\|_{n,p})$ であるので、右辺の被覆数は
 $N(d, \mathcal{F}, \|\cdot\|_{n,p})$ ($p > 1$) であるからである。 ($p=2$ は \pm の場合)

Thm (Dudley の χ^2 法)

ある $x_1, \dots, x_n \in \mathbb{R}^d$ の組をとり、 $\hat{\rho} := \sup_{f \in \mathcal{F}} \|f\|_{n,2} < \infty$ とおくと

$$\hat{R}_n(\mathcal{F}) \leq \inf_{d > 0} \left\{ 4d + \frac{12}{\sqrt{n}} \int_d^{\hat{\rho}} \sqrt{\log 2 N(\mathcal{F}, \varepsilon, \|\cdot\|_{n,2})} d\varepsilon \right\}$$

$$\hat{R}_n(\mathcal{F}) \leq \frac{2}{\sqrt{n}} \left(\sum_{i=1}^{\infty} 2^{i/2} e_{2^i}(\mathcal{F} \cup \{0\}, \|\cdot\|_{n,2}) + \sup_{f \in \mathcal{F}} \|f\|_{n,2} \right)$$

※ \mathcal{F} を有限個の代表点で近似し、その近似精度を粗いものから
 密なものまで、順番に中々という操作に対応。

Proof (エントロピー数の方をけずる)

$$T = \{ (f(x_1), \dots, f(x_n)) \mid f \in \mathcal{F} \} \cup \{ (0, \dots, 0) \} \subseteq \mathbb{R}^n$$

$$d(t, t') := \sqrt{\frac{1}{n} \sum_{i=1}^n (t_i - t'_i)^2} \quad \text{とおく。 } (T, d) \text{ は可公距離空間。}$$

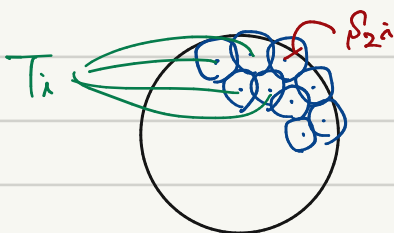
以後、 $\forall \varepsilon > 0$ に対し $e_n(T, \varepsilon)$ は有限であること一般性を失わずに
 特例として T は有界とし、 $n \rightarrow \infty$ として $e_n(T, \varepsilon) \rightarrow 0$ とする。

任意の $\varepsilon > 0$ に対し

$$\rho_1 := \sup_{t \in T} d(t, 0)$$

$$\rho_{2^i} := (1/2)^i e_{2^i}(T, d) \quad (i=1, 2, \dots)$$

とすると、 $\forall \varepsilon > 0$ に対し、 $\rho_{2^i} - \rho_{2^{i+1}} \leq \varepsilon$ となる i は
 $2^{i-1} \leq \rho_{2^i} \leq 2^i \rho_1$ であるから、 $2^{i-1} \leq \rho_1 / \varepsilon$ となる i は
 $i \leq \log_2(\rho_1 / \varepsilon) + 1$ であるから、 i は有限である。



$$\bigcup_{t \in T_i} B_d(t, \rho_{2^i}) \supseteq T \text{ となる}$$

(エントロピー数の
 定義より存在)

$T_0 = \{0\}$ とする。 好きと, $e_j(T, d) \rightarrow 0$ ($j \rightarrow \infty$) する。 $\mathcal{U} := \bigcup_{j=0}^{\infty} T_j$ となる。

$$\sup_{t \in \mathcal{U}} \frac{1}{n} \sum_{i=1}^n \sigma_i t_i = \sup_{t \in T} \frac{1}{n} \sum_{i=1}^n \sigma_i t_i$$

が成り立つ。 (\mathcal{U} は T の ϵ -網) として, T の代わり \mathcal{U} を考えればよい。

$j \geq 1$, $t \in T$ となる。 $\pi_j(t) \in T_j$ となる。 $d(t, \pi_j(t)) \leq \rho_{2^j}$ なるもの \mathcal{U} に取った。 好きと, $t \in T_j$ なら, $\pi_j(t) = t$ となる。

$j \geq 1$ として固定して, γ_n ($n \leq j$) として

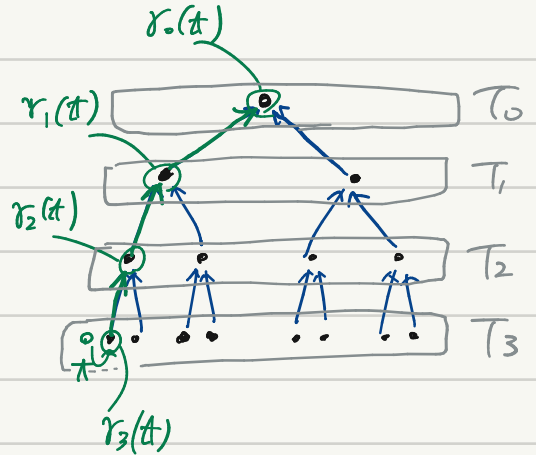
$$\gamma_j(t) = \pi_j(t)$$

$$\gamma_{n-1}(t) = \pi_{n-1} \circ \gamma_n(t) \quad (n=1, \dots, j)$$

とする。

当然 $d(\gamma_n(t), \gamma_{n-1}(t)) \leq \rho_{2^{n-1}}$ となる。

$$\text{今, } h(t) := \frac{1}{n} \sum_{i=1}^n \sigma_i t_i \quad (t \in T)$$



とすると, $t \in T_j$ となる

$$|h(t)| = |h(t) - 0|$$

$$= \left| \sum_{i=1}^j (h(\gamma_i(t)) - h(\gamma_{i-1}(t))) \right| \leq \sum_{i=1}^j \max_{t \in T_i} |h(t) - h(\pi_{i-1}(t))|$$

となる。 として, $\mathcal{U}_k := T_0 \cup \dots \cup T_k$ となる

$$\max_{t \in \mathcal{U}_k} |h(t)| \leq \sum_{j=1}^k \max_{t \in T_j} |h(t) - h(\pi_{j-1}(t))|$$

が成り立つ。 好きと, 有限集合の Rademacher 複雑度より

$$\begin{aligned} E_{\sigma_i} \left[\max_{t \in \mathcal{U}_k} |h(t)| \right] &\leq \sum_{j=1}^k E_{\sigma_i} \left[\max_{t \in T_j} |h(t) - h(\pi_{j-1}(t))| \right] \\ &\leq \sum_{j=1}^k \sqrt{\frac{2 \log(2 \cdot 2^{2^j} - 1)}{n}} \cdot \rho_{2^{j-1}} \end{aligned}$$

$\mathcal{U} = \bigcup_{j=0}^{\infty} T_j$ となる。 $k \rightarrow \infty$ となる。 単調増大性より

$$E_{\sigma_i} \left[\sup_{t \in \mathcal{U}} h(t) \right] \leq \sum_{j=0}^{\infty} \sqrt{\frac{4 \cdot 2^j}{n}} \rho_{2^j} \quad \text{ここで } \rho_1 = \sup_t d(t, 0)$$

(= 任意の t に対する)

① 深層 NN のカーネル化

スパース NN

$$\mathcal{F}_{L,S,B} := \left\{ (W^{(L)} \circ b^{(L)}) \circ \eta(W^{(L-1)} \circ b^{(L-1)}) \circ \dots \circ \eta(W^{(1)} x) \mid \right.$$

$$\sum_{\ell=1}^L \|W^{(\ell)}\|_0 + \sum_{\ell=1}^L \|b^{(\ell)}\|_0 \leq S,$$

$$\max_{\ell} \max \{ \|W^{(\ell)}\|_{\infty}, \|b^{(\ell)}\|_{\infty} \} \leq B,$$

$$\left. W^{(\ell)} \in \mathbb{R}^{m_{\ell} \times m_{\ell-1}}, b_{\ell} \in \mathbb{R}^{m_{\ell}} \right\}$$

(たとえば $\|W\|_0$ = W の非ゼロ要素の数, $\|W\|_{\infty} = \max_{i,j} |W_{ij}|$)

$$\mathcal{F}_{L,S,B} := \left\{ f(x) = \max \left\{ -R, \min \left[R, f(x) \right] \right\} \mid f \in \mathcal{F}_{L,S,B} \right\}$$

↑
clipping $\pm R$ の間におさまる.

Prop (Schmidt-Hieber, 2017)

$$\max_{\ell} m_{\ell} \leq W \text{ とおくと}$$

$$\log N(\epsilon, \mathcal{F}_{L,S,B}, \|\cdot\|_{\infty}) \leq 2SL \log(L(BWL)W \cdot \epsilon^{-1}) //$$

$$\Rightarrow \hat{R}_n(\mathcal{F}_{L,S,B}) \leq 4R \sqrt{\frac{SL \log(L(BWL)W \cdot n^{-1})}{n}}$$

$$(L, W, B = O(n) \text{ だと} \rightarrow O\left(\sqrt{\frac{SL \log(n)}{n}}\right))$$

より大抵 (2層 スパース NN) 汎化できる)

⑫ Fast learning rate と局所 Rademacher 複雑度

- 推定量 \hat{f} は f^* の "近く" に u する \rightarrow f^* の "局所的" 存在性 (用 u する)
- \rightarrow 結果的に $\frac{1}{n}$ 利速 u し f^* が示せる。

仮定

- L は 強凸: ある $\alpha > 0$ が存在し、 $\forall f \in \mathcal{F}$ に対し $\frac{\alpha}{2} \|f - f^*\|_{L^2(P_X)}^2 \leq L(f) - L(f^*)$
- L は 1- η フォーム連続, $\|L \cdot f\|_{\infty} \leq M$ ($\forall f \in \mathcal{F}$)

例 $l(y, f) = (y - f)^2$ とき

$\|f\|_{\infty} \leq \frac{1}{4}$ ($\forall f \in \mathcal{F}$), $\|f^*\|_{\infty} \leq \frac{1}{4}$, $|y - f^*(x)| \leq \frac{1}{4}$ (a.s.)
 なるが成り立つ。

← (大抵 L は l の 2 乗なので、
 適当にスเกลールすれば一般の大 L に対しても同様のことが成り立つ)

$$G_r := \{g = L \cdot f - L \cdot f^* \mid P g = L(f) - L(f^*) \leq r, f \in \mathcal{F}\}$$

$$G := \{g = L \cdot f - L \cdot f^* \mid f \in \mathcal{F}\}$$

と持。

$R_n(G_r)$ を 局所 Rademacher 複雑度 と書く。

* L フォーム連続性と強凸性より

$$\frac{\alpha}{2} \|g\|_{L^2(P_X)}^2 \leq \frac{\alpha}{2} \|f - f^*\|_{L^2(P_X)}^2 \leq P g$$

たの α . $P g \leq r$ ならば

$$\left\{ \begin{aligned} \|g\|_{L^2(P_X)} &\leq \sqrt{\frac{2r}{\alpha}} \\ \|f - f^*\|_{L^2(P_X)} &\leq \frac{2r}{\alpha} \end{aligned} \right.$$

Excess risk $P(f - f^*)^2$ が小になる \rightarrow f は f^* に 近 u

Thm (Peeling)

ある $\phi: [0, \infty) \rightarrow [0, \infty)$ なる関数が存在し、ある $\hat{r}^* > 0$ に対し
 $\forall r > \hat{r}^*$ に於て

$$\phi(4r) \leq 2\phi(r)$$

$$R_n(G_r) \leq \phi(r)$$

が成り立つならば、 $\forall r > \hat{r}^*$ に於て

$$E_{G_n, z_n} \left[\sup_{g \in G} \frac{\frac{1}{n} \sum_{i=1}^n G_i g(z_i)}{\rho g + r} \right] \leq \frac{4\phi(r)}{r}$$

(Proof)

$$\begin{aligned} \text{左辺の } E[\cdot] \text{ の中身} &\leq \sup_{g \in G_r} \frac{\frac{1}{n} \sum_{i=1}^n G_i g(z_i)}{r} \\ &+ \sum_{j=0}^{\infty} \sup_{g \in G_{r4^{j+1}} \setminus G_{r4^j}} \frac{\frac{1}{n} \sum_{i=1}^n G_i g(z_i)}{r4^j + r} \end{aligned}$$

両辺期待値を取ると、

$$\begin{aligned} E[\text{左辺}] &\leq \frac{R_n(G_r)}{r} + \sum_{j=0}^{\infty} \frac{R_n(G_{r4^{j+1}})}{r4^j + r} \\ &\leq \frac{\phi(r)}{r} + \frac{1}{r} \sum_{j=0}^{\infty} \frac{\phi(4^{j+1}r)}{4^j + 1} \\ &\leq \frac{\phi(r)}{r} + \frac{1}{r} \sum_{j=0}^{\infty} \frac{2^{j+1}\phi(r)}{4^j + 1} \leq \frac{4\phi(r)}{r} \end{aligned}$$

★ Thm (Talagrand の集中不等式)

(Z, \mathcal{A}, P) : 確率空間

\tilde{G} : (Z, \mathcal{A}) 上の可測関数の集合,

$$E[g] = 0, E[g^2] \leq \nu, \|g\|_\infty \leq B \quad (\forall g \in \tilde{G})$$

かつ \tilde{G} は $\|\cdot\|_\infty$ に閉じて可分.

すると $\forall t > 0$ に対し

$$P_Z \left[\sup_{g \in \tilde{G}} \frac{1}{n} \sum_{i=1}^n g(z_i) \geq 2 E_Z \left[\sup_{g \in \tilde{G}} \frac{1}{n} \sum_{i=1}^n g(z_i) \right] + \sqrt{\frac{2\pi\nu}{n}} + \frac{2tB}{n} \right] \leq e^{-t}$$

2. $\tilde{G} \subset \mathcal{C}$. $\tilde{G} = \left\{ \tilde{g} = \frac{Pg - g(z)}{Pg + r} \mid g \in \tilde{G} \right\} \subset \mathcal{A}$ かつ

$$\|\tilde{g}\|_\infty = \left\| \frac{Pg - g}{Pg + r} \right\|_\infty \leq \frac{M}{r} = B$$

$$\|\tilde{g}\|_{L^2}^2 = \frac{\|Pg - g\|_{L^2}^2}{(Pg + r)^2} \leq \frac{Pg^2}{2r \frac{1}{2} P(g^2)} \leq \frac{1}{\alpha r} = \nu$$

$\because \frac{1}{2} P(g^2) \leq Pg$

かつ ν : Peeling の議論 あり

$$E \left[\sup_{\tilde{g} \in \tilde{G}} (P - P_n) \tilde{g} \right] \leq 2 R_n(\tilde{G}) \leq \frac{2\phi(r)}{r}$$

$\phi(r)$ は $R_n(G_r) \leq \phi(r)$ となる

よって $\tilde{G} = \tilde{g}$ かつ

$$(P - P_n) \left(\frac{P\hat{f} - P_n\hat{f}}{P\hat{g} + r} \right) \leq \frac{16\phi(r)}{r} + \sqrt{\frac{2\pi}{\alpha r n}} + \frac{2tM}{rn}$$

$=: \psi_n(r) \quad \text{w.p. } 1 - e^{-t}$

↑
Ratio type Empirical process

$$\hat{f} := \operatorname{argmin}_{f \in \mathcal{F}} \hat{L}(f) \quad \text{と} \text{お} \text{す}.$$

$$\begin{aligned} \text{よ} \text{う} \text{に} \text{す} \text{。} \quad L(\hat{f}) - L(f^0) &= L(\hat{f}) - \hat{L}(\hat{f}) + \hat{L}(\hat{f}) - \hat{L}(f^*) + \hat{L}(f^*) - L(f^*) \\ &\quad + L(f^*) - L(f^0) + \hat{L}(f^0) - \hat{L}(f^0) \\ &= (P - P_n) \hat{g} + P_n (\ell \circ f^* - \ell \circ f^0) \\ &= (P - P_n) \hat{g} + (P_n - P) (\ell \circ f^* - \ell \circ f^0) + P (\ell \circ f^* - \ell \circ f^0) \\ &\leq (P \hat{g} + r) \gamma_n(r) + (P_n - P) (\ell \circ f^* - \ell \circ f^0) + L(f^*) - L(f^0) \end{aligned}$$

• $(P_n - P) (\ell \circ f^* - \ell \circ f^0) = o_p(1)$.

$r^* = L(f^*) - L(f^0)$ とおす。 $\frac{\alpha}{2} P \hat{g}^2 \leq r^*$ ($\hat{g} = \ell \circ f^* - \ell \circ f^0$) と。

Bernstein の不等式より

$$(P_n - P) \hat{g} = O_p \left(\sqrt{\frac{r^*}{n}} + \frac{1}{n} \right) \leq O_p \left(\frac{1}{n} \right) + r^*$$

よって、

• $(P \hat{g} + r) \gamma_n(r) = o_p(1)$.

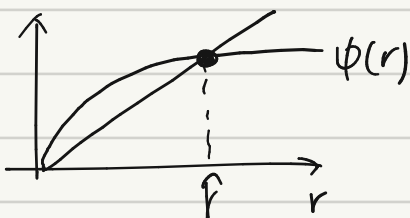
$$\psi_n(\hat{r}) \leq \frac{1}{2} \quad \text{よ} \text{う} \text{に} \text{す} \text{。} \quad \hat{r} \in \text{適} \text{当} \text{な} \text{値} \text{と} \text{す} \text{。}$$

$$P(\hat{g}) = L(\hat{f}) - L(f^0) \leq 2\hat{r} \gamma_n(\hat{r}) + 2r^* + O_p \left(\frac{1}{n} \right)$$

$$\leq \hat{r} + 2r^* + O_p \left(\frac{1}{n} \right)$$

よって、 $\gamma_n(\hat{r}) \leq \frac{1}{2}$ かつ $\max \left\{ \phi(\hat{r}), \frac{1}{\alpha n}, \frac{M \hat{g}}{n}, r^* \right\} \leq \hat{r}$ となる。

$$\Rightarrow L(\hat{f}) - L(f^0) \leq \hat{r} + 2r^* + O_p \left(\frac{1}{n} \right) \quad \text{w.p. } 1 - e^{-t}$$



$\phi(r)$ は大体凹関数 (2重)

$\phi(r) \equiv r$ となる点 r は ϕ の固定点 (fixed point)

例:

$\log N(\mathcal{F}, \varepsilon, \|\cdot\|_{n,2}) \lesssim \varepsilon^{-2p}$ (a.s.) のとき.

$$R_n(G_r) \lesssim \left(\frac{r^{\frac{1-p}{2}}}{T_n}, n^{-\frac{1}{4p}} \right)$$

$\Rightarrow R_n(G_r) \cong r$ if $\hat{r} = n^{-\frac{1}{4p}}$ で達成.

また, $f^* = f^\circ$ のとき.

$$L(\hat{f}) - L(f^\circ) \lesssim n^{-\frac{1}{4p}} //$$

○ 深層 NN の局所 Rademacher 複雑度

次の 深層 NN の Rademacher 複雑度 Bound

$$\log N(\mathcal{F}_{L,S,B}, \|\cdot\|_{\infty}) \leq 2SL \log(L(BV) \cdot W \cdot \varepsilon^{-1})$$

f). $\forall r > r^*$ におう2. $\phi(r) \leq O\left(\sqrt{\frac{SLr}{n} \log(L(BV)W \cdot n)}\right)$

(略証)

$$E[\hat{R}_n(G_r)] \leq E_2[\hat{R}_n(\mathcal{F}_r)] \quad (\text{ただし } \mathcal{F}_r := \{f - f^0 \mid L(f) - L(f^0) \leq r\})$$

↑
contraction ineq.

$$\leq E_2 \left[\inf_{d>0} \left\{ 4d + \frac{12}{rn} \int_d^{\delta} \sqrt{\log 2N(\varepsilon, \mathcal{F}_r, \|\cdot\|_{n,2})} d\varepsilon \right\} \right]$$

$$\text{ただし } \delta = \sup_{f \in \mathcal{F}_r} \|f\|_{n,2} \lesssim \sqrt{\frac{2r}{d}}$$

↑
これは示さないと存在するの2.
この δ のバリエーションは
局所 Rademacher 複雑度 ε と r のバリエーション
必要がある。

定数

$$\downarrow c\sqrt{\frac{r}{2}}$$

$$\lesssim \frac{1}{n} + \frac{1}{rn} \int_{\frac{1}{n}}^{\frac{c\sqrt{r}}{2}} \sqrt{2SL \log(L(BV)W \varepsilon^{-1}) + \log(2)} d\varepsilon$$

$$\lesssim \sqrt{\frac{SLr}{n} \log(L(BV)W \cdot n)}$$

1次元線形回帰第1応用:

$$y_i = f^*(x_i) + \varepsilon_i$$

ε_i は i.i.d. 雑音, 平均0, $|\varepsilon_i| \leq C$ (a.s.)

$$\hat{f} = \arg \min_{f \in \mathcal{F}_{L,S,B}} \sum_{i=1}^n (y_i - f(x_i))^2$$

$$L: \mathcal{R}^n \rightarrow \mathcal{R}$$

Lem (Yarotsky, 2016; Schmidt-Tieber, 2017)

$f^0 \in C^s([0,1]^d)$, $\|f^0\|_{C^s} \leq 1$ のとき.

$N \gg 1$ 任意整数 t に対し.

$\exists L = O(\log(N))$, $\exists S = O(N \log(N))$,

$\exists W = O(N)$, $\exists B = O(1)$ 2.

$$\inf_{f \in \mathcal{F}_{L,S,B}} \|f^0 - f\|_{\infty} \leq O(N^{-\frac{s}{d}})$$

$C^s([0,1]^d)$: $1 \leq s < \infty$

$$= \{f \mid \sum_{|d| \leq m} \|D^d f\|_{\infty} + \sum_{|d|=m} \sup \frac{\|D^d f(x) - D^d f(y)\|_{\infty}}{\|x-y\|_{\infty}^{s-m}} < \infty\}$$

($m = \lfloor s \rfloor$)

よ.2. $r^* = O(N^{-\frac{2s}{d}})$

$$\hat{r} = O\left(\frac{N \log(N)^2 (\log(N) + \log(n))}{n}\right)$$

← $1/n$ 分 2 ← t L t W
← $1/n$ 分 3 分 2

と r^* と \hat{r} の 2-項 2.

$$\|\hat{f} - f^0\|_{L^2(\mathcal{P}_n)}^2 = O_p(\hat{r} + r^*)$$

2-項. $N \asymp n^{\frac{d}{2s+d}}$ とおけば r^* .

$$\|\hat{f} - f^0\|_{L^2(\mathcal{P}_n)}^2 = O_p\left(n^{-\frac{2s}{2s+d}} \log(n)^3\right)$$

(注) 最近凸性非凸 L^2 -2 項 ϵ を (L^2) -最近 L^1 を達成する解析の理論的 2 項.

(注) 一般化) ξ_i が \mathcal{G}_i 分布の場合, 有界性 $|\xi_i| \leq C$ が成り立たない. この場合は Talagrand の集中不等式が使えない. 代わりに カーブ集中不等式 を用いるべき. (2017 214)

$(\xi_i)_{i=1}^n$: i.i.d. $N(0, \sigma^2)$, $(x_i)_{i=1}^n \subset \mathcal{X}$: 固定, $F: \|\cdot\|_{\infty} = \text{max}$ 可分

$$P\left(\sup_{f \in F} \left| \frac{1}{n} \sum_{i=1}^n \xi_i f(x_i) \right| \geq E\left[\sup_{f \in F} \left| \frac{1}{n} \sum_{i=1}^n \xi_i f(x_i) \right| \right] + \sqrt{\frac{\sigma^2 \|F\|_n^2}{n}} t\right) \leq \exp\left(-\frac{t^2}{2}\right) \quad (t > 0) \quad t \geq L. \quad \|F\|_n^2 = \sup_{f \in F} \frac{1}{n} \sum_{i=1}^n f(x_i)^2$$