

ニューラルネットワークの汎化誤差

汎化誤差解析の一般論

$x \in X \mapsto y \in Y$ の入出力関数を学習したい。

F : モデル ($X \rightarrow \mathbb{R}$ の関数からなる集合)

例: ニューラルネットワーク

再生核ヒルベルト空間の単位球

目標: F から汎化誤差の期待値を最小にするものを推定したい。

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

例: 二乗損失 = $l(y, f(x)) = (y - f(x))^2$

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

$L(f) = E_{x,y} [l(Y, f(x))]$ 汎化誤差, 期待損失

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

$Pf = E_{x,y} [f(x, Y)]$, $P_n f = \frac{1}{n} \sum_{i=1}^n f(x_i, y_i)$ とおく。

($L(f) = P(pf)$, $\hat{L}(f) = P_n(pf)$)

$f^* := \operatorname{argmin}_{f \in F} L(f)$ とおくと

推定量 \hat{f} に対し $L(\hat{f}) - L(f^*)$ を Excess risk (残余誤差) と言います。
これを汎化誤差とこのことも多い

※ $L(\hat{f}) - \hat{L}(\hat{f})$ のことを汎化誤差と言ったことも多いが。

判正確率は汎化率 $\rightarrow \rho$ と言ったことある。

→ 基本的考え方。

$$L(\hat{f}) - L(f^*) = \underbrace{L(\hat{f}) - \hat{L}(\hat{f})}_{\text{汎化率} \rightarrow \rho} + \underbrace{\hat{L}(\hat{f}) - \hat{L}(f^*)}_{\substack{\uparrow \text{訓練誤差を} \\ \text{最小化したときの} \\ \text{0以下}}} + \underbrace{\hat{L}(f^*) - L(f^*)}_{\substack{\downarrow \text{大数の法則} \\ L(\frac{1}{n})}}$$

よって、基本的には汎化率 $\rightarrow \rho$ をおさえれば汎化誤差をおさえられる。

\hat{f} が訓練誤差を最小化したとき $\hat{L}(\hat{f}) - \hat{L}(f^*) \leq 0$ である。よって

$$0 \leq L(\hat{f}) - L(f^*) \leq L(\hat{f}) - \hat{L}(f^*) - (\hat{L}(\hat{f}) - \hat{L}(f^*)) \\ = (P - P_n)(\hat{f} - f^*) < ?$$

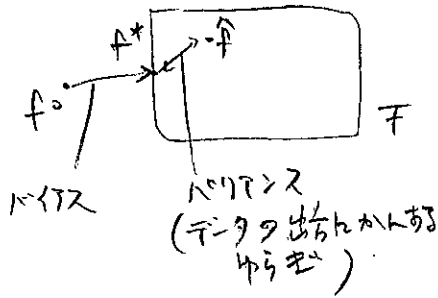
- バイアス-バリアンスのトレードオフ

$$f^0 := \operatorname{argmin} L(f)$$

f : 可測関数

f^0 の汎化誤差の差:

$$L(\hat{f}) - L(f^0) = \underbrace{L(\hat{f}) - L(f^*)}_{\text{バリアンス}} + \underbrace{L(f^*) - L(f^0)}_{\text{(モデル)バイアス}}$$



F が大きくなる \rightarrow バイアス小
 バリアンス大
 F が小さくなる \rightarrow バイアス大
 バリアンス小

適度を大とした F を用いて、バイアスとバリアンスのバランスを取ることがある。

- バリアンスのバウンズ

よりあえてバイアスは無視して、バリアンスをあたえる。

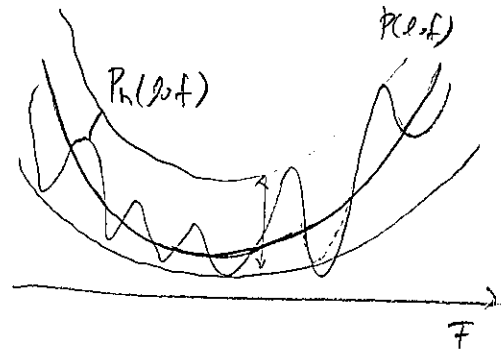
$$L(\hat{f}) - L(f^*) \leq (P - P_n)(\hat{f} - f^*)$$

この右辺を

$$\sup_{A \in \mathcal{F}} (P - P_n)(\hat{f} - f^*)$$

と表す。

② f は P -0 に依存する \Rightarrow 単純な大数の法則ではダメ。



$z = (x, \delta)$ と $g(z) = g(x, \delta)$ とする。
 $g := \hat{f} - f^*$, $G := \{g | \hat{f} - f^*, f \in \mathcal{F}\}$ と表す。

$\sup_{g \in G} (P - P_n)g$ を抑える問題になる。

* ここで $\sup_{g \in G} P_n g \neq \sup_{g \in G} \frac{1}{n} \sum_{i=1}^n \varepsilon_i g(z_i)$ ← $\{0, 1\}$ -value

これは以後可測] である仮定が、可測性からなる場合は以下の議論は成立している (Fubini の定理が成立している $E_x E_\varepsilon [\cdot]^+ \neq E_\varepsilon E_x [\cdot]^+$)

- van der Vaart & Wellner: Weak convergence and Empirical processes (1996)

Sec. 2-3

- Dudley "Uniform Central Limit Theorems": Sustlin image admissible.

Rademacher complexity

Def

G_n ($n=1, \dots, n$) ε $P(G_n=1) = P(G_n=-1) = \frac{1}{2}$ と δ i.i.d. sequence とす。
(Rademacher 変数)

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

(経験 Rademacher 複雑さ)

$$R_n(G) := E_{z_i} [\hat{R}_n(G)]$$

(Rademacher 複雑さ)

Thm (対称化)

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

(直接 $\sup (P - P_n)g$ を対称化して $R_n(G)$ の高々 2 倍 (かす))

Proof

$z_i \cup G$ i.i.d. sequence

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

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

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

($z_i' \cup z_i$ 対称 i.i.d. 分布)

$G_i \cup -G_i$ 対称 i.i.d. 分布

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

$$= 2 R_n(G)$$

(\therefore Fubini の定理) \leftarrow 解法 4
可測性 = 必要
可積分性 = 必要

Thm

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

Proof

LEM: (McDiarmid の不等式)

$X_1, \dots, X_n \in X$ と独立な (同一の分布を持つ) 確率変数とす。

$f: X^n \rightarrow \mathbb{R}$ と可測関数とす。

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

$\forall x_1, \dots, x_n, x_n' \in X$ に対して成り立つとす。

$$P(f(x_1, \dots, x_n) - E[f(x_1, \dots, x_n)] \geq \varepsilon) \leq \exp\left(-\frac{2\varepsilon^2}{\sum_{i=1}^n C_i^2}\right) \quad (\forall \varepsilon > 0) //$$

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

(9)

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

とすれば

$$|f(z_1, \dots, z_n) - f(z'_1, \dots, z'_n)| \leq \frac{2}{n} M$$

となり、これは Mc Diarnid の不等式 と 先の 対称化 により示される。 //

Thm (Rademacher 複素数の性質)

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

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

$$E_{z_i, G_i} \left[\sup_{g \in G} \frac{1}{n} \sum_{i=1}^n G_i \phi_n(g(z_i)) \right] \leq D R_n(G) \quad (\text{contraction ineq.})$$

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

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

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

Ex.

(1) 半規開区間モデル

$$\left. \begin{array}{l} - F = \{ f(x) = x^T w \mid \|w\|_2 \leq 1 \} \\ - |l(y, f(x)) - l(y, f(x'))| \leq |f(x) - f(x')| \end{array} \right\} \quad - E[\|x\|_2^2] \leq 1$$

ならば $R_n(G) \leq R_n(F)$ (contraction ineq.)

$$= E_{z_i, G_i} \left[\sup_{w: \|w\|_2 \leq 1} \frac{1}{n} \sum_{i=1}^n G_i w^T x_i \right]$$

$$\leq E_{z_i, G_i} \left[\sup_{\|w\|_2 \leq 1} \left\| \frac{1}{n} \sum_{i=1}^n G_i x_i \right\|_2 \cdot \|w\|_2 \right]$$

$$\leq \sqrt{E_{z_i, G_i} \left[\frac{1}{n^2} \sum_{i,j} G_i G_j (x_i^T x_j) \right]}$$

$$= \sqrt{\frac{1}{n^2} \sum_{i=1}^n E[\|x_i\|_2^2]} \leq \sqrt{\frac{1}{n}}$$

(2) F が有限集合のとき

$F = \{f_1, \dots, f_M\}; \|f_j\|_\infty \leq R$, ならば

$$R_n(F) \leq R \sqrt{\frac{2 \log(M)}{n}} \quad (\text{Hoeffding の不等式 により示される。})$$

(証明は略) (Massart の定理)

深層二層ネットワークの Rademacher 複素値

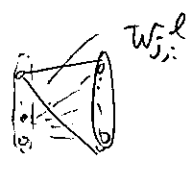
= 1/2 n 個の独立変数

$$F = \{ w^{(L)} \eta(w^{(L-1)} \eta(\dots \eta(w^{(1)} x) \dots)) \mid w^{(l)} \in \mathbb{R}^{m_l \times m_{l-1}}, \|w^{(l)}\|_{1, \infty} \leq B_l \}$$

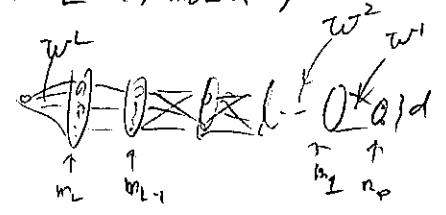
ここで、 $\eta(x) = \max(x, 0)$ (ReLU), $x \in \mathbb{R}^d$

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

とある。



($m_L = 1, m_0 = d$)



$$F_n = \{ f(x) = \sum_{j=1}^{m_{i-1}} w_j \eta(f_j(x)) \mid f_j \in F_{i-1}, \|w\|_1 \leq B_i \}$$

$$F_1 = \{ f(x) = w^T x \mid \|w\|_1 \leq B_1 \}$$

と定義すると、 $F = F_L$ とある。

Lem

$$\hat{R}_n(F_n) \leq 2 B_i \hat{R}_n(F_{i-1})$$

//

Proof

$$\hat{R}_n(F_n) = E_{G_i} \left[\sup_{\substack{\|w\|_1 \leq B_i \\ f_j \in F_{i-1}}} \frac{1}{n} \sum_{k=1}^n G_k \sum_j w_j \eta(f_j(x_k)) \right]$$

$$= E_{G_i} \left[\sup_{w, f_j} \sum_j w_j \frac{1}{n} \sum_{k=1}^n G_k \eta(f_j(x_k)) \right]$$

$$\leq E_{G_i} \left[\sup_{w, f_j} \underbrace{\|w\|_1}_{\leq B_i} \cdot \max_{1 \leq j \leq m_{i-1}} \left| \frac{1}{n} \sum_{k=1}^n G_k \eta(f_j(x_k)) \right| \right]$$

$$\leq B_i E_{G_i} \left[\sup_{f \in F_{i-1}} \left| \frac{1}{n} \sum_{k=1}^n G_k \eta(f(x_k)) \right| \right]$$

$$\begin{aligned} &= \text{ここで、} E_{G_i} \left[\sup \left| \frac{1}{n} \sum_{k=1}^n G_k g(x_k) \right| \right] \\ &= E_{G_i} \left[\max \left\{ \sup \frac{1}{n} \sum G_k g(x_k), - \inf \frac{1}{n} \sum G_k g(x_k) \right\} \right] \\ &= E_{G_i} \left[\max \left\{ \sup \frac{1}{n} \sum G_k g(x_k), \sup \frac{1}{n} \sum (-G_k) g(x_k) \right\} \right] \\ &\leq 2 E_{G_i} \left[\sup \frac{1}{n} \sum G_k g(x_k) \right] \quad (g_1 - g_2 \in \mathcal{F} \text{ 同分布}) \end{aligned}$$

つまり

$$\leq 2 B_i E_{G_i} \left[\sup_{f \in F_{i-1}} \frac{1}{n} \sum_{k=1}^n G_k \eta(f(x_k)) \right]$$

$$\leq 2 B_i \hat{R}_n(F_{i-1}) \quad (\because \text{Contraction ineq.})$$

$$\leq B_i E_{G_i} \left[\max_{j=1, \dots, d} \left| \frac{1}{n} \sum_{k=1}^n G_k x_{k,j} \right| \right] \leq \sqrt{\frac{2 \log(d)}{n}}$$

ここで、 $\|x\|_{\infty} \leq 1$ (a.s.) とする。

$$\hat{R}_n(F_1) = E_{G_n} \left[\sup_w \frac{1}{n} \sum_{k=1}^n G_k w^T x_k \right] \leq E_{G_n} \left[\sup_w \underbrace{\|w\|_1}_{\leq B_1} \max_{j=1, \dots, d} \left| \frac{1}{n} \sum_{k=1}^n G_k x_{k,j} \right| \right]$$

以上をFとの比で上下を分る。

Thm

$\|x\|_0 \leq 1$ なら

$$R_n(F) \leq 2^L \cdot \prod_{i=1}^L B_{\lambda_i} \cdot \sqrt{\frac{2B_{\lambda_i}}{n}}$$

//

よって、適度に各層のノルムを制御すれば、汎化誤差は $O_p(\frac{1}{\sqrt{n}})$ で減少するようになる。

このハイパーパラメータはノルムに依存せずノルムにおまかせしている \rightarrow 横幅が広くなる。

他のハイパー

- Spectrally Normalized Bound (Bartlett et al., 17)

$$\frac{L}{n} \prod_{k=1}^L \|W^k\|_{op} \left(\sum_{j=1}^L \frac{\|W^j - M^j\|_{1,2}^2}{\|W^j\|_{op}^2} \right)^{1/2} = \frac{1}{\sqrt{n}}$$

たとえば、 $\|W^j\|_{1,2} = \sum_{i=1}^{m_j} \|W_{i,:}^j\|_2$ であり、 M^j は任意の行列である。

・ 特に、 M^j を初期値とすれば、最適化によってノルムが小さく重みがある状態、汎化性能が良いという結果が得られる。

・ $\|W_{i,:}^j\|_2$ を用いることで、スライスごとのノルムが大きい層にのみ

・ 実験結果をよく説明する。

- (Golowich et al. 18)

$$\min \left\{ \frac{L}{n} \|W^k\|_F \cdot \frac{1}{\sqrt{n}}, \frac{L}{n} \|W^k\|_F \cdot \sqrt{\frac{L}{n}} \right\}$$

$$\text{たとえば、} \|W^k\|_F = \sqrt{\sum_{i,j} (W_{i,j}^k)^2}$$

・ 2^L である。

・ $\|W^k\|_F$ は横幅の影響を受けられる。

実際のネットワークは、各層の固有値が速く減衰する状態である。

$\|W^k\|_F$ が小さくなる。