

Neural Network Learning: Theoretical Foundations
Chapter 16-17
Anthony and Bartlett

Presented by Kuhwan Jeong ¹

¹Department of Statistics, Seoul National University, South Korea

October, 2017

- ① **Pattern Classification with Binary-Output Neural Networks**
(Chap. 2-8)
- ② **Pattern Classification with Real-Output Neural Networks**
(Chap. 9-15)
- ③ **Learning Real-Valued Functions**
(Chap. 16-21)

Outline

- ① 16. Learning Classes of Real Functions
 - 16.1 Introduction
 - 16.2 Learning Framework for Real Estimation
 - 16.3 Learning Finite Classes of Real Functions
 - 16.4 A Substitute for Finiteness

- ② 17. Uniform Convergence Results for Real Function Classes
 - 17.1 Uniform Convergence for Real Functions
 - 17.2 Remarks

16.1 Introduction

- This part examines supervised learning problems.
- Data are generated by a probability distribution P on $X \times \mathbb{R}$.
- To measure how accurately $f(x)$ approximates y , we use the *quadratic loss*.
- We define the *error* of a function $f : X \rightarrow \mathbb{R}$ with respect to P as

$$\text{er}_P(f) = \mathbb{E}(f(x) - y)^2.$$

16.2 Learning Framework for Real Estimation

- We assume that the random variable y falls in $[0, 1]$.
- F : a set of functions mapping from X into $[0, 1]$
- A learning algorithm L for F is a function

$$L : \bigcup_{m=1}^{\infty} (X \times \mathbb{R})^m \rightarrow F$$

with the following property:

for any $\epsilon \in (0, 1)$, $\delta \in (0, 1)$, $\exists m_0(\epsilon, \delta)$ s.t.

if $m \geq m_0(\epsilon, \delta)$, then for any probability distribution P on $X \times [0, 1]$,

$$P^m \{ \text{er}_P(L(z)) < \inf_{f \in \mathcal{F}} \text{er}_P(f) + \epsilon \} \geq 1 - \delta.$$

16.2 Learning Framework for Real Estimation

- *sample complexity* $m_L(\epsilon, \delta)$: the least possible value $m_0(\epsilon, \delta)$
- *estimation error* $\epsilon_L(m, \delta)$: the least possible value ϵ s.t.

$$P^m \left\{ \text{er}_P(L(z)) < \inf_{f \in \mathcal{F}} \text{er}_P(f) + \epsilon \right\} \geq 1 - \delta$$

16.3 Learning Finite Classes of Real Functions

- *sample error* $\hat{e}r_z(f)$:

$$\hat{e}r_z(f) = \frac{1}{m} \sum_{i=1}^m (f(x_i) - y_i)^2$$

- *sample error minimization* (SEM) algorithm L for a finite class F :
a function from $\bigcup_{m=1}^{\infty} Z^m$ to F s.t.

$$\hat{e}r_z(L(z)) = \min_{f \in \mathcal{F}} \hat{e}r_z(f)$$

16.3 Learning Finite Classes of Real Functions

Theorem 16.2 SEM algorithm L is a learning algorithm for finite F , whose estimation error satisfies

$$\epsilon_L(m, \delta) \leq \left(\frac{2}{m} \log \left(\frac{2|F|}{\delta} \right) \right)^{1/2}$$

Proof

$$\begin{aligned} P^m \left\{ |\hat{\text{er}}_z(f) - \text{er}_P(f)| \geq \frac{\epsilon}{2} \right\} &\leq 2e^{-\epsilon^2 m/2} \\ P^m \left\{ \exists f \in \mathcal{F} \text{ s.t. } |\hat{\text{er}}_z(f) - \text{er}_P(f)| \geq \frac{\epsilon}{2} \right\} &\leq 2|F|e^{-\epsilon^2 m/2} =: \delta \end{aligned} \quad (1)$$

Let f^* be a function in F satisfying

$$\text{er}_P(f^*) = \min_{f \in \mathcal{F}} \text{er}_P(f).$$

With probability at least $1 - \delta$,

$$\text{er}_P(L(z)) \stackrel{(1)}{<} \hat{\text{er}}_z(L(z)) + \frac{\epsilon}{2} \leq \hat{\text{er}}_z(f^*) + \frac{\epsilon}{2} \stackrel{(1)}{<} \left(\text{er}_P(f^*) + \frac{\epsilon}{2} \right) + \frac{\epsilon}{2} = \min_{f \in \mathcal{F}} \text{er}_P(F) + \epsilon. \quad \square$$

16.3 Learning Finite Classes of Real Functions

Theorem 16.3 N is a neural network with arbitrary activation functions whose output takes values in $[0, 1]$. F is the set of functions computable by N when each weight and threshold is represented using k bits. Then SEM algorithm L for F is a learning algorithm with

$$m_L(\epsilon, \delta) \leq \frac{2}{\epsilon^2} \left(kW \log 2 + \log \left(\frac{2}{\delta} \right) \right),$$

where W is the total number of weights and thresholds.

Proof

It follows from Theorem 16.2 with $|F| \leq 2^{kW}$.

16.4 A Substitute for Finiteness

- F : a set of functions mapping from X into $[0, 1]$
- metric d_{L_∞} on F :

$$d_{L_\infty}(f, g) = \sup_{x \in X} |f(x) - g(x)|$$

- Approximate-SEM algorithm \mathcal{A} for F :
a function from $\bigcup_{m=1}^{\infty} Z^m \times \mathbb{R}^+$ to F s.t.

$$\hat{e}_z(\mathcal{A}(z, \epsilon)) < \inf_{f \in \mathcal{F}} \hat{e}_z(f) + \epsilon$$

16.4 A Substitute for Finiteness

Theorem 16.5 F is totally bounded w.r.t. d_{L_∞} , and L satisfies $L(z) = \mathcal{A}(z, 1/\sqrt{2m})$ for $z \in Z^m$. Then L is a learning algorithm for F with

$$m_L(\epsilon, \delta) \leq \frac{72}{\epsilon^2} \log \left(\frac{2\mathcal{N}(\epsilon/12, F, d_{L_\infty})}{\delta} \right).$$

Proof

Lemma If $d_{L_\infty}(f, g) < \alpha$, then $|\text{er}_Q(f) - \text{er}_Q(g)| < 2\alpha$ for any probability distribution Q on Z .

Let \mathcal{C} be a $\epsilon/12$ -cover for F of cardinality $\mathcal{N}(\epsilon/12, F, d_{L_\infty})$. For any $f \in F$, $\exists \hat{f} \in \mathcal{C}$ s.t. $d_{L_\infty}(f, \hat{f}) < \epsilon/12$, and so

$$|\text{er}_P(f) - \text{er}_P(\hat{f})| < \frac{\epsilon}{6}, \quad (2)$$

$$|\hat{\text{er}}_z(f) - \hat{\text{er}}_z(\hat{f})| < \frac{\epsilon}{6}. \quad (3)$$

Let f^* be s.t.

$$\text{er}_P(f^*) < \inf_{f \in F} \text{er}_P(f) + \frac{\epsilon}{12}. \quad (4)$$

16.4 A Substitute for Finiteness

From the proof of Theorem 16.2, we have

$$P^m \left\{ \exists f \in \mathcal{C} \text{ s.t. } |\hat{\text{er}}_z(f) - \text{er}_P(f)| \geq \frac{\epsilon}{12} \right\} \leq 2\mathcal{N}(\epsilon/12, F, d_{L_\infty}) e^{-\epsilon^2 m/72} =: \delta. \quad (5)$$

For convenience, we denote $L(z)$ by f_z . With probability at least $1 - \delta$,

$$\begin{aligned} \text{er}_P(f_z) &\stackrel{(2)}{<} \text{er}_P(\hat{f}_z) + \frac{\epsilon}{6} \stackrel{(5)}{<} \left(\hat{\text{er}}_z(\hat{f}_z) + \frac{\epsilon}{12} \right) + \frac{\epsilon}{6} \stackrel{(3)}{<} \left(\hat{\text{er}}_z(f_z) + \frac{\epsilon}{6} \right) + \frac{3\epsilon}{12} \\ &\stackrel{\text{aSEM}}{<} \left(\inf_{f \in \mathcal{F}} \hat{\text{er}}_z(f) + \frac{\epsilon}{12} \right) + \frac{5\epsilon}{12} \\ &\leq \hat{\text{er}}_z(f^*) + \frac{\epsilon}{2} \stackrel{(3)}{<} \left(\hat{\text{er}}_z(\hat{f}^*) + \frac{\epsilon}{6} \right) + \frac{\epsilon}{2} \stackrel{(5)}{<} \left(\text{er}_P(\hat{f}^*) + \frac{\epsilon}{12} \right) + \frac{2\epsilon}{3} \\ &\stackrel{(2)}{<} \left(\text{er}_P(f^*) + \frac{\epsilon}{6} \right) + \frac{3\epsilon}{4} \stackrel{(4)}{<} \left(\inf_{f \in F} \text{er}_P(f) + \frac{\epsilon}{12} \right) + \frac{11\epsilon}{12} \\ &= \inf_{f \in F} \text{er}_P(f) + \epsilon. \quad \square \end{aligned}$$

17.1 Uniform Convergence for Real Functions

- Given $x = (x_1, \dots, x_k) \in X^k$,

$$F_{|x} = \{(f(x_1), \dots, f(x_k)) : f \in F\} \subset \mathbb{R}^k.$$

- uniform ϵ -covering number $\mathcal{N}_1(\epsilon, F, k)$:

$$\mathcal{N}_1(\epsilon, F, k) = \max\{\mathcal{N}(\epsilon, F_{|x}, d_1) : x \in X^k\}.$$

Theorem 17.1

$$P^m\{\exists f \in F \text{ s.t. } |\text{er}_P(f) - \hat{\text{er}}_z(f)| \geq \epsilon\} \leq 4\mathcal{N}_1(\epsilon/16, F, 2m)e^{-\epsilon^2 m/32}$$

17.1 Uniform Convergence for Real Functions

Lemma 17.2 (Symmetrization)

Let

$$Q = \{z \in Z^m : \exists f \in \mathcal{F} \text{ s.t. } |er_P(f) - \hat{er}_z(f)| \geq \epsilon\},$$
$$R = \left\{ (r, s) \in Z^m \times Z^m : \exists f \in \mathcal{F} \text{ s.t. } |\hat{er}_r(f) - \hat{er}_s(f)| \geq \frac{\epsilon}{2} \right\}.$$

Then, $P^m(Q) \leq 2P^{2m}(R)$ for $m \geq 4/\epsilon^2$.

Proof Since $|er_P(f) - \hat{er}_r(f)| \geq \epsilon$ and $|er_P(f) - \hat{er}_s(f)| < \epsilon/2$ implies $|\hat{er}_r(f) - \hat{er}_s(f)| \geq \epsilon/2$, we have

$$\begin{aligned} P^{2m}(R) &\geq P^{2m}\{\exists f \in \mathcal{F} \text{ s.t. } |er_P(f) - \hat{er}_r(f)| \geq \epsilon \ \& \ |er_P(f) - \hat{er}_s(f)| < \epsilon/2\} \\ &= \int_Q P^m\{s : \exists f \in \mathcal{F} \text{ s.t. } |er_P(f) - \hat{er}_s(f)| < \epsilon/2\} dP^m(r) \\ &\geq P^m(Q)/2 \text{ provided } m \geq 4/\epsilon^2. \square \end{aligned}$$

17.1 Uniform Convergence for Real Functions

Permutation

- Γ_m : the set of permutations on $\{1, \dots, 2m\}$ in which i can be only switched with $m + i$.
- By Lemma 4.5,

$$P^{2m}(R) \leq \max_{z \in Z^{2m}} \Pr_{\sigma}(\sigma z \in R),$$

where \Pr_{σ} is the uniform probability on Γ_m .

17.1 Uniform Convergence for Real Functions

- $l_f : Z \rightarrow [0, 1]$ is given by $l_f(x, y) = (f(x) - y)^2$.
- $l_F = \{l_f : f \in F\}$

Lemma 17.4

For all $z \in Z^m$,

$$\mathcal{N}(\epsilon, (l_F)_{|z}, d_1) \leq \mathcal{N}_1(\epsilon/2, F, m).$$

Proof Let $z = ((x_1, y_1), \dots, (x_m, y_m))$.

$$\begin{aligned} d_1(l_{f|z}, l_{g|z}) &= \frac{1}{m} \sum_{i=1}^m |(f(x_i) - y_i)^2 - (g(x_i) - y_i)^2| \\ &= \frac{1}{m} \sum_{i=1}^m |(f(x_i) - y_i)(f(x_i) + g(x_i) - 2y_i)| \leq \frac{2}{m} \sum_{i=1}^m |f(x_i) - g(x_i)| \end{aligned}$$

$$\mathcal{N}(\epsilon, (l_F)_{|z}, d_1) \leq \mathcal{N}(\epsilon/2, F_{|z}, d_1) \leq \mathcal{N}_1(\epsilon/2, F, m) \quad \square$$

17.1 Uniform Convergence for Real Functions

Lemma 17.5

$$\max_{z \in \mathcal{Z}^{2m}} \Pr_{\sigma}(\sigma z \in R) \leq 2\mathcal{N}_1(\epsilon/16, F, 2m)e^{-\epsilon^2 m/32}$$

Proof Choose $G \subset F$ s.t. $(l_G)_z$ is a minimal $\epsilon/8$ -cover for $(l_F)_z$ w.r.t. d_1 . By

Lemma 17.4,

$$|G| = \mathcal{N}(\epsilon/8, (l_F)_z, d_1) \leq \mathcal{N}_1(\epsilon/16, F, 2m).$$

If $\sigma z = (r, s) \in R$, then $\exists f \in G$ s.t. $|\hat{e}_r(f) - \hat{e}_s(f)| \geq \epsilon/4$.

$$\begin{aligned} \Pr_{\sigma}(\sigma z \in R) &\leq \Pr_{\sigma} \left(\exists f \in G : \left| \frac{1}{m} \sum_{i=1}^m (l_f(z_{\sigma(i)}) - l_f(z_{\sigma(m+i)})) \right| \geq \epsilon/4 \right) \\ &\leq |G| \max_{f \in G} \Pr_{\sigma} \left(\left| \frac{1}{m} \sum_{i=1}^m (l_f(z_{\sigma(i)}) - l_f(z_{\sigma(m+i)})) \right| \geq \epsilon/4 \right) \\ &= |G| \max_{f \in G} \Pr_{\beta} \left(\left| \frac{1}{m} \sum_{i=1}^m |l_f(z_i) - l_f(z_{m+i})| |\beta_i| \right| \geq \epsilon/4 \right) \leq 2|G|e^{-\epsilon^2 m/32}, \end{aligned}$$

where β_i is independently and uniformly drawn from $\{-1, 1\}$. \square

17.2 Remarks

- The results do not depend significantly on the choice of the quadratic loss.
- Let $Y = \mathbb{R}$ and $B \geq 1$.
- Suppose the loss function $l : [0, 1] \times Y \rightarrow [0, B]$ satisfies

$$l(y_1, y_0) - l(y_2, y_0) \leq L|y_1 - y_2|$$

for all y_0, y_1, y_2 .

- Define $l_f(z) = l(f(x), y)$.
- F is a class of $[0, 1]$ -valued functions.

Lemma 17.6 (Generalization of Lemma 17.4)

$$\max_{z \in Z^m} \mathcal{N}(\epsilon, (l_F)_{|z}, d_1) \leq \mathcal{N}_1(\epsilon/L, F, m)$$

Theorem 17.7 (Generalization of Theorem 17.1)

$$P^m \left\{ \exists f \in \mathcal{F} \text{ s.t. } \left| \mathbb{E} l_f - \frac{1}{m} \sum_{i=1}^m l_f(z_i) \right| \geq \epsilon \right\} \leq 4\mathcal{N}_1(\epsilon/(8L), F, 2m) e^{-\epsilon^2 m / (32B^4)}$$