

Generalization errors for Deep Learning

김지수 (Jisu KIM)

딥러닝의 통계적 이해 (Deep Learning: Statistical Perspective), 2025 2nd Semester (fall)

This lecture note is a combination of Prof. Joong-Ho Won's "Deep Learning: Statistical Perspective" with other lecture notes. Main references are:

Tong Zhang, Mathematical Analysis of Machine Learning Algorithms, <https://tongzhang-ml.org/lt-book.html>

Matus Telgarsky, Deep learning theory lecture notes, <https://mjt.cs.illinois.edu/dlt/>

Weinan E, Chao Ma, Stephan Wojtowytsch, Lei Wu, Towards a Mathematical Understanding of Neural Network-Based Machine Learning: what we know and what we don't, <https://arxiv.org/abs/2009.10713/>

1 Review

1.1 Basic Model for Supervised Learning

- Input(입력) / Covariate(설명 변수) : $x \in \mathbb{R}^d$, so $x = (x_1, \dots, x_d)$.
- Output(출력) / Response(반응 변수) : $y \in \mathcal{Y}$. If y is categorical, then supervised learning is "classification", and if y is continuous, then supervised learning is "regression".
- Model(모형) :

$$y \approx f(x).$$

If we include the error ϵ to the model, then it can be also written as

$$y = \phi(f(x), \epsilon).$$

For many cases, we assume additive noise, so

$$y = f(x) + \epsilon.$$

- Assumption(가정): f belongs to a family of functions \mathcal{M} . This is the assumption of a model: a model can be still used when the corresponding assumption is not satisfied in your data.
- Loss function(손실 함수): $\ell(y, a)$. A loss function measures the difference between estimated and true values for an instance of data.
- Training data(학습 자료): $\mathcal{T} = \{(y_i, x_i), i = 1, \dots, n\}$, where (y_i, x_i) is a sample from a probability distribution P_i . For many cases we assume i.i.d., or x_i 's are fixed and y_i 's are i.i.d..
- Goal(목적): we want to find f that minimizes the expected prediction error,

$$f^0 = \arg \min_{f \in \mathcal{F}} \mathbb{E}_{(Y, X) \sim P} [\ell(Y, f(X))].$$

Here, \mathcal{F} can be different from \mathcal{M} ; \mathcal{F} can be smaller than \mathcal{M} .

- Prediction model(예측 모형): f^0 is unknown, so we estimate f^0 by \hat{f} using data. For many cases we minimize on the empirical prediction error, that is taking the expectation on the empirical distribution $P_n = \frac{1}{n} \sum_{i=1}^n \delta_{(Y_i, X_i)}$.

$$\hat{f} = \arg \min_{f \in \mathcal{F}} \mathbb{E}_{P_n} [\ell(Y, f(X))] = \arg \min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \ell(Y_i, f(X_i)).$$

- Prediction(예측): if \hat{f} is a predicted function, and x is a new input, then we predict unknown y by $\hat{f}(x)$.

1.2 Rademacher complexity

Random variables ξ_1, \dots, ξ_n are called *Rademacher random variables* if they are independent, identically distributed and $\mathbb{P}(\xi_i = 1) = \mathbb{P}(\xi_i = -1) = 1/2$. Define the *Rademacher complexity* of \mathcal{F} by

$$\text{Rad}_n(\mathcal{F}) = \mathbb{E} \left(\sup_{f \in \mathcal{F}} \left(\frac{1}{n} \sum_{i=1}^n \xi_i f(Z_i) \right) \right).$$

Some authors use a slightly different definition, namely,

$$\text{Rad}_n(\mathcal{F}) = \mathbb{E} \left(\sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n \xi_i f(Z_i) \right| \right).$$

You can use either one. They lead to essentially the same results. In fact, under mild condition, two Rademacher complexity are closely related, as below:

Lemma 1. *Let $\xi = (\xi_1, \dots, \xi_n)$ be i.i.d. Rademacher. Suppose that for any $\xi \in \{\pm 1\}^n$, $\sup_{f \in \mathcal{F}} \sum_{i=1}^n \xi_i f(Z_i) \geq 0$. Then*

$$\mathbb{E}_\xi \left[\sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n \xi_i f(Z_i) \right| \middle| Z \right] \leq \mathbb{E}_\xi \left[\sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \xi_i f(Z_i) \middle| Z \right].$$

Proof. Left as HW. □

Intuitively, $\text{Rad}_n(\mathcal{F})$ is large if we can find functions $f \in \mathcal{F}$ that “look like” random noise, that is, they are highly correlated with $\sigma_1, \dots, \sigma_n$. Here are some properties of the Rademacher complexity.

Lemma. 1. *If $\mathcal{F} \subset \mathcal{G}$ then $\text{Rad}_n(\mathcal{F}, Z^n) \leq \text{Rad}_n(\mathcal{G}, Z^n)$.*

2. *Let $\text{conv}(\mathcal{F})$ denote the convex hull of \mathcal{F} . Then $\text{Rad}_n(\mathcal{F}, Z^n) = \text{Rad}_n(\text{conv}(\mathcal{F}), Z^n)$.*

3. *For any $c \in \mathbb{R}$, $\text{Rad}_n(c\mathcal{F}, Z^n) = |c| \text{Rad}_n(\mathcal{F}, Z^n)$.*

4. *Let $g : \mathbb{R} \rightarrow \mathbb{R}$ be such that $|g(y) - g(x)| \leq L|x - y|$ for all x, y . Then $\text{Rad}_n(g \circ \mathcal{F}, Z^n) \leq L \text{Rad}_n(\mathcal{F}, Z^n)$.*

5. *Suppose $\{\mathcal{F}_i\}_{i \in I}$ satisfies $0 \in \mathcal{F}_i$ for each $i \in I$. Then $\text{Rad}_n(\bigcup_{i \in I} \mathcal{F}_i, Z^n) \leq \sum_{i \in I} \text{Rad}_n(\mathcal{F}_i, Z^n)$.*

1.3 Two Layer Neural Networks

A two-layer neural network takes an input vector of d variables $x = (x_1, x_2, \dots, x_d)$ and builds a nonlinear function $f(x)$ to predict the response $y \in \mathbb{R}^D$. What distinguishes neural networks from other nonlinear methods is the particular structure of the model:

$$f(x) = f_\theta(x) = g \left(\beta_0 + \sum_{j=1}^m \beta_j \sigma(b_j + w_j^\top x) \right),$$

where $x \in \mathbb{R}^d, b_j \in \mathbb{R}, w_j \in \mathbb{R}^d, \beta_0 \in \mathbb{R}^D, \beta_j \in \mathbb{R}^D$. See Figure 1.

- $\theta = \{[\beta, a_j, b_j, w_j] : j = 1, \dots, m\}$ denotes the set of model parameters.
- x_1, \dots, x_d together is called an input layer.
- $A_j := \sigma_j(x) = \sigma(b_j + w_j^\top x)$ is called an activation.
- A_1, \dots, A_m together is called a hidden layer or hidden unit; m is the number of hidden nodes.
- $f(x)$ is called an output layer.
- g is an output function. Examples are:
 - softmax $g_i(x) = \exp(x_i) / \sum_{l=1}^D \exp(x_l)$ for classification. The softmax function estimates the conditional probability $g_i(x) = P(y = i|x)$.

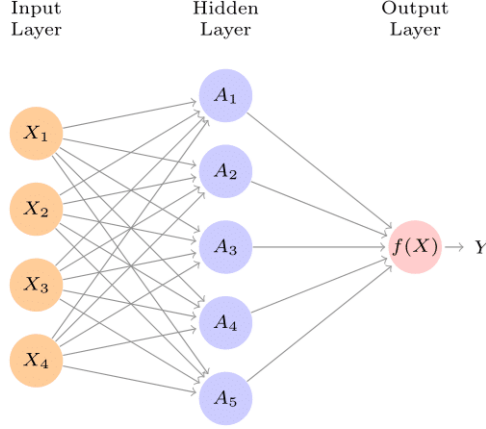


Figure 1: Neural network with a single hidden layer. The hidden layer computes activations $A_j = \sigma_j(x)$ that are nonlinear transformations of linear combinations of the inputs x_1, \dots, x_d . Hence these A_j are not directly observed. The functions σ_j are not fixed in advance, but are learned during the training of the network. The output layer is a linear model that uses these activations A_j as inputs, resulting in a function $f(x)$. Figure 10.1 from [2].

- identity/linear $g(x) = x$ for regression.
- threshold $g_i(x) = I(x_i > 0)$
- σ is called an activation function. Examples are:
 - sigmoid $\sigma(x) = 1/(1 + e^{-x})$ (see Figure 2)
 - rectified linear (ReLU) $\sigma(x) = \max\{0, x\}$ (see Figure 2)
 - identity/linear $\sigma(x) = x$
 - threshold $\sigma(x) = I(x > 0)$, threshold gives a direct multi-layer extension of the perceptron (as considered by Rosenblatt).

Activation functions in hidden layers are typically nonlinear, otherwise the model collapses to a linear model. So the activations are like derived features - nonlinear transformations of linear combinations of the features.

1.4 Multi Layer Neural Networks

Modern neural networks typically have more than one hidden layer, and often many units per layer. In theory a single hidden layer with a large number of units has the ability to approximate most functions. However, the learning task of discovering a good solution is made much easier with multiple layers each of modest size.

A deep neural network refers to the model allowing to have more than 1 hidden layers: given input $x \in \mathbb{R}^d$ and response $y \in \mathbb{R}^D$, to predict the response y . K -layer fully connected deep neural network is to build a nonlinear function $f(x)$ as

- Let $m^{(0)} = d$ and $m^{(K)} = D$
- Define recursively

$$\begin{aligned}
x^{(0)} &= x, \quad (x \in \mathbb{R}^{m^{(0)}}), \\
x_j^{(k)} &= \sigma(b_j^{(k)} + (w_j^{(k)})^\top x^{(k-1)}), \quad w_j^{(k)}, x^{(k-1)} \in \mathbb{R}^{m^{(k-1)}}, b_j \in \mathbb{R}^{m^{(k)}}, \quad k = 1, \dots, K. \\
f(x) &= g(x^{(K)}).
\end{aligned}$$

- $\theta = \{[b_j^{(k)}, w_j^{(k)}] : k = 1, \dots, K, j = 1, \dots, m^{(k)}\}$ denotes the set of model parameters.
- $m^{(k)}$ is the number of hidden units at layer k .

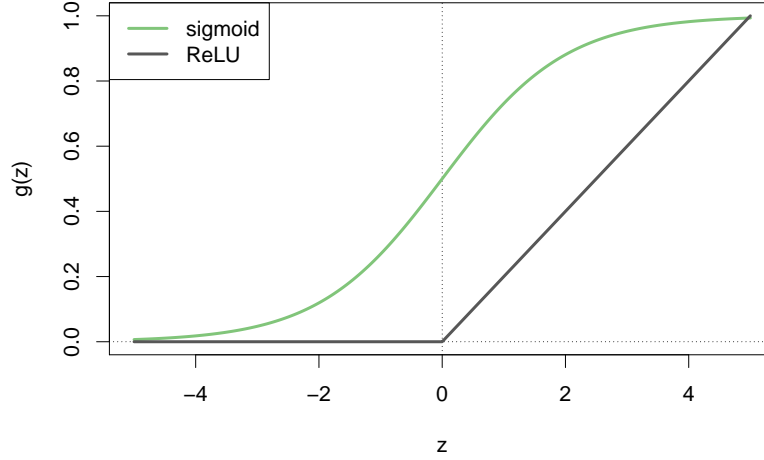


Figure 2: Activation functions. The piecewise-linear ReLU function is popular for its efficiency and computability. We have scaled it down by a factor of five for ease of comparison. Figure 10.2 from [2].

2 Notation and Goal

From here, we only consider regression problem, so $g(x) = x$. We assume $\beta_0 = 0$ and $b_j = 0$.

For the two-layer neural network with the width of the hidden layer m and activation function σ , the function space we consider is

$$\mathcal{F}_{m,\sigma} = \left\{ f_\theta : f_\theta(x) = \sum_{j=1}^m \beta_j \sigma(w_j^\top x) \right\},$$

and if we consider all two-layer neural network with arbitrary width, then

$$\mathcal{F}_\sigma = \bigcup_{m=1}^{\infty} \mathcal{F}_{m,\sigma} = \left\{ f_\theta : f_\theta(x) = \sum_{j=1}^m \beta_j \sigma(w_j^\top x), m \in \mathbb{N} \right\}.$$

Suppose the true regression function f_* is in a function class \mathcal{M} , so

$$y \approx f_*(x), \quad f_* \in \mathcal{M}.$$

Suppose are using the ℓ_2 -loss, so we find f among deep neural network class \mathcal{F} that minimizes the expected risk (평균위험),

$$f^0 = \arg \min_{f \in \mathcal{F}} \mathbb{E}_{(Y,X) \sim P} [(y - f(x))^2].$$

f_0 is the expected risk minimizing function (평균위험최소함수). And we estimate f^0 by \hat{f} using data by minimizes on the empirical risk (경험위험) on training dataset, so

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

\hat{f} is the empirical risk minimizing function (경험위험최소함수). And we set \tilde{f} be the approximation of \hat{f} by optimization (최적화); \tilde{f} is the learned function (학습된 함수).

So there are three sources of errors: approximation error, generalization error, and optimization error. See Figure 3.

$$f_* - \tilde{f} = \underbrace{f_* - f^0}_{\text{approximation error}} + \underbrace{f^0 - \hat{f}}_{\text{generalization error}} + \underbrace{\hat{f} - \tilde{f}}_{\text{optimization error}}.$$

We focus on approximation error and generalization error. What we would like to achieve is that:

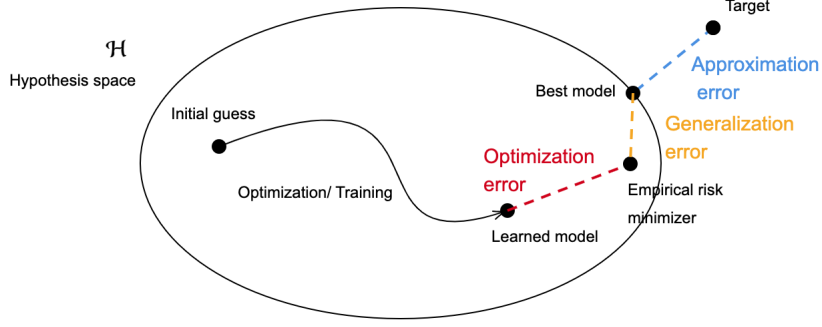


Figure 3: Diagram representing the learning procedure, the three main paradigms and their corresponding errors. Figure 2 from <https://dcn.nat.fau.eu/breaking-the-curse-of-dimensionality-with-barron-spaces/>.

For the approximation error: we would like to control $\|f_* - f^0\|_{L^2(P)}$ appropriately in terms of the width of the neural network m . Ideally, we would like to restrict the function class \mathcal{M} where f_* comes from, and define an appropriate norm $\|f_*\|_*$, so that

$$\inf_{f \in \mathcal{F}_{m,h}} \|f_* - f^0\|_{L^2(P)}^2 \lesssim \frac{\|f_*\|_*^2}{m}.$$

For the generalization error: we have seen from the concentration lecture note that with probability at least $1 - \delta$,

$$\sup_{f \in \mathcal{F}_{m,h}} \left| \frac{1}{n} \sum_{i=1}^n f(x_i) - \mathbb{E}[f] \right| \leq 2\text{Rad}(\mathcal{F}_{m,\sigma}) + \sqrt{\frac{1}{2n} \log \left(\frac{2}{\delta} \right)}.$$

Hence we would like to see that, with appropriate norm $\|f\|_{**}$ for $f \in \mathcal{F}_{m,\sigma}$, define $\mathcal{F}_{m,\sigma,Q} := \{f \in \mathcal{F}_{m,\sigma} : \|f\|_{**} \leq Q\}$, and then

$$\text{Rad}(\mathcal{F}_{m,\sigma,Q}) \lesssim \frac{Q}{\sqrt{n}}.$$

If both holds, then

$$\|f_* - \hat{f}\|_{L^2(P)}^2 = O_P \left(\frac{\|f_*\|_*^2}{m} + \frac{Q}{\sqrt{n}} \right).$$

3 Generalization error: two layer network

We now compute a bound for the Rademacher complexity of two-layer neural networks.

Theorem 2. For some constants $B_w > 0$ and $B_\beta > 0$, let

$$\mathcal{F}_{m,\sigma,B} = \{f_\theta \in \mathcal{F}_{m,\sigma} : \|\beta\|_2 \leq B_\beta, \|w_j\|_2 \leq B_w, j = 1, \dots, m\},$$

and suppose $\|Z_i\|_2 \leq C$ for all $i = 1, \dots, n$. Let σ be 1-Lipschitz. Then,

$$\text{Rad}(\mathcal{F}_{m,\sigma,B}; Z^n) \leq 2B_\beta B_w C \sqrt{\frac{m}{n}}.$$

This bound is not ideal as it depends on the number of neurons m . Empirically, it has been found that the generalization error does not increase monotonically with m . As more neurons are added to the model, thereby giving it more expressive power, studies have shown that generalization is improved [Belkin et al., 2019]. This contradicts the bound above, which states that more neurons leads to worse generalization.

Next, we look at a finer bound that results from dening a new complexity measure. A recurring theme in subsequent proofs will be the functional invariance of two-layer neural networks under a class of rescaling transformations. The key ingredient will be the positive homogeneity of the ReLU function, i.e.,

$$\alpha \sigma(x) = \sigma(\alpha x), \quad \forall \alpha > 0.$$

This implies that for any $\lambda_i > 0$ ($i = 1, \dots, m$), the transformation $\theta = \{(\beta_j, w_j)\}_{1 \leq j \leq m} \mapsto \theta' = \{(\lambda_j \beta_j, w_j / \lambda_j)\}_{1 \leq j \leq m}$ has no net effect on the neural network's functionality (i.e., $f_\theta = f_{\theta'}$) since

$$\beta_j \cdot \phi(w_j^\top x) = (\lambda_j \beta_j) \cdot \phi\left(\left(\frac{w_j}{\lambda_j}\right)^\top x\right).$$

In light of this, we devise a new complexity measure $\|\cdot\|_1$ that is also invariant under such transformations and use it to prove a better bound for the Rademacher complexity. For $f_\theta \in \mathcal{F}_{m,\sigma}$, we can write $f_\theta(x) = \sum_{j=1}^m \beta_j \sigma(w_j^\top x)$. Define a complexity of θ as

$$C(\theta) := \sum_{j=1}^m |\beta_j| \|w_j\|_2.$$

Theorem 3. *For some constant $B > 0$ consider the function class*

$$\mathcal{F}_{m,\sigma,B} = \{f_\theta \in \mathcal{F}_{m,\sigma} : C(\theta) \leq B\}. \quad (1)$$

If $\|Z_i\|_2 \leq C$ for all $i = 1, \dots, n$. Let σ be ReLU function, then

$$\text{Rad}(\mathcal{F}_{m,\sigma,B}; Z^n) \leq \frac{2BC}{\sqrt{n}}.$$

Remark 4. Compared to Theorem 2, this bound does not explicitly depend on the number of neurons m . Thus, it is possible to use more neurons and still maintain a tight bound if the value of the new complexity measure $\|\theta\|_1$ is reasonable. In contrast, the bound of Theorem 2 explicitly grows with the total number of neurons.

Moreover, Theorem 3 is stronger as we have more neurons - this is because the function class $\mathcal{F}_{m,\sigma,B}$ as defined in 1 is bigger as m increases. Because of this, it's possible to obtain a generalization guarantee that decreases as m increases, as we will see later.

Proof. Due to the positive homogeneity of the ReLU function, it will be useful to define the ℓ_2 -normalized weight vector $\bar{w}_j := w_j / \|w_j\|_2$ so that $\phi(w_j^\top x) = \|w_j\|_2 \phi(\bar{w}_j^\top x)$. Let $\xi_i = \pm 1$ being i.i.d. with probability 1/2 be Rademacher variables, then the empirical Rademacher complexity satisfies

$$\begin{aligned} \text{Rad}(\mathcal{F}_{m,\sigma,B}; Z^n) &= \mathbb{E}_\xi \left[\sup_{f_\theta \in \mathcal{F}_{m,\sigma,B}} \frac{1}{n} \sum_{i=1}^n \xi_i f_\theta(Z_i) \middle| Z \right] \\ &= \mathbb{E}_\xi \left[\sup_{\theta: C(\theta) \leq B} \frac{1}{n} \sum_{i=1}^n \xi_i \sum_{j=1}^m \beta_j \sigma(w_j^\top Z_i) \middle| Z \right] \\ &= \frac{1}{n} \mathbb{E}_\xi \left[\sup_{\theta: C(\theta) \leq B} \sum_{i=1}^n \xi_i \sum_{j=1}^m \beta_j \|w_j\|_2 \sigma(\bar{w}_j^\top Z_i) \middle| Z \right] \quad (\text{by positive homogeneity of } \sigma) \\ &= \frac{1}{n} \mathbb{E}_\xi \left[\sup_{\theta: C(\theta) \leq B} \sum_{j=1}^m \beta_j \|w_j\|_2 \sum_{i=1}^n \xi_i \sigma(\bar{w}_j^\top Z_i) \middle| Z \right] \\ &\leq \frac{1}{n} \mathbb{E}_\xi \left[\sup_{\theta: C(\theta) \leq B} \sum_{j=1}^m \beta_j \|w_j\|_2 \max_{1 \leq k \leq m} \left| \sum_{i=1}^n \xi_i \sigma(\bar{w}_k^\top Z_i) \right| \middle| Z \right], \end{aligned}$$

since $\sum_j \alpha_j \beta_j \leq \sum_j |\alpha_j| \max_k |\beta_k|$. Then from $C(\theta) \leq B$, we can further bound as

$$\begin{aligned}
\text{Rad}(\mathcal{F}_{m,\sigma,B}; Z^n) &\leq \frac{B}{n} \mathbb{E}_\xi \left[\sup_{\theta: C(\theta) \leq B} \max_{1 \leq k \leq m} \left\| \sum_{i=1}^n \xi_i \sigma(\bar{w}_k^\top Z_i) \right\| \middle| Z \right] \\
&= \frac{B}{n} \mathbb{E}_\xi \left[\sup_{\bar{w}: \|\bar{w}\|_2 = 1} \left\| \sum_{i=1}^n \xi_i \sigma(\bar{w}^\top Z_i) \right\| \middle| Z \right] \\
&\leq \frac{B}{n} \mathbb{E}_\xi \left[\sup_{\bar{w}: \|\bar{w}\|_2 \leq 1} \left\| \sum_{i=1}^n \xi_i \sigma(\bar{w}^\top Z_i) \right\| \middle| Z \right] \\
&\leq \frac{2B}{n} \mathbb{E}_\xi \left[\sup_{\bar{w}: \|\bar{w}\|_2 \leq 1} \sum_{i=1}^n \xi_i \sigma(\bar{w}^\top Z_i) \middle| Z \right] \\
&= 2B \text{Rad}(\mathcal{H}; Z^n),
\end{aligned}$$

where the last inequality is from Lemma 1, and $\mathcal{H} = \{x \mapsto \sigma(\bar{w}^\top x) : \bar{w} \in \mathbb{R}^d, \|\bar{w}\|_2 \leq 1\}$. Since the ReLU function σ is 1-Lipschitz, $\text{Rad}(\mathcal{H}; Z^n) \leq \text{Rad}(\mathcal{H}'; Z^n)$, where $\mathcal{H}' = \{x \mapsto \bar{w}^\top x : \bar{w} \in \mathbb{R}^d, \|\bar{w}\|_2 \leq 1\}$. Then $\text{Rad}(\mathcal{H}'; Z^n) \leq \frac{C}{\sqrt{n}}$ from below concludes the proof. \square

Proposition 5. Let $\mathcal{F} = \{x \mapsto \langle w, x \rangle : w \in \mathbb{R}^d, \|w\|_2 \leq B\}$ for some constant $B > 0$, and suppose $\|Z_i\|_2 \leq C$ for all $i = 1, \dots, n$. Then

$$\text{Rad}(\mathcal{F}; Z^n) \leq \frac{BC}{\sqrt{n}}.$$

Proof. Left as HW. \square

Then as suggested in the concentration lecture note, with probability at least $1 - \delta$,

$$\begin{aligned}
\sup_{f \in \mathcal{F}_{m,\sigma,B}} \left| \frac{1}{n} \sum_{i=1}^n f(x_i) - \mathbb{E}[f] \right| &\leq 2\text{Rad}(\mathcal{F}_{m,\sigma,B}) + \sqrt{\frac{1}{2n} \log \left(\frac{2}{\delta} \right)} \\
&\leq \frac{4BC}{\sqrt{n}} + \sqrt{\frac{1}{2n} \log \left(\frac{2}{\delta} \right)}.
\end{aligned}$$

There is a direct result with Barron class as well.

Theorem ([1, Theorem 15]). Let $\mathcal{F}_{m,\sigma,Q} := \{f \in \mathcal{F}_{m,\sigma} : \|f\|_B \leq Q\}$. Then we have

$$\text{Rad}(\mathcal{F}_{m,\sigma,Q}; Z^n) \leq 2Q \sqrt{\frac{2 \log(2d)}{n}}.$$

Instead of minimizing the training error, we can also consider the regularized term. For $f_\theta \in \mathcal{F}_{m,\sigma}$, we can write $f_\theta(x) = \sum_{j=1}^m \beta_j \sigma(w_j^\top x)$. Define the 1-norm of θ as

$$\|\theta\|_1 := \frac{1}{m} \sum_{j=1}^m |\beta_j| \|w_j\|_1.$$

And consider the minimization problem

$$\mathcal{L}(\theta) = \frac{1}{n} \sum_{i=1}^n (y_i - f_\theta(x_i))^2 + \lambda \sqrt{\frac{\log(2d)}{n}} \|\theta\|_1,$$

and let $\hat{\theta}^{(1)}$ be its minimizer.

Theorem ([1, Theorem 16]). Suppose $\mathcal{X} \subset \mathbb{R}^d$ is compact, and assume $f_* : \mathcal{X} \rightarrow [0, 1]$. There exists some $\lambda_0 > 0$ such that for $\lambda \geq \lambda_0$, with probability $1 - \delta$,

$$\frac{1}{n} \sum_{i=1}^n (y_i - f_{\hat{\theta}^{(1)}}(x_i))^2 \lesssim \frac{\|f_*\|_B^2}{m} + \lambda \|f_*\|_B \sqrt{\frac{\log(2d)}{n}} + \sqrt{\frac{\log(n/\delta)}{n}}.$$

References

- [1] Weinan E, Chao Ma, Stephan Wojtowytsch, and Lei Wu. Towards a mathematical understanding of neural network-based machine learning: what we know and what we don't. *CoRR*, abs/2009.10713, 2020.
- [2] Gareth James, Daniela Witten, Trevor Hastie, and Robert Tibshirani. *An introduction to statistical learning— with applications in R*. Springer Texts in Statistics. Springer, New York, [2021] ©2021. Second edition [of 3100153].