

# Non-parametric Estimation via Neural Network

Suh, Ko, Huo

August 14, 2020

## 1 Preliminary notations

In this document, we consider a non-parametric regression model,  $Y_i = f_0(X_i) + \varepsilon_i$ , where  $\varepsilon_i \sim \mathcal{N}(0, \sigma^2)$ . We are interested in finding the bound of prediction risk

$$R(\hat{f}_n, f_0) := E_{f_0} \left[ \left( \hat{f}_n(X) - f_0(X) \right)^2 \right],$$

with  $X \stackrel{D}{=} X_1$  being independent of sample  $(X_i, Y_i)$ . Let  $\hat{f}_n$  denote any local or global minimizer of empirical risk  $\frac{1}{n} \sum_{i=1}^n (Y_i - f(X_i))^2$  over the neural network class  $\mathcal{F}(L, p, s, F)$ . Note also that the subscript  $f_0$  in  $E_{f_0}$  indicates that the expectation is taken with respect to a sample generated from  $f_0$ . The sequence  $\Delta_n(\hat{f}_n, f_0)$  measures the difference between the expected empirical risk of  $\hat{f}_n$  and the global minimum over all networks in the class  $\mathcal{F}(L, p, s, F)$ , which is defined as follows:

$$\Delta_n(\hat{f}_n, f_0) := E_{f_0} \left[ \frac{1}{n} \sum_{i=1}^n (Y_i - \hat{f}_n(X_i))^2 - \inf_{f \in \mathcal{F}(L, p, s, F)} \frac{1}{n} \sum_{i=1}^n (Y_i - f(X_i))^2 \right].$$

Note that  $\Delta_n(\hat{f}_n, f_0) \geq 0$  and  $\Delta_n(\hat{f}_n, f_0) = 0$  if  $\hat{f}_n$  is an empirical minimizer. The goal of this document is to obtain both lower and upper bounds of the prediction risk in terms of  $\Delta_n(\hat{f}_n, f_0)$ .

## 2 Risk Bound w.r.t. $\Delta_n$

In this Section, we provide the Theorem on upper and lower bound of prediction risk  $R(\hat{f}_n, f_0)$ .

**Theorem 2. (Hieber, 20.)** Consider the d-variate non-parametric regression model with unknown regression function  $f_0$  satisfying  $\|f_0\|_\infty \leq F$  for some  $F \geq 1$ . Let  $\hat{f}_n$  be any estimator taking values in the class  $\mathcal{F}(L, p, s, F)$  and let  $\Delta_n(\hat{f}_n, f_0)$  be the quantity defined above. For any  $\varepsilon \in (0, 1]$ , there exists a constant  $C_\varepsilon$ , only depending on  $\varepsilon$ , such that with

$$\tau_{\varepsilon, n} := C_\varepsilon F^2 \frac{(s+1) \log(n(s+1)^L) p_0 p_{L+1}}{n},$$

$$\begin{aligned} (1 - \varepsilon)^2 \Delta_n(\hat{f}_n, f_0) - \tau_{\varepsilon, n} &\leq R(\hat{f}_n, f_0) \\ &\leq (1 + \varepsilon)^2 \left( \inf_{f \in \mathcal{F}(L, p, s, F)} E_{f_0} [\|f - f_0\|_n^2] + \Delta_n(\hat{f}_n, f_0) \right) + \tau_{\varepsilon, n}. \end{aligned}$$

**Proof of the Theorem 2. is involved with following two Lemmas : Lemma 4 and Lemma 5.**

**Lemma 4. (Hieber, 20.)** Consider the d-variate non-parametric regression model with unknown regression function  $f_0$  satisfying  $\|f_0\|_\infty \leq F$  for some  $F \geq 1$ . Let  $\widehat{f}_n$  be any estimator taking values in the class  $\mathcal{F}(L, p, s, F)$  and let  $\Delta_n(\widehat{f}_n, f_0)$  be the quantity defined above. Assume  $\{f_0\} \cup \mathcal{F} \subset \{f : [0, 1]^d \rightarrow [-F, F]\}$  for some  $F \geq 1$ . If  $\mathcal{N}_n := \mathcal{N}(\delta, \mathcal{F}, \|\cdot\|_\infty) \geq 3$ , then,

$$(1 - \varepsilon)^2 \Delta_n - F^2 \frac{18 \log \mathcal{N}_n + 76}{n\varepsilon} - 38\delta \leq R(\widehat{f}, f_0) \leq (1 + \varepsilon)^2 \left[ \inf_{f \in \mathcal{F}} E_{f_0}[\|f - f_0\|_n^2] + F^2 \frac{18 \log \mathcal{N}_n + 72}{n\varepsilon} + 32\delta F + \Delta_n \right]$$

for all  $\delta, \varepsilon \in (0, 1]$ .

**Proof sketch of Lemma 4.** Proof of Lemma 4 can be divided into several substeps. Here, we present key ideas for each substep, and all detailed proof procedures are deferred in next Section.

1. First key idea for the proof of Lemma 4 is to construct an empirical counterpart of prediction risk, where we denote it as  $\widehat{R}_n(\widehat{f}_n, f_0) := E_X[\frac{1}{n} \sum_{j=1}^n (\widehat{f}_n(X_i) - f_0(X_i))^2]$ , and relate it with the risk  $R_n(\widehat{f}_n, f_0)$  via certain inequality. Note that since  $\widehat{f}_n$  is a stochastic term,  $R(\widehat{f}_n, f_0) \neq \widehat{R}_n(\widehat{f}_n, f_0)$ . Generate i.i.d. random variables  $X'_1, \dots, X'_n$  with the same distribution as  $X$  and independent of  $(X_i)_{i=1, \dots, n}$ . Then, we can obtain a following inequality :

$$\begin{aligned} |R(\widehat{f}_n, f_0) - \widehat{R}_n(\widehat{f}_n, f_0)| &= \left| E \left[ \frac{1}{n} \sum_{i=1}^n \left\{ (\widehat{f}_n(X'_i) - f_0(X'_i))^2 - (\widehat{f}_n(X_i) - f_0(X_i))^2 \right\} \right] \right| \\ &\leq \frac{F}{n} R(\widehat{f}, f_0)^{1/2} (36n \log \mathcal{N}_n + 2^8 n) + \frac{F}{n} (6 \log \mathcal{N}_n + 11) + 26\delta F. \end{aligned}$$

Bernstein's inequality and concept of  $\delta$ -covering number are used for derivation.

2. Second key idea is to obtain both the lower and upper bound of risk  $R(\widehat{f}_n, f_0)$  in terms of  $\widehat{R}_n(\widehat{f}_n, f_0)$  by using a following inequality : Let  $a, b, c, d$  be positive real numbers, such that  $|a - b| \leq 2\sqrt{ac} + d$ . It can be easily shown for any  $\varepsilon \in (0, 1]$ ,

$$(1 - \varepsilon)b - d - \frac{c^2}{\varepsilon} \leq a \leq (1 + \varepsilon)(b + d) + \frac{(1 + \varepsilon)^2}{\varepsilon} c^2.$$

In our case, we set  $a = R(\widehat{f}, f_0)$ ,  $b = \widehat{R}_n(\widehat{f}, f_0)$ ,  $c = F(9n \log \mathcal{N}_n + 64n)^{1/2}/n$ , and  $d = F(6 \log \mathcal{N}_n + 11)/n + 26\delta F$ .

3. Lastly, a remaining task for obtaining the bounds in Lemma 4. is to get the upper and lower bound of empirical version of prediction risk,  $\widehat{R}_n(\widehat{f}_n, f_0)$ , with respect to  $\Delta_n$ . Here, we use the definition of  $\Delta_n$ . For any fixed  $f \in \mathcal{F}$ ,

$$E_{f_0} \left[ \frac{1}{n} \sum_{j=1}^n (\widehat{f}(X_i) - f_0(X_i))^2 \right] \leq E_{f_0} \left[ \frac{1}{n} \sum_{j=1}^n \left( f_0(X_i) - f(X_i) \right)^2 \right] + \Delta_n + \frac{2}{n} E_{f_0} \left[ \sum_{i=1}^n \varepsilon_i \widehat{f}(X_i) \right].$$

Observe that  $\delta$ -entropy of class  $\mathcal{F}$  (i.e.,  $\log \mathcal{N}_n$ ) appears both at the upper and lower bound of the risk  $R(\widehat{f}, f_0)$  in Lemma 4. In order to derive the bounds in Theorem 2, an upper bound on the covering number,  $\mathcal{N}_n$ , should be obtained.

**Lemma 5. (Hieber, 20.)** If  $V := \prod_{\ell=0}^{L+1} (p_\ell + 1)$ , then for any  $\delta > 0$ ,

$$\log \mathcal{N}\left(\delta, \mathcal{F}(L, p, s, \infty), \|\cdot\|_\infty\right) \leq (s+1) \log\left(2\delta^{-1}(L+1)V^2\right).$$

**Proof sketch of Lemma 5.** The key idea for the proof of Lemma 5 is to discretize the parameter space  $(\mathbf{v}_k, \mathbf{W}_k)_k$  of neural network  $f$  with certain grid size  $\varepsilon$ . Let  $(\mathbf{v}_k^*, \mathbf{W}_k^*)_k$  be the network parameter on the constructed grid and  $f^*$  be the network value evaluated at the grid point,  $(\mathbf{v}_k^*, \mathbf{W}_k^*)_k$ . Fix  $\varepsilon > 0$ . Let  $f, f^* \in \mathcal{F}$  be two networks, such that all parameters are  $\varepsilon$  away from each other. Then, the proof of Lemma 5 can be divided into following four steps :

1. Calculate the upper-bound of  $|f(X) - f^*(X)|$  in terms of  $\varepsilon, V$  and  $L$ . Actually, we can get the bound  $|f(X) - f^*(X)| \leq \varepsilon V(L+1)$ . By setting  $\varepsilon = \frac{\delta}{2V(L+1)}$ , for any  $\delta > 0$ , we can get  $|f(X) - f^*(X)| \leq \frac{\delta}{2} < \delta$ .
2. Note that the total number of parameters in the network can be bounded by  $V$  as follows:

$$\sum_{j=0}^L p_j p_{j+1} + \sum_{j=1}^L p_j = \sum_{j=0}^L (p_j + 1) p_{j+1} - p_{L+1} \leq \prod_{\ell=0}^{L+1} (p_\ell + 1) := V.$$

3. Since the absolute value of parameters in the network are bounded by 1, it is obvious that the total number of grid points in parameter space  $(\mathbf{v}_k, \mathbf{W}_k)_k$  is bounded by  $\frac{2V^2(L+1)}{\delta}$ .
4. Since  $\mathcal{F}(L, p, s, F)$  is a class of neural network functions with at most  $s$  parameters, we can obtain a bound for the covering number of class  $\mathcal{F}$  as follows,

$$\mathcal{N}\left(\delta, \mathcal{F}(L, p, s, \infty), \|\cdot\|_\infty\right) \leq \sum_{s^* \leq s} \left(\frac{2V^2(L+1)}{\delta}\right)^{s^*} \leq \left(\frac{2V^2(L+1)}{\delta}\right)^{s+1}.$$

All we need to do is to derive the bound of  $|f(X) - f(X^*)| \leq \varepsilon V(L+1)$  for which we defer the proof in Section 4.

### 3 Proof of Lemma 4.

#### 3.1 Step 0.

First, we divide the cases into two: 1) when  $\log \mathcal{N}_n \geq n$  and 2) when  $\log \mathcal{N}_n \leq n$ . In this step, when  $\log \mathcal{N}_n \geq n$  holds, we will show that the upper and lower bounds in Lemma 4. hold trivially. For the upper bound, note that the risk  $R(\hat{f}, f_0)$  is bounded by  $4F^2$ . The assumptions  $\log \mathcal{N}_n \geq n$  and  $\varepsilon \in (0, 1]$  make the claimed upper bound in Lemma 4. already greater than  $18F^2$ . For the lower bound, first observe that  $R(\hat{f}, f_0) \geq 0$ , then it is enough to show that when  $\log \mathcal{N}_n \geq n$ , the claimed lower bound in Lemma 4 is less than or equal to 0. It is easy to see that the following is true, for any fixed empirical risk minimizer,  $\hat{f}$ , we have :

$$\begin{aligned} \hat{R}_n(\hat{f}_n, f_0) - \hat{R}_n(\tilde{f}_n, f_0) &= E_{f_0} \left[ \frac{1}{n} \sum_{i=1}^n (\hat{f}(X_i) - f_0(X_i))^2 - (\tilde{f}(X_i) - f_0(X_i))^2 \right] \\ &= E_{f_0} \left[ \frac{1}{n} \sum_{i=1}^n (\hat{f}(X_i) - Y_i + \varepsilon_i)^2 - (\tilde{f}(X_i) - Y_i + \varepsilon_i)^2 \right] \\ &= \Delta_n + E_{f_0} \left[ \frac{2}{n} \sum_{i=1}^n \varepsilon_i \hat{f}(X_i) \right] - E_{f_0} \left[ \frac{2}{n} \sum_{i=1}^n \varepsilon_i \tilde{f}(X_i) \right]. \end{aligned}$$

From this equation, we can see that  $\Delta_n$  is bounded by  $8F^2$ , since  $F \geq 1$ . The lower bound of Lemma 4 is bounded by  $(1 - \varepsilon)^2 \Delta_n - 18F^2 \leq -10F^2 \leq 0$ . So the lower bound holds trivially.

### 3.2 Step 1.

From this step, we assume  $\log \mathcal{N}_n \leq n$ . Recall that the definition of  $\delta$ -covering number  $\mathcal{N}(\delta, \mathcal{F}(L, p, s, F), \|\cdot\|_\infty)$  is the minimal number of  $\|\cdot\|_\infty$ -balls with radius  $\delta$  that covers  $\mathcal{F}(L, p, s, F)$ . By construction there exists a (random)  $j^*$  such that  $\|\hat{f} - f_{j^*}\|_\infty \leq \delta$  for any  $\hat{f}$  achieved by any optimization methods. Generate i.i.d. random variables  $X'_1, \dots, X'_n$  with the same distribution as  $X$  and independent of  $(X_i)_{i=1, \dots, n}$ , then we can get a following inequality under the assumptions:  $\|f_j\|_\infty, \|f_0\|_\infty, \delta \leq F$ .

$$\begin{aligned}
\left| R(\hat{f}, f_0) - R_n(\hat{f}, f_0) \right| &= \left| E_{f_0} \left[ \frac{1}{n} \sum_{i=1}^n (\hat{f}(X'_i) - f_0(X'_i))^2 - (\hat{f}(X_i) - f_0(X_i))^2 \right] \right| \\
&\leq E_{f_0} \left| \frac{1}{n} \sum_{i=1}^n \underbrace{(f_{j^*}(X'_i) - f_0(X'_i))^2 - (f_{j^*}(X_i) - f_0(X_i))^2}_{:= g_{j^*}(X_i, X'_i)} \right| + 9\delta F \\
&= E_{f_0} \left[ \left| \frac{1}{n} \sum_{i=1}^n \frac{g_{j^*}(X_i, X'_i)}{r_{j^*} F} (r_{j^*} F) \right| \right] + 9\delta F \\
&\leq \frac{F}{n} E_{f_0} \left[ \left| \max_j \sum_{i=1}^n \left( \frac{g_j(X_i, X'_i)}{r_j F} \right) \right| r_{j^*} \right] + 9\delta F
\end{aligned}$$

For second inequality, we use a following relation :

$$\begin{aligned}
&(\hat{f}(X'_i) - f_0(X'_i))^2 - (\hat{f}(X_i) - f_0(X_i))^2 \\
&= (\hat{f}(X'_i) - f_{j^*}(X'_i) + f_{j^*}(X'_i) - f_0(X'_i))^2 - (\hat{f}(X_i) - f_{j^*}(X_i) + f_{j^*}(X_i) - f_0(X_i))^2 \\
&= (f_{j^*}(X'_i) - f_0(X'_i))^2 - (f_{j^*}(X_i) - f_0(X_i))^2 + \underbrace{(\hat{f}(X'_i) - f_{j^*}(X'_i))^2}_{\leq \delta^2 \leq \delta F} - \underbrace{(\hat{f}(X_i) - f_{j^*}(X_i))^2}_{\geq 0} \\
&\quad - 2 \underbrace{(\hat{f}(X'_i) - f_{j^*}(X'_i))}_{\leq \delta} \underbrace{(f_{j^*}(X'_i) - f_0(X'_i))}_{\geq -2F} - 2 \underbrace{(\hat{f}(X_i) - f_{j^*}(X_i))}_{\leq \delta} \underbrace{(f_{j^*}(X_i) - f_0(X_i))}_{\geq -2F} \\
&\leq (f_{j^*}(X'_i) - f_0(X'_i))^2 - (f_{j^*}(X_i) - f_0(X_i))^2 + 9\delta F.
\end{aligned}$$

Now we set  $r_{j^*} := \sqrt{n^{-1} \log \mathcal{N}_n} \vee E^{1/2}[(f_{j^*}(X) - f_0(X))^2 | (X_i, Y_i)_i]$ . Note that we have a following bound for  $r_{j^*}$  by using the notion of  $\delta$ -covering number once again:

$$r_{j^*} \leq \sqrt{\frac{\log \mathcal{N}_n}{n}} + E^{1/2} \left[ (\hat{f}(X) - f_0(X))^2 | (X_i, Y_i)_i \right] + \delta$$

Denoting  $T := \left| \max_j \sum_{i=1}^n g_j(X_i, X'_i) / r_j F \right|$  and  $U := E^{1/2}[(\hat{f}(X) - f_0(X))^2 | (X_i, Y_i)_i]$ , and using Cauchy-Schwartz inequality for random variables  $U$  and  $T$  (i.e.,  $E[UT] \leq E^{1/2}[U^2] E^{1/2}[T^2]$ ), we have:

$$\begin{aligned}
\left| R(\hat{f}, f_0) - R_n(\hat{f}, f_0) \right| &\leq \frac{F}{n} E_{f_0} \left[ TU + \left( \delta + \sqrt{\frac{\log \mathcal{N}_n}{n}} \right) T \right] + 9\delta F \\
&\leq \frac{F}{n} R(\hat{f}, f_0)^{1/2} E^{1/2}[T^2] + \frac{F}{n} \left( \delta + \sqrt{\frac{\log \mathcal{N}_n}{n}} \right) E[T] + 9\delta F. \tag{1}
\end{aligned}$$

Now, the first and second moment of  $T$  need to be controlled. (i.e.,  $E[T]$  and  $E[T^2]$ ). First observe that  $Eg_j(X_i, X'_i)/r_j = 0$ ,  $\text{Var}(g_j(X_i, X'_i)/r_j) \leq 8F^2$  and  $|g_j(X_i, X'_i)/r_j| \leq 4F^2/r_j$ . We can use Bernstein's inequality which states that for independent and centered random variables  $U_1, \dots, U_n$ , satisfying  $|U_i| \leq M$ ,

$$P(|\sum_{i=1}^n U_i| \geq t) \leq 2 \exp(-t^2/[2Mt/3 + 2 \sum_{i=1}^n \text{Var}(U_i)]).$$

$$\begin{aligned} P(T \geq t) &= P\left(\left|\max_j \sum_{i=1}^n \left(\frac{g_j(X_i, X'_i)}{r_j F}\right)\right| \geq t\right) \leq P\left(\max_j \left|\sum_{i=1}^n \left(\frac{g_j(X_i, X'_i)}{r_j}\right)\right| \geq Ft\right) \\ &\leq \mathcal{N}_n \max_j P\left(\left|\sum_{i=1}^n \left(\frac{g_j(X_i, X'_i)}{r_j}\right)\right| \geq Ft\right) \\ &\leq 2\mathcal{N}_n \max_j \exp\left(-\frac{(Ft)^2}{2 \cdot 4F^2t/(3r_j) + 2 \cdot 8F^2n}\right) \\ &= 2\mathcal{N}_n \max_j \exp\left(-\frac{t^2}{8t/(3r_j) + 16n}\right) \wedge 1 \end{aligned}$$

The first term in the denominator dominates for large enough  $t$ : Specifically, for  $t \geq 6\sqrt{n \log \mathcal{N}_n}$ , we have

$$\exp\left(-\frac{t^2}{8t/(3r_j) + 16n}\right) \leq \exp\left(-\frac{t}{16/(3r_j)}\right) \leq \exp\left(-\frac{3t\sqrt{\log \mathcal{N}_n}}{16\sqrt{n}}\right),$$

where we use the inequality  $r_j \geq \sqrt{n^{-1} \log \mathcal{N}_n}$  in the last inequality. We therefore find following inequality for the two moments of  $T$ :

$$\begin{aligned} E[T] &= \int_0^\infty P(T \geq t) dt \leq 6\sqrt{n \log \mathcal{N}_n} + \int_{6\sqrt{n \log \mathcal{N}_n}}^\infty 2\mathcal{N}_n \exp\left(-\frac{3t\sqrt{\log \mathcal{N}_n}}{16\sqrt{n}}\right) dt \\ &\leq 6\sqrt{n \log \mathcal{N}_n} + \frac{32}{3} \sqrt{\frac{n}{\log \mathcal{N}_n}}. \end{aligned}$$

Here, in second inequality, for  $3 \leq \mathcal{N}_n$ , we use :

$$2\mathcal{N}_n \int_{6\sqrt{n \log \mathcal{N}_n}}^\infty \exp\left(-\frac{3t\sqrt{\log \mathcal{N}_n}}{16\sqrt{n}}\right) dt = \left(\mathcal{N}_n^{-\frac{7}{2}}\right) \frac{32}{3} \sqrt{\frac{n}{\log \mathcal{N}_n}} \leq \frac{32}{3} \sqrt{\frac{n}{\log \mathcal{N}_n}}.$$

Similarly  $E[T^2]$  can be controlled as follows:

$$\begin{aligned} E[T^2] &= \int_0^\infty P(T^2 \geq u) du = \int_0^\infty P(T \geq \sqrt{u}) du \\ &\leq 36n \log \mathcal{N}_n + \int_{36n \log \mathcal{N}_n}^\infty 2\mathcal{N}_n \exp\left(-\frac{3\sqrt{u}\sqrt{\log \mathcal{N}_n}}{16\sqrt{n}}\right) du \\ &\leq 36n \log \mathcal{N}_n + 2^8 n. \end{aligned}$$

In the third inequality, we use the identity :  $\int_{b^2}^\infty e^{-\sqrt{u}a} du = 2(ab+1)e^{-ab}/b^2$ . By setting  $a = \frac{3\sqrt{\log \mathcal{N}_n}}{16\sqrt{n}}$  and  $b = 6\sqrt{n \log \mathcal{N}_n}$  and using  $3 \leq \mathcal{N}_n$  once again, we have :

$$\int_{36n \log \mathcal{N}_n}^\infty 2\mathcal{N}_n \exp\left(-\frac{3\sqrt{u}\sqrt{\log \mathcal{N}_n}}{16\sqrt{n}}\right) du = \left(\frac{1}{2} + \frac{4}{9 \log \mathcal{N}_n}\right) \mathcal{N}_n^{-\frac{1}{8}} \cdot 2^8 n \leq 2^8 n.$$

Lastly, plugging the bound of  $E[T]$  and  $E[T^2]$  in (1) and using the relation  $1 \leq \log \mathcal{N}_n \leq n$ , we get :

$$\left|R(\hat{f}, f_0) - R_n(\hat{f}, f_0)\right| \leq \frac{F}{n} R(\hat{f}, f_0)^{1/2} (36n \log \mathcal{N}_n + 2^8 n)^{1/2} + \frac{F}{n} (6 \log \mathcal{N}_n + 11) + 26\delta F. \quad (2)$$

### 3.3 Step 2.

Let  $a, b, c, d$  be positive real numbers, such that  $|a - b| \leq 2\sqrt{ac} + d$ . It can be easily shown for any  $\varepsilon \in (0, 1]$ ,

$$(1 - \varepsilon)b - d - \frac{c^2}{\varepsilon} \leq a \leq (1 + \varepsilon)(b + d) + \frac{(1 + \varepsilon)^2}{\varepsilon} c^2. \quad (3)$$

Upper bound of (3) can be obtained as follows:

$$\begin{aligned} a - b &\leq 2\sqrt{ac} + d = 2\sqrt{\frac{\varepsilon}{1+\varepsilon}}a\sqrt{\frac{1+\varepsilon}{\varepsilon}}c + d \\ &\leq \frac{\varepsilon}{1+\varepsilon}a + \frac{1+\varepsilon}{\varepsilon}c^2 + d. \end{aligned}$$

Lower bound of (3) can be obtained as follows:

$$\begin{aligned} b - a &\leq 2\sqrt{ac} + d = 2\sqrt{\frac{\varepsilon}{1-\varepsilon}}a\sqrt{\frac{1-\varepsilon}{\varepsilon}}c + d \\ &\leq \frac{\varepsilon}{1-\varepsilon}a + \frac{1-\varepsilon}{\varepsilon}c^2 + d. \end{aligned}$$

Rearranging the terms and using the fact  $\varepsilon \in (0, 1]$  yield the claim. Now, set  $a = R(\hat{f}, f_0)$ ,  $b = \hat{R}_n(\hat{f}, f_0)$ ,  $c = F(9n \log \mathcal{N}_n + 64n)^{1/2}/n$ , and  $d = F(6 \log \mathcal{N}_n + 11)/n + 27\delta F$ . Finally, we relate the risk  $R(\hat{f}, f_0) = E[(\hat{f}(X) - f_0(X))^2]$  to its empirical counterpart  $\hat{R}_n(\hat{f}, f_0)$  via the inequalities,

$$\begin{aligned} (1 - \varepsilon)\hat{R}_n(\hat{f}, f_0) - \frac{F^2}{n\varepsilon}(15 \log \mathcal{N}_n + 75) - 26\delta F &\leq R(\hat{f}, f_0) \\ &\leq (1 + \varepsilon)\left(\hat{R}_n(\hat{f}, f_0) + (1 + \varepsilon)\frac{F^2}{n\varepsilon}(12 \log \mathcal{N}_n + 70) + 26\delta F\right). \end{aligned}$$

Note that we have used  $2 \leq (1 + \varepsilon)/\varepsilon$  for obtaining the upper-bound.

### 3.4 Step 3.

Given an estimator  $\tilde{f}$  taking values in  $\mathcal{F}$ , let  $j'$  be such that  $\|\tilde{f} - f_{j'}\|_\infty \leq \delta$ . First, using the independence assumption of  $\varepsilon_i$  and  $X_i$ , it can be easily seen that  $E[\varepsilon_i f_0(X_i)] = 0$  for deterministic  $f_0$ . Secondly, it is also easy to see that  $|E[\sum_{i=1}^n \varepsilon_i(\tilde{f}(X_i) - f_{j'}(X_i))]| \leq n\delta$  by using the fact  $E[|\varepsilon_i|] = \sqrt{\frac{2}{\pi}}$ . So the bound of can be obtained as follows :

$$\begin{aligned} \left| E \left[ \frac{2}{n} \sum_{i=1}^n \varepsilon_i \tilde{f}(X_i) \right] \right| &= \left| E \left[ \frac{2}{n} \sum_{i=1}^n \varepsilon_i \left( \tilde{f}(X_i) - f_0(X_i) \right) \right] \right| \\ &\leq \left| E \left[ \frac{2}{n} \sum_{i=1}^n \varepsilon_i \left( f_{j'}(X_i) - f_0(X_i) \right) \right] \right| + 2\delta \\ &= \frac{2}{\sqrt{n}} \left| E \left[ \frac{\sum_{i=1}^n \varepsilon_i (f_{j'}(X_i) - f_0(X_i))}{\sqrt{n} \|f_{j'} - f_0\|_n} \|f_j - f_0\|_n \right] \right| + 2\delta \\ &\leq \frac{2}{\sqrt{n}} E \left[ \left| \frac{\sum_{i=1}^n \varepsilon_i (f_{j'}(X_i) - f_0(X_i))}{\sqrt{n} \|f_{j'} - f_0\|_n} \right| \|f_{j'} - f_0\|_n \right] + 2\delta \\ &\leq \frac{2}{\sqrt{n}} E \left[ \underbrace{\left| \frac{\sum_{i=1}^n \varepsilon_i (f_{j'}(X_i) - f_0(X_i))}{\sqrt{n} \|f_{j'} - f_0\|_n} \right|}_{:= \xi_{j'}} (\|\tilde{f} - f_0\|_n + \delta) \right] + 2\delta \end{aligned}$$

Here, note that random variable  $\xi_{j'}$  follows standard normal (i.e.,  $\xi_{j'} \sim \mathcal{N}(0, 1)$ ). Next, we further need to control the term  $E \left[ |\xi_{j'}| (\|\tilde{f} - f_0\|_n + \delta) \right]$  as follows:

$$\begin{aligned}
E \left[ |\xi_{j'}| (\|\tilde{f} - f_0\|_n + \delta) \right] &= E \left[ |\xi_{j'}| \|\tilde{f} - f_0\|_n \right] + \delta E[|\xi_{j'}|] \\
&\leq E \left[ |\xi_{j'}| \|\tilde{f} - f_0\|_n \right] + \delta \\
&\leq \hat{R}_n^{1/2}(\tilde{f}, f_0) E^{1/2}[\xi_{j'}^2] + \delta \\
&\leq \hat{R}_n^{1/2}(\tilde{f}, f_0) \sqrt{3 \log \mathcal{N}_n + 1} + \delta \\
&\leq \left( \hat{R}_n^{1/2}(\tilde{f}, f_0) + \delta \right) \sqrt{3 \log \mathcal{N}_n + 1}
\end{aligned}$$

In first inequality, the fact  $E[|\xi_{j'}|] = \sqrt{\frac{2}{\pi}} \leq 1$  is used. For second inequality, Cauchy-Schwarz inequality is employed. For third inequality, the fact  $E[\xi_{j'}^2] \leq 3 \log \mathcal{N}_n + 1$  is used. In the last inequality, we use  $1 \leq \log \mathcal{N}_n$ . Here, we provide a simple proof for the fact :

$$E[\xi_{j'}^2] \leq 3 \log \mathcal{N}_n + 1, \quad \xi_{j'} \sim \mathcal{N}(0, 1).$$

*Proof.* Let  $Z = \max_{j=1, \dots, \mathcal{N}_n} \xi_{j'}^2$ . Since  $Z \leq \sum_{j=1}^{\mathcal{N}_n} \xi_{j'}^2$ , we have  $E[Z] \leq \mathcal{N}_n$ . For  $\mathcal{N}_n \in \{1, 2, 3\}$ , it can be checked that  $\mathcal{N}_n \leq 3 \log(\mathcal{N}_n) + 1$  and the result holds in this case. Now, we consider the case where  $\mathcal{N}_n \geq 4$ . Mill's ratio gives  $P(|\xi_1| \geq \sqrt{t}) = 2P(\xi_1 \geq \sqrt{t}) \leq \frac{2}{\sqrt{2\pi t}} e^{-t/2}$ . For any  $T \geq 0$  and by using union bound,

$$\begin{aligned}
E[Z] &= \int_0^\infty P(Z \geq t) dt \\
&\leq T + \int_T^\infty P(Z \geq t) dt \leq T + \mathcal{N}_n \int_T^\infty P(\xi_1^2 \geq t) dt \\
&\leq T + \mathcal{N}_n \int_T^\infty P(\xi_1 \geq \sqrt{t}) dt \leq T + \mathcal{N}_n \int_T^\infty \frac{2}{\sqrt{2\pi t}} e^{-t/2} dt \\
&\leq T + \frac{2\mathcal{N}_n}{\sqrt{2\pi T}} \int_T^\infty e^{-t/2} dt = T + \frac{4\mathcal{N}_n}{\sqrt{2\pi T}} e^{-T/2}.
\end{aligned}$$

Set  $T = 2 \log(\mathcal{N}_n)$ , we find

$$E[Z] \leq 2 \log(\mathcal{N}_n) + \frac{2}{\sqrt{\pi \log(\mathcal{N}_n)}}.$$

This yields the claim for  $\mathcal{N}_n \geq 4$ . □

Further, because of  $\log \mathcal{N}_n \leq n$ , we have  $2n^{-1/2} \delta \sqrt{3 \log \mathcal{N}_n + 1} \leq 4\delta$ . By combining the facts, we finally get a following inequality:

$$\left| E \left[ \frac{2}{n} \sum_{i=1}^n \varepsilon_i \tilde{f}(X_i) \right] \right| \leq 2 \sqrt{\frac{\hat{R}_n(\tilde{f}, f_0) (3 \log \mathcal{N}_n + 1)}{n}} + 6\delta.$$

### 3.5 Step 4.

By the definition of  $\Delta_n$ , for any fixed  $f \in \mathcal{F}$ , we have  $E[\frac{1}{n} \sum_{i=1}^n (Y_i - \hat{f}(X_i))^2] \leq E[\frac{1}{n} \sum_{i=1}^n (Y_i - f(X_i))^2] + \Delta_n$ . We want to get a bound for  $\hat{R}_n(\hat{f}, f_0)$  for any fixed  $f \in \mathcal{F}$ ,

$$\begin{aligned}
\hat{R}_n(\hat{f}, f_0) &= E_{f_0} \left[ \frac{1}{n} \sum_{j=1}^n (\hat{f}(X_i) - f_0(X_i))^2 \right] \\
&= E_{f_0} \left[ \frac{1}{n} \sum_{j=1}^n (\hat{f}(X_i) + \varepsilon_i - Y_i)^2 \right] \\
&= E_{f_0} \left[ \frac{1}{n} \sum_{j=1}^n \left( (Y_i - \hat{f}(X_i))^2 + 2\varepsilon_i(Y_i - \hat{f}(X_i)) + \varepsilon_i^2 \right) \right] \\
&\leq E_{f_0} \left[ \frac{1}{n} \sum_{j=1}^n (Y_i - f(X_i))^2 \right] + \Delta_n - \frac{2}{n} E_{f_0} \left[ \sum_{i=1}^n \varepsilon_i (Y_i - \hat{f}(X_i)) \right] + E_{f_0}[\varepsilon_1^2] \\
&\leq E_{f_0} \left[ \frac{1}{n} \sum_{j=1}^n \left( f_0(X_i) - f(X_i) + \varepsilon_i \right)^2 \right] + \Delta_n + \frac{2}{n} E_{f_0} \left[ \sum_{i=1}^n \varepsilon_i (\hat{f}(X_i) - f_0(X_i)) \right] - E_{f_0}[\varepsilon_1^2] \\
&\leq E_{f_0} \left[ \frac{1}{n} \sum_{j=1}^n \left( f_0(X_i) - f(X_i) \right)^2 \right] + \Delta_n + \frac{2}{n} E_{f_0} \left[ \sum_{i=1}^n \varepsilon_i \hat{f}(X_i) \right].
\end{aligned}$$

Note that we have used the relation  $E[\varepsilon_i f(X_i)] = 0$  at second to the last inequality. Finally, combining the upper bound we get on  $\frac{2}{n} E_{f_0} \left[ \sum_{i=1}^n \varepsilon_i \hat{f}(X_i) \right]$  at Step 3., we obtain the upper-bound for  $\hat{R}_n(\hat{f}, f_0)$  as follows:

$$\hat{R}_n(\hat{f}, f_0) \leq E_{f_0}[\|f - f_0\|_n^2] + 2\sqrt{\frac{\hat{R}_n(\tilde{f}, f_0)(3 \log \mathcal{N}_n + 1)}{n}} + 6\delta + \Delta_n.$$

Lastly using the upper bound of (3) and by setting  $a = \hat{R}_n(\hat{f}, f_0)$ ,  $b = 0$ ,  $c = \sqrt{\frac{3 \log \mathcal{N}_n + 1}{n}}$ , and  $d = E_{f_0}[\|f - f_0\|_n^2] + 6\delta + \Delta$ , we get the following upper bound of  $\hat{R}_n(\hat{f}, f_0)$  for any  $\varepsilon \in (0, 1]$ :

$$\hat{R}_n(\hat{f}, f_0) \leq (1 + \varepsilon) \left( \inf_{f \in \mathcal{F}} E_{f_0}[\|f - f_0\|_n^2] + F^2 \frac{6 \log \mathcal{N}_n + 2}{n\varepsilon} + 6\delta + \Delta_n \right).$$

As for the lower bound of  $\hat{R}_n(\hat{f}, f_0)$ , first it is trivial to see that

$$\begin{aligned}
\hat{R}_n(\hat{f}, f_0) - \hat{R}_n(\tilde{f}, f_0) &= \Delta_n + E_{f_0} \left[ \frac{2}{n} \sum_{i=1}^n \varepsilon_i \hat{f}(X_i) \right] - E_{f_0} \left[ \frac{2}{n} \sum_{i=1}^n \varepsilon_i \tilde{f}(X_i) \right] \\
&\geq \Delta_n - 2\sqrt{\frac{\hat{R}_n(\hat{f}, f_0)(3 \log \mathcal{N}_n + 1)}{n}} - 2\sqrt{\frac{\hat{R}_n(\tilde{f}, f_0)(3 \log \mathcal{N}_n + 1)}{n}} - 12\delta \\
&\geq \Delta_n - \frac{\varepsilon}{1 - \varepsilon} \hat{R}_n(\hat{f}, f_0) - \frac{3 \log \mathcal{N}_n + 1}{\varepsilon n} - \hat{R}_n(\tilde{f}, f_0) - 12\delta.
\end{aligned}$$

Here in the third inequality, we have used following inequalities, for any  $\varepsilon \in (0, 1]$ ,

$$\begin{aligned}
2\sqrt{\frac{\hat{R}_n(\hat{f}, f_0)(3 \log \mathcal{N}_n + 1)}{n}} &\leq \frac{\varepsilon}{1 - \varepsilon} \hat{R}_n(\hat{f}, f_0) + \frac{1 - \varepsilon}{\varepsilon} \left( \frac{3 \log \mathcal{N}_n + 1}{n} \right), \\
2\sqrt{\frac{\hat{R}_n(\tilde{f}, f_0)(3 \log \mathcal{N}_n + 1)}{n}} &\leq \hat{R}_n(\tilde{f}, f_0) + \frac{3 \log \mathcal{N}_n + 1}{n}.
\end{aligned}$$



Rearranging the term, for any  $\varepsilon \in (0, 1]$ ,

$$\widehat{R}_n(\widehat{f}, f_0) \geq (1 - \varepsilon) \left( \Delta_n - \frac{3 \log \mathcal{N}_n + 1}{n\varepsilon} - 12\delta \right).$$