

# Summary of Theory

Miaoyan Wang, Oct 4, 2020

## 1 Set-up

Consider the linear function class

$$\mathcal{F}(d, r, s) = \{\mathbf{X} \mapsto \langle \mathbf{X}, \mathbf{B} \rangle \mid \text{rank}(\mathbf{B}) \leq r, \text{Supp}(\mathbf{B}) \leq s, \mathbf{B} \in \mathbb{R}^{d \times d}\}. \quad (1)$$

For any function  $f \in \mathcal{F}(d, r, s)$ , we define  $\|f\|_F = \|\mathbf{B}\|_F$ . (do I need to care about  $\mathbf{X}$ ?? consistent with RHKS?). Distinguish to  $L_2$  norm of  $f$ ?

Let  $\{(\mathbf{X}_i, y_i) \in \mathbb{R}^{d \times d} \times \{\pm 1\} : i = 1, \dots, n\}$  denote the i.i.d. training sample from an unknown distribution  $\mathbb{P}(\mathbf{X}, y)$ . We are interested in the high-dimensional regime as  $n, d \rightarrow \infty$ , while holding  $s, r$  as fixed constants.

Define the restricted eigenvalue as

$$\lambda_{\max}(\mathbf{X}) = \max_{\mathbf{x} \in \mathcal{S}^d, \|\mathbf{x}\|_0 \leq s} \mathbf{a}^T \mathbf{X} \mathbf{a}.$$

Assume there exists a constant  $C > 0$  such that  $\lambda_{\max}(\mathbf{X}) \leq C$ , a.s. as  $d \rightarrow \infty$ . The restricted eigenvalue condition incorporates (i) bounded feature  $\|\mathbf{X}\|_F \leq C$ , and (ii) Gaussian feature  $\mathbf{X}$  with i.i.d.  $N(0, 1)$  entries. In the later case, the feature space is unbounded,  $\|\mathbf{X}\|_F \asymp \mathcal{O}(d)$  as  $d \rightarrow \infty$ , but spectral norm is bounded  $\lambda_{\max}(\mathbf{X}) \asymp s = \mathcal{O}(1)$ . We also need  $\|\mathbf{B}\|_F \leq \max_{\mathbf{X} \in \mathbb{R}^{d \times d}} \langle \mathbf{B}, \mathbf{X} \rangle$  (this is natural) do we need assumptions on  $\mathbf{X}$  in the probability estimation?? We need  $\|f\|_\infty$  bounded in the  $L_2$  entropy. We require  $\mathbf{X}$  spread out, in particular, cover the set of  $\mathbf{B}$

## 2 Theory

**Definition 1** (Classification risk and surrogate risk). Let  $f(\cdot) : \mathbb{R}^{d \times d} \mapsto \mathbb{R}$  be the decision function of interest,  $\ell(\cdot) : \mathbb{R} \mapsto \mathbb{R}_{\geq 0}$  be a surrogate loss function in terms of the margin  $yf(\mathbf{X})$ . We define the 0/1 classification risk and surrogate risk,

$$R(f) = \mathbb{P}(y \neq \text{sign}f(\mathbf{X})), \quad R_\ell(f) = \mathbb{E}\ell(yf(\mathbf{X})).$$

**Assumption 1** (Surrogate loss). Assume that the surrogate loss  $\ell$  satisfies the following two conditions

- i.  $\ell$  is a  $L$ -Lipschitz function and  $\ell$  entrywise dominates the 0/1 loss. This assumption implies that  $R(f) \leq R_\ell(f)$  for all functions  $f$ .

ii. The loss is Fisher consistent,

$$R(f_{\text{bayes}}) = R_\ell(f_{\text{bayes}}), \quad \text{and} \quad \arg \min_{\text{all possible } f} R_\ell(f) = \arg \min_{\text{all possible } f} R(f). \quad (2)$$

That is, replacing 0/1 loss by surrogate loss does not change the minimal risk and minimizer. Note that the right hand side of (2) is obtained at  $f_{\text{bayes}}(\cdot): \mathbf{X} \mapsto \text{sign}\{\mathbb{P}(y = 1|\mathbf{X}) - 1/2\}$  a.s. except for the decision boundary  $\{\mathbf{X}: \mathbb{P}(y = 1|\mathbf{X}) = 1/2\}$ . (where do we use the global minimum assumption?)

We denote the empirical risks calculated from the training sample,

$$\hat{R}(f) = \frac{1}{n} \sum_{i=1}^n (y_i \neq \text{sign} f(\mathbf{X}_i)), \quad \hat{R}_\ell(f) = \frac{1}{n} \sum_{i=1}^n \ell(y_i f(\mathbf{X}_i)).$$

**Proposition 1** (Generalization). Consider the function class  $\mathcal{F}(d, r, s)$  in (1) and a surrogate loss  $\ell$  under Assumption 1i. With very high probability over  $\{(\mathbf{X}_i, y_i)\}$ , we have

$$\sup_{f \in \mathcal{F}(d, r, s)} |R_\ell(f) - \hat{R}_\ell(f)| \leq CLs \log d \sqrt{\frac{r}{n}} \sup_{\mathbf{B} \in \mathcal{F}(d, r, s)} \|\mathbf{B}\|_F.$$

**Corollary 2** (Excess risk). Consider the same assumptions as in Proposition 1. Let  $\hat{f} \in \mathcal{F}(d, r, s)$  denote empirical surrogate risk minimizer constrained to the function class  $\mathcal{F}(d, r, s)$ ,

$$\hat{f} = \arg \min_{f \in \mathcal{F}(d, r, s)} \frac{1}{n} \sum_{i=1}^n \ell(y_i f(\mathbf{X}_i)).$$

Then with very high probability over the training set, the surrogate risk error satisfies

$$R_\ell(\hat{f}) - \inf_{f \in \mathcal{F}(d, r, s)} R_\ell(f) \leq 2CLs \log d \sqrt{\frac{r}{n}} \sup_{\mathbf{B} \in \mathcal{F}(d, r, s)} \|\mathbf{B}\|_F.$$

and the classification error satisfies

$$R(\hat{f}) - \inf_{f \in \mathcal{F}(d, r, s)} R(f) \leq 2CLs \log d \sqrt{\frac{r}{n}} \sup_{\mathbf{B} \in \mathcal{F}(d, r, s)} \|\mathbf{B}\|_F.$$

Write in terms of Bayes error?

**Remark 1** (Two shortcomings). First, the convergence in sample size is of order  $\mathcal{O}(1/\sqrt{n})$ . This can be improved to  $\mathcal{O}(1/n)$  upon mean-variance (low noise) conditions. Second, the generalization bound depends on  $\sup_{\mathbf{B} \in \mathcal{F}(d, r, s)} \|\mathbf{B}\|_F$ , which may depends on  $d$ . In fact, we can use an adaptive penalization to replace this term to  $\|\mathbf{B}^*\|_F$ , where  $\mathbf{B}^*$  is the matrix that induces  $f^*$ ; i.e.,  $f^*(\mathbf{X}) = \langle \mathbf{X}, \mathbf{B}^* \rangle$ .

To overcome the above pitfalls, we propose to use the following penalized empirical surrogate risk

minimizer

$$\hat{f}_\lambda = \arg \min_{f \in \mathcal{F}(d, r, s)} \left\{ \frac{1}{n} \sum_{i=1}^n \ell(y_i f(\mathbf{X}_i)) + \lambda \|f\|_F^2 \right\}, \quad (3)$$

where  $\mathcal{F}(d, r, s) = \{\mathbf{X} \mapsto \langle \mathbf{X}, \mathbf{B} \rangle \mid \text{rank}(\mathbf{B}) \leq r, \text{Supp}(\mathbf{B}) \leq s, \mathbf{B} \in \mathbb{R}^{d \times d}\}$ .

**Assumption 2.** Consider the following assumptions:

- i. As  $n, d \rightarrow \infty$ , there exists a sequence  $f_n^* \in \mathcal{F}(d, r, s)$  that is (1) bounded  $\|f_n^*\|_F^2 \leq J^*$  for some constant  $C > 0$ , and (2)  $R_\ell(f_n^*) - R_\ell(f_{\text{bayes}}) \leq a_n \rightarrow 0$ . The upper bound  $J^*$  is allowed to depend on  $d = d(n)$ .
- ii. There exist constants  $a > 0$  and  $\rho \in [0, 1]$ , such that, for any sufficient small  $\delta > 0$ ,

$$\text{Var}[\ell(yf(\mathbf{X})) - \ell(yf_{\text{bayes}}(\mathbf{X}))] \leq a \mathbb{E}[\ell(yf(\mathbf{X})) - \ell(yf_{\text{bayes}}(\mathbf{X}))]^\rho$$

holds for all  $f \in \{f \in \mathcal{F}(d, r, s) : R_\ell(f) - R_\ell(f_{\text{bayes}}) \leq \delta\}$  in a  $\delta$ -neighborhood of  $f_{\text{bayes}}$ .

- iii. There exist constants  $b > 0$  and  $\alpha \geq 0$ , such that, for any sufficient small  $\delta > 0$ ,

$$\mathbb{E}|\text{sign}f(\mathbf{X}) - \text{sign}f_{\text{bayes}}(\mathbf{X})| \leq b [R_\ell(f) - R_\ell(f_{\text{bayes}})]^\alpha$$

holds for all  $f \in \{f \in \mathcal{F}(d, r, s) : R_\ell(f) - R_\ell(f_{\text{bayes}}) \leq \delta\}$  in a  $\delta$ -neighborhood of  $f_{\text{bayes}}$ .

**Theorem 3** (Main result for classification accuracy). Suppose Assumption 1 and Assumption 2i-2ii hold. Consider a penalized empirical surrogate risk minimizer (3) with the regularity parameter  $\lambda \asymp \frac{1}{J^*} \left( \frac{rs \log d}{n} \right)^{1/(2-\rho)}$ . We have that, with very high probability, the surrogate error and classification error satisfy

$$R_\ell(\hat{f}_\lambda) - R_\ell(f_{\text{bayes}}) \leq C \left( \frac{rs \log d}{n} + a_n \right)^{1/(2-\rho)}, \quad R(\hat{f}_\lambda) - R(f_{\text{bayes}}) \leq C \left( \frac{rs \log d}{n} + a_n \right)^{1/(2-\rho)}. \quad (4)$$

Furthermore, suppose Assumption 2iii holds. Then, Assumption 2ii holds with  $\rho = \alpha \wedge 1$ , and the estimation error for level sets satisfies

$$\mathbb{P}(\hat{\mathbf{X}} \Delta \mathbf{X}_{\text{bayes}}) \leq C \left( \frac{rs \log d}{n} + a_n \right)^{\alpha/(2-\alpha \wedge 1)},$$

where  $\hat{\mathbf{X}} = \{\mathbf{X} : \hat{f}(\mathbf{X}) \geq 0\}$  and  $\mathbf{X}_{\text{bayes}} = \{\mathbf{X} : f_{\text{bayes}}(\mathbf{X}) \geq 0\}$  are estimated and true level sets, respectively.

**Remark 2.** The bound (4) implies the consistency of risk estimator in the ultra high-dimensional regime. In particular, the dimension of feature matrices  $d$  is allowed to grow sub-exponentially in sample size  $n$ ; i.e.,  $d = o(e^n)$ .

**Theorem 4** (Main result for probability estimation). Consider the same assumptions of Theorem 3. For the estimator  $\hat{p}: \mathbb{R}^{d \times d} \rightarrow [0, 1]$  obtained from level-set estimation,

$$\hat{p}(\mathbf{X}) = \frac{1}{H} \sum_{h=1}^H \mathbb{1}\{\mathbf{X}: \hat{f}_h(\mathbf{X}) \geq 0\}, \quad \text{for all } \mathbf{X} \in \mathbb{R}^{d \times d}.$$

With probability at least  $1 - C \exp(-an\lambda^{2-\alpha})$ ,

$$\mathbb{E}|\hat{p}(\mathbf{X}) - p(\mathbf{X})| \geq \underbrace{\frac{1}{2H}}_{\text{discretization error}} + \frac{a(H+1)}{2} \mathcal{O} \left( \underbrace{\frac{rs \log d}{n}}_{\text{statistical error}} + \underbrace{a_n}_{\text{approximation error}} \right)^{\alpha^2}.$$

**Corollary 5.** Assume  $a_n \leq \frac{rs \log d}{n}$ . Choosing  $H \asymp (\frac{rs \log d}{n})^{-\alpha^2/2}$  gives the estimation error,

$$\mathbb{E}|\hat{p}(\mathbf{X}) - p(\mathbf{X})| \leq C \left( \frac{rs \log d}{n} \right)^{\rho/2}.$$

### 3 Algorithm

### 4 Numerical experiments

Three goals

1. Assess the 0/1 and hinge loss errors
2. Assess the level-set estimation
3. Assess the probability estimation

### 5 Proofs

*Proof of Proposition 1.* It suffices to bound the Radamecher complexity

$$\begin{aligned} \text{Rad}(f) &= \frac{1}{n} \mathbb{E} \max_{f \in \mathcal{F}(d,r,s)} \sum_{i \in [n]} \sigma_i \langle \mathbf{X}_i, \mathbf{B} \rangle \\ &\leq \frac{1}{n} \|\mathbf{B}\|_* \left\| \sum_{i \in [n]} \sigma_i \mathbf{X}_i \right\|_{\text{sp}} \\ &\leq \frac{1}{n} \mathbb{E} \max_{\text{possible choice } \mathcal{S} \text{ from } [d] \text{ given set } \mathcal{S}} \max_i \langle \mathbf{B}, \sum_i \sigma_i \mathbf{X}_i \rangle \\ &\leq \frac{1}{n} \log \binom{d}{s} \mathbb{E} \max_{\{\mathbf{a}_r, \mathbf{b}_r, \lambda_r\}} \left\langle \sum_r \lambda_r \mathbf{a}_r \mathbf{b}_r^T, \sum_i \sigma_i \mathbf{X}_i \right\rangle \end{aligned}$$

$$\begin{aligned}
&\leq \frac{s}{n} \log d \sum_r \lambda_r \langle \mathbf{a}_r \mathbf{b}_r, \sum_i \sigma_i \mathbf{X}_i \rangle \\
&\leq \frac{s}{n} \log d \sqrt{n} \lambda_{\max}(\mathbf{X}) \sqrt{r} \|\mathbf{B}\|_F \\
&\leq Cs \sqrt{\frac{r}{n}} \log d \|\mathbf{B}\|_F.
\end{aligned}$$

□

*Proof of Corollary 2.* Consider the decomposition

$$\begin{aligned}
R_\ell(\hat{f}) - \inf_{f \in \mathcal{F}(d,r,s)} R_\ell(f) &\leq R_\ell(\hat{f}) - \hat{R}_\ell(\hat{f}) + \hat{R}_\ell(\hat{f}) - \hat{R}_\ell(f^*) + \hat{R}_\ell(f^*) - R_\ell(f^*) \\
&\leq |R_\ell(\hat{f}) - \hat{R}_\ell(\hat{f})| + \hat{R}_\ell(f^*) - R_\ell(f^*) \\
&\leq 2\text{Rad}(f) \\
&\leq 2CLs \sqrt{\frac{r}{n}} \log d \max_{f \in \mathcal{F}(d,r,s)} \|\mathbf{B}\|_F
\end{aligned}$$

□

*Proof of Theorem 3.* We set  $\delta_n = \mathcal{O}(\frac{rs \log d}{n})^\rho$ . Then

$$\begin{aligned}
&\mathbb{P} \left\{ R_\ell(\hat{f}) - R_\ell(f_{\text{bayes}}) \geq \delta \right\} \\
&\leq \mathbb{P} \left\{ \sup_{\{f \in \mathcal{F}(d,r,s) : R_\ell(f) - R_\ell(f_{\text{bayes}}) \geq \delta\}} \left( \hat{R}(f^*) + \lambda \|f^*\|_F^2 - \hat{R}(f) - \lambda \|f\|_F^2 \right) \geq 0 \right\} \\
&\leq 3.5 \exp(-an\lambda^{2-\rho}),
\end{aligned}$$

where the last line comes from Lemma 7, by taking  $\lambda \asymp \delta_n$ . (where  $a_n$  enters?)

The classification bound follows by noting that  $R(\hat{f}) \leq R_\ell(\hat{f})$  and  $R(f_{\text{bayes}}) = R_\ell(f_{\text{bayes}})$ . □

**Definition 2** (bracketing number, uniform entropy, and bounded functions). Consider a function set  $\mathcal{F}$ , and let  $\varepsilon > 0$ . We call  $\{(f_m^l, f_m^u)\}_{m=1}^M$  an  $L_2$ -metric,  $\varepsilon$ -bracketing function set of  $\mathcal{F}$ , if for every  $f \in \mathcal{F}$ , there exists an  $m \in [M]$  such that

$$f_m^l(\mathbf{X}) \leq f(\mathbf{X}) \leq f_m^u(\mathbf{X}), \quad \text{for all } \mathbf{X} \in \mathbb{R}^{d \times d},$$

and

$$\|f_m^l - f_m^u\|_2 \stackrel{\text{def}}{=} \sqrt{\mathbb{E}|f_m^l(\mathbf{X}) - f_m^u(\mathbf{X})|^2} \leq \varepsilon, \quad \text{for all } m = 1, \dots, M.$$

The bracketing number with  $L_2$ -metric,  $\mathcal{H}_{[\cdot]}(\varepsilon, \mathcal{F}, \|\cdot\|_2)$ , is defined as the logarithm of the smallest cardinality of the  $\varepsilon$ -bracketing function set of  $\mathcal{F}$ . Furthermore, consider the set of functions with

$L_\infty$  bound no larger than  $M$ , denoted  $\mathcal{F}(M) = \{f \in \mathcal{F}: \|f\|_\infty \leq M\}$ . Then we have

$$\mathcal{H}_{[\cdot]}(\varepsilon, \mathcal{F}(M), \|\cdot\|_2) \leq \mathcal{H}(\varepsilon, \mathcal{F}(M), \|\cdot\|_\infty).$$

**Lemma 6** (Uniform entropy for bounded functions in  $\mathcal{F}(d, r, s)$ ). Let  $\mathcal{F}(d, r, s)$  denote the function class in (1). Consider the subset of functions with  $L_\infty$  bound no larger than  $M$ , denoted  $\mathcal{F}(M) = \{f \in \mathcal{F}(d, r, s): \|f\|_\infty \leq M\}$  for  $M = 1, 2, \dots$ . When  $\varepsilon$  sufficiently small, we have

$$\mathcal{H}(\varepsilon, \mathcal{F}(M), \|\cdot\|_\infty) \leq 5rs \log \frac{M\sqrt{r}\lambda_{\max}d}{\varepsilon}.$$

*Proof.* For a given matrix  $\mathbf{B}$ , the definition  $f(\mathbf{X}) = \langle \mathbf{X}, \mathbf{B} \rangle$  implies that  $\|\mathbf{B}\|_F \leq \max_{\mathbf{X} \in \mathbb{R}^{d \times d}} \langle \mathbf{X}, \mathbf{B} \rangle = \|\mathbf{B}\|_\infty \leq \|\mathbf{B}\|_F \sqrt{r}\lambda_{\max}$ . Therefore, we have

$$\mathcal{H}(\varepsilon\sqrt{r}\lambda_{\max}, \mathcal{F}(M), \|\cdot\|_\infty) = \mathcal{H}(\varepsilon, \mathcal{B}(M), \|\cdot\|_F), \quad (5)$$

where we have defined the matrix set  $\mathcal{B}(M) = \{\mathbf{B} \in \mathbb{R}^{d \times d}: \text{rank}(\mathbf{B}) \leq r, \text{Supp}(\mathbf{B}) \leq s, \|\mathbf{B}\|_F \leq M\}$ . Note that the  $F$ -norm of matrix  $\mathbf{B}$  is equivalent to the  $l_2$ -norm of the vector  $\text{vec}(\mathbf{B})$ , and the vector  $l_2$ -norm is lower bounded by vector  $l_\infty$ -norm in Euclidean space. Therefore, it suffices to bound  $\mathcal{H}(\varepsilon, \mathcal{B}(M), \|\cdot\|_\infty)$ . Now fix a subset  $S \subset [d]$  with  $|S| = s$ , and let  $\mathcal{B}_S(M) \subset \mathcal{B}(M)$  denote the subset of matrices satisfying  $\mathbf{B}(i, j) = 0$  whenever  $(i, j) \notin S^2$ . **Based on ...**, the  $\varepsilon$ -covering of  $\mathcal{B}_S(M)$  has entropy

$$\mathcal{H}(\varepsilon, \mathcal{B}_S(M), \|\cdot\|_\infty) \leq r(2s + 1) \log \left( \frac{M}{\varepsilon} \right).$$

In view of  $\mathcal{B}(M) \subset \bigcup_{S \subset [d], |S|=s} \mathcal{B}_S(M)$ , an  $\varepsilon$ -covering set  $\mathcal{B}(M)$  is then given by the union of  $\varepsilon$ -covering set of  $\mathcal{B}_S(M)$ . Using Stirling's bound, we derive that

$$\mathcal{H}(\varepsilon, \mathcal{B}(M), \|\cdot\|_\infty) \leq 2s \log \frac{d}{s} + r(2s + 1) \log \frac{M}{\varepsilon} \leq 5rs \log \frac{Md}{\varepsilon}.$$

Substituting  $\varepsilon$  by  $\varepsilon/\sqrt{r}\lambda_{\max}$  into (5) concludes the proof.  $\square$

**Lemma 7** (Metric of local  $\mathcal{F}(d, r, s)$ ). Let  $\delta > 0$  be the solution to the following inequality,

$$\max_{M \geq 2} \left\{ \frac{1}{\delta + \lambda(M/2 - 1)} \int_{\delta + \lambda(M/2 - 1)}^{(\delta + \lambda(M/2 - 1))^{\rho/2}} \sqrt{\mathcal{H}(\varepsilon, \mathcal{F}(M), \|\cdot\|_\infty)} d\varepsilon \right\} \leq n^{1/2}.$$

Then we have  $\delta = \mathcal{O} \left( \frac{rs \log d}{n} \right)^\rho$  provided that  $\lambda \leq 4\delta$  (??).

**Theorem 8.** Let  $\mathcal{F}$  be a class of functions. Let  $T, V \in (0, \infty)$  denote the upper bound of functions in  $\mathcal{F}$  in  $L_\infty$  and  $L_2$  norms; that is,  $\sup_{f \in \mathcal{F}} \|f\|_\infty \leq T$  and  $\sup_{f \in \mathcal{F}} \text{Var}(f) \leq v$ . Let  $E_n(f) = \frac{1}{n} \sum_{i=1}^n (f(Y_i))$  be the empirical process. Define  $x_n^*$  be the solution of the equation to the following

equation

$$\frac{1}{x} \int_x^{\sqrt{V}} \sqrt{\mathcal{H}_{[\cdot]}(\varepsilon, \mathcal{F}, \|\cdot\|_2)} d\varepsilon = \sqrt{n}.$$

Suppose the  $\sqrt{V} \leq T$ , and

$$x_n^* \lesssim \frac{V}{T}, \quad \text{and} \quad \mathcal{H}_{[\cdot]}(\sqrt{V}, \mathcal{F}, \|\cdot\|_2) \lesssim \frac{nx_n^*}{T}.$$

Then we have

$$\mathbb{P} \left( \sup_{f \in \mathcal{F}} E_n(f) \geq \mathbb{E}f(Y) + x_n^* \right) \lesssim \exp(-cnx_n^*).$$

**Remark 3.** The function set  $\mathcal{F}$ , and bounds  $T$ ,  $v$  are allowed to depend on  $n$ .

We view  $\mathcal{Y}_\Omega = \{\bar{\mathcal{Y}}(\omega) : \omega \in \Omega\}$  as a collection of  $n$  i.i.d. random variables where the randomness is induced from both  $\bar{\mathcal{Y}}$  and  $\omega \sim \Pi$ , and view the tensor  $\mathcal{Z}$  as a function that maps  $\bar{\mathcal{Y}}_\Omega$  to  $L(\mathcal{Z}, \bar{\mathcal{Y}}_\Omega)$ .

Specifically, the data takes the form  $\{y_i : i = 1, 2, \dots, n\}$ , where for each  $i \in [n]$ ,  $y_i$  is i.i.d. sampled from all entries of  $\mathcal{Y}$  based on  $\omega \sim \Pi$ . We denote  $y_i$  i.i.d. random variables where the randomness is induced from both  $\Pi$  and noise in the tensor model.

The loss function then takes the form

$$L(\mathcal{Z}, \mathcal{Y}_\Omega) = \frac{1}{n} \sum_{i=1}^n \ell_i(y_i),$$

where  $\ell_i(y_i) = |y_i| |\text{sign}(z_i) - \text{sign}(y_i)|$ . The collection of function  $\{\ell_i : i \in [d_1] \times \dots \times [d_K]\}$  is one-to-one  $\mathcal{Z}$ . For notational simplicity, we write  $L(\mathcal{Z})$  in place of  $L(\mathcal{Z}, \mathcal{Y}_\Omega)$ . The relevant probability statements, such as  $\mathbb{E}$  and  $\text{Var}$ , are taken with respect to  $\mathcal{Y}(\omega)$ .

Because  $\bar{\Theta}$  is the global minimizer of  $\text{Risk}(\cdot)$ , and by definition,  $L(\hat{\mathcal{Z}}) \leq L(\bar{\Theta})$ , we have the following inclusion of the event

$$\{\text{Risk}(\hat{\mathcal{Z}}) - \text{Risk}(\bar{\Theta}) \geq T_n\} \subset \left\{ \sup_{\mathcal{Z} \in \mathcal{F}} (\text{Risk}(\mathcal{Z}) - \text{Risk}(\bar{\Theta}) + L(\bar{\Theta}) - L(\mathcal{Z})) \geq T_n \right\}.$$

Therefore,

$$\mathbb{P} \left( \text{Risk}(\hat{\mathcal{Z}}) - \text{Risk}(\bar{\Theta}) \geq T_n \right) \leq \mathbb{P} \left( \sup_{\mathcal{Z} \in \mathcal{F}} |L(\mathcal{Z}) - \text{Risk}(\mathcal{Z}) - L(\bar{\Theta}) + \text{Risk}(\bar{\Theta})| \geq T_n \right).$$

We then use the empirical process to uniformly bound the stochastic residual

$$E_n(\mathcal{Z}) := L(\mathcal{Z}) - \text{Risk}(\mathcal{Z}) - L(\bar{\Theta}) + \text{Risk}(\bar{\Theta}).$$

To show this, we notice that the stochastic residual is a sum of i.i.d. r.v.'s

$$\begin{aligned} E_n(\mathcal{Z}) &= \frac{1}{n} \sum_{i=1}^n \underbrace{[\ell_{\mathcal{Z}}(y_i) - \ell_{\bar{\Theta}}(y_i) + \mathbb{E}\ell_{\bar{\Theta}}(y_i) - \mathbb{E}\ell_{\mathcal{Z}}(y_i)]}_{\text{mean-zero, i.i.d. r.v.'s}} \\ &= \frac{1}{n} \sum_{i=1}^n e_i + \text{Risk}(\bar{\Theta}) - \text{Risk}(\mathcal{Z}) \end{aligned}$$

where

$$\text{Var}[e_i] \leq \mathbb{E}[e_i]^\alpha + \frac{1}{\rho} \mathbb{E}[e_i].$$

With high probability

$$\begin{aligned} \max_{\mathcal{Z}} \frac{1}{n} \sum_{i=1}^n e_i &\leq T_n + \text{Risk}(\mathcal{Z}) - \text{Risk}(\bar{\Theta}) \\ \sup_{\mathcal{Z}} E_n(\mathcal{Z}) &\leq T_n \end{aligned}$$

**Lemma 9.** Let  $\mathcal{F}$  be a class of functions, and  $(y_i)_{i \in [n]}$  be an i.i.d. sample from random variable  $y$ . Suppose the variance-to-mean relationship holds uniformly over  $\mathcal{F}$ ,

$$\text{Var}f(y) = [\mathbb{E}f(y)]^\beta + \frac{1}{\rho} \mathbb{E}f(y), \quad \text{for all } f \in \mathcal{F},$$

where  $\beta \in [0, 1]$  is a constant. Then

$$\mathbb{P} \left( \sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n f(y_i) \geq \mathbb{E}f(y) + T_n \right) \leq \exp(-T_n).$$

To bound the right-hand side, we partition  $\{\mathcal{Z} \in \mathcal{F} : \text{Risk}(\mathcal{Z}) - \text{Risk}(\bar{\Theta}) \geq L_n\}$  into a union of  $A_s = \{\mathcal{Z} \in \mathcal{F} : 2^{s-1}L_n \leq \text{Risk}(\mathcal{Z}) - \text{Risk}(\bar{\Theta}) < 2^s L_n\}$  for  $s = 1, 2, \dots$ . Then it suffices to bounding the corresponding probability over each  $A_s$ . Towards this end, we need to bound the first and second moment of  $\Delta_n(\mathcal{Z}, \mathcal{Y}_\Omega)$ .

For the first moment we have

$$\inf_{\mathcal{Z} \in A_s} \mathbb{E}\Delta_n(\mathcal{Z}, \mathcal{Y}) = \inf_{\mathcal{Z} \in A_s} \mathbb{E}[L(\mathcal{Z}, \mathcal{Y}_\Omega) - L(\bar{\Theta}, \mathcal{Y}_\Omega)] \geq \inf_{\mathcal{Z} \in A_s} [\text{Risk}(\mathcal{Z}) - \text{Risk}(\bar{\Theta})] \geq \underbrace{2^{s-1}L_n}_{=:M(s)}$$

for any  $s = 1, 2, \dots$ . For the second moment, it follows from Lemma ?? that

$$\begin{aligned} \sup_{\mathcal{Z} \in A_s} \text{Var}\Delta_n(\mathcal{Z}, \mathcal{Y}) &= \sup_{\mathcal{Z} \in A_s} \text{Var}[L(\mathcal{Z}, \bar{\mathcal{Y}}_\Omega) - L(\bar{\Theta}, \bar{\mathcal{Y}}_\Omega)] \\ &\leq \underbrace{M^{\frac{\alpha}{1+\alpha}}(s) + \frac{1}{\rho} M(s)}_{=:V(s)} \end{aligned}$$



We now apply Shen & Wong for each of the set  $\{\mathcal{Z} \in A_s\}$

$$\begin{aligned}
\mathbb{P}\left(\text{Risk}(\hat{\mathcal{Z}}) - \text{Risk}(\bar{\Theta}) \geq L_n\right) &\leq \sum_{s=1}^{\infty} \mathbb{P}\left(\sup_{A_s}(\mathbb{E}\Delta_n - \Delta_n) \geq M_n\right) \leq \sum_{s=1}^{\infty} \exp\left[-\frac{ncM^2(s)}{V(s) + TM(s)}\right] \quad (6) \\
&\leq \sum_{s=1}^{\infty} \exp\left[-\frac{nM^2(s)}{M^{\frac{\alpha}{1+\alpha}}(s) + \frac{1}{\rho}M(s)}\right] \\
&\lesssim \exp(-d^{\frac{\alpha}{\alpha+1}}n^{\frac{1}{\alpha+1}}) \\
&\lesssim \exp(-d), \text{ provided } n \geq d.
\end{aligned}$$

where the convergence rate  $L_n > 0$  is determined by the solution to the following inequality,

$$\frac{1}{L_n} \int_{L_n}^{\sqrt{L_n^{\alpha/(\alpha+1)} + \frac{L_n}{\rho}}} \sqrt{\mathcal{H}(\varepsilon, \mathcal{F}, \|\cdot\|_2)} d\varepsilon \leq \sqrt{n}. \quad (7)$$

In particular, the smallest  $L_n$  satisfying (7) yields the best upper bound of the error rate. Here  $\mathcal{H}(\varepsilon, \mathcal{F}, \|\cdot\|_2)$  denotes the  $L_2$ -metric,  $\varepsilon$ -bracketing number (c.f. Definition 2) of family  $\mathcal{F}$ .

It remains to solve for the smallest possible  $L_n$  in (7). Based on Lemma 7, the inequality (7) is satisfied with

$$L_n \asymp t_n^{(\alpha+1)/(\alpha+2)} + \frac{1}{\rho}t_n, \quad \text{where } t_n = \frac{Kd_{\max}r}{n}.$$

Combining (??) and (6) gives

$$\text{Risk}(\hat{\mathcal{Z}}) - \text{Risk}(\bar{\Theta}) \geq \left(\frac{Kd_{\max}r}{n}\right)^{(\alpha+1)/(\alpha+2)} + \frac{1}{\rho}\left(\frac{Kd_{\max}r}{n}\right).$$

with probability at least  $1 - \exp(-\sqrt{Knd_{\max}r})$ .