

# Contents

<b>Appendices</b>	<b>1</b>
<b>A Additional theoretical results</b>	<b>2</b>
A.1 Sign rank and matrix rank . . . . .	2
A.2 Adjusting for additional covariates . . . . .	5
A.3 Extension to sub-Gaussian noise . . . . .	5
A.4 Extension to unbounded number of mass points . . . . .	6
A.5 Connection to structured matrix model with functional coefficients . . . . .	7
<b>B Proofs</b>	<b>9</b>
B.1 Main notation . . . . .	9
B.2 Proof of Theorem 1 . . . . .	10
B.3 Proof of Theorem 2 . . . . .	10
B.4 Proofs of Theorem 3, Part (a) in Theorems 5 and 7 . . . . .	12
B.5 Proofs of Theorem 4, Part (b) in Theorems 5 and 7 . . . . .	21
B.6 Proofs of Theorem 6 and Theorem A.3.1 . . . . .	23
<b>C Auxiliary lemmas</b>	<b>25</b>

## A Additional theoretical results

### A.1 Sign rank and matrix rank

In the main paper, we have provided several examples with high matrix rank but low sign rank. This section provides more examples and their proofs.

**Example 1** (Max graphon). Suppose the matrix  $\Theta \in \mathbb{R}^{d \times d}$  takes the form

$$\Theta(i, j) = \log \left( 1 + \frac{1}{d} \max(i, j) \right), \text{ for all } (i, j) \in [d]^2.$$

Then

$$\text{rank}(\Theta) = d, \quad \text{and} \quad \text{srnk}(\Theta - \pi) \leq 2 \text{ for all } \pi \in \mathbb{R}.$$

*Proof.* The full-rankness of  $\Theta$  is verified from elementary row operations as follows

$$\begin{pmatrix} (\Theta_2 - \Theta_1)/(\log(1 + \frac{2}{d}) - \log(1 + \frac{1}{d})) \\ (\Theta_3 - \Theta_2)/(\log(1 + \frac{3}{d}) - \log(1 + \frac{2}{d})) \\ \vdots \\ (\Theta_d - \Theta_{d-1})/(\log(1 + \frac{d}{d}) - \log(1 + \frac{d-1}{d})) \\ \Theta_d/\log(1 + \frac{d}{d}) \end{pmatrix} = \begin{pmatrix} 1 & 0 & \ddots & \ddots & 0 \\ 1 & 1 & \ddots & \ddots & \ddots \\ \vdots & \vdots & \ddots & \ddots & \ddots \\ 1 & 1 & 1 & 1 & 0 \\ 1 & 1 & 1 & 1 & 1 \end{pmatrix},$$

where  $\Theta_i$  denotes the  $i$ -th row of  $\Theta$ . Now it suffices to show  $\text{srnk}(\Theta - \pi) \leq 2$  for  $\pi$  in the feasible range  $(\log(1 + \frac{1}{d}), \log 2)$ . In this case, there exists an index  $i^* \in \{2, \dots, d\}$ , such that  $\log(1 + \frac{i^*-1}{d}) < \pi \leq \log(1 + \frac{i^*}{d})$ . By definition, the sign matrix  $\text{sgn}(\Theta - \pi)$  takes the form

$$\text{sgn}(\Theta(i, j) - \pi) = \begin{cases} -1, & \text{both } i \text{ and } j \text{ are smaller than } i^*; \\ 1, & \text{otherwise.} \end{cases}$$

Therefore, the matrix  $\text{sgn}(\Theta - \pi)$  is a rank-2 block matrix, which implies  $\text{srnk}(\Theta - \pi) = 2$ .  $\square$

In fact, Example 1 is a special case of the following proposition.

**Proposition 1** (Min/Max graphon). Let  $g: \mathbb{R} \rightarrow \mathbb{R}$  be a continuous function such that  $g(z) = 0$  has at most  $r \geq 1$  distinct real roots. For given numbers  $x_i, y_j \in [0, 1]$  all  $(i, j) \in [d]^2$ , define a matrix  $\Theta \in \mathbb{R}^{d \times d}$  with entries

$$\Theta(i, j) = g(\max(x_i, y_j)), \quad (i, j) \in [d]^2. \tag{1}$$

Then, the sign rank of  $\Theta$  satisfies

$$\text{srnk}(\Theta) \leq 2r.$$

The same conclusion holds if we use min in place of max in (1).

*Proof.* Without loss of generality, assume  $x_1 \leq \dots \leq x_d$  and  $y_1 \leq \dots \leq y_d$ . Based on the construction of  $\Theta$ , the reordering does not change the rank of  $\Theta$ . Let  $z_1 < \dots < z_r$  be the  $r$  distinct real roots for the equation  $g(z) = 0$ . We separate the proof for two cases,  $r = 1$  and  $r \geq 2$ .

- When  $r = 1$ . The continuity of  $g(\cdot)$  implies that the function  $g(z)$  has at most one sign change point. Based on the similar argument as in Example 1, the matrix  $\text{sgn}(\Theta)$  is a rank-2 block matrix; i.e.,

$$\text{sgn}(\Theta) = 1 - 2\mathbf{a} \otimes \mathbf{b} \quad \text{or} \quad \text{sgn}(\Theta) = 2\mathbf{a} \otimes \mathbf{b} - 1,$$

where  $\mathbf{a}, \mathbf{b}$  are binary vectors defined by

$$\mathbf{a} = (\underbrace{1, \dots, 1}_{\text{positions for which } x_i < z_1}, 0, \dots, 0)^T, \quad \mathbf{b} = (\underbrace{1, \dots, 1}_{\text{positions for which } y_j < z_1}, 0, \dots, 0)^T$$

Therefore,  $\text{srnk}(\Theta) \leq \text{rank}(\text{sgn}(\Theta)) = 2$ .

- When  $r \geq 2$ . By continuity, the function  $g(z)$  is non-zero and remains an unchanged sign in each of the intervals  $(z_s, z_{s+1})$  for  $1 \leq s \leq r - 1$ . Define the index set

$$\mathcal{I} = \{s \in \mathbb{N}_+ : \text{the interval } (z_s, z_{s+1}) \text{ in which } g(z) < 0\}.$$

We now prove that the sign matrix  $\text{sgn}(\Theta)$  has rank bounded by  $2r - 1$ . To see this, consider the matrix indices for which  $\text{sgn}(\Theta) = -1$ ,

$$\begin{aligned} \{(i, j) : \Theta(i, j) < 0\} &= \{(i, j) : g(\max(x_i, y_j)) < 0\} \\ &= \cup_{s \in \mathcal{I}} \{(i, j) : \max(x_i, y_j) \in (z_s, z_{s+1})\} \\ &= \cup_{s \in \mathcal{I}} \left( \{(i, j) : x_i < z_{s+1}, y_j < z_{s+1}\} \cap \{(i, j) : x_i \leq z_s, y_j \leq z_{s+1}\}^c \right). \end{aligned} \quad (2)$$

The equation (2) is equivalent to

$$\mathbb{1}(\Theta(i, j) < 0) = \sum_{s \in \mathcal{I}} (\mathbb{1}(x_i < z_{s+1})\mathbb{1}(y_j < z_{s+1}) - \mathbb{1}(x_i \leq z_s)\mathbb{1}(y_j \leq z_s)), \quad (3)$$

for all  $(i, j) \in [d]^2$ , where  $\mathbb{1}(\cdot) \in \{0, 1\}$  denotes the indicator function. The equation (3) implies the low-rank representation of  $\text{sgn}(\Theta)$ ,

$$\text{sgn}(\Theta) = 1 - 2 \sum_{s \in \mathcal{I}} (\mathbf{a}_{s+1} \otimes \mathbf{b}_{s+1} - \bar{\mathbf{a}}_s \otimes \bar{\mathbf{b}}_s), \quad (4)$$

where  $\mathbf{a}_{s+1}, \bar{\mathbf{a}}_s$  are binary vectors defined

$$\mathbf{a}_{s+1} = (\underbrace{1, \dots, 1}_{\text{positions for which } x_i < z_{s+1}}, 0, \dots, 0)^T, \quad \text{and} \quad \bar{\mathbf{a}}_s = (\underbrace{1, \dots, 1}_{\text{positions for which } x_i \leq z_s}, 0, \dots, 0)^T,$$

and  $\mathbf{b}_{s+1}, \bar{\mathbf{b}}_s$  are binary vectors defined similarly by using  $y_j$  in place of  $x_i$ . Therefore, by (4) and

the assumption  $|\mathcal{I}| \leq r - 1$ , we conclude that

$$\text{srnk}(\Theta) \leq 1 + 2(r - 1) = 2r - 1.$$

Combining two cases yields that  $\text{srnk}(\Theta) \leq 2r$  for any  $r \geq 1$ .  $\square$

**Example 2** (Banded matrices). Let  $\mathbf{a} = (1, 2, \dots, d)^T$  be a  $d$ -dimensional vector, and define a  $d$ -by- $d$  banded matrix  $\mathbf{M} = |\mathbf{a} \otimes \mathbf{1} - \mathbf{1} \otimes \mathbf{a}|$ . Then

$$\text{rank}(\mathbf{M}) = d, \quad \text{and} \quad \text{srnk}(\mathbf{M} - \pi) \leq 3, \quad \text{for all } \pi \in \mathbb{R}.$$

*Proof.* Note that  $\mathbf{M}$  is a banded matrix with entries

$$\mathbf{M}(i, j) = |i - j|, \quad \text{for all } (i, j) \in [d]^2.$$

Elementary row operation directly shows that  $\mathbf{M}$  is full rank as follows,

$$\begin{pmatrix} (\mathbf{M}_1 + \mathbf{M}_d)/(d-1) \\ \mathbf{M}_1 - \mathbf{M}_2 \\ \mathbf{M}_2 - \mathbf{M}_3 \\ \vdots \\ \mathbf{M}_{d-1} - \mathbf{M}_d \end{pmatrix} = \begin{pmatrix} 1 & 1 & 1 & \cdots & 1 & 1 \\ -1 & 1 & 1 & \cdots & 1 & 1 \\ -1 & -1 & 1 & \cdots & 1 & 1 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ -1 & -1 & -1 & \cdots & -1 & 1 \end{pmatrix}.$$

We now show  $\text{srnk}(\mathbf{M} - \pi) \leq 3$  by construction. Define two vectors  $\mathbf{b} = (2^{-1}, 2^{-2}, \dots, 2^{-d})^T \in \mathbb{R}^d$  and  $\text{rev}(\mathbf{b}) = (2^{-d}, \dots, 2^{-1})^T \in \mathbb{R}^d$ . We construct the following matrix

$$\mathbf{A} = \mathbf{b} \otimes \text{rev}(\mathbf{b}) + \text{rev}(\mathbf{b}) \otimes \mathbf{b}. \quad (5)$$

The matrix  $\mathbf{A} \in \mathbb{R}^{d \times d}$  is banded with entries

$$\mathbf{A}(i, j) = \mathbf{A}(j, i) = \mathbf{A}(d - i, d - j) = \mathbf{A}(d - j, d - i) = 2^{-d-1} (2^{j-i} + 2^{i-j}), \quad \text{for all } (i, j) \in [d]^2.$$

Furthermore, the entry value  $\mathbf{A}(i, j)$  decreases with respect to  $|i - j|$ ; i.e.,

$$\mathbf{A}(i, j) \geq \mathbf{A}(i', j'), \quad \text{for all } |i - j| \geq |i' - j'|. \quad (6)$$

Notice that for a given  $\pi \in \mathbb{R}$ , there exists  $\pi' \in \mathbb{R}$  such that  $\text{sgn}(\mathbf{A} - \pi') = \text{sgn}(\mathbf{M} - \pi)$ . This is because both  $\mathbf{A}$  and  $\mathbf{M}$  are banded matrices satisfying monotonicity (6). By definition (5),  $\mathbf{A}$  is a rank-2 matrix. Henceforce,  $\text{srnk}(\mathbf{M} - \pi) = \text{srnk}(\mathbf{A} - \pi') \leq 3$ .  $\square$

**Example 3** (Identity matrices). Let  $\mathbf{I}$  be a  $d$ -by- $d$  identity matrix. Then

$$\text{rank}(\mathbf{I}) = d, \quad \text{and} \quad \text{srnk}(\mathbf{I} - \pi) \leq 3 \quad \text{for all } \pi \in \mathbb{R}.$$

*Proof.* Depending on the value of  $\pi$ , the sign matrix  $\text{sgn}(\mathbf{I} - \pi)$  falls into one of the two cases:

1.  $\text{sgn}(\mathbf{I} - \pi)$  is a matrix of all 1, or of all  $-1$ ;
2.  $\text{sgn}(\mathbf{I} - \pi) = 2\mathbf{I} - \mathbf{1}_d \otimes \mathbf{1}_d$ .

The former case is trivial, so it suffices to show  $\text{srnk}(\mathbf{I} - \pi) \leq 3$  in the second case. Based on Example 2, the rank-2 matrix  $\mathbf{A}$  in (5) satisfies

$$A(i, j) \begin{cases} = 2^{-d}, & i = j, \\ \geq 2^{-d} + 2^{-d-2}, & i \neq j. \end{cases}$$

Therefore,  $\text{sgn}(2^{-d} + 2^{-d-3} - \mathbf{A}) = 2\mathbf{I} - \mathbf{1}_d \otimes \mathbf{1}_d$ . We conclude that  $\text{srnk}(\mathbf{I} - \pi) \leq \text{rank}(2^{-d} + 2^{-d-3} - \mathbf{A}) = 3$ .  $\square$

## A.2 Adjusting for additional covariates

Our method allows a mixture of matrix-valued predictors and usual vector-valued predictors, i.e., classifiers of the type  $f(\mathbf{X}) = \langle \mathbf{X}, \mathbf{B} \rangle + \mathbf{W}^T \mathbf{C}$ , where  $\mathbf{X} \in \mathbb{R}^{d_1 \times d_2}$  represents the matrix-valued predictor of our interest, and  $\mathbf{W} \in \mathbb{R}^p$  represents the additional covariate including intercept. In our neuroimaging analysis (see Section 7.1 of the main paper), for example, we have used  $\mathbf{W}$  to capture covariate such as age, gender, etc in the prediction model. Our algorithm can be easily extended to allow additional covariates in the model. Specifically, we allow classifiers of the type  $f(\mathbf{X}, \mathbf{W}) = \langle \mathbf{X}, \mathbf{B} \rangle + \mathbf{W}^T \mathbf{C}$ , where  $\mathbf{W} \in \mathbb{R}^p$  is the usual vector-valued covariate, and  $\mathbf{C} \in \mathbb{R}^p$  is the unconstrained coefficient. The only change is the primal update in the algorithm (Line 4 in Algorithm 1 of main paper). The decision variables now consist of  $(\mathbf{B}, \mathbf{C})$  and we solve them simultaneously. Because both  $\mathbf{B}$  and  $\mathbf{C}$  are unconstrained decision variables, the algorithm lends itself well to this context.

## A.3 Extension to sub-Gaussian noise

In the main paper, we have assumed the bounded noise (and thus bounded response) in the regression model. Here we extend the results to unbounded response with sub-Gaussian noise. For notational simplicity, we state the results for the matrix completion problem with  $d_1 = d_2 = d$ . The results extend similarly to general nonparamatrix matrix regression; we omit the elaboration but only state the difference in the remark.

Consider the signal plus noise model on matrix  $\mathbf{Y} \in \mathbb{R}^{d \times d}$ ,

$$\mathbf{Y} = \mathbf{\Theta} + \mathbf{E},$$

where  $\mathbf{E}$  consists of zero-mean, independent noise entries, and  $\mathbf{\Theta} \in \mathcal{M}_{\text{sgn}}(r)$  is an  $\alpha$ -smooth matrix. Theoretical results in Section 4 of the main paper are based on bounded observation  $\|\mathbf{Y}\|_\infty \leq 1$ . Here, we extend the results to unbounded observation with the following assumption.

**Assumption 1** (Sub-Gaussian noise).

1. There exists a constant  $\beta > 0$ , independent of matrix dimension, such that  $\|\Theta\|_\infty \leq \beta$ . Without loss of generality, we set  $\beta = 1$ .
2. The noise entries  $\mathbf{E}(\omega)$  are independent zero-mean sub-Gaussian random variables with variance proxy  $\sigma^2 > 0$ ; i.e.,  $\mathbb{P}(|\mathbf{E}(\omega)| \geq B) \leq 2e^{-B^2/2\sigma^2}$  for all  $B > 0$ .

We say that an event  $E$  occurs “with high probability” if  $\mathbb{P}(E)$  tends to 1 as the dimension  $d \rightarrow \infty$ . The following result show that the sub-Gaussian noise incurs an additional  $\log d$  factor compared to the bounded case.

**Theorem A.3.1** (Extension of Theorem 6 to sub-Gaussian noise). *Consider the same conditions of Theorem 6. Under Assumption 1, with high probability over training data  $\mathbf{Y}_\Omega$ , we have*

(a) (Sign matrix estimation). For all  $\pi \in [-1, 1]$  except for a finite number of levels,

$$\text{MAE}(\text{sgn}\hat{Z}_\pi, \text{sgn}(\Theta - \pi)) \lesssim \left( \frac{r\sigma^2 d \log d}{|\Omega|} \right)^{\frac{\alpha+1}{\alpha+2}} + \frac{1}{\rho(\pi, \mathcal{N})} \left( \frac{r\sigma^2 d \log d}{|\Omega|} \right). \quad (7)$$

(b) (Signal matrix estimation) Set  $H \asymp \left( \frac{|\Omega|}{r\sigma^2 d} \right)^{1/2}$ . We have

$$\text{MAE}(\hat{\Theta}, \Theta) \lesssim \tilde{O} \left\{ \left( \frac{r\sigma^2 d \log d}{|\Omega|} \right)^{\min(\frac{\alpha}{\alpha+2}, \frac{1}{2})} \right\}.$$

The proof is provided in Section B.6.

**Remark 1** (Extending to general non-parametric matrix regression). For matrix nonparametric regression (Theorem 4 of the main paper), the extension of bounded noise to sub-Gaussian noise incurs an additional  $\log n$  factor, where  $n$  is the sample size. The techniques of handling sub-Gaussian noise is identical to the above extension, and is thus omitted in the paper.

#### A.4 Extension to unbounded number of mass points

Theorem 6 of our main paper assumes the bounded  $|\mathcal{N}|_{\text{cover}} < c < \infty$  for some constant  $c > 0$ , where  $|\mathcal{N}|_{\text{cover}}$  is defined as the covering number of  $\mathcal{N}$  with  $2\Delta s$ -bin’s. Recall that  $\mathcal{N}$  corresponds to regions of jumps greater than  $\Delta s = 1/d^2$  in the CDF  $G(\pi) = \mathbb{P}_{\omega \sim \Pi}(\Theta(\omega) \leq \pi)$ . This setup gives a cleaner exposition of our results but may be restricted in some cases. For example, the high-rank matrices in Example 5 and Figure 1(b) are excluded, because  $\alpha = \infty$  and  $|\mathcal{N}|_{\text{cover}} = d$  in this setup. Fortunately, our framework still applies to this family of matrices with a little amendment.

We now extend the setup to allow for more general structured matrices including those in Example 5. Redefine  $\Delta s = 1/d$ . Correspondingly, redefine the smoothness index  $\alpha$  and the set  $\mathcal{N}$  for the psudo density of  $\Theta(\omega)$  with new bin width  $2\Delta s$ . Let  $|\mathcal{N}|_{\text{cover}}$  be the covering number of  $\mathcal{N}$  with new

$2\Delta s$ -bin's. Under this new setup, the signal matrix in Example 5 has  $|\mathcal{N}|_{\text{cover}} = 0$  and  $\alpha < \infty$ . Following the same line as in Theorem 6 and use the fact that  $\Delta s \lesssim t_n$ , we obtain that

$$\text{MAE}(\hat{\Theta}, \Theta) \lesssim t_d^{\alpha/(\alpha+2)} \log H + \frac{1}{H} + t_d H \log H, \quad \text{with } t_d = \frac{dr}{|\Omega|}.$$

Therefore, setting  $H \asymp \sqrt{\frac{1+|\mathcal{N}|_{\text{cover}}}{t_d}}$  yields the error bound

$$\text{MAE}(\hat{\Theta}, \Theta) \leq \tilde{O} \left\{ \left( \frac{dr}{|\Omega|} \right)^{\min(\frac{\alpha}{2+\alpha}, \frac{1}{2})} \right\}. \quad (8)$$

The result (8) applies to cases when the signal matrices belong to  $\mathcal{M}_{\text{sgn}}(r)$  and have at most  $d$  distinct entries with repetition patterns.

### A.5 Connection to structured matrix model with functional coefficients

In Section 6.2 of the main paper, we simulate data  $(\mathbf{X}_i, Y_i)_{i=1}^n$  from latent variable model  $(\mathbf{X}, Y)|\pi$  based on the following scheme,

$$\pi \stackrel{\text{i.i.d.}}{\sim} \text{Unif}[0, 1] \xrightarrow{\text{conditional on } \pi} \begin{cases} Y \sim \text{Ber}(\pi), \quad Y \perp \mathbf{X} | \pi, \\ \mathbf{X} = \llbracket \mathbf{X}_{pq} \rrbracket, \quad \text{where } \mathbf{X}_{pq} \stackrel{\text{indep.}}{\sim} \mathcal{N}(g_{pq}(\pi) \mathbb{1}(\text{edge } (p, q) \text{ is active}), \sigma^2). \end{cases}$$

Notice that, for any given  $\pi$ ,  $\mathbf{X}$  is a rank- $r$ ,  $(s_1, s_2)$  matrix as shown in Fig 6 of the main paper.

Here we provide justification to this simulation. We will show that the, in the absence of noise  $\sigma = 0$ , the conditional expectation  $\mathbb{E}(Y|\mathbf{X}) = f(\mathbf{X})$  from the above simulation falls into the low-rank sign-representable function family of our interest.

Specifically, we consider a structured matrix model with functional coefficients

$$\mathbf{X}_\pi \stackrel{\text{def}}{=} \mathbf{B}_0 + \sum_{s=1}^r g_s(\pi) \mathbf{B}_s + \sigma \mathbf{E}, \quad Y_\pi \sim \text{Ber}(\pi), \quad \mathbf{X}_\pi \perp Y_\pi | \pi, \quad (9)$$

where  $\pi \in [0, 1]$  is drawn from  $\text{Unif}[0, 1]$ ;  $\mathbf{E}$  is a noise matrix consisting of i.i.d. entries in  $N(0, 1)$ ;  $\sigma$  is the noise level;  $\mathbf{B}_0$  is an arbitrary baseline matrix;  $(\mathbf{B}_s)_{s=1}^r$  is a set of rank-1 matrices in  $\{0, 1\}^{d_1 \times d_2}$  that satisfy three conditions:

1. non-overlapping supports, i.e.,  $\langle \mathbf{B}_s, \mathbf{B}_{s'} \rangle = 0$  for all  $s \neq s'$
2. bounded total support, i.e.,  $\sum_{s \in [r]} \text{supp}(\mathbf{B}_s) \leq (s_1, s_2)$ ;
3. At least one of the functions  $(g_s)_{s=1}^r$  is strictly monotonic with respect to  $\pi$  for all  $s \in [r]$ .

**Proposition 2** (Connection to structured matrix model with functional coefficients). Let  $\mathbb{P}_{\mathbf{X}, Y}$  denote the joint distribution induced by  $(\mathbf{X}_\pi, Y_\pi)_{\pi \in [0, 1]}$  drawn from (9). In the noiseless case  $\sigma = 0$ , let  $f(\mathbf{X}) = \mathbb{E}(Y|\mathbf{X})$  denote the regression function based on  $\mathbb{P}_{\mathbf{X}, Y}$ . Then  $f \in \mathcal{F}_{\text{sgn}}(r, s_1, s_2)$ .

*Proof.* We restrict ourselves to the noiseless case with  $\sigma = 0$  in (9). Let

$$\mathcal{X} = \{\mathbf{X}_\pi : \mathbf{X}_\pi \text{ has structure specified in (9) for } \pi \in [0, 1]\}$$

denote the predictor space. The mapping between  $\pi$  and  $\mathbf{X} \in \mathcal{X}$  is one-to-one based on the construction of  $\mathbf{X}_\pi$ . We use  $\Pi: [0, 1] \rightarrow \mathcal{X}$  to denote the mapping and  $\Pi^{-1}$  the inverse. Based on the property 3, without loss of generality, assume  $g_1$  is a strictly increasing function.

For any given  $\pi \in [0, 1]$ , we have

$$\mathbb{E}_{Y|\pi}[Y|\pi] = \pi = \Pi^{-1}(\mathbf{X}).$$

This implies the regression function  $f = \Pi^{-1}$ . To show  $f \in \mathcal{F}_{\text{sgn}}(r, s_1, s_2)$ , it suffices to show  $\Pi^{-1} \in \mathcal{F}_{\text{sgn}}(r, s_1, s_2)$ . For any given  $\pi' \in [0, 1]$ , write

$$\begin{aligned} \{\mathbf{X} \in \mathcal{X} : \text{sgn}(\Pi^{-1} - \pi') = 1\} &= \{\mathbf{X} \in \mathcal{X} : \Pi^{-1}(\mathbf{X}) \geq \pi'\} \\ &= \{\mathbf{X} \in \mathcal{X} : g_1(\Pi^{-1}(\mathbf{X})) \geq g_1(\pi')\} \\ &= \{\mathbf{X} \in \mathcal{X} : \langle \mathbf{X}, \mathbf{B}_1 \rangle \geq g_1(\pi')\langle \mathbf{B}_1, \mathbf{B}_1 \rangle + \langle \mathbf{B}_0, \mathbf{B}_1 \rangle\}, \end{aligned}$$

where the second line uses the fact that  $g_1$  is strictly increasing.

Therefore, the sign function  $\text{sgn}(\Pi^{-1} - \pi')$  can be expressed as the sign of trace function,

$$\text{sgn}(\Pi^{-1} - \pi') = \text{sgn}\left(\underbrace{\langle \mathbf{X}, \mathbf{B}_1 \rangle}_{\text{trace}} - \underbrace{g_1(\pi')\langle \mathbf{B}_1, \mathbf{B}_1 \rangle - \langle \mathbf{B}_0, \mathbf{B}_1 \rangle}_{\text{intercept}}\right), \quad \text{for all } \mathbf{X} \in \mathcal{X},$$

where  $\mathbf{B}_1$  is a rank-1,  $\text{supp-}(s_1, s_2)$  matrix coefficient. The proof is complete.  $\square$

**Remark 2.** The above result shows the connection of our method to joint matrix model (9)  $(\mathbf{X}_\pi, Y_\pi)_{\pi \in [0, 1]}$ . We should point out, despite of the seeming similarity, a fundamental challenge arises in our setting when the latent index  $\pi$  is unobserved. Our level-set approach essentially learns the right ordering of  $\mathbf{X}_\pi$  against the index  $\pi \in [0, 1]$  (see Figure 2 of the main paper), thereby facilitating the estimation of regression function  $f$ .



## B Proofs

### B.1 Main notation

Notation	Definition
$(\mathbf{X}, Y)$	matrix predictor and univariate response
$(\mathbf{X}_i, Y_i)_{i=1}^n$	a sample of size $n$
$\mathcal{X}$	predictor space
$\bar{Y}_\pi = Y - \pi$	shifted response
$f: \mathbf{X} \mapsto \mathbb{E}(Y \mathbf{X})$	ground truth regression function
$\hat{f}: \mathbf{X} \mapsto \mathbb{R}$	estimated regression function
$f_{\text{bayes}, \pi} = \text{sgn}(f - \pi)$	Bayes classifier at level $\pi$
$S_{\text{bayes}}(\pi) = \{\mathbf{X} \in \mathcal{X}: f(\mathbf{X}) \geq \pi\}$	Indicator set corresponding to $f_{\text{bayes}, \pi}$
$r$	matrix rank
$(s_1, s_2)$	support parameter
$\mathcal{F}_{\text{sgn}}(r)$	set of $r$ -sign representable functions
$\Phi(r, s_1, s_2)$	rank- $r$ , supp- $(s_1, s_2)$ trace functions
$\Phi(r)$	family of rank- $r$ trace functions
$B$	rank- $r$ , supp- $(s_1, s_2)$ matrix in trace function
$\alpha$	smoothness index of $G(\pi)$
$\mathcal{N}$	set of mass points associated with CDF $G(\pi) = \mathbb{P}_{\mathbf{X}}[f(\mathbf{X}) \leq \pi]$
$\rho(\pi, \mathcal{N})$	distance from $\pi$ to nearest point in $\mathcal{N}$
$H$	resolution parameter in sign aggregation
$\phi$	an arbitrary classifier function from $\mathcal{X}$ to $\mathbb{R}$
$S_\phi = \{\mathbf{X} \in \mathcal{X}: \phi(\mathbf{X}) \geq 0\}$	Indicator set corresponding to $\phi$
$F$	surrogate large-margin loss function from $\mathbb{R}$ to $\mathbb{R}_{\geq 0}$
$\hat{\phi}_{\pi, F}$	estimated classifier function based on regularized empirical $F$ -risk
$\ell_{\pi, F}$	weighted $F$ -loss function, i.e., $\ell_{\pi, F}(\phi; (\mathbf{X}, Y)) =  \bar{Y}_\pi  F(\phi(\mathbf{X}) \text{sgn} \bar{Y}_\pi)$
$\text{Risk}_\pi$	weighted 0-1 risk
$\text{Risk}_{\pi, F}$	weighted surrogate $F$ -risk
$\widehat{\text{Risk}}_{\pi, F}$	empirical weighted $F$ -risk
$S, S_1, S_2$	subsets in $\mathcal{X}$
$d_\Delta(S_1, S_2)$	probability set difference, equal to $\mathbb{P}_{\mathbf{X}}(\mathbf{X} \in \mathcal{X}: \mathbf{X} \in S_1/S_2 \text{ or } S_2/S_1)$
$d_\pi(S_1, S_2)$	risk difference, equal to $\text{Risk}_\pi(\text{sgn} S_1) - \text{Risk}_\pi(\text{sgn} S_2)$
$\mathbf{Y}$	data matrix with complete observation
$\Omega \subset [d_1] \times [d_2]$	index set of observations
$\mathbf{Y}_\Omega$	data matrix with incomplete observation
$\mathcal{M}_{\text{sgn}}(r)$	family of rank- $r$ sign representable matrices
$\Theta \in \mathcal{M}_{\text{sgn}}(r)$	signal matrix in matrix completion problem
$E$	noise matrix
$Z$	an arbitrary matrix

## B.2 Proof of Theorem 1

*Proof.* Fix  $\pi \in [-1, 1]$ . For any arbitrary function  $\phi \in \Phi(r)$ , we evaluate the excess risk between  $\text{sgn}(f - \pi)$  and  $\text{sgn}\phi$ ,

$$\begin{aligned} & \text{Risk}_\pi(\text{sgn}\phi) - \text{Risk}_\pi(\text{sgn}(f - \pi)) \\ &= \frac{1}{2} \mathbb{E}_{\mathbf{X}} \underbrace{\mathbb{E}_{Y|\mathbf{X}} \{ |Y - \pi| [ |\text{sgn}(Y - \pi) - \text{sgn}\phi| - |\text{sgn}(Y - \pi) - \text{sgn}(f - \pi)| ] \}}_{\stackrel{\text{def}}{=} I}. \end{aligned} \quad (10)$$

Here,  $I = I(\mathbf{X})$  is a function of  $\mathbf{X}$ , and its expression can be simplified as

$$\begin{aligned} I &= \mathbb{E}_{Y|\mathbf{X}} [(Y - \pi)(\text{sgn}(f - \pi) - \text{sgn}\phi)\mathbf{1}(Y \geq \pi) + (\pi - Y)(\text{sgn}\phi - \text{sgn}(f - \pi))\mathbf{1}(Y < \pi)] \\ &= \mathbb{E}_{Y|\mathbf{X}} [(\text{sgn}(f - \pi) - \text{sgn}\phi)(Y - \pi)] \\ &= [\text{sgn}(f - \pi) - \text{sgn}\phi] [f - \pi] \\ &= |\text{sgn}(f - \pi) - \text{sgn}\phi| |f - \pi|, \end{aligned} \quad (11)$$

where the third line uses the fact  $\mathbb{E}_{Y|\mathbf{X}} Y = f(\mathbf{X})$ . Combining (11) with (10), we conclude that, for all  $\phi \in \Phi(r)$ ,

$$\text{Risk}_\pi(\text{sgn}\phi) - \text{Risk}_\pi(\text{sgn}(f - \pi)) = \frac{1}{2} \mathbb{E}_{\mathbf{X}} |\text{sgn}(f - \pi) - \text{sgn}\phi| |f - \pi| \geq 0,$$

where the last line equals to zero when  $\text{sgn}\phi = \text{sgn}(f - \pi)$  or  $f \equiv \pi$  is a constant function. Note that  $(f - \pi)$  is  $r$ -sign representable by assumption. Therefore,

$$\text{Risk}_\pi(\text{sgn}(f - \pi)) = \inf\{\text{Risk}_\pi(\text{sgn}\phi) : \phi \in \Phi(r)\}.$$

Based on the definition of 0-1 classification loss, the  $\text{Risk}_\pi(\cdot)$  relies only on the sign of the argument function. Therefore, for all functions  $\bar{f} : \mathcal{X} \rightarrow \mathbb{R}$  that have the same sign as  $\text{sgn}(f - \pi)$ , we have

$$\text{Risk}_\pi(\bar{f}) = \inf\{\text{Risk}_\pi(\text{sgn}\phi) : \phi \in \Phi(r)\} = \inf\{\text{Risk}_\pi(\phi) : \phi \in \Phi(r)\}.$$

□

## B.3 Proof of Theorem 2

*Proof.* Fix  $\pi \in [-1, 1]$ . For ease of notation, we drop the dependence of  $\pi$  in  $S_{\text{bayes}}(\pi)$  and simply write  $S_{\text{bayes}}$ . Based on (11) in the proof of Theorem 1, we have

$$\begin{aligned} d_\pi(S, S_{\text{bayes}}) &\stackrel{\text{def}}{=} \text{Risk}_\pi(\text{sgn}(S)) - \text{Risk}_\pi(\text{sgn}(S_{\text{bayes}})) \\ &= \frac{1}{2} \mathbb{E}_{\mathbf{X}} (|\text{sgn}(S) - \text{sgn}(S_{\text{bayes}})| |\pi - f|) \\ &= \int_{\mathbf{X} \in S \Delta S_{\text{bayes}}} |f(\mathbf{X}) - \pi| d\mathbb{P}_{\mathbf{X}}. \end{aligned} \quad (12)$$

We divide the proof into two cases:  $\alpha > 0$  and  $\alpha = \infty$ .

Case 1:  $\alpha > 0$ .

Consider an arbitrary set  $S \subset \mathbb{R}^{d_1 \times d_2}$ . Let  $t$  be an arbitrary number in the interval  $[0, 1]$ , and define the set  $A = \{\mathbf{X} \in \mathcal{X} : |f(\mathbf{X}) - \pi| > t\}$ .

$$\begin{aligned} \int_{\mathbf{X} \in S \Delta S_{\text{bayes}}} |f(\mathbf{X}) - \pi| d\mathbb{P}_{\mathbf{X}} &\geq t [\mathbb{P}_{\mathbf{X}}((S \Delta S_{\text{bayes}}) \cap A)] \\ &\geq t (\mathbb{P}_{\mathbf{X}}(S \Delta S_{\text{bayes}}) - \mathbb{P}_{\mathbf{X}}(A^c)) \\ &\geq t (\mathbb{P}_{\mathbf{X}}(S \Delta S_{\text{bayes}}) - Ct^\alpha), \quad \text{for all } 0 \leq t < \rho(\pi, \mathcal{N}), \end{aligned}$$

where the last inequality is from  $\alpha$ -globally smoothness condition. Combining the above inequality with the identity (12) yields

$$d_\pi(S, S_{\text{bayes}}) \geq t (\mathbb{P}_{\mathbf{X}}(S \Delta S_{\text{bayes}}) - Ct^\alpha), \quad \text{for all } 0 \leq t < \rho(\pi, \mathcal{N}). \quad (13)$$

We maximize the lower bound of (13) with respect to  $t$ , and obtain the optimal  $t_{\text{opt}}$ ,

$$t_{\text{opt}} = \begin{cases} \rho(\pi, \mathcal{N}), & \text{if } \mathbb{P}_{\mathbf{X}}(S \Delta S_{\text{bayes}}) > C(1 + \alpha)\rho^\alpha(\pi, \mathcal{N}), \\ \left[ \frac{1}{2C(1 + \alpha)} \mathbb{P}_{\mathbf{X}}(S \Delta S_{\text{bayes}}) \right]^{1/\alpha}, & \text{if } \mathbb{P}_{\mathbf{X}}(S \Delta S_{\text{bayes}}) \leq C(1 + \alpha)\rho^\alpha(\pi, \mathcal{N}). \end{cases}$$

The corresponding lower bound of the inequality (13) becomes

$$d_\pi(S, S_{\text{bayes}}) \geq \begin{cases} c_1 \rho(\pi, \mathcal{N}) \mathbb{P}_{\mathbf{X}}(S \Delta S_{\text{bayes}}), & \text{if } \mathbb{P}_{\mathbf{X}}(S \Delta S_{\text{bayes}}) > C(1 + \alpha)\rho^\alpha(\pi, \mathcal{N}), \\ c_2 [\mathbb{P}_{\mathbf{X}}(S \Delta S_{\text{bayes}})]^{\frac{1+\alpha}{\alpha}}, & \text{if } \mathbb{P}_{\mathbf{X}}(S \Delta S_{\text{bayes}}) \leq C(1 + \alpha)\rho^\alpha(\pi, \mathcal{N}), \end{cases}$$

where  $c_1, c_2 > 0$  are two constants independent of  $S$ . Combining both cases gives

$$d_\Delta(S, S_{\text{bayes}}) \stackrel{\text{def}}{=} \mathbb{P}_{\mathbf{X}}(S \Delta S_{\text{bayes}}) \lesssim [d_\pi(S, S_{\text{bayes}})]^{\frac{\alpha}{1+\alpha}} + \frac{1}{\rho(\pi, \mathcal{N})} d_\pi(S, S_{\text{bayes}}), \quad (14)$$

where we have absorbed the constants into the relationship  $\lesssim$ .

Case 2:  $\alpha = \infty$ .

The inequality (13) now becomes

$$d_\pi(S, S_{\text{bayes}}) \geq t \mathbb{P}_{\mathbf{X}}(S \Delta S_{\text{bayes}}) = t d_\Delta(S, S_{\text{bayes}}), \quad \text{for all } 0 \leq t < \rho(\pi, \mathcal{N}). \quad (15)$$

The conclusion (14) follows by taking  $t = \frac{\rho(\pi, \mathcal{N})}{2}$  in the inequality (15).

□

**Remark 3** (Bounding  $L_1$  distance by classification risk). The bound controls the  $L_1$  distance to  $f_{\text{bayes}, \pi} = \text{sgn}(f - \pi)$  using the classification excess risk to  $\text{Risk}_\pi(f_{\text{bayes}, \pi})$ . The result applies

uniformly to  $\pi \in [-1, 1]$  if  $f$  is globally- $\alpha$  smooth; i.e., the bound

$$\|\text{sgn}\phi - f_{\text{bayes},\pi}\|_1 \lesssim [\text{Risk}_\pi(\phi) - \text{Risk}_\pi(f_{\text{bayes},\pi})]^{\frac{\alpha}{1+\alpha}} + \frac{1}{\rho(\pi, \mathcal{N})} [\text{Risk}_\pi(\phi) - \text{Risk}_\pi(f_{\text{bayes},\pi})]$$

holds for all functions  $\phi: \mathcal{X} \rightarrow \mathbb{R}$  and for all  $\pi \in [-1, 1]$  except for a finite number of points. In fact, the similar inequality holds by replacing the 0-1 risk to hinge risk or  $T$ -truncated hinge risk. Specifically, the following bound holds for all functions  $\phi: \mathcal{X} \rightarrow \mathbb{R}$  and all  $\pi \in [-1, 1]$  except for a finite number of points.

- For hinge loss  $F(z) = (1 - z)_+$ ,

$$\|\phi - f_{\text{bayes},\pi}\|_1 \lesssim [\text{Risk}_{\pi,F}(\phi) - \text{Risk}_{\pi,F}(f_{\text{bayes},\pi})]^{\frac{\alpha}{1+\alpha}} + \frac{1}{\rho(\pi, \mathcal{N})} [\text{Risk}_{\pi,F}(\phi) - \text{Risk}_{\pi,F}(f_{\text{bayes},\pi})]$$

- For  $T$ -truncated hinge loss  $F(z) = \min((1 - z)_+, T)$  with  $T \geq 2$ ,

$$\|\phi^T - f_{\text{bayes},\pi}\|_1 \lesssim [\text{Risk}_{\pi,F}(\phi) - \text{Risk}_{\pi,F}(f_{\text{bayes},\pi})]^{\frac{\alpha}{1+\alpha}} + \frac{1}{\rho(\pi, \mathcal{N})} [\text{Risk}_{\pi,F}(\phi) - \text{Risk}_{\pi,F}(f_{\text{bayes},\pi})],$$

where  $\phi^T$  is a truncation of function  $\phi$  (see formal definition in (22)).

See Lemma 4.

## B.4 Proofs of Theorem 3, Part (a) in Theorems 5 and 7

We provide a unified framework that incorporates Theorem 3, Part (a) in Theorems 5 and 7 in the main paper. For any  $\pi \in [-1, 1]$ , write  $\bar{Y}_\pi = Y - \pi$ , and let  $\ell_{\pi,F}(\phi; (\mathbf{X}, Y))$  denote the weighted  $F$ -loss

$$\ell_{\pi,F}(\phi; (\mathbf{X}, Y)) \stackrel{\text{def}}{=} |\bar{Y}_\pi| F(\text{sgn}(Y - \pi)\phi(\mathbf{X})),$$

where the loss function  $F$  could be either standard 0-1 loss  $F(z) = \mathbf{1}(z > 0)$  or surrogate loss satisfying Assumption 1. Consider the large-margin estimate

$$\hat{\phi}_{\pi,F} = \arg \min_{\phi \in \Phi(r, s_1, s_2)} \left\{ \frac{1}{n} \sum_{i=1}^n \ell_{\pi,F}(\phi; (\mathbf{X}_i, Y_i)) + \lambda \|\phi\|_F^2 \right\}. \quad (16)$$

The following theorem states the accuracy for sign function estimate  $\text{sgn}(\hat{\phi}_{\pi,F}): \mathcal{X} \rightarrow \{-1, 1\}$ . We say that an event  $E$  occurs “with high probability” if  $\mathbb{P}(E)$  tends to 1 as the dimension  $d \rightarrow \infty$ .

**Theorem B.4.1** (Large-margin sign estimation). *Fix  $\pi \notin \mathcal{N}$ . Suppose the regression function  $f \in \mathcal{F}_{\text{sgn}}(r, s_1, s_2)$  is  $(\pi, \alpha)$ -smooth over  $\mathcal{X}$ . Then, with high probability at least  $1 - \exp(-nt_n)$  over training data  $(\mathbf{X}_i, Y_i)_{i \in [n]}$ , the estimate (16) satisfies*

$$\|\text{sgn}\hat{\phi}_{\pi,F} - f_{\text{bayes},\pi}\|_1 \lesssim t_n^{\alpha/(2+\alpha)} + \frac{1}{\rho^2(\pi, \mathcal{N})} t_n, \quad (17)$$

under the following three specifications:

(a) (Theorem 3) 0-1 loss  $F(z) = \mathbb{1}(z > 0)$ , no penalization  $\lambda = 0$ ,  $(s_1, s_2) = (d_1, d_2)$ , and  $t_n = \frac{1}{n} r d_{\max}$ .

(b) (Theorem 5) 0-1 loss  $F(z) = \mathbb{1}(z > 0)$ , no penalization  $\lambda = 0$ , constant  $(s_1, s_2)$ , and  $t_n = \frac{1}{n} r (s_1 + s_2) \log d_{\max}$ .

(c) (Theorem 7) Surrogate loss satisfying Assumption 1, constant  $(s_1, s_2)$ ,  $t_n = \frac{1}{n} r (s_1 + s_2) \log d_{\max}$ , penalization  $\lambda \asymp (t_n^{(\alpha+1)/(\alpha+2)} + \rho t_n)$ , approximation error  $a_n^{(\alpha+1)/(\alpha+2)} \leq t_n$ .

Here, the constants absorbed in the  $\lesssim$  of (17) are independent of  $\pi$ .

**Remark 4** (Ridge penalization). The estimation under 0-1 loss requires no penalization, because only the sign, but not the magnitude, of  $\phi$  affects the 0-1 risk. One can constrain  $\|\phi\|_F = 1$  in the empirical 0-1 risk minimization without altering the solution. In contrast, the surrogate loss such as hinge loss is scale-sensitive, rendering the possible unboundedness of  $\phi$ . We impose penalization to control the magnitude of the  $\|\phi\|_F$  and thus the local complexity. The resulting estimation enjoys the fast convergence as in sieve estimate [Shen and Wong, 1994] under well tuned  $\lambda$ .

We provide the proof after introducing two main lemmas. There are two key ingredients in the proof. The first step is to quantify the convergence of  $\hat{\phi}_{\pi,F}$ 's excess  $F$ -risk using Lemmas 1 and 2. The second step is to relate the excess  $F$ -risk to excess 0-1 risk using Lemma 1, and then establish the sign function accuracy using Theorem 2.

Recall that  $\hat{\phi}_{\pi,F}$  is the minimizer of empirical  $F$ -risk. To quantify the  $\hat{\phi}_{\pi,F}$ 's excess  $F$ -risk, we notice that

$$\begin{aligned} & \text{Risk}_{\pi,F}(\hat{\phi}_{\pi,F}) - \inf_{\text{all } \phi} \text{Risk}_{\pi,F}(\phi) \\ &= \underbrace{\text{Risk}_{\pi,F}(\hat{\phi}_{\pi,F}) - \inf_{\phi \in \Phi(r, s_1, s_2)} \text{Risk}_{\pi,F}(\phi^*)}_{\text{estimation error}} + \underbrace{\inf_{\phi \in \Phi(r, s_1, s_2)} \text{Risk}_{\pi,F}(\phi) - \inf_{\text{all } \phi} \text{Risk}_{\pi,F}(\phi)}_{\text{approximation error}}, \end{aligned}$$

The simplest way to bound  $\hat{\phi}_{\pi,F}$ 's excess risk is to use a uniform convergence of excess risk over classifiers  $\Phi(r, s_1, s_2)$ ; however, this approach ignores the local complexity around  $\hat{\phi}_{\pi,F}$  and yields a suboptimal rate. Here we adopt the local iterative techniques of Wang et al. [2008, Theorem 3] to obtain a better rate. The improvement stems from the fact that, under considered assumptions, the variance of the excess loss is bounded in terms of its expectation. Because the variance decreases as we approach the optimal  $\phi_\pi^*$ , the risk of the empirical minimizer converges more quickly to the optimal risk than the simple uniform converge results would suggest.

The following result summarizes the key properties of four common losses: 0-1 loss, hinge loss,  $T$ -truncated hinge loss, and psi-loss. Here, the  $T$ -truncated hinge loss is defined as  $F(z) = \min((1 - z)_+, T)$  for a given  $T \geq 2$ . We will use  $T$ -truncated hinge loss to facilitate the proofs of Lemma 2 and Theorem B.4.1.

**Lemma 1** (Conversion inequalities). *Suppose the regression function  $f$  is  $(\pi, \alpha)$ -smooth, and denote*

$f_{\text{bayes},\pi} = \text{sgn}(f - \pi)$  for  $\pi \in [-1, 1]$ . Let  $F$  be 0-1 loss, hinge loss,  $T$ -truncated hinge loss, or psi-loss. Then, the following three properties hold for all  $\pi \in [-1, 1]$ .

(a) *Optimality*:  $\inf_{\text{all } \phi} \text{Risk}_{\pi,F}(\phi) = \text{Risk}_{\pi,F}(f_{\text{bayes},\pi})$ .

(b) *Excess risk bound*: for all classifiers  $\phi: \mathcal{X} \rightarrow \mathbb{R}$ ,

$$\text{Risk}_{\pi}(\phi) - \text{Risk}_{\pi}(f_{\text{bayes},\pi}) \leq C [\text{Risk}_{\pi,F}(\phi) - \text{Risk}_{\pi,F}(f_{\text{bayes},\pi})], \quad (18)$$

where  $C = 1$  for 0-1, hinge loss or  $T$ -truncated loss, and  $C = 2$  for psi-loss.

(c) *Variance-to-mean relationship*: Suppose  $F$  is 0-1 loss,  $T$ -truncated loss, or psi-loss. Then, for all classifiers  $\phi: \mathcal{X} \rightarrow \mathbb{R}$ ,

$$\begin{aligned} & \text{Var} [\ell_{\pi,F}(\phi; (\mathbf{X}, Y)) - \ell_{\pi,F}(f_{\text{bayes},\pi}; (\mathbf{X}, Y))] \\ & \lesssim [\text{Risk}_{\pi,F}(\phi) - \text{Risk}_{\pi,F}(f_{\text{bayes},\pi})]^{\alpha/(1+\alpha)} + \frac{1}{\rho(\pi, \mathcal{N})} [\text{Risk}_{\pi,F}(\phi) - \text{Risk}_{\pi,F}(f_{\text{bayes},\pi})]. \end{aligned} \quad (19)$$

**Remark 5.** The property (c) holds only for bounded loss functions, i.e, excluding hinge loss.

*Proof.* Case 1:  $F(z) = \mathbb{1}(z < 0)$  is 0-1 loss.

Properties (a) and (b) directly follow from Theorem 1. To prove (c), we expand the variance by

$$\begin{aligned} \text{Var} [\ell_{\pi}(\phi; (\mathbf{X}, Y)) - \ell_{\pi}(f_{\text{bayes},\pi}, (\mathbf{X}, Y))] & \lesssim \mathbb{E} |\ell_{\pi}(\phi; (\mathbf{X}, Y)) - \ell_{\pi}(f_{\text{bayes},\pi}, (\mathbf{X}, Y))|^2 \\ & \lesssim \mathbb{E} |\ell_{\pi}(\phi; (\mathbf{X}, Y)) - \ell_{\pi}(f_{\text{bayes},\pi}, (\mathbf{X}, Y))| \\ & \lesssim \mathbb{E} |\text{sgn} \bar{Y}_{\pi} - \text{sgn} \phi(\mathbf{X})| - |\text{sgn} \bar{Y}_{\pi} - f_{\text{bayes},\pi}(\mathbf{X})| \\ & \leq \mathbb{E} |\text{sgn} \phi - f_{\text{bayes},\pi}|, \end{aligned} \quad (20)$$

where the second line comes from the boundedness of 0-1 loss, and the third line comes from the boundedness of weight  $|\bar{Y}_{\pi}|$ , and fourth line comes from the inequality  $||a - b| - |c - b|| \leq |a - c|$  for  $a, b, c \in \{-1, 1\}$ . Here we have absorbed the constant multipliers in  $\lesssim$ . Therefore, the conclusion (19) then directly follows by applying Remark 3 to (20).

Case 2:  $F(z) = (1 - z)_+$  is hinge loss.

Property (a) was firstly introduced in Wang et al. [2008, Lemma 1], and here we provide an alternative proof. A direct calculation (see Lemma 4) shows that

$$\text{Risk}_{\pi,F}(\phi) - \text{Risk}_{\pi,F}(f_{\text{bayes},\pi}) \geq \mathbb{E} |\phi - f_{\text{bayes},\pi}| |f - \pi| \geq 0,$$

Therefore,  $\inf_{\text{all } \phi} \text{Risk}_{\pi,F}(\phi) = \text{Risk}_{\pi,F}(f_{\text{bayes},\pi})$ . Property (18) is from Scott [2011, Corollary 1].

Case 3: When  $F(z) = 2 \min(1, (1 - z)_+)$  is psi-loss.

Again, the property (a) follows from Wang et al. [2008, Lemma 1]. For the property (18), we

use [Scott \[2011, Theorem 1\]](#) to find the transformation function  $\psi$  that relates 0-1 risk to F-risk:

$$\psi(\text{Risk}_\pi(\phi) - \text{Risk}_\pi(f_{\text{bayes},\pi})) \leq \text{Risk}_{\pi,F}(\phi) - \text{Risk}_{\pi,F}(f_{\text{bayes},\pi}).$$

To put our problem in the context of [Scott \[2011\]](#), we need additional notation. For any function measurable  $g: x \mapsto g(x)$ , we write  $g = g^+ - g^-$ , where  $g^+$  and  $g^-$  are two non-negative functions given by

$$g^+(x) = \max\{g(x), 0\} = \begin{cases} g(x), & \text{if } g(x) > 0, \\ 0, & \text{otherwise,} \end{cases} \quad g^-(x) = \max\{-g(x), 0\} = \begin{cases} -g(x), & \text{if } g(x) < 0, \\ 0, & \text{otherwise.} \end{cases}$$

Under this notation, we have  $|g| = g^+ + g^-$ .

Define the conditional  $F$ -risk

$$C_{\pi,F}(\mathbf{X}, t) := F(t)\mathbb{E}_{Y|\mathbf{X}}(Y - \pi)^+ + F(-t)\mathbb{E}_{Y|\mathbf{X}}(Y - \pi)^-.$$

A direct calculation shows that

$$C_{\pi,F}(\mathbf{X}, t) = \begin{cases} 2\mathbb{E}_{Y|\mathbf{X}}(Y - \pi)^-, & \text{if } t \geq 1, \\ 2\mathbb{E}_{Y|\mathbf{X}}|Y - \pi| - 2t\mathbb{E}_{Y|\mathbf{X}}(Y - \pi)^+, & \text{if } t \in [0, 1), \\ 2\mathbb{E}_{Y|\mathbf{X}}|Y - \pi| + 2t\mathbb{E}_{Y|\mathbf{X}}(Y - \pi)^-, & \text{if } t \in [-1, 0), \\ 2\mathbb{E}_{Y|\mathbf{X}}(Y - \pi)^+, & \text{if } t < -1. \end{cases}$$

Therefore, following the notation of [Scott \[2011\]](#), we have

$$H_{\pi,F}(\mathbf{X}) := \inf_{t \in \mathbb{R}: t(f(\mathbf{X}) - \pi) \leq 0} C_{\pi,F}(\mathbf{X}, t) - \inf_{t \in \mathbb{R}} C_{\pi,F}(\mathbf{X}, t) = 2|f(\mathbf{X}) - \pi|.$$

Applying [Scott \[2011, Theorem 1\]](#) to the above setup gives the excess risk transformation rule:  $\psi: z \rightarrow 2|z|$ . Therefore, the property (18) is proved.

To prove (19), notice that

$$\begin{aligned} & \text{Var} \left\{ |\bar{Y}_\pi| \left[ F(\phi(\mathbf{X})\text{sgn}\bar{Y}_\pi) - F(f_{\text{bayes},\pi}(\mathbf{X})\text{sgn}\bar{Y}_\pi) \right] \right\} \\ & \lesssim \mathbb{E} |\bar{Y}_\pi| |F(\phi(\mathbf{X})\text{sgn}\bar{Y}_\pi) - F(f_{\text{bayes},\pi}(\mathbf{X})\text{sgn}\bar{Y}_\pi)| \\ & \lesssim \underbrace{\mathbb{E} |1 - \text{sgn}(\phi(\mathbf{X})\bar{Y}_\pi) - F(f_{\text{bayes},\pi}(\mathbf{X})\text{sgn}\bar{Y}_\pi)|}_{=:(i)} + \underbrace{\mathbb{E} |\bar{Y}_\pi| |F(\phi(\mathbf{X})\text{sgn}\bar{Y}_\pi) - (1 - \text{sgn}(\phi(\mathbf{X})\bar{Y}_\pi))|}_{=:(ii)} \end{aligned} \quad (21)$$

The first term (i) is bounded as follows

$$\begin{aligned} (i) &= \mathbb{E} |\text{sgn}(\phi(\mathbf{X})\bar{Y}_\pi) - \text{sgn}(f_{\text{bayes},\pi}(\mathbf{X})\bar{Y}_\pi)| \lesssim d_\Delta(S_\phi, S_{\text{bayes}}(\pi)) \\ &\lesssim d_\pi^\alpha(S_\phi, S_{\text{bayes}}(\pi)) + \frac{1}{\rho(\pi, \mathcal{N})} d_\pi(S_\phi, S_{\text{bayes}}(\pi)), \end{aligned}$$

where the first line uses the fact that  $F(1) = 0$  and  $F(-1) = 2$ , and last inequality is from Theorem 2. Here we define indicator set corresponding  $\phi$  as  $S_\phi = \{\mathbf{X} \in \mathcal{X} : \phi(\mathbf{X}) \geq 0\}$ . The second term (ii) is bounded as follows

$$\begin{aligned}
\text{(ii)} &= \mathbb{E} [|\bar{Y}_\pi| F(\phi(\mathbf{X}) \text{sgn} \bar{Y}_\pi) - |\bar{Y}_\pi| (1 - \text{sgn}(\phi(\mathbf{X}) \bar{Y}_\pi))] \\
&= \mathbb{E} [|\bar{Y}_\pi| F(\phi(\mathbf{X}) \text{sgn} \bar{Y}_\pi) - |\bar{Y}_\pi| F(f_{\text{bayes}, \pi}(\mathbf{X}) \text{sgn} \bar{Y}_\pi)] \\
&\quad + \mathbb{E} [|\bar{Y}_\pi| (1 - \text{sgn}(f_{\text{bayes}, \pi} \bar{Y}_\pi)) - |\bar{Y}_\pi| (1 - \text{sgn}(\phi(\mathbf{X}) \bar{Y}_\pi))] \\
&\leq [\text{Risk}_{\pi, F}(\phi) - \text{Risk}_{\pi, F}(f_{\text{bayes}, \pi})] + d_\pi(S_\phi, S_{\text{bayes}(\pi)}),
\end{aligned}$$

where the first equality is based on  $F(z) = 1 - \text{sgn}(z)$  if  $z = 1$  or  $-1$ , and the last inequality is from definition of  $d_\pi(\cdot, \cdot)$ . Notice we have  $d_\pi(S_\phi, S_{\text{bayes}(\pi)}) = \text{Risk}_\pi(\phi) - \text{Risk}_\pi(f_{\text{bayes}, \pi})$  by definition. Therefore, the proof is complete by combining (21), (18) and bounds (i)-(ii).

Case 4:  $F(z) = \min((1 - z)_+, T)$  for  $T$ -truncated hinge loss, for given  $T \geq 2$ . A direct calculation (c.f. Remark 6 after Lemma 4) shows that

$$\text{Risk}_{\pi, F}(\phi) - \text{Risk}_{\pi, F}(f_{\text{bayes}, \pi}) \geq \mathbb{E} |\phi^T - f_{\text{bayes}, \pi}| |f - \pi| \geq 0,$$

where  $\phi^T : \mathcal{X} \rightarrow [-(T - 1), (T - 1)]$  denotes the  $(T - 1)$ -truncation of  $\phi$ ,

$$\phi^T = \begin{cases} T - 1 & \text{if } \phi > T - 1, \\ \phi, & \text{if } |\phi| \leq T - 1, \\ -(T - 1), & \text{if } \phi < -(T - 1). \end{cases} \quad (22)$$

Therefore,  $\inf_{\text{all } \phi} \text{Risk}_{\pi, F}(\phi) = \text{Risk}_{\pi, F}(f_{\text{bayes}, \pi})$ . To show property (18), we again use Scott [2011, Theorem 1] to find the transformation function  $\psi$  that relates 0-1 risk to F-risk:

$$\psi(\text{Risk}_\pi(\phi) - \text{Risk}_\pi(f_{\text{bayes}, \pi})) \leq \text{Risk}_{\pi, F}(\phi) - \text{Risk}_{\pi, F}(f_{\text{bayes}, \pi}).$$

Using similar arguments as in Case 3, we obtain the conditional  $F$ -risk

$$C_{\pi, F}(\mathbf{X}, t) = \begin{cases} \min \{T, (1 + t) \mathbb{E}_{Y|\mathbf{X}}(Y - \pi)^-\}, & \text{if } t \geq 1, \\ \mathbb{E}_{Y|\mathbf{X}} |Y - \pi| - t(f(\mathbf{X}) - \pi), & \text{if } t \in [0, 1), \\ \mathbb{E}_{Y|\mathbf{X}} |Y - \pi| + t(f(\mathbf{X}) - \pi), & \text{if } t \in [-1, 0), \\ \min \{T, (1 - t) \mathbb{E}_{Y|\mathbf{X}}(Y - \pi)^+\}, & \text{if } t < -1. \end{cases}$$

Therefore, following the notation of Scott [2011], we have

$$H_{\pi, F}(\mathbf{X}) := \inf_{t \in \mathbb{R} : t(f(\mathbf{X}) - \pi) \leq 0} C_{\pi, F}(\mathbf{X}, t) - \inf_{t \in \mathbb{R}} C_{\pi, F}(\mathbf{X}, t) = |f(\mathbf{X}) - \pi|.$$

Applying Scott [2011, Theorem 1] to the above setup gives the excess risk transformation rule:  $\psi : z \rightarrow |z|$ . Therefore, the property (18) is proved.



To prove (19), we use Lemma 4 and the boundedness condition of  $\|F\|_\infty \leq T$ . Specifically, we bound the variance using the  $L$ -1 distance between  $\phi$  and  $f_{\text{bayes},\pi}$ ,

$$\begin{aligned} & \text{Var} \left\{ |\bar{Y}_\pi| \left[ F(\phi(\mathbf{X})\text{sgn}\bar{Y}_\pi) - F(f_{\text{bayes},\pi}(\mathbf{X})\text{sgn}\bar{Y}_\pi) \right] \right\} \\ & \leq 4\mathbb{E} |F(\phi(\mathbf{X})\text{sgn}\bar{Y}_\pi) - F(f_{\text{bayes},\pi}(\mathbf{X})\text{sgn}\bar{Y}_\pi)|^2 \\ & \lesssim T\mathbb{E} |F(\phi(\mathbf{X})\text{sgn}\bar{Y}_\pi) - F(f_{\text{bayes},\pi}(\mathbf{X})\text{sgn}\bar{Y}_\pi)| \\ & \lesssim T\mathbb{E} |\phi^T - f_{\text{bayes},\pi}| \end{aligned}$$

where  $T > 0$  is the upper bound of truncated hinge loss, the second inequality comes from the boundedness of the  $T$ -truncated hinge loss, and the last line comes from 1-Lipschitz continuity and the definition of  $F$ . Applying Remark 6 in Lemma 4 on the last inequality complete the proof.  $\square$

Below we establish the estimation convergence rate for  $\hat{\phi}_{\pi,F}$ 's excess F-risk. The variance-to-mean relationship in Lemma 1 plays a key role in determining the convergence rate based on empirical process theory [Shen and Wong, 1994]. The proof adopts the local iterative techniques from Wang et al. [2008, Theorem 3]. Similar techniques have been used in Bartlett et al. [2006, Theorem 4] for similar estimate but without ridge penalization.

**Lemma 2** (Excess  $F$ -risk error). *Consider the set-up as in Theorem B.4.1. Then, with high probability and  $t_n$  specified in Theorem B.4.1, the following holds for all  $\pi \notin \mathcal{N}$ .*

(a) *If  $F$  is 0-1 loss or psi-loss, then*

$$\text{Risk}_\pi(\hat{\phi}_{\pi,F}) - \text{Risk}_\pi(f_{\text{bayes},\pi}) \lesssim \text{Risk}_{\pi,F}(\hat{\phi}_{\pi,F}) - \text{Risk}_{\pi,F}(f_{\text{bayes},\pi}) \lesssim t_n^{(\alpha+1)/(\alpha+2)} + \frac{1}{\rho(\pi, \mathcal{N})} t_n.$$

(b) *If  $F$  is hinge loss, then*

$$\text{Risk}_\pi(\hat{\phi}_{\pi,F}) - \text{Risk}_\pi(f_{\text{bayes},\pi}) \lesssim \text{Risk}_{F'}(\hat{\phi}_{\pi,F}) - \text{Risk}_{F'}(f_{\text{bayes},\pi}) \lesssim t_n^{(\alpha+1)/(\alpha+2)} + \frac{1}{\rho(\pi, \mathcal{N})} t_n,$$

where  $\text{Risk}_{F'}(\phi) := \mathbb{E} [|\bar{Y}_\pi| F'(\phi(\mathbf{X})\text{sgn}\bar{Y}_\pi)]$  denotes the risk evaluated under  $T$ -truncated hinge loss  $F' = \min(T, (1-z)_+)$ , and  $T = \max(2, J) \geq \max(2, \|\phi_\pi^{(n)}\|_K)$  is a constant based on Assumption 1(a).

*Proof of Lemma 2.* In Theorem 3 and 5, we have  $\lambda = a_n = 0$ . We consider the most general case  $\lambda, a_n \geq 0$  as in Theorem 7. We first consider the bounded loss (0-1 loss or psi-loss). The modification for unbounded loss (hinge loss) incurs only slight difference in the proof.

Fix  $\pi \notin \mathcal{N}$ , and write  $\rho = \rho(\pi, \mathcal{N})$ . For any function  $\phi \in \Phi(r, s_1, s_2)$  of consideration, define

empirical weighted  $F$ -risk

$$\widehat{\text{Risk}}_{\pi,F}(\phi) = \frac{1}{n} \sum_{i=1}^n \ell_{\pi,F}(\phi; (\mathbf{X}_i, Y_i)).$$

Under the notation, our estimate  $\hat{\phi}_{\pi,F}$  is the minimizer of the regularized empirical  $F$ -risk,

$$\hat{\phi}_{\pi,F} = \arg \min_{\phi \in \Phi(r, s_1, s_2)} \left\{ \widehat{\text{Risk}}_{\pi,F}(\phi) + \lambda \|\phi\|_K^2 \right\}. \quad (23)$$

We are interested in the convergence of  $\hat{\phi}_{\pi,F}$ 's excess risk,  $\text{Risk}_{\pi,F}(\hat{\phi}_{\pi,F}) - \text{Risk}_{\pi,F}(f_{\text{bayes},\pi})$ .

Let  $L_n$  denote the convergence rate to seek. By the definition of  $\hat{\phi}_{\pi,F}$ , we have

$$\widehat{\text{Risk}}_{\pi,F}(\hat{\phi}_{\pi,F}) + \lambda \|\hat{\phi}_{\pi,F}\|_K^2 \leq \widehat{\text{Risk}}_{\pi,F}(\phi_\pi^{(n)}) + \lambda J,$$

where  $\phi_\pi^{(n)}$  is a sequence of functions in Assumption 1(a). Therefore, we have the following inclusion of probability events,

$$\begin{aligned} & \left\{ (\mathbf{X}_i, Y_i)_{i \in [n]} : \text{Risk}_{\pi,F}(\hat{\phi}_{\pi,F}) - \text{Risk}_{\pi,F}(f_{\text{bayes},\pi}) \geq 2L_n \right\} \\ & \subset \left\{ (\mathbf{X}_i, Y_i)_{i \in [n]} : \exists \phi \in \Phi(r, s_1, s_2), \text{ s.t. } \text{Risk}_{\pi,F}(\phi; (\mathbf{X}, Y)) - \text{Risk}_{\pi,F}(f_{\text{bayes},\pi}) \geq 2L_n, \right. \\ & \quad \left. \text{and } \widehat{\text{Risk}}_{\pi,F}(\hat{\phi}_{\pi,F}) + \lambda \|\hat{\phi}_{\pi,F}\|_K^2 \leq \widehat{\text{Risk}}_{\pi,F}(\phi_\pi^{(n)}) + \lambda J \right\} \\ & \subset \left\{ (\mathbf{X}_i, Y_i)_{i \in [n]} : \sup_{\substack{\phi \in \Phi(r, s_1, s_2) \\ \text{Risk}_{\pi,F}(\phi; (\mathbf{X}, Y)) - \text{Risk}_{\pi,F}(f_{\text{bayes},\pi}) \geq 2L_n}} \left[ \widehat{\text{Risk}}_{\pi,F}(\phi_\pi^{(n)}) - \widehat{\text{Risk}}_{\pi,F}(\phi) \right] \geq \lambda \|\phi\|_K^2 - \lambda J \right\} \\ & \subset \bigcup_{\phi \in A_{s,k}} \left\{ (\mathbf{X}_i, Y_i)_{i \in [n]} : \sup_{\phi \in A_{s,k}} \left[ \widehat{\text{Risk}}_{\pi,F}(\phi_\pi^{(n)}) - \widehat{\text{Risk}}_{\pi,F}(\phi) \right] \geq \lambda \|\phi\|_K^2 - \lambda J \right\}. \quad (24) \end{aligned}$$

In the last line of (24), we have partitioned the set  $\{\phi \in \Phi(r, s_1, s_2) : \mathbb{E}\Delta(\phi; (\mathbf{X}, Y)) \geq L_n\}$  into a union of  $A_{s,k}$ , with

$$A_{s,k} = \{\phi \in \Phi(r, s_1, s_2) : (s+1)L_n \leq \text{Risk}_{\pi,F}(\phi) - \text{Risk}_{\pi,F}(f_{\text{bayes},\pi}) < (s+2)L_n, (k-1)J \leq |\phi| < kJ\},$$

for  $s, k = 1, 2, \dots$

Let  $\Gamma$  denote the target probability for the first line in (24). To bound  $\Gamma$ , it suffices to bound the sum of probabilities over sets  $A_{s,k}$ . For each  $A_{s,k}$ , we consider the scaled empirical process,

$$\begin{aligned} v_n(\phi) &:= \left[ \widehat{\text{Risk}}_{\pi,F}(\phi_\pi^{(n)}) - \widehat{\text{Risk}}_{\pi,F}(\phi) \right] - \left[ \text{Risk}_{\pi,F}(\phi_\pi^{(n)}) - \text{Risk}_{\pi,F}(\phi) \right] \\ &= \frac{1}{n} \sum_i \left( \ell_{\pi,F}(\phi_\pi^{(n)}; (\mathbf{X}_i, Y_i)) - \ell_{\pi,F}(\phi; (\mathbf{X}_i, Y_i)) - \mathbb{E} \left[ \ell_{\pi,F}(\phi_\pi^{(n)}; (\mathbf{X}_i, Y_i)) - \ell_{\pi,F}(\phi; (\mathbf{X}_i, Y_i)) \right] \right). \end{aligned} \quad (25)$$

Notice that

$$\begin{aligned}
\text{Risk}_{\pi,F}(\phi) - \text{Risk}_{\pi,F}(\phi_\pi^{(n)}) &= \text{Risk}_{\pi,F}(\phi) - \text{Risk}_{\pi,F}(f_{\text{bayes},\pi}) + \text{Risk}_{\pi,F}(f_{\text{bayes},\pi}) - \text{Risk}_{\pi,F}(\phi_\pi^{(n)}) \\
&\geq (s+1)L_n - a_n \\
&\geq sL_n,
\end{aligned} \tag{26}$$

where the first inequality is from the fact that  $\phi \in A_{s,k}$  and Assumption 1(a) and the last inequality uses the condition that  $a_n \lesssim L_n$ .

Combining the definition of  $v_n$  in (25) and (26) gives (24) as

$$\begin{aligned}
\Gamma &\leq \sum_{s,k=1}^{\infty} \mathbb{P} \left\{ \sup_{\phi \in A_{s,k}} v_n(\phi) \geq sL_n + \lambda \|\phi\|_K^2 - \lambda J \right\} \\
&\leq \sum_{s,k=1}^{\infty} \mathbb{P} \left\{ \sup_{\phi \in A_{s,k}} v_n(\phi) \geq sL_n + \lambda(k-2)J =: M(s,k) \right\},
\end{aligned} \tag{27}$$

where  $M(s,k) > 0$  for all  $s,k \in \mathbb{N}$  from the condition  $\lambda J \leq L_n/2$ .

The variance of the empirical process is bounded by,

$$\begin{aligned}
&\sup_{A_{s,k}} \text{Var} \left[ \ell_{\pi,F}(\phi_\pi^{(n)}; (\mathbf{X}, Y)) - \ell_{\pi,F}(\phi; (\mathbf{X}, Y)) \right] \\
&\leq \sup_{A_{s,k}} 2 \left\{ \text{Var} \left[ \ell_{\pi,F}(\phi_\pi^{(n)}; (\mathbf{X}, Y)) - \ell_{\pi,F}(f_{\text{bayes},\pi}; (\mathbf{X}, Y)) \right] + \text{Var} \left[ \ell_{\pi,F}(\phi; (\mathbf{X}, Y)) - \ell_{\pi,F}(f_{\text{bayes},\pi}; (\mathbf{X}, Y)) \right] \right\} \\
&\lesssim [M(s,k)]^{\alpha/(1+\alpha)} + \frac{M(s,k)}{\rho} =: V(s,k),
\end{aligned} \tag{28}$$

where the last inequality is from Lemma 1.

We next bound the right-hand-side of (27) choosing an  $L_n$  that satisfies the conditions in Shen and Wong [1994, Theorem 3]. (The specification of  $L_n$  is deferred to the next paragraph). Once such  $L_n$  is chosen, then it follows from Shen and Wong [1994, Theorem 3] that

$$\begin{aligned}
\Gamma &\lesssim \sum_{s,k} \exp \left( -\frac{nM^2(s,k)}{V(s,k) + M(s,k)} \right) \\
&\lesssim \sum_{s,k} \exp(-\rho n M(s,k)) \\
&= \sum_{s,k} \exp(-n\rho s L_n - n\rho J \lambda(k-2)) \\
&\leq \left( \frac{e^{-n\rho L_n}}{1 - e^{-n\rho L_n}} \right) \left( \frac{e^{n\rho \lambda J}}{1 - e^{n\rho \lambda J}} \right) \leq e^{-n\rho L_n/2}.
\end{aligned} \tag{29}$$

where the last inequality uses the condition  $\lambda J \leq L_n/2$ . Notice that this bound for  $\Gamma$  has a sub-exponential decay in  $L_n$ .

Now, we specify an  $L_n$  that satisfies the condition of [Shen and Wong \[1994, Theorem 3\]](#). The convergence rate  $L_n > 0$  is determined by the solution to the following inequality,

$$\sup_{k \geq 2} \frac{1}{x} \int_x^{\sqrt{x^{\alpha/(\alpha+1)} + \rho^{-1}x}} \sqrt{\mathcal{H}_{[\cdot]}(\varepsilon, \Phi^k, \|\cdot\|_2)} d\varepsilon \lesssim n^{1/2}, \quad \text{where } x = L_n + \lambda J(k-2). \quad (30)$$

In particular, the smallest  $L_n$  satisfying (30) yields the best upper bound of the error rate. Here  $\mathcal{H}_{[\cdot]}(\varepsilon, \Phi^k, \|\cdot\|_2)$  denotes the  $L_2$ -norm,  $\varepsilon$ -bracketing number (c.f. Definition 1) for function family  $\Phi^k$ , and, we have denoted  $\Phi^k = \{\phi \in \Phi(r, s_1, s_2) : \|\phi\|_K^2 \leq k\}$ , i.e., the subset of functions in  $\Phi(r, s_1, s_2)$  with magnitudes bounded by  $k$ , for  $k \in \mathbb{N}_+$ .

It remains to solve for the smallest possible  $L_n$  in (30). Based on Lemma 6, the inequality (30) is satisfied with

$$L_n \asymp t_n^{(\alpha+1)/(\alpha+2)} + \frac{1}{\rho} t_n, \quad \text{and} \quad \lambda \asymp \frac{L_n}{J}, \quad (31)$$

where

$$t_n = \begin{cases} \frac{r d_{\max}}{n}, & \text{low-rank model } \phi \in \Phi(r), \\ \frac{r(s_1+s_2) \log d_{\max}}{n}, & \text{low-rank and two-way sparse model } \phi \in \Phi(r, s_1, s_2). \end{cases}$$

Plugging (31) into (29) gives that,

$$\begin{aligned} \Gamma &= \mathbb{P} \left[ \text{Risk}_{\pi, F}(\hat{\phi}_{\pi, F}) - \text{Risk}_{\pi, F}(f_{\text{bayes}, \pi}) \geq L_n \right] \\ &\lesssim \exp(-n\rho L_n) \lesssim \exp(-nt_n). \end{aligned}$$

where the last inequality uses the fact that  $\rho L_n \geq t_n$ . The proof is then complete by bounding the 0-1 risk by  $F$ -risk.

For unbounded loss (hinge loss), we seek to bound the  $F'$ -risk of  $\hat{\phi}_{\pi, F}$ , where  $F'$  is  $T$ -truncated version of  $F$ . The general strategy is to evaluate  $\hat{\phi}_{\pi, F}$ 's error using  $F'$ -risk. Note that the estimate  $\hat{\phi}_{\pi, F}$  (23) is defined under unbounded loss  $F$ . Therefore, the inclusion (24) changes to

$$\begin{aligned} &\left\{ (\mathbf{X}_i, Y_i)_{i \in [n]} : \text{Risk}_{\pi, F}(\hat{\phi}_{\pi, F}) - \text{Risk}_{\pi, F}(f_{\text{bayes}, \pi}) \geq 2L_n \right\} \\ &\subset \left\{ (\mathbf{X}_i, Y_i)_{i \in [n]} : \exists \phi \in \Phi(r, s_1, s_2), \text{ s.t. } \text{Risk}_{\pi, F}(\phi; (\mathbf{X}, Y)) - \text{Risk}_{\pi, F}(f_{\text{bayes}, \pi}) \geq 2L_n, \right. \\ &\quad \left. \text{and } \widehat{\text{Risk}}_{\pi, F}(\hat{\phi}_{\pi, F}) + \lambda \|\hat{\phi}_{\pi, F}\|_K^2 \leq \widehat{\text{Risk}}_{\pi, F}(\phi_{\pi}^{(n)}) + \lambda J \right\} \\ &\subset \left\{ (\mathbf{X}_i, Y_i)_{i \in [n]} : \exists \phi \in \Phi(r, s_1, s_2), \text{ s.t. } \text{Risk}_{\pi, F}(\phi; (\mathbf{X}, Y)) - \text{Risk}_{\pi, F}(f_{\text{bayes}, \pi}) \geq 2L_n, \right. \\ &\quad \left. \text{and } \widehat{\text{Risk}}_{\pi, F'}(\hat{\phi}_{\pi, F}) + \lambda \|\hat{\phi}_{\pi, F}\|_K^2 \leq \widehat{\text{Risk}}_{\pi, F'}(\phi_{\pi}^{(n)}) + \lambda J \right\}. \end{aligned}$$

where the second line is newly added, and the last line comes from the fact that

$$\widehat{\text{Risk}}_{\pi, F'}(\phi) \leq \widehat{\text{Risk}}_{\pi, F}(\phi) \text{ for all } \phi \in \Phi(r, s_1, s_2) \quad \text{and} \quad \widehat{\text{Risk}}_{\pi, F'}(\phi_\pi^{(n)}) = \widehat{\text{Risk}}_{\pi, F}(\phi_\pi^{(n)}),$$

setting  $T = \max(2, J) > \max(2, \sup_n \|\phi_\pi^{(n)}\|_K)$ , which is a constant. The remaining proof follows a similar line of argument as the bounded case. In particular, we invoke the mean-to-variance relationship for bounded  $F'$ -loss in (28). The final conclusion follows from the excess bound inequality for  $T$ -truncated risk (c.f. Lemma 1).  $\square$

*Proof of Theorem B.4.1.* Write  $\rho = \rho(\pi, \mathcal{N})$ . Combining Theorem 2 and Lemma 2 gives

$$\begin{aligned} \|\text{sgn}\phi - f_{\text{bayes}, \pi}\|_1 &\lesssim \left[ \text{Risk}_\pi(\hat{\phi}_{\pi, F}) - \text{Risk}_\pi(f_{\text{bayes}, \pi}) \right]^{\alpha/(\alpha+1)} + \frac{1}{\rho} \left[ \text{Risk}_\pi(\hat{\phi}_{F, \pi}) - \text{Risk}_\pi(f_{\text{bayes}, \pi}) \right] \\ &\lesssim t_n^{\alpha/(\alpha+2)} + \frac{1}{\rho^{\alpha/\alpha+1}} t_n^{\alpha/(\alpha+1)} + \frac{1}{\rho} t_n^{(\alpha+1)/(\alpha+2)} + \frac{1}{\rho^2} t_n \\ &\leq 4t_n^{\alpha/(\alpha+2)} + \frac{4}{\rho^2} t_n, \end{aligned}$$

where the last line follows from the fact that  $a(b^2 + b^{(\alpha+2)/(\alpha+1)} + b + 1) \leq 4a(b^2 + 1)$  with  $a = \frac{t_n}{\rho^2}$  and  $b = \rho t_n^{-1/(\alpha+2)}$ . The proof is complete by specializing  $t_n$  in each context.  $\square$

## B.5 Proofs of Theorem 4, Part (b) in Theorems 5 and 7

*Proof of Theorem 4.* It follows from the definition of  $\hat{f}$  that

$$\begin{aligned} \|\hat{f} - f\|_1 &= \mathbb{E} \left| \frac{1}{2H+1} \sum_{\pi \in \mathcal{H}} \text{sgn}\hat{\phi}_\pi - f \right| \\ &\leq \mathbb{E} \left| \frac{1}{2H+1} \sum_{\pi \in \mathcal{H}} (\text{sgn}\hat{\phi}_\pi - \text{sgn}(f - \pi)) \right| + \mathbb{E} \left| \frac{1}{2H+1} \sum_{\pi \in \mathcal{H}} \text{sgn}(f - \pi) - f \right| \\ &\leq \frac{1}{2H+1} \sum_{\pi \in \mathcal{H}} \|\text{sgn}\hat{\phi}_\pi - \text{sgn}(f - \pi)\|_1 + \frac{1}{H}, \end{aligned} \tag{32}$$

where the last line comes from the triangle inequality and the inequality

$$\left| \frac{1}{2H+1} \sum_{\pi \in \mathcal{H}} \text{sgn}(f(\mathbf{X}) - \pi) - f(\mathbf{X}) \right| \leq \frac{1}{H}, \quad \text{for all } \mathbf{X} \in \mathcal{X}.$$

It suffices to bound the first term in (32). We shall prove that, with probability at least  $1 - \exp(-nt \log H)$ ,

$$\frac{1}{2H+1} \sum_{\pi \in \mathcal{H}} \|\text{sgn}\hat{\phi}_\pi - \text{sgn}(f - \pi)\|_1 \leq t^{\alpha/(\alpha+2)} + \frac{1}{H} + Ht, \tag{33}$$

for any  $t \geq t_n$ , where  $t_n$  is specified in Theorem B.4.1.

For any  $t > t_n$ , define the union of events

$$E = \left\{ \left\| \text{sgn} \hat{\phi}_\pi - \text{sgn}(f - \pi) \right\|_1 \leq t^{\alpha/(2+\alpha)} + \frac{t}{\rho^2(\pi, \mathcal{N})} \text{ for all } \pi \in \mathcal{H} \right\}.$$

We first bound the probability of the event  $E$ , and then show  $E$  implies (33). Based on Theorem B.4.1 and union bound over  $\pi \in \mathcal{H}$ , we have that, for any  $t \geq t_n$ ,

$$\begin{aligned} \mathbb{P}(E) &\geq 1 - \sum_{\pi \in \mathcal{H}} \mathbb{P} \left( \left\| \text{sgn} \hat{\phi}_\pi - \text{sgn}(f - \pi) \right\|_1 \geq t^{\alpha/(2+\alpha)} + \frac{t}{\rho^2(\pi, \mathcal{N})} \text{ for a given } \pi \right) \\ &\gtrsim 1 - (2H + 1) \exp(-nt) \gtrsim 1 - \exp(-nt \log H). \end{aligned}$$

One the other hand, we now show that the  $E$  implies (33). Theorem B.4.1 shows that the sign function accuracy depends on the closeness of  $\pi \in \mathcal{H}$  to the mass points in  $\mathcal{H}$ . Therefore, we partition the level set  $\pi \in \mathcal{H}$  based on their closeness to  $\mathcal{H}$ . Specifically, let  $\mathcal{N}_H \stackrel{\text{def}}{=} \bigcup_{\pi' \in \mathcal{N}} (\pi' - \frac{1}{H}, \pi' + \frac{1}{H})$  denote the set of levels at least  $\frac{1}{H}$ -close to the mass points. We expand left hand side of (33) by

$$\begin{aligned} &\frac{1}{2H + 1} \sum_{\pi \in \mathcal{H}} \left\| \text{sgn} \hat{\phi}_\pi - \text{sgn}(f - \pi) \right\|_1 \\ &= \frac{1}{2H + 1} \sum_{\pi \in \mathcal{H} \cap \mathcal{N}_H} \left\| \text{sgn} \hat{\phi}_\pi - \text{sgn}(f - \pi) \right\|_1 + \frac{1}{2H + 1} \sum_{\pi \in \mathcal{H} \cap \mathcal{N}_H^c} \left\| \text{sgn} \hat{\phi}_\pi - \text{sgn}(f - \pi) \right\|_1. \end{aligned} \quad (34)$$

By assumption, the first term of (34) involves only finite number of summands and thus can be bounded by  $4C/(2H + 1)$  where  $C > 0$  is a constant such that  $|\mathcal{N}| \leq C$ . We bound the second term using the explicit forms of  $\rho(\pi, \mathcal{N})$  in the sequence  $\pi \in \Pi \cap \mathcal{N}_H^c$ .

$$\begin{aligned} \frac{1}{2H + 1} \sum_{\pi \in \mathcal{H} \cap \mathcal{N}_H^c} \left\| \text{sgn} \hat{\phi}_\pi - \text{sgn}(f - \pi) \right\|_1 &\lesssim \frac{1}{2H + 1} \sum_{\pi \in \mathcal{H} \cap \mathcal{N}_H^c} t^{\alpha/(2+\alpha)} + \frac{t}{2H + 1} \sum_{\pi \in \mathcal{H} \cap \mathcal{N}_H^c} \frac{1}{\rho^2(\pi, \mathcal{N})} \\ &\leq t^{\alpha/(2+\alpha)} + \frac{t}{2H + 1} \sum_{\pi \in \mathcal{H} \cap \mathcal{N}_H^c} \sum_{\pi' \in \mathcal{N}} \frac{1}{|\pi - \pi'|^2} \\ &\leq t^{\alpha/(2+\alpha)} + \frac{t}{2H + 1} \sum_{\pi' \in \mathcal{N}} \sum_{\pi \in \mathcal{H} \cap \mathcal{N}_H^c} \frac{1}{|\pi - \pi'|^2} \\ &\leq t^{\alpha/(2+\alpha)} + 2CHt, \end{aligned}$$

where the first inequality uses the property of  $E$ , and the last inequality follows from Lemma 3. Combining the bounds for the last two terms in (34) completes the proof for inequality (33). Finally, plugging (33) into (32) yields

$$\mathbb{P} \left( \left\| \hat{f} - f \right\|_1 \lesssim t^{\alpha/(\alpha+2)} + \frac{1}{H} + tH \right) \geq 1 - \exp(-nt \log H),$$

for any  $t \geq t_n$ . In particular, taking  $t = t_n \log H$  gives

$$\left\| \hat{f} - f \right\|_1 \lesssim t_n^{\alpha/(\alpha+2)} \log H + \frac{1}{H} + t_n H \log H$$

with high probability at least  $1 - \exp(-nt)$ . Setting  $H \asymp t_n^{-1/2}$  yields the desired upper bound.  $\square$

**Lemma 3.** Fix  $\pi' \in \mathcal{N}$  and a sequence  $\mathcal{H} = \{-1, \dots, -1/H, 0, 1/H, \dots, 1\}$  with  $H \geq 2$ . Then,

$$\sum_{\pi \in \mathcal{H} \cap \mathcal{N}_H^c} \frac{1}{|\pi - \pi'|^2} \leq 4H^2.$$

*Proof of Lemma 3.* Notice that all points  $\pi \in \mathcal{H} \cap \mathcal{N}_H^c$  satisfy  $|\pi - \pi'| > \frac{1}{H}$  for all  $\pi' \in \mathcal{N}$ . We use this fact to compute the sum

$$\begin{aligned} \sum_{\pi \in \mathcal{H} \cap \mathcal{N}_H^c} \frac{1}{|\pi - \pi'|^2} &= \sum_{\frac{h}{H} \in \mathcal{H} \cap \mathcal{N}_H^c} \frac{1}{|\frac{h}{H} - \pi'|^2} \\ &\leq 2H^2 \sum_{h=1}^H \frac{1}{h^2} \\ &\leq 2H^2 \left\{ 1 + \int_1^2 \frac{1}{x^2} dx + \int_2^3 \frac{1}{x^2} dx + \dots + \int_{H-1}^H \frac{1}{x^2} dx \right\} \\ &= 2H^2 \left( 1 + \int_1^H \frac{1}{x^2} dx \right) \leq 4H^2, \end{aligned}$$

where the third line uses the monotonicity of  $\frac{1}{x^2}$  for  $x \geq 1$ .  $\square$

## B.6 Proofs of Theorem 6 and Theorem A.3.1

*Proof of Theorem 6.* Theorem 6 follows from the same line of proof as in Theorem 4, with slight modification to account for discrete measure space. For any matrix  $\mathbf{Z} \in \mathbb{R}^{d \times d}$ , we use  $f_{\mathbf{Z}}: [d]^2 \rightarrow \mathbb{R}$  to denote the function induced by matrix  $\mathbf{Z}$  such that  $f_{\mathbf{Z}}(\omega) = \mathbf{Z}(\omega)$  for  $\omega \in [d]^2$ . Set  $\mathcal{X} = \{\mathbf{e}_i \otimes \mathbf{e}_j: (i, j) \in [d]^2\}$  be the discrete feature space, and  $n = |\Omega|$  the sample size. Under this set up,  $\|\hat{f} - f\|_1 = \mathbb{E}_{\mathbf{X}} |\hat{f}(\mathbf{X}) - f(\mathbf{X})| = \mathbb{E}_{\omega} |\hat{\Theta}(\omega) - \Theta(\omega)| = \text{MAE}(\hat{\Theta} - \Theta)$ . Notice that the small tolerance  $\Delta s = 1/d^2$  in the pseudo density is dominated by the derived convergence rate. Applying Theorem 4 to this setting finishes the proof.  $\square$

*Proof of Theorem A.3.1.* By setting  $s = \log(d_{\max})$  in Lemma 7, we have

$$\mathbb{P}(\|\mathbf{E}\|_{\infty} \geq \sqrt{4\sigma^2 \log d_{\max}}) \leq 2d_{\max}^{-2}.$$

We divide the sample space into two exclusive events:

- Event I:  $\|\mathbf{E}\|_{\infty} \geq \sqrt{4\sigma^2 \log d_{\max}}$ ;
- Event II:  $\|\mathbf{E}\|_{\infty} < \sqrt{4\sigma^2 \log d_{\max}}$ .

Because the Event I occurs with probability tending to zero, we restrict ourselves to the Event II only by following the proof of Theorem 3. We summarize the key difference compared to Section 4.

Let  $\ell_\omega(\cdot)$  denote the 0-1 loss evaluated at the  $\omega$ -th value of matrix. We expand the variance by

$$\begin{aligned}
\text{Var} [\ell_\omega(\mathbf{Z}, \bar{\mathbf{Y}}_\Omega) - \ell_\omega(\bar{\boldsymbol{\Theta}}, \bar{\mathbf{Y}}_\Omega)] &\leq \mathbb{E} |\ell_\omega(\mathbf{Z}(\omega), \bar{\mathbf{Y}}(\omega)) - \ell_\omega(\bar{\boldsymbol{\Theta}}(\omega), \bar{\mathbf{Y}}(\omega))|^2 \\
&= \mathbb{E} |\bar{\mathbf{Y}}(\omega) - \bar{\boldsymbol{\Theta}}(\omega) + \bar{\boldsymbol{\Theta}}(\omega)|^2 |\text{sgn} \mathbf{Z}(\omega) - \text{sgn} \bar{\boldsymbol{\Theta}}(\omega)| \\
&\leq 2(4\sigma^2 \log d_{\max} + 2) \mathbb{E} |\text{sgn} \mathbf{Z} - \text{sgn} \bar{\boldsymbol{\Theta}}| \\
&\lesssim (\sigma^2 \log d_{\max}) \text{MAE}(\text{sgn} \mathbf{Z}, \text{sgn} \bar{\boldsymbol{\Theta}}),
\end{aligned} \tag{35}$$

where the third line uses the facts  $\|\bar{\boldsymbol{\Theta}}\|_\infty \leq 2$  and  $\|\bar{\mathbf{Y}} - \bar{\boldsymbol{\Theta}}\|_\infty^2 = \|\mathbf{E}\|_\infty^2 < 4\sigma^2 \log d_{\max}$  within the Event II; the last line comes from the definition of MAE and the asymptotic  $\sigma^2 \log d_{\max} \gg 1$  provided that  $\sigma > 0$  with  $d_{\max}$  sufficiently large.

Based on (35), the  $(\alpha, \pi)$ -smoothness of  $\boldsymbol{\Theta}$  implies that for all measurable functions  $f_{\mathbf{Z}}$ , we have

$$\text{Var} \Delta_i(f_{\mathbf{Z}}, \bar{\mathbf{Y}}) \lesssim (\sigma^2 \log d_{\max}) \left\{ [\mathbb{E} \Delta_i(f_{\mathbf{Z}}, \bar{\mathbf{Y}})]^{\frac{\alpha}{1+\alpha}} + \frac{1}{\rho} \mathbb{E} \Delta_i(f_{\mathbf{Z}}, \bar{\mathbf{Y}}) \right\}. \tag{36}$$

The empirical process with variance-to-mean relationship (36) gives that

$$\mathbb{P} \left( \text{Risk}(\hat{\mathbf{Z}}) - \text{Risk}(\bar{\boldsymbol{\Theta}}) \geq L_n \right) \lesssim \exp(-nt_n), \tag{37}$$

where the convergence rate  $L_n$  is obtained by the same way in the proof of Lemma 6,

$$L_n \asymp t_n^{(\alpha+1)/(\alpha+2)} + \frac{1}{\rho} t_n, \quad \text{with } t_n = \frac{r\sigma^2 d_{\max} \log d_{\max}}{n}. \tag{38}$$

Combining (37) and (38), we obtain that, with high probability,

$$\text{Risk}(\hat{\mathbf{Z}}) - \text{Risk}(\bar{\boldsymbol{\Theta}}) \lesssim \left( \frac{r\sigma^2 d_{\max} \log d_{\max}}{|\Omega|} \right)^{(\alpha+1)/(\alpha+2)} + \frac{1}{\rho(\pi, \mathcal{N})} \left( \frac{r\sigma^2 d_{\max} \log d_{\max}}{|\Omega|} \right), \tag{39}$$

where constants have been absorbed into the  $\lesssim$  relationship. Therefore, combining (39) and (12) completes the proof for sign matrix estimation error in (7). The signal estimation error follows the same proof of Theorem 4.  $\square$



## C Auxiliary lemmas

**Lemma 4** (Hinge loss and  $L$ -1 distance). *Consider the same set-up as in Theorem 7. Let  $F(z) = (1 - z)_+$  be the hinge loss. Then, the  $L$ -1 distance between  $\phi$  and  $f_{\text{bayes},\pi}$  is bounded by their excess risk; i.e.,*

$$\|\phi - f_{\text{bayes},\pi}\|_1 \lesssim [\text{Risk}_{\pi,F}(\phi) - \text{Risk}_{\pi,F}(f_{\text{bayes},\pi})]^{\frac{\alpha}{1+\alpha}} + \frac{1}{\rho(\pi, \mathcal{N})} [\text{Risk}_{\pi,F}(\phi) - \text{Risk}_{\pi,F}(f_{\text{bayes},\pi})],$$

for all functions  $\phi: \mathcal{X} \rightarrow \mathbb{R}$ .

**Remark 6.** With little modification in the proof, similar inequality also holds for  $T$ -truncated hinge loss  $F(z) = \min(T, (1 - z)_+)$  with  $T \geq 2$ . Specifically,

$$\|\phi^T - f_{\text{bayes},\pi}\|_1 \leq [\text{Risk}_{\pi,F}(\phi) - \text{Risk}_{\pi,F}(f_{\text{bayes},\pi})]^{\frac{\alpha}{1+\alpha}} + \frac{1}{\rho(\pi, \mathcal{N})} [\text{Risk}_{\pi,F}(\phi) - \text{Risk}_{\pi,F}(f_{\text{bayes},\pi})],$$

where  $\phi^T: \mathcal{X} \rightarrow [-(T-1), (T-1)]$  is the truncated  $\phi$  defined in (22).

*Proof of Lemma 4.* For ease of notation, we drop the random variable  $\mathbf{X}$  in the function expression, and simply use  $\phi, f_{\text{bayes},\pi}, f$ , to represent the trace function, Bayes rule, and the regression function, respectively. The meaning should be clear given the contexts.

We expand the excess risk using the definition of hinge loss,

$$\begin{aligned} & \text{Risk}_{\pi,F}(\phi) - \text{Risk}_{\pi,F}(f_{\text{bayes},\pi}) \\ &= \mathbb{E}[|\bar{Y}_\pi|(1 - \phi \text{sgn} \bar{Y}_\pi)_+] - \mathbb{E}[|\bar{Y}_\pi|(1 - f_{\text{bayes},\pi} \text{sgn} \bar{Y}_\pi)_+] \\ &= \int_{\mathbf{X}} (1 - \phi)_+ \int_{y>\pi} (y - \pi) dy d\mathbb{P}_{\mathbf{X}} + \int_{\mathbf{X}} (1 + \phi)_+ \int_{y\leq\pi} (\pi - y) dy d\mathbb{P}_{\mathbf{X}} \\ &\quad - \int_{\mathbf{X}} (1 - f_{\text{bayes},\pi})_+ \int_{y>\pi} (y - \pi) dy d\mathbb{P}_{\mathbf{X}} - \int_{\mathbf{X}} (1 + f_{\text{bayes},\pi})_+ \int_{y\leq\pi} (\pi - y) dy d\mathbb{P}_{\mathbf{X}}. \end{aligned} \quad (40)$$

In order to evaluate the integral, we divide the domain  $\mathbf{X}$  into four exclusive regions:

- Region I =  $\{\mathbf{X}: f < \pi \text{ and } \phi \geq -1\}$ . In this region,  $f_{\text{bayes},\pi} = -1$ , and the integrand in (40) reduces to

$$\begin{aligned} \Phi_{\text{I}} &:= [(1 - \phi)_+ - 2] \mathbb{E}_{Y|\mathbf{X}}(Y - \pi) \mathbb{1}(Y > \pi) + (\phi + 1)_+ \mathbb{E}_{Y|\mathbf{X}}(\pi - Y) \mathbb{1}(Y \leq \pi) \\ &\geq -(\phi + 1) \mathbb{E}_{Y|\mathbf{X}}(Y - \pi) \mathbb{1}(Y > \pi) - (\phi + 1) \mathbb{E}_{Y|\mathbf{X}}(Y - \pi) \mathbb{1}(Y \leq \pi) \\ &= (\phi + 1)(\pi - f) = |\phi - f_{\text{bayes},\pi}| |f - \pi|. \end{aligned}$$

- Region II =  $\{\mathbf{X}: f < \pi \text{ and } \phi < -1\}$ . In this region,  $f_{\text{bayes},\pi} = -1$ , and the integrand in (40) reduces to

$$\Phi_{\text{II}} := -(\phi + 1) \mathbb{E}_{Y|\mathbf{X}}(Y - \pi) \mathbb{1}(Y > \pi) \geq -|\phi + 1|(f - \pi) = |\phi - f_{\text{bayes},\pi}| |f - \pi|.$$

- Region III =  $\{\mathbf{X}: f \geq \pi \text{ and } \phi \leq 1\}$ . In this region,  $f_{\text{bayes},\pi} = 1$ , and the integrand in (40) reduces to

$$\begin{aligned}\Phi_{\text{III}} &:= (1 - \phi)_+ \mathbb{E}_{Y|\mathbf{X}}(Y - \pi) \mathbb{1}(Y > \pi) + [(1 + \phi)_+ - 2] \mathbb{E}_{Y|\mathbf{X}}(\pi - Y) \mathbb{1}(Y \leq \pi) \\ &\geq (1 - \phi) \mathbb{E}_{Y|\mathbf{X}}(Y - \pi) \mathbb{1}(Y > \pi) + (\phi - 1) \mathbb{E}_{Y|\mathbf{X}}(\pi - Y) \mathbb{1}(Y \leq \pi) \\ &= (1 - \phi)(f - \pi) = |\phi - f_{\text{bayes},\pi}| |f - \pi|.\end{aligned}$$

- Region IV =  $\{\mathbf{X}: f \geq \pi \text{ and } \phi > 1\}$ . In this region,  $f_{\text{bayes},\pi} = 1$ , and the integrand in (40) reduces to

$$\Phi_{\text{IV}} := (\phi - 1) \mathbb{E}_{Y|\mathbf{X}}(\pi - Y) \mathbb{1}(Y \leq \pi) \geq (\phi - 1)(f - \pi) = |\phi - f_{\text{bayes},\pi}| |f - \pi|.$$

Therefore, the integral is evaluated as

$$\begin{aligned}\text{Risk}_{\pi,F}(\phi) - \text{Risk}_{\pi,F}(f_{\text{bayes},\pi}) &= \int_{\text{I}} \Phi_{\text{I}} d\mathbb{P}_{\mathbf{X}} + \int_{\text{II}} \Phi_{\text{II}} d\mathbb{P}_{\mathbf{X}} + \int_{\text{III}} \Phi_{\text{III}} d\mathbb{P}_{\mathbf{X}} + \int_{\text{IV}} \Phi_{\text{IV}} d\mathbb{P}_{\mathbf{X}} \\ &\geq \mathbb{E}|\phi - f_{\text{bayes},\pi}| |f - \pi|.\end{aligned}\tag{41}$$

Note that the function  $|f - \pi|$  is  $\alpha$ -smooth. Using the same techniques as in Theorem 2 to the last line of (41), we conclude

$$\mathbb{E}|\phi - f_{\text{bayes},\pi}| \lesssim [\text{Risk}_{\pi,F}(f) - \text{Risk}_{\pi,F}(f_{\text{bayes},\pi})]^{\frac{\alpha}{1+\alpha}} + \frac{1}{\rho(\pi, \mathcal{N})} [\text{Risk}_{\pi,F}(f) - \text{Risk}_{\pi,F}(f_{\text{bayes},\pi})].$$

□

**Definition 1** (Bracketing number). Consider a function set  $\Phi$ , and let  $\varepsilon > 0$ . We call  $\{(f_m^l, f_m^u)\}_{m=1}^M$  an  $L_2$ -metric,  $\varepsilon$ -bracketing function set of  $\Phi$ , if for every  $f \in \Phi$ , 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

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

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

**Lemma 5** (Bracketing number for bounded functions in  $\Phi(r, s_1, s_2)$ ). *Let  $\Phi(r, s_1, s_2)$  denote the collection of trace functions,*

$$\Phi(r, s_1, s_2) = \{\phi: \mathbf{X} \mapsto \langle \mathbf{X}, \mathbf{B} \rangle + b \mid \text{rank}(\mathbf{B}) \leq r, \text{supp}(\mathbf{B}) \leq (s_1, s_2), (\mathbf{B}, b) \in \mathbb{R}^{d_1 \times d_2} \times \mathbb{R}\},$$

*Assume, for simplicity, that the intercept  $b$  is known and that the function domain satisfies  $\mathbb{P}(\|\mathbf{X}\|_F \leq 1) = 1$ . For any given  $k \in \mathbb{N}_+$ , consider the subset of functions in  $\Phi(r, s_1, s_2)$  with magnitudes bounded*

by  $k$ , denoted by  $\Phi^k = \{f \in \Phi(r, s_1, s_2) : \|f\|_F^2 \leq k\}$ . Then, there exists a constant  $C > 0$  such that

$$\mathcal{H}_{[]}(\varepsilon, \Phi^k, \|\cdot\|_2) \leq Cr(s_1 + s_2) \log \frac{kd}{\varepsilon}.$$

*Proof.* For any given  $k \in \mathbb{N}_+$ , define the matrix class

$$\mathcal{B} = \{\mathbf{B} \in \mathbb{R}^{d \times d} : \text{rank}(\mathbf{B}) \leq r, \text{supp}(\mathbf{B}) \leq (s_1, s_2), \|\mathbf{B}\|_F^2 \leq k\}.$$

Based on the assumption of known  $b$ , there is an one-to-one correspondence between functions in  $\Phi^k$  and matrices in  $\mathcal{B}$ ,

$$\Phi^k = \{f : \mathbf{X} \mapsto \langle \mathbf{X}, \mathbf{B} \rangle + b \mid \mathbf{B} \in \mathcal{B}\}.$$

Furthermore, every pair of two functions  $f_1 = \langle \mathbf{X}, \mathbf{B}_1 \rangle$ ,  $f_2 = \langle \mathbf{X}, \mathbf{B}_2 \rangle \in \Phi(r, s_1, s_1)$  satisfies the norm relationship

$$\|f_1 - f_2\|_2 \leq \|f_1 - f_2\|_\infty = \sup_{\|\mathbf{X}\|_F \leq 1} |\langle \mathbf{X}, \mathbf{B}_1 \rangle - \langle \mathbf{X}, \mathbf{B}_2 \rangle| \leq \|\mathbf{B}_1 - \mathbf{B}_2\|_F.$$

Based on [Kosorok \[2007, Theorem 9.23\]](#), the  $L_2$ -metric,  $(2\varepsilon)$ -bracketing number in  $\Phi^k$  is bounded by

$$\mathcal{H}_{[]} (2\varepsilon, \Phi^k, \|\cdot\|_2) \leq \mathcal{H}(\varepsilon, \mathcal{B}, \|\cdot\|_F),$$

where  $\mathcal{H}$  denotes the log covering number for the (non-bracketing) set. Therefore, it suffices to bound  $\mathcal{H}(\varepsilon, \mathcal{B}, \|\cdot\|_F)$ . Now fix two subsets  $S_1, S_2 \subset [d]$  with  $|S_1| = s_1$  and  $|S_2| = s_2$ , where  $|\cdot|$  denotes the cardinality of the sets. Let  $\mathcal{B}_{S_1, S_2} \subset \mathcal{B}$  denote the subset of matrices satisfying  $\mathbf{B}(i, j) = 0$  whenever  $(i, j) \notin S_1 \times S_2$ . Based on [Candes and Plan \[2011, Lemma 3.1\]](#), the log covering number for  $\mathcal{B}_{S_1, S_2}$  is

$$\mathcal{H}(\varepsilon, \mathcal{B}_{S_1, S_2}, \|\cdot\|_F) \leq r(s_1 + s_2 + 1) \log \left( \frac{9\sqrt{k}}{\varepsilon} \right).$$

In view of the construction  $\mathcal{B} \subset \bigcup \{\mathcal{B}_{S_1, S_2} : S_1 \times S_2 \subset [d_1] \times [d_2], |S_1| = s_1, |S_2| = s_2\}$ , an  $\varepsilon$ -covering set  $\mathcal{B}$  is then given by the union of  $\varepsilon$ -covering set of  $\mathcal{B}_{S_1, S_2}$ . Using Stirling's bound, we derive that

$$\begin{aligned} \mathcal{H}(\varepsilon, \mathcal{B}, \|\cdot\|_F) &\leq \log \left\{ \binom{d}{s_1} \binom{d}{s_2} \exp [\mathcal{H}(\varepsilon, \mathcal{B}_{S_1, S_2}, \|\cdot\|_F)] \right\} \\ &\leq s_1 \log \frac{d}{s_1} + s_2 \log \frac{d}{s_2} + C' r(s_1 + s_2 + 1) \log \frac{k}{\varepsilon} \\ &\leq Cr(s_1 + s_2) \log \frac{kd}{\varepsilon}, \end{aligned}$$

where  $C, C' > 0$  are constants. □

**Lemma 6** (Local complexity of  $\Phi(r, s_1, s_2)$ ). *Define  $\Phi^k = \{f \in \Phi(r, s_1, s_2) : \|f\|_F^2 \leq k\}$  for all*

$k \in \mathbb{N}_+$ ; i.e.,  $\Phi^k$  is the subset of functions in  $\Phi(r, s_1, s_2)$  with magnitudes bounded by  $k$ . Set

$$L_n \asymp \left( \frac{r(s_1 + s_2) \log d}{n} \right)^{\frac{\alpha+1}{\alpha+2}} + \frac{1}{\rho(\pi, \mathcal{N})} \left( \frac{r(s_1 + s_2) \log d}{n} \right), \quad \text{and } \lambda \asymp \frac{L_n}{J}.$$

Then, the following inequality is satisfied for all  $k \in \{2, 3, \dots\}$ ,

$$\frac{1}{x} \int_x^{\sqrt{x^{\alpha/(\alpha+1)} + \frac{x}{\rho(\pi, \mathcal{N})}}} \sqrt{\mathcal{H}_{[\cdot]}(\varepsilon, \Phi^k, \|\cdot\|_2)} d\varepsilon \leq n^{1/2}, \quad \text{where } x := L_n + \lambda J(k/2 - 1).$$

*Proof.* To simplify the notation, we write  $\rho = \rho(\pi, \mathcal{N})$ , and define

$$g(x, k) = \frac{1}{x} \int_x^{\sqrt{x^{\alpha/(\alpha+1)} + \rho^{-1}x}} \sqrt{r(s_1 + s_2) \log \left( \frac{kd}{\varepsilon} \right)} d\varepsilon, \quad \text{for all } k \in \{2, 3, \dots\},$$

where we have inserted the bracketing number based on Lemma 5. Notice that

$$\begin{aligned} g(x, k) &\leq \frac{\sqrt{r(s_1 + s_2)}}{L} \int_x^{\sqrt{x^{\alpha/(\alpha+1)} + \rho^{-1}x}} \sqrt{\log \left( \frac{kd}{x} \right)} d\varepsilon \\ &\leq \sqrt{r(s_1 + s_2)(\log k + \log d - \log x)} \left( \frac{\sqrt{x^{\alpha/(2\alpha+2)} + \sqrt{\rho^{-1}x}}}{x} - 1 \right) \\ &\leq \sqrt{r(s_1 + s_2)(\log k + \log d)} \left( \frac{1}{x^{(\alpha+2)/(2\alpha+2)} + \frac{1}{\sqrt{\rho x}}} \right), \end{aligned} \tag{42}$$

where the second line follows from  $\sqrt{a+b} \leq \sqrt{a} + \sqrt{b}$  for  $a, b > 0$ . It remains to verify that  $g(L_n + \lambda J(k/2 - 1), k) \leq n^{1/2}$  for all  $k \in \{2, 3, \dots\}$ . Notice that

$$L_n + \lambda J(k/2 - 1) = (k/2 + C) \left\{ \left( \frac{r(s_1 + s_2) \log d}{n} \right)^{(\alpha+1)/(\alpha+2)} + \frac{1}{\rho} \left( \frac{r(s_1 + s_2) \log d}{n} \right) \right\},$$

for some universal constant  $C > 0$ . Plugging the above expression into the last line of (42) gives

$$g(L_n + \lambda J(k/2 - 1), k) \leq n^{1/2} \sqrt{\frac{\log k + \log d}{(k/2)^{(\alpha+2)/(\alpha+1)} \log d}} + n^{1/2} \sqrt{\frac{\log k + \log d}{(k/2) \log d}} \leq C' n^{1/2},$$

for all  $k \in \{2, 3, \dots\}$ , where  $C' > 0$  is a constant independent of  $k$  and  $d$ . The proof is therefore complete.  $\square$

**Theorem C.0.1** (Theorem 3 in Shen and Wong [1994]). *Let  $\mathcal{F}$  be a class of functions defined on  $\mathcal{X}$  with  $\sup_{f \in \mathcal{F}} \|f\|_\infty \leq T$ . Let  $(\mathbf{X}_i)_{i=1}^n$  be i.i.d. random variables with distribution  $\mathbb{P}_{\mathbf{X}}$  over  $\mathcal{X}$ . Set  $\sup_{f \in \mathcal{F}} \text{Var} f(\mathbf{X}) = V < \infty$ . Define the empirical process  $\hat{\mathbb{E}}f = \frac{1}{n} \sum_{i=1}^n f(\mathbf{X}_i)$ . 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

$$x_n^* \lesssim \frac{V}{T}, \quad \text{and} \quad \mathcal{H}_{[\cdot]}(\sqrt{V}, \mathcal{F}, \|\cdot\|_2) \lesssim \frac{n(x_n^*)^2}{V}.$$

Then we have

$$\mathbb{P} \left( \sup_{f \in \mathcal{F}} \hat{\mathbb{E}} f - \mathbb{E} f \geq x_n^* \right) \lesssim \exp \left( -\frac{n(x_n^*)^2}{V + T x_n^*} \right).$$

**Remark 7.** Both upper bounds  $T$  and  $V$  are allowed to depend on  $n$ . By applying Theorem C.0.1 to  $\mathcal{G} = \{g: g = \mathbb{E} f - f, f \in \mathcal{F}\}$ , we obtain the other side inequality,

$$\mathbb{P} \left\{ \sup_{f \in \mathcal{F}} (\mathbb{E} f - \hat{\mathbb{E}} f) \geq x_n^* \right\} \lesssim \exp \left( -\frac{n(x_n^*)^2}{V + T x_n^*} \right).$$

**Lemma 7** (sub-Gaussian maximum). *Let  $X_1, \dots, X_n$  be independent sub-Gaussian zero-mean random variables with variance proxy  $\sigma^2$ . Then, for any  $s > 0$ ,*

$$\mathbb{P} \left\{ \max_{1 \leq i \leq n} |X_i| \geq \sqrt{2\sigma^2(\log n + s)} \right\} \leq 2e^{-s}.$$

*Proof.* The conclusion follows from

$$\mathbb{P} \left[ \max_{1 \leq i \leq n} X_i \geq u \right] \leq \sum_{i=1}^n \mathbb{P}[X_i \geq u] \leq n e^{-\frac{u^2}{2\sigma^2}} = e^{-s},$$

where we set  $u = \sqrt{2\sigma^2(\log n + s)}$ . □

## References

- Peter L Bartlett, Michael I Jordan, and Jon D McAuliffe. Convexity, classification, and risk bounds. *Journal of the American Statistical Association*, 101(473):138–156, 2006.
- Emmanuel J Candes and Yaniv Plan. Tight oracle inequalities for low-rank matrix recovery from a minimal number of noisy random measurements. *IEEE Transactions on Information Theory*, 57(4):2342–2359, 2011.
- Michael R Kosorok. *Introduction to empirical processes and semiparametric inference*. Springer Science & Business Media, 2007.
- Clayton Scott. Surrogate losses and regret bounds for cost-sensitive classification with example-dependent costs. In *International Conference on Machine Learning*, 2011.
- Xiaotong Shen and Wing Hung Wong. Convergence rate of sieve estimates. *The Annals of Statistics*, pages 580–615, 1994.
- Junhui Wang, Xiaotong Shen, and Yufeng Liu. Probability estimation for large-margin classifiers. *Biometrika*, 95(1):149–167, 2008.