

A Proofs

In this section, we provide the proofs for our main results. First we characterize the implications of our general framework for the models in §2. We then establish the statistical convergence rates of the proposed procedure and the corresponding minimax lower bounds.

A.1 Proof of Lemma 3.2

Let \mathbf{X} and \mathbf{X}' be two independent random vectors following $\mathcal{N}(\mathbf{0}, \mathbf{I}_p)$. Let Y and Y' be two binary responses that depend on \mathbf{X}, \mathbf{X}' via (1.1). Then we have

$$\mathbb{E}(\mathbf{M}) = \mathbb{E}[(Y - Y')^2(\mathbf{X} - \mathbf{X}')(\mathbf{X} - \mathbf{X}')^\top].$$

Note that $(Y - Y')^2$ is a binary random variable taking values in $\{0, 4\}$. We have

$$\begin{aligned} \mathbb{E}[(Y - Y')^2 | \mathbf{X} = \mathbf{x}, \mathbf{X}' = \mathbf{x}'] &= 4 \cdot \mathbb{P}[(Y - Y')^2 = 4 | \mathbf{X} = \mathbf{x}, \mathbf{X}' = \mathbf{x}'] \\ &= 4 \cdot \mathbb{P}(Y = 1 | \mathbf{X} = \mathbf{x}) \cdot \mathbb{P}(Y' = -1 | \mathbf{X}' = \mathbf{x}') + 4 \cdot \mathbb{P}(Y' = 1 | \mathbf{X}' = \mathbf{x}') \cdot \mathbb{P}(Y = -1 | \mathbf{X} = \mathbf{x}) \\ &= 2 - 2f(\langle \mathbf{x}, \boldsymbol{\beta}^* \rangle) f(\langle \mathbf{x}', \boldsymbol{\beta}^* \rangle). \end{aligned} \quad (\text{A.1})$$

There exists some rotation matrix $\mathbf{Q} \in \mathbb{R}^{p \times p}$ such that $\mathbf{Q}\boldsymbol{\beta}^* = \mathbf{e}_1 := [1, 0, \dots, 0]^\top$. Let $\bar{\mathbf{X}} := \mathbf{Q}\mathbf{X}$ and $\bar{\mathbf{X}}' := \mathbf{Q}\mathbf{X}'$. Then we have

$$\mathbb{E}[(Y - Y')^2 | \mathbf{X} = \mathbf{x}, \mathbf{X}' = \mathbf{x}'] = \mathbb{E}[(Y - Y')^2 | \bar{\mathbf{X}} = \mathbf{Q}\mathbf{x}, \bar{\mathbf{X}}' = \mathbf{Q}\mathbf{x}'] = 2 - 2 \cdot f(\bar{x}_1) \cdot f(\bar{x}'_1),$$

where \bar{x}_1 and \bar{x}'_1 denote the first entries of $\bar{\mathbf{x}} := \mathbf{Q}\mathbf{x}$ and $\bar{\mathbf{x}}' := \mathbf{Q}\mathbf{x}'$ respectively. Note that $\bar{\mathbf{X}}$ and $\bar{\mathbf{X}}'$ also follow $\mathcal{N}(\mathbf{0}, \mathbf{I}_p)$ since symmetric Gaussian distribution is rotation invariant. Then we have

$$\begin{aligned} \mathbb{E}(\mathbf{M}) &= \mathbb{E} \left\{ [2 - 2f(\bar{\mathbf{X}}_1)f(\bar{\mathbf{X}}'_1)] (\mathbf{X} - \mathbf{X}')(\mathbf{X} - \mathbf{X}')^\top \right\} \\ &= \mathbf{Q}^\top \mathbb{E} \left\{ [2 - 2f(\bar{\mathbf{X}}_1)f(\bar{\mathbf{X}}'_1)] (\bar{\mathbf{X}} - \bar{\mathbf{X}}')(\bar{\mathbf{X}} - \bar{\mathbf{X}}')^\top \right\} \mathbf{Q} \\ &= 4\mathbf{Q}^\top [(\mu_1^2 - \mu_0\mu_2 + \mu_0^2) \cdot \mathbf{e}_1\mathbf{e}_1^\top + (1 - \mu_0^2) \cdot \mathbf{I}_p] \mathbf{Q} = 4\phi(f) \cdot \boldsymbol{\beta}^*\boldsymbol{\beta}^{*\top} + 4(1 - \mu_0^2) \cdot \mathbf{I}_p. \end{aligned}$$

The third equality is from the definitions of μ_0, μ_1, μ_2 in (3.2) and the last equality is from (3.1).

A.2 Proof of Lemma 3.4

Flipped logistic regression. For flipped logistic regression, the link function f is defined in (2.1), where ζ is the intercept. For $\zeta = 0$, we have

$$f(z) = \frac{e^z - 1}{e^z + 1} + 2p_e \cdot \frac{1 - e^z}{1 + e^z}.$$

Note that f is odd. Hence, by (3.2) we have $\mu_0 = \mu_2 = 0$. Meanwhile, from Stein's lemma, we have

$$\mu_1 = \mathbb{E}[f'(z)] = \mathbb{E} \left[(1 - 2p_e) \cdot \frac{2e^z}{(1 + e^z)^2} \right] = (1 - 2p_e) \cdot \mathbb{E} \left[\frac{2e^z}{(1 + e^z)^2} \right].$$

We thus have $\phi(f) = \mu_1^2 \geq C(1 - 2p_e)^2$ for some constant C .

Robust one-bit compressed sensing. Recall in robust one-bit compressed sensing, we have

$$f(z) = 2 \cdot \mathbb{P}(z + \epsilon > 0) - 1,$$

where $\epsilon \sim \mathcal{N}(0, \sigma^2)$ is the noise term in (2.2). In particular, note that

$$f(z) + f(-z) = 2 \cdot [\mathbb{P}(\epsilon > z) + \mathbb{P}(\epsilon > -z)] - 2 = 0.$$

Hence, $f(z)$ is an odd function, which implies $\mu_0 = \mu_2 = 0$ by (3.2). For μ_1 defined in (3.2), we have

$$\begin{aligned} \mu_1 &= \mathbb{E}[f(z)z] = \mathbb{E} \left\{ [2 \cdot \mathbb{P}(\epsilon > -z) - 1] z \right\} = \mathbb{E}[\mathbb{P}(|\epsilon| < |z|) |z|] \geq \mathbb{E} \left\{ \left[1 - 2e^{-z^2/(2\sigma^2)} \right] |z| \right\} \\ & \quad (\text{A.2}) \end{aligned}$$

$$= \mathbb{E}(|z|) - \int_{-\infty}^{\infty} \frac{2}{\sqrt{2\pi}} e^{-\frac{u^2}{2\sigma^2}} e^{-\frac{z^2}{2}} |u| du = \mathbb{E}(|z|) \left(1 - 2 \frac{\sigma^2}{1 + \sigma^2} \right) = \mathbb{E}(|z|) \frac{1 - \sigma^2}{1 + \sigma^2}.$$

Here the inequality is from the fact that $\mathbb{P}(|\epsilon| < |z|) \geq 1 - 2e^{-\frac{z^2}{2\sigma^2}}$ since $\epsilon \sim \mathcal{N}(0, \sigma^2)$. For $\sigma^2 < 1/2$, we have

$$\phi(f) = \mu_1^2 \geq C \left(\frac{1 - \sigma^2}{1 + \sigma^2} \right)^2,$$

where $C = \mathbb{E}(|z|)$ with $z \sim \mathcal{N}(0, 1)$. For $\sigma^2 \geq 1/2$, rather than applying $\mathbb{P}(|\epsilon| < |z|) \geq 1 - 2e^{-\frac{z^2}{2\sigma^2}}$ in the inequality of (A.2), we apply $\mathbb{P}(|\epsilon| < |z|) \geq \frac{2}{\sqrt{2\pi}\sigma} e^{-\frac{z^2}{2\sigma^2}} |z|$ since $\epsilon \sim \mathcal{N}(0, \sigma^2)$. We then obtain

$$\mu_1 \geq \mathbb{E} \left[\frac{2}{\sqrt{2\pi}\sigma} e^{-\frac{z^2}{2\sigma^2}} z^2 \right] = \frac{2}{\sqrt{2\pi}\sigma} \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}} e^{-\frac{u^2}{2\sigma^2}} e^{-\frac{u^2}{2}} u^2 du \geq \frac{C'}{\sigma} \left(\frac{\sigma^2}{1 + \sigma^2} \right)^{\frac{3}{2}}.$$

Finally, for $\sigma^2 \geq 1/2$ we have

$$\phi(f) \geq \frac{C'\sigma^4}{(1 + \sigma^2)^3}.$$

One-bit phase retrieval. For the one-bit phase retrieval model, the major difference from the previous two models is that $f(z)$ is even, which results in $\mu_1 = 0$. By the definition in (3.2), we have

$$\mu_0 = \mathbb{E}[f(z)] = \mathbb{P}(|z| \geq \theta) - \mathbb{P}(|z| < \theta),$$

and

$$\mu_2 = \mathbb{E}[f(z)z^2] = \mathbb{P}(|z| \geq \theta)\mathbb{E}(z^2 \mid |z| \geq \theta) - \mathbb{P}(|z| < \theta)\mathbb{E}(z^2 \mid |z| < \theta).$$

For notational simplicity, we define $p_1 = \mathbb{P}(|z| \geq \theta)$. We have

$$\phi(f) = \mu_0(\mu_0 - \mu_2) = 2p_1(2p_1 - 1)[1 - \mathbb{E}(z^2 \mid |z| > \theta)], \quad (\text{A.3})$$

where the second equality follows from the fact that

$$\mathbb{P}(|z| \geq \theta) + \mathbb{P}(|z| < \theta) = 1, \quad (\text{A.4})$$

and

$$\mathbb{P}(|z| \geq \theta)\mathbb{E}(z^2 \mid |z| \geq \theta) + \mathbb{P}(|z| < \theta)\mathbb{E}(z^2 \mid |z| < \theta) = \mathbb{E}(z^2) = 1. \quad (\text{A.5})$$

By (A.4) and (A.5) we have $p_1 > 0$ and $\mathbb{E}(z^2 \mid |z| \geq \theta) > 1$ for $\theta > 0$. Hence, for $\theta < \theta_m$ with θ_m being the median of $|z|$ with $z \sim \mathcal{N}(0, 1)$, we have $p_1 \geq 1/2$, which further implies $\phi(f) < 0$ by (A.3). Otherwise we have $\phi(f) > 0$. Thus, we have $\text{sign}[\phi(f)] = \text{sign}(\theta - \theta_m)$.

In the following we establish a lower bound for $|\phi(f)|$. Note that

$$\mathbb{E}(z^2 \mid |z| \geq \theta) = \frac{2}{p_1} \int_{\theta}^{+\infty} \frac{1}{\sqrt{2\pi}} e^{-\frac{z^2}{2}} z^2 dz = \frac{2\theta}{p_1\sqrt{2\pi}} e^{-\frac{\theta^2}{2}} + 1. \quad (\text{A.6})$$

Plugging (A.6) into (A.3) yields

$$\phi(f) = -2(2p_1 - 1) \frac{2\theta}{\sqrt{2\pi}} e^{-\frac{\theta^2}{2}}. \quad (\text{A.7})$$

For $0 < \theta < \theta_m$, which implies $p_1 \geq 1/2$, we have

$$p_1 - \frac{1}{2} = 2 \int_{\theta}^{\theta_m} \frac{1}{\sqrt{2\pi}} e^{-\frac{z^2}{2}} dz \geq \frac{2}{\sqrt{2\pi}} e^{-\frac{\theta_m^2}{2}} (\theta_m - \theta). \quad (\text{A.8})$$

By plugging (A.6) into (A.7), we have

$$|\phi(f)| \geq \frac{8}{\sqrt{2\pi}} e^{-\frac{\theta_m^2}{2}} (\theta_m - \theta) \frac{2\theta}{\sqrt{2\pi}} e^{-\frac{\theta^2}{2}} \geq C\theta(\theta_m - \theta)e^{-\frac{\theta^2}{2}}. \quad (\text{A.9})$$

For $\theta > \theta_m$, which implies $p_1 < 1/2$, similarly to (A.8), we have

$$\frac{1}{2} - p_1 = 2 \int_{\theta_m}^{\theta} \frac{1}{\sqrt{2\pi}} e^{-\frac{z^2}{2}} dz \geq \frac{2}{\sqrt{2\pi}} e^{-\frac{\theta^2}{2}} (\theta - \theta_m). \quad (\text{A.10})$$

Thus, we conclude that

$$|\phi(f)| \geq C'\theta(\theta - \theta_m)e^{-\theta^2}.$$

A.3 Proof of Theorem 3.5

Let $\widehat{\beta}$ be the top eigenvector of \mathbf{M} and $\widehat{\lambda}_1, \widehat{\lambda}_2$ be the first and second largest eigenvalues of \mathbf{M} . We use λ_1, λ_2 to denote the first and second largest eigenvalues of $\mathbb{E}(\mathbf{M})$. From Lemma 3.2, we already know that

$$\lambda_1 = 4\phi(f) + 4(1 - \mu_0^2), \quad \text{and} \quad \lambda_2 = 4(1 - \mu_0^2).$$

By the triangle inequality, we have

$$\|\beta^t - \beta^*\|_2 \leq \|\beta^* - \widehat{\beta}\|_2 + \|\beta^t - \widehat{\beta}\|_2.$$

The first term on the right hand side is the statistical error and the second term is the optimization error. From standard analysis of the power method, we have

$$\|\beta^t - \widehat{\beta}\|_2 \leq \sqrt{\frac{1 - \alpha^2}{\alpha^2}} \cdot (\widehat{\lambda}_2 / \widehat{\lambda}_1)^t,$$

where $\alpha = \langle \beta^0, \widehat{\beta} \rangle$. By the definition in (3.4), \mathbf{M} is the sample covariance matrix of $n/2$ independent realizations of the random vector $(Y - Y')(X - X') \in \mathbb{R}^p$. Since \mathbf{X} is Gaussian and Y is bounded, $(Y - Y')(X - X')$ is sub-Gaussian. By standard concentration results (see e.g. Theorem 5.39 in [26]), there some constants C, C_1 such that for any $t \geq 0$, with probability at least $1 - 2e^{-Ct^2}$,

$$\|\mathbf{M} - \mathbb{E}(\mathbf{M})\|_2 \leq \max(\delta, \delta^2) \|\mathbb{E}(\mathbf{M})\|_2,$$

where $\delta = C_1 \sqrt{\frac{p}{n}} + \frac{t}{\sqrt{n}}$. We let $t = \sqrt{p}$, then for any $\xi \in (0, 1)$, we have that $\|\mathbf{M} - \mathbb{E}(\mathbf{M})\|_2 \leq \xi \|\mathbb{E}(\mathbf{M})\|_2$ when $n \geq C_2 p / \xi^2$ for sufficiently large constant C_2 . Conditioning on $\|\mathbf{M} - \mathbb{E}(\mathbf{M})\|_2 \leq \xi \|\mathbb{E}(\mathbf{M})\|_2$, from Weyl's inequality, we have

$$\widehat{\lambda}_1 \geq 4(1 - \xi)[\phi(f) + 1 - \mu_0^2], \quad \text{and} \quad \widehat{\lambda}_2 \leq 4\xi\phi(f) + 4(1 + \xi)(1 - \mu_0^2).$$

Furthermore, for any $\gamma \in ((1 - \mu_0^2) / [\phi(f) + 1 - \mu_0^2], 1)$, by restricting

$$\xi \leq \frac{\gamma\phi(f) + (\gamma - 1)(1 - \mu_0^2)}{(1 + \gamma)[\phi(f) + 1 - \mu_0^2]}, \quad (\text{A.11})$$

we have

$$\|\beta^t - \widehat{\beta}\|_2 \leq \sqrt{\frac{1 - \alpha^2}{\alpha^2}} \cdot \gamma^t.$$

Now we turn to the statistical error. By Wedin's sin theorem, for some positive constant $C > 0$, we have

$$\sin \angle(\beta^*, \widehat{\beta}) \leq C \cdot \frac{\xi \|\mathbb{E}(\mathbf{M})\|_2}{\lambda_1 - \lambda_2}. \quad (\text{A.12})$$

Elementary calculation yields

$$\|\widehat{\beta} - \beta^*\|_2 = 2 \sin[\angle(\beta^*, \widehat{\beta})/2] \leq \sqrt{2} \sin \angle(\beta^*, \widehat{\beta}). \quad (\text{A.13})$$

As $\xi \lesssim \sqrt{p/n}$, combining (A.12) and (A.13), we have

$$\|\widehat{\beta} - \beta^*\|_2 \lesssim \frac{\phi(f) + 1 - \mu_0^2}{\phi(f)} \cdot \sqrt{\frac{p}{n}}.$$

Putting all pieces together, we conclude that if ξ satisfies (A.11) and $n \gtrsim p/\xi^2$, then we have that with probability at least $1 - 2e^{-Cp}$,

$$\|\beta^t - \beta^*\|_2 \leq C \cdot \frac{\phi(f) + 1 - \mu_0^2}{\phi(f)} \cdot \sqrt{\frac{p}{n}} + \sqrt{\frac{1 - \alpha^2}{\alpha^2}} \cdot \gamma^t.$$

as required.

A.4 Proof of Theorem 3.6

The analysis of Algorithm 2 follows from a combination of [27] (for the initialization via convex relaxation) and [31] (for the original truncated power method). Recall that κ is defined in (3.9).

Assume the initialization β^0 is \widehat{s} -sparse with $\|\beta^0\|_2 = 1$, and satisfies

$$\|\beta^0 - \beta^*\|_2 \leq C \min \left\{ \sqrt{\kappa(1 - \kappa^{1/2})/2}, \sqrt{2\kappa/4} \right\}, \quad (\text{A.14})$$

for $\widehat{s} = C' \max \{ \lceil 1/(\kappa^{-1/2} - 1)^2 \rceil, 1 \} \cdot s$. Theorem 1 of [31] implies that

$$\|\beta^t - \beta^*\|_2 \leq C'' \cdot \frac{[\phi(f) + (1 - \mu_0^2)]^{\frac{5}{2}} (1 - \mu_0^2)^{\frac{1}{2}}}{\phi^3} \cdot \sqrt{\frac{s \log p}{n}} + \kappa^t \cdot \sqrt{\min \{ (1 - \kappa^{1/2})/2, 1/8 \}}$$

with high probability. Therefore, we only need to prove the initialization β^0 obtained in Algorithm 2 satisfies the condition in (A.14).

Corollary 3.3 of [27] shows that the minimizer to the minimization problem in line 3 of Algorithm 2 satisfies

$$\|\Pi^0 - \beta^* \cdot (\beta^*)^\top\|_2 \leq C \cdot \frac{\phi(f) + (1 - \mu_0^2)}{\phi} \cdot s \sqrt{\frac{\log p}{n}}$$

with high probability. Corollary 3.2 of [27] implies, the first eigenvector of Π^0 , denoted as $\bar{\beta}^0$, satisfies

$$\|\bar{\beta}^0 - \beta^*\|_2 \leq C' \cdot \frac{\phi(f) + (1 - \mu_0^2)}{\phi} \cdot s \sqrt{\frac{\log p}{n}}$$

with the same probability. However, $\bar{\beta}^0$ is not necessarily \widehat{s} -sparse. Using Lemma 12 of [31], we obtain that the truncate step in lines 12-15 of Algorithm 2 ensures that β^0 is \widehat{s} -sparse and also satisfies

$$\|\beta^0 - \beta^*\|_2 \leq (1 + 2\sqrt{\widehat{s}/s}) \cdot \|\bar{\beta}^0 - \beta^*\|_2 \leq 3\|\bar{\beta}^0 - \beta^*\|_2,$$

where the last inequality follows from our assumption that $\widehat{s} \geq s$. Therefore, we only have to set n to be sufficiently large such that

$$\|\beta^0 - \beta^*\|_2 \leq C' \cdot \frac{\phi(f) + (1 - \mu_0^2)}{\phi(f)} \cdot s \sqrt{\frac{\log p}{n}} \leq C \min \left\{ \sqrt{\kappa(1 - \kappa^{1/2})/2}, \sqrt{2\kappa/4} \right\},$$

which is ensured by setting $n \geq n_{\min}$ with

$$n_{\min} = C' \cdot s^2 \log p \cdot \phi(f)^2 \cdot \min \{ \kappa(1 - \kappa^{1/2})/2, \kappa/8 \} / [(1 - \mu_0^2) + \phi(f)]^2,$$

as specified in our assumption. Thus we conclude the proof.

A.5 Proof of Theorem 3.7

The proof of the minimax lower bound follows from the basic idea of reducing an estimation problem to a testing problem, and then invoking Fano's inequality to lower bound the testing error. We first introduce a finite packing set for $\mathbb{S}^{p-1} \cap \mathbb{B}_0(s, p)$.

Lemma A.1. *Consider the set $\{0, 1\}^p$ equipped with Hamming distance δ . For $s \leq p/4$, there exists a finite subset $\mathcal{Q} \subset \{0, 1\}^p$ such that*

$$\delta(\theta, \theta') > s/2, \quad \forall (\theta, \theta') \in \mathcal{Q} \times \mathcal{Q} \text{ and } \theta \neq \theta', \quad \|\theta\|_0 = s, \text{ for all } \theta \in \mathcal{Q}.$$

The cardinality of such a set satisfies

$$\log(|\mathcal{Q}|) \geq 8/3 \cdot s \log(p/s).$$

Proof. See the proof of Lemma 4.10 in [18]. □

We use $\mathcal{Q}(p, s)$ to denote the finite set specified in Lemma A.1. For $\xi < 1$, we construct a finite subset $\bar{\mathcal{Q}}(p, s, \xi) \subset \mathbb{S}^{p-1} \cap \mathbb{B}_0(s, p)$ as

$$\bar{\mathcal{Q}}(p, s, \xi) := \left\{ \beta \in \mathbb{R}^p : \beta = \left(\sqrt{1 - \xi^2}, \frac{\xi}{\sqrt{s-1}} \cdot \mathbf{w} \right), \text{ where } \mathbf{w} \in \mathcal{Q}(p-1, s-1) \right\}. \quad (\text{A.15})$$

It is easy to verify that set $\bar{\mathcal{Q}}(p, s, \xi)$ has the following properties:

- For any $\theta \in \bar{\mathcal{Q}}(p, s, \xi)$, it holds that $\|\theta\|_2 = 1$ and $\|\theta\|_0 = s$.

- For distinct $\boldsymbol{\theta}, \boldsymbol{\theta}' \in \overline{\mathcal{Q}}(p, s, \xi)$, $\|\boldsymbol{\theta} - \boldsymbol{\theta}'\|_2 \geq \sqrt{2}\xi/2$ and $\|\boldsymbol{\theta} - \boldsymbol{\theta}'\|_2 \leq \sqrt{2}\xi$.
- $\log |\overline{\mathcal{Q}}(p, s, \xi)| \geq Cs \log(p/s)$ for some positive constant C .

In order to derive lower bound of $\mathcal{R}(n, m, L, \mathcal{B})$ with $\mathcal{B} = \mathbb{S}^{p-1} \cap \mathbf{B}_0(s, p)$, we assume that the infimum over f in (3.13) is obtained for a certain $f^* \in \mathcal{F}(m, L)$, namely

$$\mathcal{R}(n, m, L, \mathcal{B}) = \inf_{\widehat{\boldsymbol{\beta}} \in \mathbb{S}^{p-1}} \sup_{\boldsymbol{\beta} \in \mathbb{S}^{p-1} \cap \mathbf{B}_0(s, p)} \mathbb{E} \|\widehat{\boldsymbol{\beta}}(\mathcal{X}_{f^*}^n) - \boldsymbol{\beta}\|_2 \geq \inf_{\widehat{\boldsymbol{\beta}} \in \mathbb{S}^{p-1}} \sup_{\boldsymbol{\beta} \in \overline{\mathcal{Q}}(p, s, \xi)} \mathbb{E} \|\widehat{\boldsymbol{\beta}}(\mathcal{X}_{f^*}^n) - \boldsymbol{\beta}\|_2.$$

Note that for any $\xi > 0$, we have $\|\boldsymbol{\beta}_1 - \boldsymbol{\beta}_2\|_2 \geq \frac{\sqrt{2}}{2}\xi$ for any two distinct vectors $(\boldsymbol{\beta}_1, \boldsymbol{\beta}_2)$ in $\overline{\mathcal{Q}}(p, s, \xi)$. Therefore, we are in a position to apply standard minimax risk lower bound. Following Lemma 3 in Yu [30], we obtain

$$\inf_{\widehat{\boldsymbol{\beta}} \in \mathbb{S}^{p-1}} \sup_{\boldsymbol{\beta} \in \overline{\mathcal{Q}}(p, s, \xi)} \mathbb{E} \|\widehat{\boldsymbol{\beta}}(\mathcal{X}_{f^*}^n) - \boldsymbol{\beta}\|_2 \geq \frac{\sqrt{2}}{4}\xi \left(1 - \frac{\max_{\boldsymbol{\beta}, \boldsymbol{\beta}' \in \overline{\mathcal{Q}}(p, s, \xi)} D_{KL}(P_{\boldsymbol{\beta}'} \| P_{\boldsymbol{\beta}}) + \log 2}{\log |\overline{\mathcal{Q}}(p, s, \xi)|} \right). \quad (\text{A.16})$$

In the following, we derive an upper bound for the term involving KL divergence on the right hand side of the above inequality. For any $\boldsymbol{\beta}, \boldsymbol{\beta}' \in \overline{\mathcal{Q}}(p, s, \xi)$, we have

$$\begin{aligned} D_{KL}(P_{\boldsymbol{\beta}'} \| P_{\boldsymbol{\beta}}) &\leq n \cdot D_{KL}[P_{\boldsymbol{\beta}'}(Y, \mathbf{X}) \| P_{\boldsymbol{\beta}}(Y, \mathbf{X})] = n \cdot \mathbb{E}_{\mathbf{X}} \{ D_{KL}[P_{\boldsymbol{\beta}'}(Y | \mathbf{X}) \| P_{\boldsymbol{\beta}}(Y | \mathbf{X})] \} \\ &= \frac{1}{2}n \cdot \mathbb{E}_{\mathbf{X}} \left\{ [1 + f^*(\mathbf{X}^\top \boldsymbol{\beta})] \log \frac{1 + f^*(\mathbf{X}^\top \boldsymbol{\beta})}{1 + f^*(\mathbf{X}^\top \boldsymbol{\beta}')} + [1 - f^*(\mathbf{X}^\top \boldsymbol{\beta})] \log \frac{1 - f^*(\mathbf{X}^\top \boldsymbol{\beta})}{1 - f^*(\mathbf{X}^\top \boldsymbol{\beta}')} \right\} \\ &\leq \frac{1}{2}n \cdot \mathbb{E}_{\mathbf{X}} \left\{ [1 + f^*(\mathbf{X}^\top \boldsymbol{\beta})] \left[\frac{1 + f^*(\mathbf{X}^\top \boldsymbol{\beta})}{1 + f^*(\mathbf{X}^\top \boldsymbol{\beta}')} - 1 \right] + [1 - f^*(\mathbf{X}^\top \boldsymbol{\beta})] \left[\frac{1 - f^*(\mathbf{X}^\top \boldsymbol{\beta})}{1 - f^*(\mathbf{X}^\top \boldsymbol{\beta}')} - 1 \right] \right\}. \end{aligned} \quad (\text{A.17})$$

In the last inequality, we utilize the fact that $\log z \leq z - 1$. Then by elementary calculation, we have

$$D_{KL}(P_{\boldsymbol{\beta}'} \| P_{\boldsymbol{\beta}}) \leq n \cdot \mathbb{E}_{\mathbf{X}} \left\{ \frac{[f^*(\mathbf{X}^\top \boldsymbol{\beta}) - f^*(\mathbf{X}^\top \boldsymbol{\beta}')]^2}{[1 + f^*(\mathbf{X}^\top \boldsymbol{\beta}')] \cdot [1 - f^*(\mathbf{X}^\top \boldsymbol{\beta}')] } \right\}. \quad (\text{A.18})$$

Using $|f(z)| \leq 1 - m$ and the Lipschitz continuity condition of f , we have

$$D_{KL}(P_{\boldsymbol{\beta}'} \| P_{\boldsymbol{\beta}}) \leq n \cdot \mathbb{E}_{\mathbf{X}} \left\{ \frac{L^2 \langle \mathbf{X}, \boldsymbol{\beta} - \boldsymbol{\beta}' \rangle^2}{m(1-m)} \right\} = \frac{nL^2 \|\boldsymbol{\beta} - \boldsymbol{\beta}'\|_2^2}{m(1-m)} \leq \frac{2nL^2 \xi^2}{m(1-m)}. \quad (\text{A.19})$$

Note that (A.17)-(A.19) hold for any $\boldsymbol{\beta}, \boldsymbol{\beta}' \in \overline{\mathcal{Q}}(p, s, \xi)$. We thus have

$$\max_{\boldsymbol{\beta}, \boldsymbol{\beta}' \in \overline{\mathcal{Q}}(p, s, \xi)} D_{KL}(P_{\boldsymbol{\beta}'} \| P_{\boldsymbol{\beta}}) \leq \frac{2nL^2 \xi^2}{m(1-m)}.$$

Now we proceed with (A.16) using the above result. The right hand side is thus lower bounded by

$$\frac{\sqrt{2}}{4}\xi \left(1 - \frac{2L^2 n \xi^2 / [m(1-m)] + \log 2}{|\overline{\mathcal{Q}}(p, s, \xi)|} \right) \geq \frac{\sqrt{2}}{4}\xi \left(1 - \frac{2L^2 n \xi^2 / [m(1-m)] + \log 2}{Cs \log(p/s)} \right),$$

where the last inequality is from $|\overline{\mathcal{Q}}(p, s, \xi)| \geq Cs \log(p/s)$. Finally, consider the case where the sample size n is sufficiently large such that

$$n \geq \frac{m(1-m)}{2L^2} \cdot [Cs \log(p/s)/2 - \log 2],$$

by choosing

$$\xi^2 = \frac{m(1-m)}{2L^2 n} \cdot [Cs \log(p/s)/2 - \log 2], \quad (\text{A.20})$$

we thus have

$$\mathcal{R}(n, m, L, \mathcal{B}) \geq C' \cdot \frac{\sqrt{m(1-m)}}{L} \cdot \sqrt{\frac{s \log(p/s)}{n}}$$

as required.