

# Stats 300A HW3 Solutions

Song Mei

October 18, 2018

## Problem 1

(a)

Define a projector operator to be the following:

$$\text{Proj}(\boldsymbol{\theta}) = \begin{cases} \boldsymbol{\theta}, & \text{if } \boldsymbol{\theta} \in \Theta^\varepsilon, \\ \arg \min_{\boldsymbol{\theta}' \in \Theta} \|\boldsymbol{\theta} - \boldsymbol{\theta}'\|_2^2, & \text{if } \boldsymbol{\theta} \notin \Theta^\varepsilon. \end{cases}$$

Since  $\Theta$  is a convex compact set, the minimizer  $\arg \min_{\boldsymbol{\theta}' \in \Theta} \|\boldsymbol{\theta} - \boldsymbol{\theta}'\|_2^2$  is unique, so that Proj operator is well defined.

Given an estimator  $\hat{\boldsymbol{\theta}} : \mathcal{X} \rightarrow \mathbb{R}^d$  such that  $P_{\boldsymbol{\theta}}(\hat{\boldsymbol{\theta}}(\mathbf{X}) \notin \Theta^\varepsilon) > \delta$ , we take  $\tilde{\boldsymbol{\theta}} = \text{Proj}(\hat{\boldsymbol{\theta}})$ . Then for any  $\boldsymbol{\theta} \in \Theta$ , we have

$$\|\tilde{\boldsymbol{\theta}}(\mathbf{x}) - \boldsymbol{\theta}\|_2^2 \leq \|\hat{\boldsymbol{\theta}}(\mathbf{x}) - \boldsymbol{\theta}\|_2^2 - \eta \mathbf{1}\{\hat{\boldsymbol{\theta}}(\mathbf{x}) \notin \Theta^\varepsilon\},$$

where

$$\eta = \min_{\boldsymbol{\theta} \in \Theta, \boldsymbol{\theta}' \in \partial \Theta^\varepsilon} \|\boldsymbol{\theta}' - \boldsymbol{\theta}\|_2^2 - \|\text{Proj}(\boldsymbol{\theta}') - \boldsymbol{\theta}\|_2^2.$$

Since  $\Theta$  is a convex compact set, and  $\partial \Theta^\varepsilon$  is a compact set, we have  $\eta > 0$ .

As a result, we have for any  $\boldsymbol{\theta} \in \Theta$ ,

$$E_{\boldsymbol{\theta}}[\|\tilde{\boldsymbol{\theta}}(\mathbf{X}) - \boldsymbol{\theta}\|_2^2] \leq E_{\boldsymbol{\theta}}[\|\hat{\boldsymbol{\theta}}(\mathbf{X}) - \boldsymbol{\theta}\|_2^2] - \eta P_{\boldsymbol{\theta}}(\hat{\boldsymbol{\theta}}(\mathbf{X}) \notin \Theta^\varepsilon) \leq E_{\boldsymbol{\theta}}[\|\hat{\boldsymbol{\theta}}(\mathbf{X}) - \boldsymbol{\theta}\|_2^2] - \eta \delta,$$

and

$$\sup_{\boldsymbol{\theta} \in \Theta} E_{\boldsymbol{\theta}}[\|\tilde{\boldsymbol{\theta}}(\mathbf{X}) - \boldsymbol{\theta}\|_2^2] \leq \sup_{\boldsymbol{\theta} \in \Theta} E_{\boldsymbol{\theta}}[\|\hat{\boldsymbol{\theta}}(\mathbf{X}) - \boldsymbol{\theta}\|_2^2] - \eta \delta.$$

That means  $\tilde{\boldsymbol{\theta}}$  has strictly better worst risk than  $\hat{\boldsymbol{\theta}}$ , so that  $\hat{\boldsymbol{\theta}}$  is not minimax optimal on  $\Theta$ .

(b)

First we consider the case when  $\mathbf{M} \neq \mathbf{0}$ . Since  $\Theta$  is a compact convex set, we take  $R$  large enough so that  $\Theta^\varepsilon \subseteq B(\mathbf{0}, R)$  for some small  $\varepsilon > 0$ . The estimator  $\hat{\boldsymbol{\theta}}(\mathbf{y}) = \mathbf{M}\mathbf{y} + \boldsymbol{\theta}_0 \stackrel{d}{=} \mathbf{M}\mathbf{D}\boldsymbol{\theta} + \boldsymbol{\theta}_0 + \sigma \mathbf{M}\mathbf{g}$ , where  $\mathbf{g} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_n)$ . Note  $\sigma \mathbf{M}\mathbf{g}$  is not identically  $\mathbf{0}$  when  $\mathbf{M} \neq \mathbf{0}$ , and  $\Theta$  is a compact set, we have

$$\inf_{\boldsymbol{\theta} \in \Theta} P_{\boldsymbol{\theta}}(\|\mathbf{M}\mathbf{D}\mathbf{y} + \boldsymbol{\theta}_0\|_2 \geq R) \equiv \delta > 0.$$

By problem (a), we conclude that  $\hat{\boldsymbol{\theta}}$  cannot be minimax optimal on  $\Theta$ .

**Remark 1.** To show  $\hat{\boldsymbol{\theta}} = \boldsymbol{\theta}_0$  is not minimax optimal, we need to make the additional assumption that  $\mathbf{D} \in \mathbb{R}^{n \times d}$  has full column rank, otherwise this conclusion doesn't hold. In the following, we prove this conclusion under this additional assumption.

Then we consider the case when  $\mathbf{M} = \mathbf{0}$ . That means,  $\hat{\boldsymbol{\theta}} = \boldsymbol{\theta}_0$ . If  $\boldsymbol{\theta}_0 \notin \Theta$ , it is obvious  $\hat{\boldsymbol{\theta}}$  is not minimax optimal on  $\Theta$ . Hence we consider the case when  $\hat{\boldsymbol{\theta}} = \boldsymbol{\theta}_0 \in \Theta$ .

We claim that the  $\hat{\boldsymbol{\theta}} = \boldsymbol{\theta}_0$  cannot be the Bayes estimator for any prior except the prior  $\delta(\boldsymbol{\theta}_0)$ . Suppose this claim holds, the Bayes risk  $R_B(\hat{\boldsymbol{\theta}}, \delta(\boldsymbol{\theta}_0)) = 0$ . Since  $\Theta$  contains at least two points, it is easy to see that the minimax risk should be large than 0, hence  $\delta(\boldsymbol{\theta}_0)$  is not the least favorable prior. By minimax theorem, the minimax estimator should be the Bayes estimator for least favorable prior. Therefore,  $\hat{\boldsymbol{\theta}} = \boldsymbol{\theta}_0$  cannot be the minimax estimator.

Now suffice to show the claim above. Suppose  $Q$  is a prior probability distribution on  $\Theta$  and  $Q(\Theta \setminus \{\boldsymbol{\theta}_0\}) > 0$ , then the Bayes estimator under prior  $Q$  and square loss should be the posterior expectation  $\hat{\boldsymbol{\theta}}_Q(\mathbf{x}) = \mathbb{E}_Q[\boldsymbol{\theta}|\mathbf{x}]$ . We would like to show  $\mathbb{E}_Q[\boldsymbol{\theta}|\mathbf{x}] \neq \boldsymbol{\theta}_0$ . The intuition why  $\mathbb{E}_Q[\boldsymbol{\theta}|\mathbf{x}] \neq \boldsymbol{\theta}_0$  can be explained by the following: when  $\|\mathbf{x}\|_2 \rightarrow \infty$ , the posterior expectation  $\mathbb{E}_Q[\boldsymbol{\theta}|\mathbf{x}]$  should be at the boundary of the support of  $Q$ . In the following we show the above intuition rigorously.

By the fact that  $Q(\Theta \setminus \{\boldsymbol{\theta}_0\}) > 0$ , there exists a neighborhood  $\mathcal{B}(\boldsymbol{\theta}_*, \delta)$  such that  $Q(\mathcal{B}(\boldsymbol{\theta}_*, \delta)) \equiv \eta > 0$  and  $\|\boldsymbol{\theta}_* - \boldsymbol{\theta}_0\|_2 \geq 2\delta$ . Now we take  $\mathbf{x}_k = \mathbf{D}[\boldsymbol{\theta}_0 + k(\boldsymbol{\theta}_* - \boldsymbol{\theta}_0)]$ , then we have (denoting  $\varphi_n(\mathbf{x}) = (1/(2\pi)^{n/2}) \exp\{-\|\mathbf{x}\|_2^2/2\}$  to be the standard Gaussian density function on  $\mathbb{R}^n$ )

$$\langle \mathbb{E}_Q[\boldsymbol{\theta}|\mathbf{x}_k] - \boldsymbol{\theta}_0, \boldsymbol{\theta}_* - \boldsymbol{\theta}_0 \rangle = \frac{\int_{\Theta} \langle \boldsymbol{\theta} - \boldsymbol{\theta}_0, \boldsymbol{\theta}_* - \boldsymbol{\theta}_0 \rangle \varphi_n(\mathbf{D}(\boldsymbol{\theta}_0 - \boldsymbol{\theta} + k(\boldsymbol{\theta}_* - \boldsymbol{\theta}_0))/\sigma) Q(d\boldsymbol{\theta})}{\int_{\Theta} \varphi_n(\mathbf{D}(\boldsymbol{\theta}_0 - \boldsymbol{\theta} + k(\boldsymbol{\theta}_* - \boldsymbol{\theta}_0))/\sigma) Q(d\boldsymbol{\theta})}.$$

The integration in the numerator above can be decomposed into the integration in  $\mathcal{B}(\boldsymbol{\theta}_*, \delta)$  and the integration outside  $\mathcal{B}(\boldsymbol{\theta}_*, \delta)$ ,

$$\begin{aligned} & \int_{\Theta} \langle \boldsymbol{\theta} - \boldsymbol{\theta}_0, \boldsymbol{\theta}_* - \boldsymbol{\theta}_0 \rangle \varphi_n(\mathbf{D}(\boldsymbol{\theta}_0 - \boldsymbol{\theta} + k(\boldsymbol{\theta}_* - \boldsymbol{\theta}_0))/\sigma) Q(d\boldsymbol{\theta}) \\ & \geq \|\boldsymbol{\theta}_* - \boldsymbol{\theta}_0\|_2 (\|\boldsymbol{\theta}_* - \boldsymbol{\theta}_0\| - \delta) \frac{1}{(2\pi\sigma^2)^{n/2}} \exp\{-\|\mathbf{D}[\boldsymbol{\theta}_* - \boldsymbol{\theta}_0] - \mathbf{u}\|_2^2/(2\sigma^2)\} \eta \\ & \quad - \|\boldsymbol{\theta}_* - \boldsymbol{\theta}_0\|_2 \text{Diam}(\Theta) \frac{1}{(2\pi\sigma^2)^{n/2}} \exp\{-\|\mathbf{D}(\boldsymbol{\theta}_* - \boldsymbol{\theta}_0)\|_2^2/(2\sigma^2)\} (1 - \eta), \end{aligned}$$

where  $\text{Diam}(\Theta)$  gives the diameter of  $\Theta$ , and  $\mathbf{u} = [(\|\boldsymbol{\theta}_* - \boldsymbol{\theta}_0\|_2 - \delta)/\|\boldsymbol{\theta}_* - \boldsymbol{\theta}_0\|_2](\boldsymbol{\theta}_* - \boldsymbol{\theta}_0)$ . Note (we already assumed  $\mathbf{D}$  has full column rank)

$$\lim_{k \rightarrow \infty} \frac{\exp\{-\|\mathbf{D}[\boldsymbol{\theta}_* - \boldsymbol{\theta}_0] - \mathbf{u}\|_2^2/(2\sigma^2)\}}{\exp\{-\|\mathbf{D}(\boldsymbol{\theta}_* - \boldsymbol{\theta}_0)\|_2^2/(2\sigma^2)\}} = \infty,$$

hence for large  $k$ , we have

$$\langle \mathbb{E}_Q[\boldsymbol{\theta}|\mathbf{x}_k] - \boldsymbol{\theta}_0, \boldsymbol{\theta}_* - \boldsymbol{\theta}_0 \rangle > 0.$$

That means, we have  $\mathbb{E}_Q[\boldsymbol{\theta}|\mathbf{x}_k] \neq \boldsymbol{\theta}_0$  for large  $k$ . This proves the claim.

**(c)**

Let  $\Theta = \{-1, 1\}$ ,  $\mathsf{P}_1 = \mathsf{P}_0 = \delta(0)$  (no matter what  $\theta$  is, the data  $X$  is deterministically 0). Hence we only need to consider the estimator that is a constant mapping (Rao-Blackwell theorem tells us that we don't need to consider randomized estimator). The risk function for any constant estimator is  $R(\hat{\theta} = a; \Theta) = \sup\{(1-a)^2, (-1-a)^2\}$ . Minimizing this over  $a$ , the minimax estimator is  $\hat{\theta} = 0$ . For this estimator, for  $\varepsilon < 1/4$ ,  $\mathsf{P}_0(\hat{\theta} \notin \{-1, 1\}^\varepsilon) = \mathsf{P}_1(\hat{\theta} \notin \{-1, 1\}^\varepsilon) = 1$ .

**(d)**

Consider the estimator  $\tilde{\theta} = \text{Proj}(\hat{\theta})$ , where Proj operator enjoy the same definition of Problem (a), then we have

$$L(\tilde{\theta}(x), \theta) \leq L(\hat{\theta}(x), \theta) - \eta \mathbf{1}\{\hat{\theta}(x) \notin \Theta^\varepsilon\},$$

where

$$\eta = \min_{\theta \in \Theta, \theta' \in \partial\Theta^\varepsilon} L(\theta', \theta) - L(\text{Proj}(\theta'), \theta).$$

Since  $L$  is strictly decreasing for  $a < \theta$  and strictly increasing for  $a > \theta$ , and  $\Theta$  and  $\partial\Theta^\varepsilon$  are compact sets, we have  $\eta > 0$ .

As a result, we have for any  $\theta \in \Theta$ ,

$$R(\tilde{\theta}, \theta) = \mathbb{E}_\theta[L(\tilde{\theta}(X), \theta)] \leq \mathbb{E}_\theta[L(\hat{\theta}(X), \theta)] - \eta \mathbb{P}_\theta(\hat{\theta}(X) \notin \Theta^\varepsilon) \leq R(\hat{\theta}, \theta) - \eta\delta.$$

Since  $R(\hat{\theta}, \theta)$  is continuous in  $\theta$ ,  $R(\hat{\theta}, \theta)$  can attain the maximum, and we have

$$\sup_{\theta \in \Theta} R(\tilde{\theta}, \theta) \leq \sup_{\theta \in \Theta} R(\hat{\theta}, \theta) - \eta\delta.$$

That means  $\hat{\theta}$  is not minimax optimal on  $\Theta$ .

## Problem 2

(a)

Let  $\theta_1 = 1/2 - 1/(2\sqrt{2})$ ,  $\theta_2 = 1/2 + 1/(2\sqrt{2})$ , and let  $q = 1/2$ . Under the square loss, the Bayes optimal estimator for  $Q$  is given by the conditional expectation

$$\begin{aligned}\hat{\theta}_B(x) &= \mathbb{E}[\theta|X=x] \\ &= \begin{cases} \frac{\theta_1^2 + \theta_2^2}{\theta_1 + \theta_2} & \text{if } x=1 \\ \frac{\theta_1(1-\theta_1) + \theta_2(1-\theta_2)}{1-\theta_1+1-\theta_2} & \text{if } x=0 \end{cases} \\ &= \begin{cases} \frac{3}{4} & \text{if } x=1 \\ \frac{1}{4} & \text{if } x=0 \end{cases} \\ &= \frac{x}{2} + \frac{1}{4}.\end{aligned}\tag{1}$$

The above implies that  $\hat{\theta}_B(Q) = \hat{\theta}_{MM}$ .

(b)

As suggested in the hint, there exists an integer  $m$ , such that choosing  $q_i \geq 0$  for  $i = 0, 1, \dots, m$  such that (here  $Q$  is the measure induced by a  $\text{Beta}(\sqrt{n}/2, \sqrt{n}/2)$  random variable)

$$\sum_{i=0}^m q_i \left(\frac{i}{m}\right)^k = \int \theta^k Q(d\theta) \quad \text{for all } k = 0, 1, \dots, n+1.\tag{2}$$

Then the above implies that, for any polynomial  $p$  of degree at most  $n+1$ , we have

$$\sum_{i=0}^m q_i p\left(\frac{i}{m}\right) = \int p(\theta) Q(d\theta).\tag{3}$$

Consider the prior distribution:

$$Q_1 = \sum_{i=0}^{n+1} q_i \delta\left(\frac{i}{m}\right)\tag{4}$$

The Bayes optimal estimator is given by the conditional expectation

$$\begin{aligned}\hat{\theta}_{Q_1}(X) &= \mathbb{E}_{Q_2}[\theta|X] \\ &= \frac{\sum_{i=0}^{n+1} q_i (i/m)^{X+1} (1-i/m)^{n-X}}{\sum_{i=0}^{n+1} q_i (i/m)^X (1-i/m)^{n-X}}.\end{aligned}\tag{5}$$

On the other hand, the Bayes estimator with respect to  $\text{Beta}(\sqrt{n}/2, \sqrt{n}/2)$  is given by

$$\begin{aligned}\hat{\theta}_{MM}(X) &= \frac{\sqrt{n}}{1+\sqrt{n}} \cdot \frac{X}{n} + \frac{1}{1+\sqrt{n}} \cdot \frac{1}{2} \\ &= \mathbb{E}_Q[\theta|X] \\ &= \frac{\int \theta^{X+1} (1-\theta)^{n-X} Q(d\theta)}{\int \theta^X (1-\theta)^{n-X} Q(d\theta)}.\end{aligned}\tag{6}$$

Let  $p_1(t; X) = t^{X+1} (1-t)^{n-X}$ ,  $p_2(t; X) = t^X (1-t)^{n-X}$ , then it clear that both  $p_1$  and  $p_2$  as a function of  $t$  are polynomial of degree at most  $n+1$ . Hence by (3) we have

$$\hat{\theta}_{Q_1}(X) = \frac{\sum_{i=0}^m p_1\left(\frac{i}{m}, X\right) q_i}{\sum_{i=0}^m p_2\left(\frac{i}{m}, X\right) q_i} = \frac{\int p_1(\theta, X) Q(d\theta)}{\int p_2(\theta, X) Q(d\theta)} = \hat{\theta}_{MM}(X).\tag{7}$$

Therefore,

$$\hat{\theta}_{Q_1}(X) = \hat{\theta}_{MM}(X) = \frac{\sqrt{n}}{1+\sqrt{n}} \cdot \frac{X}{n} + \frac{1}{1+\sqrt{n}} \cdot \frac{1}{2}.\tag{8}$$

### Problem 3

(a)

Since  $L$  is upper bounded by  $L_0$ ,  $R(A, \boldsymbol{\theta})$  is also bounded from above by  $L_0$  for all  $A \in \mathcal{A}$  and  $\boldsymbol{\theta} \in \Theta$ . Given  $Q$ , for any statistical procedure  $A$ , we have

$$\begin{aligned} R(A, Q) &= \int_{\mathbb{R}^d} R(A, \boldsymbol{\theta}) Q(d\boldsymbol{\theta}) = \int_{\Theta} R(A, \boldsymbol{\theta}) Q(d\boldsymbol{\theta}) + \int_{\Theta^c} R(A, \boldsymbol{\theta}) Q(d\boldsymbol{\theta}) \\ &\leq \sup_{\boldsymbol{\theta} \in \Theta} R(A, \boldsymbol{\theta}) + L_0 Q(\Theta^c). \end{aligned} \quad (9)$$

Hence

$$R_B(Q) - L_0 Q(\Theta^c) \leq R(A, Q) - L_0 Q(\Theta^c) \leq \sup_{\boldsymbol{\theta} \in \Theta} R(A, \boldsymbol{\theta}). \quad (10)$$

Since the above is true for all  $A$ , taking the infimum over  $A \in \mathcal{A}$  gives

$$R_M(\Theta) \geq R_B(Q) - L_0 Q(\Theta^c). \quad (11)$$

(b)

Let  $\hat{\boldsymbol{\theta}}$  be any estimator, and let  $\tilde{\boldsymbol{\theta}}$  be the projection of  $\hat{\boldsymbol{\theta}}$  onto  $\mathbb{B}^d(\mathbf{0}, M\sqrt{k})$ . That is

$$\tilde{\boldsymbol{\theta}} = \min \left\{ \frac{M\sqrt{k}}{\|\hat{\boldsymbol{\theta}}\|_2}, 1 \right\} \hat{\boldsymbol{\theta}}. \quad (12)$$

Then it is clear that  $L(\tilde{\boldsymbol{\theta}}, \boldsymbol{\theta}) \leq L(\hat{\boldsymbol{\theta}}, \boldsymbol{\theta})$  with probability 1 for all  $\boldsymbol{\theta} \in \Theta(d, k, M) \subset \mathbb{B}^d(\mathbf{0}, M\sqrt{k})$ . Since  $\tilde{\boldsymbol{\theta}} \in \mathbb{B}^d(\mathbf{0}, M\sqrt{k})$ , it is sufficient to only consider estimators taking values in  $\mathbb{B}^d(\mathbf{0}, M\sqrt{k})$ . In this case, since both  $\tilde{\boldsymbol{\theta}}$  and  $\boldsymbol{\theta}$  are in a ball with radius  $M\sqrt{k}$ , there distance square is upper bounded by the diameter square of the ball. That is, for all  $\boldsymbol{\theta} \in \Theta$  and  $\tilde{\boldsymbol{\theta}}$  in the above form, we have

$$L(\tilde{\boldsymbol{\theta}}, \boldsymbol{\theta}) \leq 4M^2k. \quad (13)$$

Therefore it is also sufficient to replace the square loss by  $\tilde{L}(\hat{\boldsymbol{\theta}}, \boldsymbol{\theta}) = \min\{\|\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}\|_2^2, 4M^2k\}$ .

(c)

Let  $G = \Pi_d \times \Sigma_d$  be a group, where  $\Pi_d$  is the permutation group on  $\{1, \dots, d\}$ , and  $\Sigma_d = \{+1, -1\}^d$  is the sign changing group. For any  $g = [\pi, \boldsymbol{\sigma}] \in G$  ( $\pi$  is a permutation, where  $\{\pi(1), \dots, \pi(d)\} = \{1, \dots, d\}$  as a set;  $\boldsymbol{\sigma} = [\sigma_1, \dots, \sigma_d]^\top \in \{+1, -1\}^d$ ), the action of  $\varphi_g$  on  $\mathbf{x} = (x_1, \dots, x_d)^\top \in \mathbb{R}^d$  gives  $\varphi_g(\mathbf{x}) = (\sigma_1 x_{\pi(1)}, \dots, \sigma_d x_{\pi(d)})^\top$ . We would like to show our statistical model is invariant under this group. First we have  $L(a, \boldsymbol{\theta}) = \|a - \boldsymbol{\theta}\|_2^2 = \|\varphi_g(a) - \varphi_g(\boldsymbol{\theta})\|_2^2 = L(\varphi_g(a), \varphi_g(\boldsymbol{\theta}))$ . Next we have  $\mathbb{P}_{g(\boldsymbol{\theta})}(\mathbf{X} \in S) = \mathbb{P}_{\mathbf{Z} \sim \mathcal{N}(0, \sigma^2 \mathbf{I}_d)}(\varphi_g(\boldsymbol{\theta}) + \mathbf{Z} \in S) = \mathbb{P}_{\mathbf{Z} \sim \mathcal{N}(0, \sigma^2 \mathbf{I}_d)}(\varphi_g(\boldsymbol{\theta}) + \varphi_g(\mathbf{Z}) \in S) = \mathbb{P}_{\mathbf{Z} \sim \mathcal{N}(0, \sigma^2 \mathbf{I}_d)}(\varphi_g(\boldsymbol{\theta} + \mathbf{Z}) \in S) = \mathbb{P}_{\boldsymbol{\theta}}(\varphi_g(\mathbf{X}) \in S) = (\varphi_g)_{\#} \mathbb{P}_{\boldsymbol{\theta}}(\mathbf{X} \in S)$ . Hence our model is invariant under this group. Since minimax theorem holds for this model, there exists a least favorable prior. According to invariant least favorable prior theorem, there exists a least favorable prior that is invariant under the group action. This invariant least favorable prior can only be written in the form  $Q = \sum_{\ell=0}^k p_\ell Q_\ell$ .

(d)

By part (b) we know that  $R_M(d, k; M) = \tilde{R}_M(d, k; M)$ , and we can replace the loss  $L$  by  $\tilde{L}$ , which is bounded from above by  $4M^2k$ . By part (a) we have

$$R_M(d, k; M) = \tilde{R}_M(d, k; M) \geq \tilde{R}_B(Q_{M,\epsilon}) - 4M^2k Q_{M,\epsilon}(\Theta^c). \quad (14)$$

Let  $\mathbf{X} \in \mathbb{R}^d$  be a random variable whose induced measure is  $Q_{M,\epsilon}$ , then it is clear that  $Q_{M,\epsilon}(\Theta^c)$  is equal to  $\mathbb{P}(\|\mathbf{X}\|_0 > k)$ . Since the coordinates of  $\mathbf{X}$  are independent and  $\mathbf{1}(X_i \neq 0)$  has Bernoulli( $\epsilon$ ) distribution,  $\|\mathbf{X}\|_0$  has Binomial( $d, \epsilon$ ) distribution. Therefore, (14) becomes

$$R_M(d, k; M) \geq \tilde{R}_B(Q_{M,\epsilon}) - 4M^2k \mathbb{P}(\text{Binom}(d, \epsilon) > k). \quad (15)$$

(e)

Note  $\boldsymbol{\theta} = (\theta_1, \dots, \theta_d) \sim Q_{M,\epsilon} = q_{M,\epsilon}^{\otimes d}$ , and  $\mathbf{X} \sim \mathcal{N}(\boldsymbol{\theta}, \sigma^2 \mathbf{I}_d)$ . We have  $(X_i, \theta_i)$  for  $i \in [d]$  are mutually independent. Hence the Bayes estimator which is the posterior mean gives

$$(\hat{\boldsymbol{\theta}}_B(\mathbf{x}))_j = \mathbb{E}[\theta_j | \mathbf{X} = \mathbf{x}] = \mathbb{E}[\theta_j | X_j = x_j].$$

Hence

$$\begin{aligned} R_B(Q_{M,\epsilon}) &= \mathbb{E}_{Q_{M,\epsilon}}[\|\hat{\boldsymbol{\theta}}_B - \boldsymbol{\theta}\|_2^2] \\ &= \sum_{j \in [d]} \mathbb{E}_{Q_{M,\epsilon}}[((\hat{\boldsymbol{\theta}}_B)_j - \theta_j)^2] \\ &= \sum_{j \in [d]} \mathbb{E}_{Q_{M,\epsilon}}[(\mathbb{E}(\theta_j | X_j) - \theta_j)^2] \\ &= \sum_{j \in [d]} \mathbb{E}_{q_{M,\epsilon}}[(\mathbb{E}(\theta_j | X_j) - \theta_j)^2] \\ &= dR_B(q_{M,\epsilon}). \end{aligned} \tag{16}$$

Since we have

$$P(\text{Binom}(d, \epsilon) > k) \leq e^{-k\eta^2/4}, \tag{17}$$

which implies that  $kP(\text{Binom}(d, \epsilon) > k) = o_\eta(k)$ , using (15), (16) and (7) in the question gives

$$R_M(d, k; M) \geq dR_B(q_{M,\epsilon}) - (M^2 + 1)o_\eta(k) - 4M^2o_\eta(k). \tag{18}$$

Since a constant times  $o_\eta(k)$  is still  $o_\eta(k)$ , the  $-4M^2o_\eta(k)$  above can be merged with the first  $M^2o_\eta(k)$ , so it can be simplifies to

$$R_M(d, k; M) \geq dR_B(q_{M,\epsilon}) - (M^2 + 1)o_\eta(k). \tag{19}$$