

## A Asymptotic Lower Bound of Shtarkov Complexity for Standard Normal Location Models

We show an asymptotic lower bound of the Shtarkov complexity of standard normal location models.

**Lemma 8** *Consider the  $d$ -dimensional standard normal location model, given by  $f_X(\theta) = \frac{1}{2} \|X - \theta\|_2^2 + \frac{d}{2} \ln 2\pi$ , where  $X \in \mathcal{X} = \mathbb{R}^d$ . Let  $\gamma = \lambda \|\theta\|_1$  for  $\lambda \geq 0$ . Then we have*

$$S(\gamma) \geq d \ln \left( 1 + \frac{e^{-\lambda^2/2}}{\sqrt{2\pi}\lambda^3} (1 + o(1)) \right).$$

**Proof** By definition of  $S(\gamma)$ , we have

$$\begin{aligned} S(\gamma) &= \ln \int e^{-m(f_X + \gamma)} \nu(dX) \\ &= d \ln \int_{-\infty}^{\infty} \frac{\sup_{t \in \mathbb{R}} \exp \left[ -\frac{1}{2}(x-t)^2 - \lambda|t| \right]}{\sqrt{2\pi}} dx \\ &= d \ln \frac{1}{\sqrt{2\pi}} \left[ \int_{-\infty}^{-\lambda} e^{-\lambda(-\lambda-x) - \frac{\lambda^2}{2}} dx + \int_{-\lambda}^{\lambda} e^{-\frac{x^2}{2}} dx + \int_{\lambda}^{\infty} e^{-\lambda(x-\lambda) - \frac{\lambda^2}{2}} dx \right] \\ &= d \ln \left[ 2\Phi(\lambda) - 1 + \frac{2e^{-\lambda^2/2}}{\sqrt{2\pi}} \int_0^{\infty} e^{-\lambda x} dx \right] \\ &= d \ln \left[ 2\Phi(\lambda) - 1 + \sqrt{\frac{2}{\pi}} \frac{e^{-\lambda^2/2}}{\lambda} \right], \end{aligned}$$

where  $\Phi(\lambda)$  denotes the standard normal distribution function. Now, by Komatu (1955),  $\Phi(\lambda)$  is bounded below with  $\Phi(\lambda) > 1 - 2\phi(\lambda)/(\sqrt{2} + x^2 + x)$  for  $\phi(\lambda)$  being the standard normal density, which yields the lower bound of interest after a few lines of elementary calculation.  $\blacksquare$

## B Lower Bound on Minimax Regret of Smooth Models

We describe how we adopt the minimax risk lower bound as to show the minimax-regret lower bound.

The story of the proof is based on Donoho and Johnstone (1994). First, the so-called three-point prior is constructed to approximate the least favorable prior. Then, since the approximate prior violates the  $\ell_1$ -constraint, the degree of the violation is shown to be appropriately bounded to derive a valid lower bound.

The goal of our proof is to establish a lower bound on the minimax regret with respect to logarithmic losses,

whereas their proof is about the minimax risk with respect to  $\ell_q$ -loss. Therefore, below we present the proof highlighting (i) an approximate least favorable prior for *logarithmic losses* over  $\ell_1$ -balls and (ii) the way to bound *regrets* on the basis of risk bounds.

Let  $\mathcal{H} = \{\theta \in \mathbb{R}^d \mid \|\theta\|_1 \leq B\}$  be a  $\ell_1$ -ball. Let  $X \sim \mathcal{N}_d[\theta, I_d/L]$  be a  $d$ -dimensional normal random variable with mean  $\theta \in H$  and precision  $L > 0$ . We denote the distribution just by  $X \sim \theta$  where any confusion is unlikely. Let  $h \in \mathcal{H}$  be a predictor associated with any sub-probability distribution  $P(\cdot|h) \in \mathcal{M}_+(\mathbb{R}^d)$ . For notational simplicity, we may write  $f_X(\theta) = \frac{L}{2} \|X - \theta\|_2^2 + \frac{d}{2} \ln \frac{2\pi}{L}$  and  $f_X(h) = \ln \frac{dP(X|h)}{d\nu}$  where  $\nu$  is the Lebesgue measure over  $\mathbb{R}^d$ .

Consider the risk function

$$R_d(h, \theta) \stackrel{\text{def}}{=} \mathbb{E}_{X \sim \theta} [f_X(h) - f_X(\theta)],$$

and the Bayes risk function

$$R_d(h, \pi) \stackrel{\text{def}}{=} \mathbb{E}_{\theta \sim \pi} [R_d(h, \theta)],$$

where  $\pi \in \mathcal{P}(\mathcal{H})$  denotes prior distributions on  $\mathcal{H}$ . Then, the minimax Bayes risk bounds below the minimax regret,

$$\begin{aligned} \text{REG}^*(\mathcal{H}) &= \inf_{h \in \mathcal{H}} \sup_{\theta \in \mathcal{H}} \sup_{X \in \mathbb{R}^d} f_X(h) - f_X(\theta) \\ &\geq \inf_{h \in \mathcal{H}} \sup_{\pi \in \mathcal{P}(\mathcal{H})} \mathbb{E}_{\theta \sim \pi} \mathbb{E}_{X \sim \theta} [f_X(h) - f_X(\theta)] \\ &= \inf_{h \in \mathcal{H}} \sup_{\pi \in \mathcal{P}(\mathcal{H})} R_d(h, \pi). \end{aligned}$$

The minimax theorem states that there exists a saddle point  $(h^*, \pi_*)$  such that

$$\begin{aligned} R_d(h^*, \pi_*) &= \inf_{h \in \mathcal{H}} \sup_{\pi \in \mathcal{P}(\mathcal{H})} R_d(h, \pi) \\ &= \sup_{\pi \in \mathcal{P}(\mathcal{H})} \inf_{h \in \mathcal{H}} R_d(h, \pi) \stackrel{\text{def}}{=} \sup_{\pi \in \mathcal{P}(\mathcal{H})} R_d(\pi), \end{aligned}$$

and  $\pi_*$  is referred to as the least favorable prior. We want to approximate  $\pi_*$  to give an analytic approximation of  $R_d(\pi_*)$ , which is a lower bound of  $\text{REG}^*(\mathcal{H})$ .

Let  $F_{\epsilon, \mu} \in \mathcal{P}(\mathbb{R})$  be the three-point prior defined by

$$F_{\epsilon, \mu} = (1 - \epsilon)\delta_0 + \frac{\epsilon}{2}(\delta_{-\mu} + \delta_{\mu})$$

for  $\epsilon, \mu > 0$ . We show that the corresponding achievable Bayes risk  $R_1(F_{\epsilon, \mu})$  tends to be the entropy of the prior  $F_{\epsilon, \mu}$  in some limit of small  $\epsilon$ .

**Lemma 9** *Take  $\mu = \mu(\epsilon) = \sqrt{2L^{-1} \ln \epsilon^{-1}}$ . Let  $H_{\epsilon} = H(F_{\epsilon, \mu}) = (1 - \epsilon) \ln(1 - \epsilon)^{-1} + \epsilon \ln 2\epsilon^{-1}$  be the entropy of the prior. Then we have*

$$R_1(F_{\epsilon, \mu}) \sim H_{\epsilon} \sim \epsilon \ln \frac{1}{\epsilon}$$

as  $\epsilon \rightarrow 0$ . Here,  $x \sim y$  denotes the asymptotic equality such that  $x/y \rightarrow 1$ .

**Proof** First, we show the famous inequality on the entropy given by  $R_1(F_{\epsilon,\mu}) \leq H_\epsilon$ . Let  $\hat{P}(\cdot|h) = \mathbb{E}_{\theta \sim F_{\epsilon,\mu}} P(\cdot|\theta) = (1-\epsilon)P(\cdot|0) + \frac{\epsilon}{2}(P(\cdot|-\mu) + P(\cdot|\mu))$  be the Bayes marginal distribution with respect to  $F_{\epsilon,\mu}$ . Then we have

$$\begin{aligned} H_\epsilon - R_1(F_{\epsilon,\mu}) &= H_\epsilon - R_1(h, F_{\epsilon,\mu}) \\ &= H_\epsilon - \mathbb{E}_{\theta \sim F_{\epsilon,\mu}} \mathbb{E}_{X \sim \theta} \ln \frac{dP(X|\theta)}{dP(X|h)} \\ &= H_\epsilon - (1-\epsilon) \mathbb{E}_{P(X|0)} \ln \frac{dP(X|0)}{dP(X|h)} \\ &\quad - \epsilon \mathbb{E}_{P(X|\mu)} \ln \frac{dP(X|\mu)}{dP(X|h)} \\ &= (1-\epsilon) \mathbb{E}_{P(X|0)} \ln \left( 1 + \frac{\epsilon}{1-\epsilon} \frac{dP(X|\mu) + dP(X|-\mu)}{2dP(X|0)} \right) \\ &\quad + \epsilon \mathbb{E}_{P(X|\mu)} \ln \left( 1 + \frac{1-\epsilon}{\epsilon} \frac{2dP(X|0) + dP(X|-\mu)}{dP(X|\mu)} \right) \\ &\geq 0. \end{aligned}$$

Now, we show that, with the specific value of  $\mu = \mu(\epsilon)$ , the gap is negligible compared to the entropy itself. Applying Jensen's inequality, we have

$$\begin{aligned} H_\epsilon - R_1(F_{\epsilon,\mu}) &\leq \epsilon + \epsilon \mathbb{E}_{P(X|\mu)} \ln \left( 1 + (1-\epsilon) \left( 2e^{-L\mu X} + \epsilon^3 e^{-2L\mu X} \right) \right) \\ &\leq \epsilon(1 + \ln 4 + \mathbb{E}_{P(X|\mu)} \max\{0, -2L\mu X\}) \\ &= \epsilon \left( 1 + \ln 4 + \mathbb{E}_{Z \sim \mathcal{N}[0,1]} \max\{0, 2\sqrt{L}\mu(Z - \sqrt{L}\mu)\} \right) \\ &\quad (\because -\sqrt{L}(X - \mu) = Z) \\ &\leq \epsilon \left( 1 + \ln 4 + 2\sqrt{L}\mu\epsilon \right) \\ &= \epsilon \left( 1 + \ln 4 + 2\epsilon\sqrt{2\ln \frac{1}{\epsilon}} \right) = o(H_\epsilon). \end{aligned}$$

Thus we get  $H_\epsilon \sim R_1(F_{\epsilon,\mu})$ .  $\blacksquare$

Now we show that the  $d$ -th Kronecker product of  $F_{\epsilon,\mu}$ ,  $F_{\epsilon,\mu}^d$ , can be used to bound the Bayes minimax risk  $R_d(\pi_*)$  with an appropriate choice of  $\epsilon$  and  $\mu$ . To this end, let  $\pi_+ = F_{\epsilon,\mu}^d | \mathcal{H}$  be the conditional prior restricted over the  $\ell_1$ -ball  $\mathcal{H}$ .

**Lemma 10** Take  $\epsilon\mu = (1-c)B/d$  and  $\mu = \sqrt{2L^{-1} \ln \epsilon^{-1}}$  for  $0 < c < 1$ . Then, if  $\epsilon \rightarrow 0$  and  $d\epsilon \rightarrow \infty$ , we have

$$R_d(\pi_*) \geq R_d(\pi_+) \sim R_d(F_{\epsilon,\mu}^d) \sim d\epsilon \ln \frac{1}{\epsilon}.$$

**Proof** First of all, the inequality is trivial from the definition of  $R_d(\pi)$ . Moreover, the second asymptotic equality immediately follows from Lemma 9.

Now we consider the first asymptotic equality. Let  $h$  be the Bayesian predictor with respect to the prior  $F_{\epsilon,\mu}$  and  $h^+$  be the one with respect to the conditional prior  $\pi_+$ . Then we have

$$\begin{aligned} R_d(F_{\epsilon,\mu}^d) &= R_d(h, F_{\epsilon,\mu}^d) \\ &= \mathbb{E}_{\theta \sim F_{\epsilon,\mu}^d} [R_d(h, \theta)] \\ &= F_{\epsilon,\mu}^d(\mathcal{H}) R_d(h, \pi_+) + \mathbb{E}_{\theta \sim F_{\epsilon,\mu}^d} [R_d(h, \theta) \cdot \mathbf{1}\{\theta \notin \mathcal{H}\}] \\ &\geq F_{\epsilon,\mu}^d(\mathcal{H}) \cdot R_d(\pi_+) \end{aligned}$$

and

$$\begin{aligned} R_d(F_{\epsilon,\mu}^d) &\leq R_d(h^+, F_{\epsilon,\mu}^d) \\ &= \mathbb{E}_{\theta \sim F_{\epsilon,\mu}^d} [R_d(h^+, \theta)] \\ &= F_{\epsilon,\mu}^d(\mathcal{H}) \cdot R_d(\pi_+) + \\ &\quad \mathbb{E}_{\theta \sim F_{\epsilon,\mu}^d} [R_d(h^+, \theta) \cdot \mathbf{1}\{\theta \notin \mathcal{H}\}]. \end{aligned}$$

Let  $N$  be the number of nonzero elements in  $\theta \sim F_{\epsilon,\mu}^d$ . Then  $N$  is subjects to the Binomial distribution  $\text{Bin}(d, \epsilon)$ . On the other hand, the event  $\theta \in \mathcal{H}$  is equal to  $\{\|\theta\|_1 \leq B\} = \{N \leq B/\mu = \mathbb{E}N/(1-c)\}$ . Therefore, applying the Chebyshev's inequality, we get

$$\begin{aligned} P_d &\stackrel{\text{def}}{=} F_{\epsilon,\mu}^d(\mathcal{H}^c) = \Pr \left\{ \frac{N - \mathbb{E}N}{\mathbb{E}N} > \frac{c}{1-c} \right\} \\ &\leq \frac{(1-c)^2}{c^2 d \epsilon} \rightarrow 0. \end{aligned}$$

Similarly, we have  $\mathbb{E}|N - \mathbb{E}N|/\mathbb{E}N \rightarrow 0$ . Now observe that

$$\begin{aligned} &\mathbb{E}_{\theta \sim F_{\epsilon,\mu}^d} [R_d(h^+, \theta) \cdot \mathbf{1}\{\theta \notin \mathcal{H}\}] \\ &\leq \mathbb{E}_{\theta \sim F_{\epsilon,\mu}^d} \mathbb{E}_{\varphi \sim \pi_+} [R_d(\varphi, \theta) \cdot \mathbf{1}\{\theta \notin \mathcal{H}\}] \\ &\leq 2L \mathbb{E}_{\theta \sim F_{\epsilon,\mu}^d} \mathbb{E}_{\varphi \sim \pi_+} \left[ \left( \|\varphi\|_2^2 + \|\theta\|_2^2 \right) \cdot \mathbf{1}\{\theta \notin \mathcal{H}\} \right] \\ &\leq 2L\mu^2 \mathbb{E} [P_d N + N \cdot \mathbf{1}\{N > B/\mu\}] \\ &\quad (\because \|\theta\|_2^2 = \mu^2 N) \\ &\leq 2L\mu^2 \mathbb{E}N \left( 2P_d + \frac{\mathbb{E}|N - \mathbb{E}N|}{\mathbb{E}N} \right) \\ &= 4d\epsilon \ln \frac{1}{\epsilon} \left( 2P_d + \frac{\mathbb{E}|N - \mathbb{E}N|}{\mathbb{E}N} \right) \\ &= o(R_d(F_{\epsilon,\mu}^d)). \end{aligned}$$

Thus, combining all above, we get

$$\begin{aligned} &(1 + o(1))R_d(\pi_+) \\ &= (1 - P_d)R_d(\pi_+) \\ &\leq R_d(F_{\epsilon,\mu}^d) \\ &\leq (1 - P_d) \cdot R_d(\pi_+) + \\ &\quad \mathbb{E}_{\theta \sim F_{\epsilon,\mu}^d} [R_d(h^*, \theta) \cdot \mathbf{1}\{\theta \notin \mathcal{H}\}]. \\ &= (1 - o(1))R_d(\pi_*) + o(R_d(F_{\epsilon,\mu}^d)), \end{aligned}$$

which implies the desired asymptotic equality  $R_d(F_{\epsilon,\mu}) \sim R_d(\pi_+)$ . ■

Summing these up, we have an asymptotic lower bound on the minimax regret which is the same as the upper bound given by the ST prior within a factor of two (see Theorem 7). This implies that both the regret of the ST prior and the Bayes risk of the prior  $\pi_+$  are tight with respect to the minimax-regret rate except with a factor of two.

**Theorem 11 (Minimax lower bound)** *Suppose that  $\omega(1) = \ln(d/\sqrt{L}) = o(L)$ . Then we have*

$$\text{REG}^*(\mathcal{H}) \gtrsim \frac{B}{2} \sqrt{2L \ln \frac{d}{\sqrt{L}}},$$

where  $x \gtrsim y$  means that there exists  $y' \sim y$  such that  $x \geq y'$ .

**Proof** The assumptions of Lemma 10 are satisfied for all  $0 < c < 1$  since

$$\begin{aligned} \epsilon &\lesssim \epsilon \sqrt{\ln \frac{1}{\epsilon}} = \frac{1-c}{d} \sqrt{\frac{L}{2}} \rightarrow 0, \\ d\epsilon &= (1-c) \sqrt{\frac{L}{2 \ln \frac{1}{\epsilon}}} \sim (1-c) \sqrt{\frac{L}{2 \ln \frac{d}{\sqrt{L}}}} \rightarrow \infty. \end{aligned}$$

Thus, we have

$$\text{REG}^*(\mathcal{H}) \geq R_d(\pi_*) \gtrsim d\epsilon \ln \frac{1}{\epsilon} \sim (1-c) \frac{B}{2} \sqrt{2L \ln \frac{d}{\sqrt{L}}}$$

for all  $0 < c < 1$ . Slowly moving  $c$  toward zero completes the theorem. ■

## C Existence of Gap between LREG\* and LREG<sup>Bayes</sup> under $\ell_1$ -Penalty

Below we show that, under standard normal location models, the Bayesian luckiness minimax regret is strictly larger than the non-Bayesian luckiness minimax regret if  $\gamma$  is nontrivial and has a non-differentiable point. Here we refer to  $\gamma$  as *trivial* when there exists  $\theta_0$  such that  $\gamma(\theta) = \infty$  for all  $\theta \neq \theta_0$ .

**Lemma 12** *Let  $f_X(\theta) = \frac{1}{2}(X - \theta)^2 + \frac{1}{2} \ln 2\pi$  for  $X \in \mathbb{R}$  and  $\theta \in \mathbb{R}$ . Then, for all nontrivial, convex and non-differentiable penalties  $\gamma : \mathbb{R} \rightarrow \overline{\mathbb{R}}$ ,*

$$\text{LREG}^*(\gamma) < \text{LREG}^{\text{Bayes}}(\gamma).$$

**Proof** Let  $\mathcal{F} = \{f_X \mid X \in \mathbb{R}\}$  and recall that  $\text{LREG}^{\text{Bayes}}(\gamma) = \inf_{w \in \mathcal{E}(\mathcal{F}_\gamma)} \ln w [e^{-\gamma}]$  by Theorem 1. Let  $\|\cdot\|_\gamma$  be the metric of pre-priors  $w \in \mathcal{M}_+(\mathbb{R})$  given by  $\|w\|_\gamma = w [e^{-\gamma}]$ . Owing to the continuity of  $w \mapsto \ln w [e^{-\gamma}]$  and the completeness of  $\mathcal{E}(\mathcal{F}_\gamma) \subset \mathcal{M}_+(\mathbb{R})$ , it suffices to show that there exists no pre-prior  $w \in \mathcal{E}(\mathcal{F}_\gamma)$  such that  $\ln w [e^{-\gamma}] = S(\gamma)$ . Let us prove this by contradiction. Now, assume that  $\ln w [e^{-\gamma}] = S(\gamma)$ . Observe that

$$\begin{aligned} 0 &= w [e^{-\gamma}] - \exp S(\gamma) \\ &= w \left[ \int e^{-f_X - \gamma} \nu(dX) \right] - \int e^{-m(f_X + \gamma)} \nu(dX) \\ &= \int \left\{ w [e^{-f_X - \gamma}] - e^{-m(f_X + \gamma)} \right\} \nu(dX), \end{aligned}$$

which means  $w [e^{-f_X - \gamma}] = e^{-m(f_X + \gamma)}$  for almost every  $X$  since  $w \in \mathcal{E}(\mathcal{F}_\gamma)$ . Note that  $f_X(\theta)$  is continuous with respect to  $X$ , and then we have  $w [e^{-f_X - \gamma}] = e^{-m(f_X + \gamma)}$  for all  $X$ . After some rearrangement and differentiation, we have

$$\begin{aligned} 0 &= \frac{d}{dX} w [e^{-f_X - \gamma + m(f_X + \gamma)}] \\ &= w \left[ \frac{d e^{-f_X - \gamma + m(f_X + \gamma)}}{dX} \right] \\ &= w_\theta \left[ (\theta - \theta_X^*) e^{-f_X - \gamma + m(f_X + \gamma)} \right], \quad (13) \end{aligned}$$

where  $\theta_X^* = \arg m(f_X + \gamma)$ . Here we exploited Dan-skin's theorem at the last equality. One more differentiation gives us

$$\begin{aligned} 0 &= \frac{d}{dX} w_\theta \left[ (\theta - \theta_X^*) e^{-f_X - \gamma + m(f_X + \gamma)} \right], \\ &= w_\theta \left[ \left\{ (\theta - \theta_X^*)^2 - \frac{d\theta_X^*}{dX} \right\} e^{-f_X - \gamma + m(f_X + \gamma)} \right] \end{aligned}$$

for all  $X \in \mathbb{R}$ .

Note that we have  $\frac{d\theta_X^*}{dX} \Big|_{X=t} = 0$  for any non-differentiable points  $t$  of  $\gamma$ . Then it implies that  $w = c\delta_{\theta_t^*}$  where  $\delta_s$  denotes the Kronecker delta measure. Then, according to (13), we have

$$\begin{aligned} 0 &= w_\theta \left[ (\theta - \theta_X^*) e^{-f_X - \gamma + m(f_X + \gamma)} \right], \\ &= c(\theta_t^* - \theta_X^*) e^{-f_X(\theta_t^*) - \gamma(\theta_t^*) + m(f_X + \gamma)}, \end{aligned}$$

which means that  $\theta_X^* = \theta_t^*$  is a constant independent of  $X$ . However, this contradicts to the assumption that  $\gamma$  is nontrivial. ■

As a remark, we note that this lemma is easily extended to multidimensional exponential family of distributions.