

---

## Robust Bayesian Recourse (Supplementary Material)

---

Tuan-Duy H. Nguyen<sup>1</sup>

Ngoc Bui<sup>1</sup>

Duy Nguyen<sup>1</sup>

Man-Chung Yue<sup>2</sup>

Viet Anh Nguyen<sup>1</sup>

<sup>1</sup>VinAI Research, Vietnam

<sup>2</sup>The University of Hong Kong

### A PROOFS OF SECTION 5

**Lemma 5.3 (re-stated).** *There exists a distribution  $\mathbb{Q}_0^*$  that solves (3) and is a mixture of at most  $N_0$  Gaussian components. Moreover, problem (3) is equivalent to a separable problem of the form*

$$\max \{L(x, \mathbb{Q}_0) : \mathbb{Q}_0 \in \mathbb{B}_{\varepsilon_0}(\widehat{\mathbb{P}}_0^\sigma)\} = \begin{cases} \max & \frac{1}{N_0} \sum_{i \in \mathcal{I}_0} f(x|\mu_i, \Sigma_i) \\ \text{s. t.} & (\mu_i, \Sigma_i) \in \mathbb{R}^p \times \mathbb{S}_{\geq \sigma}^p \\ & c((\mu_i, \Sigma_i), (\widehat{x}_i, \sigma^2 \bar{I})) \leq \varepsilon_0 \quad \forall i \in \mathcal{I}_0. \end{cases}$$

An analogous result holds for problem (4) with the corresponding subscript  $y = 1$ .

*Proof of Lemma 5.3.* There exists a distribution  $\mathbb{Q}_0^*$  that solves (3) and is a mixture of at most  $N_0$  Gaussian components. Moreover, problem (3) is equivalent to a separable problem of the form

$$\begin{aligned} & \max \{L(x, \mathbb{Q}_0) : \mathbb{Q}_0 \in \mathbb{B}_{\varepsilon_0}(\widehat{\mathbb{P}}_0^\sigma)\} \\ = & \begin{cases} \max & \frac{1}{N_0} \sum_{i \in \mathcal{I}_0} f(x|\mu_i, \Sigma_i) \\ \text{s. t.} & (\mu_i, \Sigma_i) \in \mathbb{R}^p \times \mathbb{S}_{\geq \sigma}^p \\ & c((\mu_i, \Sigma_i), (\widehat{x}_i, \sigma^2 \bar{I})) \leq \varepsilon_0 \quad \forall i \in \mathcal{I}_0. \end{cases} \end{aligned}$$

We use  $\forall i$  implies  $\forall i \in \mathcal{I}_0$ , and  $\sum_i$  is also taken over the same set. Given any  $x$ , the likelihood of  $x$  under any Gaussian mixture  $\mathbb{Q}_0$  can be written using the corresponding measure  $\nu_0$  as

$$L(x, \mathbb{Q}_0) = \int_{\mathbb{R}^p \times \mathbb{S}_+^p} f(x|\mu, \Sigma) \nu_0(d\mu, d\Sigma).$$

Recall that  $\Xi = \mathbb{R}^p \times \mathbb{S}_{\geq \sigma}^p$ . Using the definition of the type- $\infty$  Wasserstein, we find

$$\begin{aligned} & W_c(\nu_0, \widehat{\nu}_0) \leq \varepsilon_0 \\ \Leftrightarrow & \exists \lambda \in \Lambda(\nu_0, \widehat{\nu}_0) \text{ such that} \\ & \text{ess sup}_\lambda \{c((\mu, \Sigma), (\mu', \Sigma')) : (\mu, \Sigma, \mu', \Sigma') \in \Xi \times \Xi\} \leq \varepsilon_0 \\ \Leftrightarrow & \forall i \exists \lambda_i \in \mathcal{P}(\Xi) \text{ such that} \\ & \text{ess sup}_{\lambda_i} \{c((\mu, \Sigma), (\widehat{x}_i, \sigma I)) : (\mu, \Sigma) \in \Xi\} \leq \varepsilon_0 \\ \Leftrightarrow & \forall i \exists \lambda_i \in \mathcal{P}(\Xi) \text{ such that} \\ & c((\mu, \Sigma), (\widehat{x}_i, \sigma I)) \leq \varepsilon_0 \quad (\mu, \Sigma) \in \text{supp}(\lambda_i), \end{aligned}$$

where the second equivalence follows from that  $\hat{\nu}_0 = \frac{1}{N_0} \sum_i \delta_{(\hat{x}_i, \sigma^2 I)}$  and hence any  $\lambda \in \Lambda(\nu_0, \hat{\nu}_0)$  takes the form  $\frac{1}{N_0} \sum_i \lambda_i \otimes \delta_{(\hat{x}_i, \sigma^2 I)}$  for some probability measures  $\lambda_i \in \mathcal{P}(\Xi)$ , and the third equivalence follows from Lemma A.1. Hence, problem (3) is equivalent to

$$\begin{aligned} & \begin{cases} \max & \int_{\mathbb{R}^p \times \mathbb{S}_{\geq \sigma}^p} f(x|\mu, \Sigma) \nu_0(d\mu, d\Sigma) \\ \text{s. t.} & \nu_0 \in \mathcal{P}(\mathbb{R}^p \times \mathbb{S}_{\geq \sigma}^p) \\ & \mathbb{W}_c(\nu_0, \hat{\nu}_0) \leq \varepsilon_0 \end{cases} \\ = & \begin{cases} \max & \frac{1}{N_0} \sum_i \int_{\mathbb{R}^p \times \mathbb{S}_{\geq \sigma}^p} f(x|\mu_i, \Sigma_i) \lambda_i(d\mu_i, d\Sigma_i) \\ \text{s. t.} & \lambda_i \in \mathcal{P}(\mathbb{R}^p \times \mathbb{S}_{\geq \sigma}^p) \quad \forall i \\ & c((\mu_i, \Sigma_i), (\hat{x}_i, \sigma^2 I)) \leq \varepsilon_0 \quad \forall (\mu_i, \Sigma_i) \in \text{supp}(\lambda_i) \quad \forall i. \end{cases} \end{aligned}$$

It is easy now to employ a greedy argument to show that the optimal solution for  $\lambda_i$  should be a Dirac delta distribution supported on one point in the space of  $\mathbb{R}^p \times \mathbb{S}_{\geq \sigma}^p$ . This leads to the conclusion regarding the maximization problem (3).

An similar argument can be applied for the minimization problem (4), the detailed proof is omitted.  $\square$

**Lemma A.1.** *For any  $\lambda \in \mathcal{P}(\Xi)$ ,  $\hat{x} \in \mathbb{R}^p$ ,  $\sigma, \varepsilon > 0$  and any function  $c : \Xi \times \Xi \rightarrow \mathbb{R}$  such that the map  $(\mu, \Sigma) \mapsto c((\mu, \Sigma), (\hat{x}, \sigma^2 I))$  is continuous, we have  $\text{ess sup}_\lambda c((\mu, \Sigma), (\hat{x}, \sigma^2 I)) \leq \varepsilon$  if and only if  $c((\mu, \Sigma), (\hat{x}, \sigma^2 I)) \leq \varepsilon$  for any  $(\mu, \Sigma) \in \text{supp}(\lambda)$ .*

*Proof of Lemma A.1.* We first prove the ‘‘only if’’ direction. Suppose that there exists  $(\mu', \Sigma') \in \text{supp}(\lambda)$  such that

$$c((\mu', \Sigma'), (\hat{x}, \sigma^2 I)) > \varepsilon.$$

By continuity of the map  $(\mu, \Sigma) \mapsto c((\mu, \Sigma), (\hat{x}, \sigma^2 I))$ , there exists an open neighbourhood  $U \subseteq \Xi$  containing  $(\mu', \Sigma')$  such that

$$c((\mu, \Sigma), (\hat{x}, \sigma^2 I)) > \varepsilon \quad \forall (\mu, \Sigma) \in U.$$

By the definition of support,  $\lambda(U) > 0$ . Therefore,

$$\Pr_\lambda(c((\mu, \Sigma), (\hat{x}, \sigma^2 I)) \leq \varepsilon) = 1 - \Pr_\lambda(c((\mu, \Sigma), (\hat{x}, \sigma^2 I)) > \varepsilon) \leq 1 - \lambda(U) < 1,$$

which contradicts to that  $\text{ess sup}_\lambda c((\mu, \Sigma), (\hat{x}, \sigma^2 I)) \leq \varepsilon$ .

We next prove the ‘‘if’’ direction. By the law of total probability and the fact that  $c((\mu, \Sigma), (\hat{x}, \sigma^2 I)) \leq \varepsilon$  for any  $(\mu, \Sigma) \in \text{supp}(\lambda)$ ,

$$\begin{aligned} & \Pr_\lambda(c((\mu, \Sigma), (\hat{x}, \sigma^2 I)) \leq \varepsilon) \\ = & \Pr_\lambda(c((\mu, \Sigma), (\hat{x}, \sigma^2 I)) \leq \varepsilon | (\mu, \Sigma) \in \text{supp}(\lambda)) \lambda(\text{supp}(\lambda)) \\ & + \Pr_\lambda(c((\mu, \Sigma), (\hat{x}, \sigma^2 I)) \leq \varepsilon | (\mu, \Sigma) \notin \text{supp}(\lambda)) (1 - \lambda(\text{supp}(\lambda))) \\ = & 1 \cdot 1 + \Pr_\lambda(c((\mu, \Sigma), (\hat{x}, \sigma^2 I)) \leq \varepsilon | (\mu, \Sigma) \notin \text{supp}(\lambda)) \cdot 0 = 1, \end{aligned}$$

which completes the proof.  $\square$

**Proposition 5.4 (re-stated).** *Fix any index  $i \in \mathcal{I}_0$ . For any  $\hat{x}_i \in \mathbb{R}^p$ ,  $x \in \mathbb{R}^p$  and  $\varepsilon_0 \in \mathbb{R}_+$ , we have*

$$\frac{\exp(-\alpha_i)}{(2\pi)^{p/2}} = \begin{cases} \max & f(x|\mu_i, \Sigma_i) \\ \text{s. t.} & (\mu_i, \Sigma_i) \in \mathbb{R}^p \times \mathbb{S}_{\geq \sigma}^p \\ & c((\mu_i, \Sigma_i), (\hat{x}_i, \sigma^2 I)) \leq \varepsilon_0, \end{cases}$$

where  $\alpha_i$  is the optimal value of the two-dimensional optimization problem

$$\min_{\substack{a \in \mathbb{R}_+, d_p \in [\sigma, +\infty) \\ a^2 + (d_p - \sigma)^2 \leq \varepsilon_0^2}} \log d_p + \frac{(\|x - \hat{x}_i\|_2 - a)^2}{2d_p^2} + (p-1) \log \sigma.$$

*Proof of Proposition 5.4.* Let  $\alpha_i$  be the optimal value of the negative log-likelihood minimization problem

$$\alpha_i = \begin{cases} \min & \frac{1}{2} \log \det \Sigma_i + \frac{1}{2} (x - \mu_i)^\top \Sigma_i^{-1} (x - \mu_i) \\ \text{s. t.} & \mu_i \in \mathbb{R}^p, \Sigma_i \in \mathbb{S}_+^p \\ & \|\mu_i - \hat{x}_i\|_2^2 + \text{Tr} [\Sigma_i + \sigma^2 I - 2((\sigma^2 I)^{\frac{1}{2}} \Sigma_i (\sigma^2 I)^{\frac{1}{2}})^{\frac{1}{2}}] \leq \varepsilon_0^2 \\ & \Sigma_i \succeq \sigma^2 I. \end{cases}$$

It is easy to see that

$$\max\{f(x|\mu_i, \Sigma_i) : (\mu_i, \Sigma_i) \in \mathbb{R}^p \times \mathbb{S}_{\geq \sigma}^p, c((\mu_i, \Sigma_i), (\hat{x}_i, \sigma^2 I)) \leq \varepsilon_0\} = \frac{1}{\sqrt{(2\pi)^p}} \exp(-\alpha_i).$$

It remains to provide a simpler formulation to determine  $\alpha_i$ . To simplify the notation, we omit the index  $i$  on all variables and parameters. We reparameterize  $\Sigma = V \text{diag}(d^2) V^\top$  for a vector  $d \in \mathbb{R}_+^p$ , where  $\text{diag}(d^2)$  denotes a  $\mathbb{R}^{p \times p}$  diagonal matrix with its  $j$ -th diagonal entries equals to  $d_j^2$ , and  $\text{O}(p)$  is the set of  $p$ -dimensional orthogonal matrices

$$\text{O}(p) = \{V \in \mathbb{R}^{p \times p} : V^\top V = I_p\}.$$

The negative log-likelihood minimization problem is further equivalent to

$$\begin{aligned} \min & \sum_{j=1}^p \log d_j + \frac{1}{2} (V^\top (x - \mu))^\top \text{diag}(d^{-2}) (V^\top (x - \mu)) \\ \text{s. t.} & d \in \mathbb{R}_+^p, V \in \text{O}(p), \mu \in \mathbb{R}^p \\ & \|\mu - \hat{x}\|_2^2 + \sum_{j=1}^p (d_j - \sigma)^2 \leq \varepsilon_0^2 \\ & d \geq \sigma, \end{aligned}$$

where  $d \geq \sigma$  implies the element-wise constraints  $d_j \geq \sigma$  for any  $j = 1, \dots, p$ . We introduce an auxiliary variable  $a \in \mathbb{R}_+$  and rewrite the optimization problem in an equivalent way as

$$\min_{\substack{a \in \mathbb{R}_+, d \in \mathbb{R}_+^p, d \geq \sigma \\ a^2 + \sum_j (d_j - \sigma)^2 \leq \varepsilon_0^2}} \min_{\substack{\mu \in \mathbb{R}^p \\ \|\mu - \hat{x}\|_2^2 = a^2}} \min_{V \in \text{O}(p)} \sum_{j=1}^p \log d_j + \frac{1}{2} (V^\top (x - \mu))^\top \text{diag}(d^{-2}) (V^\top (x - \mu)).$$

Notice that the above optimization problem is invariant to the ordering of the entries of  $d$ . As a consequence, without any loss of generality, we can assume that  $d_p$  is the maximum value across all  $d_j$ . By Lemma B.1, the above optimization problem becomes

$$\min_{\substack{a \in \mathbb{R}_+, d \in \mathbb{R}_+^p, d \geq \sigma \\ a^2 + \sum_j (d_j - \sigma)^2 \leq \varepsilon_0^2 \\ d_p = \max\{d\}}} \min_{\mu \in \mathbb{R}^p} \sum_{j=1}^p \log d_j + \frac{1}{2d_p^2} \|x - \mu\|_2^2.$$

Using Lemma B.2, we obtain the equivalent optimization problem

$$\min_{\substack{a \in \mathbb{R}_+, d_p \in \mathbb{R}_+, d_p \geq \sigma \\ a^2 + \sum_j (d_j - \sigma)^2 \leq \varepsilon_0^2 \\ d_p = \max\{d\}}} \sum_{j=1}^p \log d_j + \frac{1}{2d_p^2} (\|x - \hat{x}\|_2 - a)^2.$$

Rewriting the above problem into a two-layer optimization problem

$$\min_{\substack{a \in \mathbb{R}_+, d_p \in \mathbb{R}_+, d_p \geq \sigma \\ a^2 + (d_p - \sigma)^2 \leq \varepsilon_0^2}} \left\{ \log d_p + \frac{1}{2d_p^2} (\|x - \hat{x}\|_2 - a)^2 + \min_{\substack{d_j \in \mathbb{R}_+, d_j \geq \sigma \forall j=1, \dots, p-1 \\ \sum_{j=1}^{p-1} (d_j - \sigma)^2 \leq \varepsilon_0^2 - a^2 - (d_p - \sigma)^2 \\ d_j \leq d_p \forall j=1, \dots, p-1}} \sum_{j=1}^{p-1} \log d_j \right\}. \quad (1)$$

Notice that for any  $d_p$  that is feasible for the outer minimization problem, the inner minimization problem over  $d_j$ ,  $\forall j = 1, \dots, p-1$  admits a non-empty feasible set. Indeed, because  $d_p \geq \sigma$ , the value  $d_j = \sigma$ ,  $j = 1, \dots, p-1$  is a

feasible solution for the inner problem. We now focus on solving the inner minimization problem. As  $\log(\cdot)$  is an increasing function, for any  $s \geq 0$ , we find

$$\min_{\substack{d_p \geq d_j \geq \sigma \forall j=1, \dots, p-1 \\ \sum_{j=1}^{p-1} (d_j - \sigma)^2 \leq s}} \sum_{j=1}^{p-1} \log d_j = (p-1) \log \sigma,$$

which holds because the optimization problem on the left hand side admits the optimal solution  $d_j^* = \sigma$  for all  $j = 1, \dots, p-1$ . This completes the proof.  $\square$

**Proposition 5.5 (re-stated).** Fix any index  $i \in \mathcal{I}_1$ . For any  $\hat{x}_i \in \mathbb{R}^p$ ,  $x \in \mathbb{R}^p$  and  $\varepsilon_1 \in \mathbb{R}_+$ , we have

$$\frac{\exp(\alpha_i)}{(2\pi)^{p/2}} = \begin{cases} \min & f(x|\mu_i, \Sigma_i) \\ \text{s. t.} & (\mu_i, \Sigma_i) \in \mathbb{R}^p \times \mathbb{S}_{\geq \sigma}^p \\ & c((\mu_i, \Sigma_i), (\hat{x}_i, \sigma^2 I)) \leq \varepsilon_1, \end{cases}$$

where  $\alpha_i$  is the optimal value of the two-dimensional optimization problem

$$\min_{\substack{a \in \mathbb{R}_+, d_1 \in [\sigma, +\infty) \\ a^2 + p(d_1 - \sigma)^2 \leq \varepsilon_1^2}} \left\{ -\log d_1 - \frac{(\|x - \hat{x}_i\|_2 + a)^2}{2d_1^2} - (p-1) \log \left( \sigma + \sqrt{\frac{\varepsilon^2 - a^2 - (d_1 - \sigma)^2}{p-1}} \right) \right\}.$$

*Proof of Proposition 5.5.* Let  $\alpha_i$  be the optimal value of the log-likelihood minimization problem

$$\alpha_i = \begin{cases} \min & -\frac{1}{2} \log \det \Sigma_i - \frac{1}{2} (x - \mu_i)^\top \Sigma_i^{-1} (x - \mu_i) \\ \text{s. t.} & \mu_i \in \mathbb{R}^p, \Sigma_i \in \mathbb{S}_+^p \\ & \|\mu_i - \hat{x}_i\|_2^2 + \text{Tr} [\Sigma_i + \sigma^2 I - 2((\sigma^2 I)^{\frac{1}{2}} \Sigma_i (\sigma^2 I)^{\frac{1}{2}})^{\frac{1}{2}}] \leq \varepsilon_1^2 \\ & \Sigma_i \succeq \sigma^2 I. \end{cases}$$

It is easy to see that

$$\min \{f(x|\mu_i, \Sigma_i) : (\mu_i, \Sigma_i) \in \mathbb{R}^p \times \mathbb{S}_{\geq \sigma}^p, c((\mu_i, \Sigma_i), (\hat{x}_i, \sigma^2 I)) \leq \varepsilon_1\} = \frac{1}{(2\pi)^{p/2}} \exp(\alpha_i).$$

It remains to provide the computational routine to determine  $\alpha_i$ . To simplify the notation, we omit the index  $i$  on all variables and parameters. We reparameterize  $\Sigma = V \text{diag}(d^2) V^\top$  for a vector  $d \in \mathbb{R}_+^p$ , where  $\text{diag}(d^2)$  denotes a  $\mathbb{R}^{p \times p}$  diagonal matrix with its  $j$ -th diagonal entries equals to  $d_j^2$ , and  $O(p)$  is the set of  $p$ -dimensional orthogonal matrices

$$O(p) = \{V \in \mathbb{R}^{p \times p} : V^\top V = I_p\}.$$

The log-likelihood minimization problem is further equivalent to

$$\begin{aligned} \min & -\sum_{j=1}^p \log d_j - \frac{1}{2} (V^\top (x - \mu))^\top \text{diag}(d^{-2}) (V^\top (x - \mu)) \\ \text{s. t.} & d \in \mathbb{R}_+^p, V \in O(p), \mu \in \mathbb{R}^p \\ & \|\mu - \hat{x}\|_2^2 + \sum_{j=1}^p (d_j - \sigma)^2 \leq \varepsilon_1^2 \\ & d \geq \sigma, \end{aligned}$$

where  $d \geq \sigma$  implies the element-wise constraints  $d_j \geq \sigma$  for any  $j = 1, \dots, p$ . We introduce an auxiliary variable  $a \in \mathbb{R}_+$  and rewrite the optimization problem in an equivalent way as

$$\min_{\substack{a \in \mathbb{R}_+, d \in \mathbb{R}_+^p, d \geq \sigma \\ a^2 + \sum_j (d_j - \sigma)^2 \leq \varepsilon_1^2}} \min_{\substack{\mu \in \mathbb{R}^p \\ \|\mu - \hat{x}\|_2^2 = a^2}} \min_{V \in O(p)} -\sum_{j=1}^p \log d_j - \frac{1}{2} (V^\top (x - \mu))^\top \text{diag}(d^{-2}) (V^\top (x - \mu)).$$

Notice that the above optimization problem is invariant to the ordering of the entries of  $d$ . As a consequence, without any loss of generality, we can assume that  $d_1$  is the minimum value across all  $d_j$ . By Lemma B.1, the above optimization problem becomes

$$\min_{\substack{a \in \mathbb{R}_+, d \in \mathbb{R}_+^p, d \geq \sigma \\ a^2 + \sum_j (d_j - \sigma)^2 \leq \varepsilon_1^2 \\ d_1 = \min\{d\}}} \min_{\mu \in \mathbb{R}^p} -\sum_{j=1}^p \log d_j - \frac{1}{2d_1^2} \|x - \mu\|_2^2.$$

Using Lemma B.2, we obtain the equivalent optimization problem

$$\min_{\substack{a \in \mathbb{R}_+, d \in \mathbb{R}_+^p, d \geq \sigma \\ a^2 + \sum_j (d_j - \sigma)^2 \leq \varepsilon_1^2 \\ d_1 = \min\{d\}}} - \sum_{j=1}^p \log d_j - \frac{1}{2d_1^2} (\|x - \hat{x}\|_2 + a)^2.$$

Notice that the constraint  $\sigma \leq d_1 = \min\{d\}$  implies that  $p(d_1 - \sigma)^2 \leq \sum_j (d_j - \sigma)^2$ . As a consequence, any feasible value for  $d_1$  should satisfy  $a^2 + p(d_1 - \sigma)^2 \leq \varepsilon_1^2$ . Separating the variable  $d$  into two groups  $d_1$  and  $d_2, \dots, d_p$  leads to a two-layer optimization problem

$$\min_{\substack{a \in \mathbb{R}_+, d_1 \in \mathbb{R}_+, d_1 \geq \sigma \\ a^2 + p(d_1 - \sigma)^2 \leq \varepsilon_1^2}} \left\{ -\log d_1 - \frac{1}{2d_1^2} (\|x - \hat{x}\|_2 + a)^2 + \min_{\substack{d_j \in \mathbb{R}_+, d_j \geq d_1 \forall j=2, \dots, p \\ \sum_{j=2}^p (d_j - \sigma)^2 \leq \varepsilon_1^2 - a^2 - (d_1 - \sigma)^2}} - \sum_{j=2}^p \log d_j \right\}. \quad (2)$$

Consider momentarily the minimization problem

$$\min_{\substack{d_j \in \mathbb{R}_+ \\ \sum_{j=2}^p (d_j - \sigma)^2 \leq \varepsilon_1^2 - a^2 - (d_1 - \sigma)^2}} - \sum_{j=2}^p \log d_j,$$

where the constraints  $d_j \geq d_1$  have been intentionally omitted. Proposition B.3 asserts that this optimization problem has the optimal value

$$-(p-1) \log \left( \sigma + \sqrt{\frac{\varepsilon_1^2 - a^2 - (d_1 - \sigma)^2}{p-1}} \right)$$

at the optimal solution  $d_j^* = \sigma + \sqrt{\frac{\varepsilon_1^2 - a^2 - (d_1 - \sigma)^2}{p-1}}$ , which also by the outer constraint  $a^2 + p(d_1 - \sigma)^2 \leq \varepsilon_1^2$  satisfies  $d_j \geq d_1 \forall j = 2, \dots, p$ . Thus it is indeed the optimal solution to the inner minimization problem in (2). As a consequence, problem (2) is equivalent to

$$\min_{\substack{a \in \mathbb{R}_+, d_1 \in \mathbb{R}_+, d_1 \geq \sigma \\ a^2 + p(d_1 - \sigma)^2 \leq \varepsilon_1^2}} \left\{ -\log d_1 - \frac{1}{2d_1^2} (\|x - \hat{x}\|_2 + a)^2 - (p-1) \log \left( \sigma + \sqrt{\frac{\varepsilon_1^2 - a^2 - (d_1 - \sigma)^2}{p-1}} \right) \right\}.$$

This completes the proof.  $\square$

## B AUXILIARY RESULTS

The following preparatory results are necessary to prove Propositions 5.4 and 5.5.

**Lemma B.1** (Eigenbasis solution). *Let  $E \in \mathbb{R}^{p \times p}$  be a diagonal matrix satisfying  $E_{11} \leq \dots \leq E_{pp}$ . Then, for any  $w \in \mathbb{R}^p$ , we have*

$$\max_{V \in O(p)} w^\top V E V^\top w = E_{pp} \|w\|_2^2.$$

*Proof of Lemma B.1.* The claim holds trivially when  $w = 0$ . Consider now any  $w \in \mathbb{R}^p \setminus \{0\}$ . Since  $V E V^\top \preceq E_{pp} \cdot I_p$ , we find

$$\max_{V \in O(p)} w^\top V E V^\top w \leq \max_{V \in O(p)} w^\top V (E_{pp} \cdot I_p) V^\top w = E_{pp} \|w\|_2^2.$$

On the other hand, taking  $V^* = [v_1^*, \dots, v_p^*] \in O(p)$  with  $v_p^* = \frac{w}{\|w\|_2}$ , and using the orthogonality of the columns of  $V^*$ , we have

$$w^\top V^* E V^{*\top} w = E_{pp} \|w\|_2^2.$$

This shows that  $V^*$  is an optimal solution and completes the proof.  $\square$

**Lemma B.2** (Quadratic optimization). *For any  $x \in \mathbb{R}^p$ ,  $\hat{x} \in \mathbb{R}^p$  and  $a \in \mathbb{R}_+$ , the following assertions hold.*

- *Convex quadratic minimization:*

$$\min_{\mu \in \mathbb{R}^p: \|\mu - \hat{x}\|_2^2 = a^2} \|x - \mu\|_2^2 = (\|x - \hat{x}\|_2 - a)^2,$$

where the minimum is attained at  $\mu^* = \frac{a}{\|x - \hat{x}\|_2} x + (1 - \frac{a}{\|x - \hat{x}\|_2}) \hat{x}$ .

- *Convex quadratic maximization:*

$$\max_{\mu \in \mathbb{R}^p: \|\mu - \hat{x}\|_2^2 = a^2} \|x - \mu\|_2^2 = (\|x - \hat{x}\|_2 + a)^2,$$

where the maximum is attained at  $\mu^* = -\frac{a}{\|x - \hat{x}\|_2} x + (1 + \frac{a}{\|x - \hat{x}\|_2}) \hat{x}$ .

The results in Lemma B.2 are dispersed in the literature. An elementary proof is provided here for completeness.

*Proof of Lemma B.2.* By the triangle inequality, for any  $\mu$  such that  $\|\mu - \hat{x}\|_2 = a$ , we have

$$\|x - \mu\|_2 \geq \| \|x - \hat{x}\| - \|\mu - \hat{x}\| \| = \| \|x - \hat{x}\| - a \|,$$

where the lower bound can be attained by taking  $\mu = \frac{a}{\|x - \hat{x}\|_2} x + (1 - \frac{a}{\|x - \hat{x}\|_2}) \hat{x}$ . Therefore,

$$\min_{\mu \in \mathbb{R}^p: \|\mu - \hat{x}\|_2^2 = a^2} \|x - \mu\|_2^2 = (\|x - \hat{x}\|_2 - a)^2$$

Similarly, by the triangle inequality we have

$$\|x - \mu\|_2 \leq \| \|x - \hat{x}\| + \|\hat{x} - \mu\| \| = \|x - \hat{x}\| + a,$$

and the upper bound can be attained by  $\mu = -\frac{a}{\|x - \hat{x}\|_2} x + (1 + \frac{a}{\|x - \hat{x}\|_2}) \hat{x}$ . This completes the proof.  $\square$

**Proposition B.3** (Logarithm maximization). *For any  $s, \sigma \geq 0$  and positive integer  $k$ , we have*

$$k \log \left( \sqrt{\frac{s}{k}} + \sigma \right) = \begin{cases} \max_{e \in \mathbb{R}_+^k} \sum_{j=1}^k \log e_j \\ \text{s. t. } \sum_{j=1}^k (\sigma - e_j)^2 \leq s. \end{cases} \quad (3)$$

Moreover, the optimal solution  $e^*$  satisfies  $e_j^* = \sqrt{\frac{s}{k}} + \sigma$  for any  $j = 1, \dots, k$ .

*Proof of Proposition B.3.* Let  $e^* \in \mathbb{R}_+^k$  be an optimal solution to the maximization problem (3). Suppose there exist two indices  $m$  and  $n$  such that  $e_m^* \neq e_n^*$ . Consider  $e'$  defined by

$$e'_j = \begin{cases} \frac{1}{2}(e_m^* + e_n^*), & \text{if } j \in \{m, n\}, \\ e_j^*, & \text{otherwise.} \end{cases}$$

By the convexity of the function  $x \mapsto (x - \sigma)^2$ ,

$$(e'_m - \sigma)^2 + (e'_n - \sigma)^2 = 2 \left( \frac{e_m^* + e_n^*}{2} - \sigma \right)^2 \leq (e_m^* - \sigma)^2 + (e_n^* - \sigma)^2,$$

which implies that  $e'$  is a feasible solution to problem (3). Furthermore, since  $e_m^* \neq e_n^*$ , by the concavity of the function  $x \mapsto \log x$ , we have that

$$\log e_m^* + \log e_n^* < 2 \log \left( \frac{e_m^* + e_n^*}{2} \right) = \log e'_m + \log e'_n,$$

which violates the optimality of  $e^*$ . Therefore, any optimal solution  $e^*$  must have all entries identical. Using this, we get from the constraint that

$$|e_j^* - \sigma| \leq \sqrt{\frac{s}{k}} \quad \forall j = 1, \dots, k.$$

By continuity of the objective and constraint functions, we must have

$$|e_j^* - \sigma| = \sqrt{\frac{s}{k}} \quad \forall j = 1, \dots, k.$$

Since the objective function is increasing in  $e_j^*$ , the optimal solution is given by

$$e_j^* = \sigma + \sqrt{\frac{s}{k}} \quad \forall j = 1, \dots, k.$$

The optimal value can then be obtained by direct computation. This completes the proof.  $\square$

## C FIRST-ORDER ALGORITHMS

### C.1 OPTIMISTIC LIKELIHOOD PROBLEM

For the optimistic likelihood problem, Theorem 5.1 reduces the task to solving the 2-dimensional problem

$$\min_{\substack{a \in \mathbb{R}_+, d_p \in [\sigma, +\infty) \\ a^2 + (d_p - \sigma)^2 \leq \varepsilon_0^2}} \log d_p + \frac{(\|x - \hat{x}_i\|_2 - a)^2}{2d_p^2} + (p-1) \log \sigma.$$

By letting

$$d_p = v_2 + \sigma, \quad \text{and} \quad a = v_1,$$

we can obtain the equivalent form

$$\min_{\substack{v_1, v_2 \geq 0 \\ v_1^2 + v_2^2 \leq \varepsilon_0^2}} F(v), \tag{4}$$

where the objective function is given by

$$F(v) = \log(v_2 + \sigma) + \frac{(\|x - \hat{x}_i\|_2 - v_1)^2}{2(v_2 + \sigma)^2} + (p-1) \log \sigma.$$

If we denote by  $\mathcal{V} = \{v \in \mathbb{R}^2 : v_1, v_2 \geq 0, v_1^2 + v_2^2 \leq \varepsilon_0^2\}$  the feasible region of the above minimization problem, then the projection  $\text{Proj}_{\mathcal{V}}(v)$  can be computed in closed-form via

$$\text{Proj}_{\mathcal{V}}(v) = \begin{cases} v, & \text{if } v_1, v_2 \geq 0, v_1^2 + v_2^2 \leq \varepsilon_0^2, \\ \frac{\varepsilon_0}{\|v\|_2} v, & \text{if } v_1, v_2 \geq 0, v_1^2 + v_2^2 > \varepsilon_0^2, \\ (0, \varepsilon_0)^\top, & \text{if } v_1 < 0, v_2 > \varepsilon_0, \\ (0, v_2)^\top, & \text{if } v_1 < 0, 0 \leq v_2 \leq \varepsilon_0, \\ (\varepsilon_0, 0)^\top, & \text{if } v_1 > \varepsilon_0, v_2 < 0, \\ (v_1, 0)^\top, & \text{if } 0 \leq v_1 \leq \varepsilon_0, v_2 < 0, \\ (0, 0)^\top, & \text{if } v_1, v_2 < 0. \end{cases}$$

Algorithm 1 is a projected gradient descent routine to solve problem (4). The convergence guarantee for Algorithm 1 follows from Beck [2017, Theorem 10.15].

---

**Algorithm 1** Projected Gradient Descent Algorithm with Backtracking Line-Search

---

**Algorithm parameters:** Line search parameters  $\theta \in (0, 1), \beta > 0$

**Initialization:** Set  $v^0 \leftarrow 0$

**for**  $t = 0, 1, \dots$  **do**

Find the smallest integer  $k \geq 0$  such that

$$F(\text{Proj}_{\mathcal{V}}(v^t - \theta^k \beta \nabla F(v^t))) \leq F(v^t) - \frac{1}{2\theta^k \beta} \|v^t - \text{Proj}_{\mathcal{V}}(v^t - \theta^k \beta \nabla F(v^t))\|_2^2$$

Set  $s^t = \theta^k \beta$  and set  $v^{t+1} = \text{Proj}_{\mathcal{V}}(v^t - s^t \nabla F(v^t))$ .

**end for**

---

## C.2 PESSIMISTIC LIKELIHOOD PROBLEM

For the pessimistic likelihood problem, Theorem 5.2 reduces the task to solving the 2-dimensional problem

$$\min_{\substack{a \in \mathbb{R}_+, d_1 \in [\sigma, +\infty) \\ a^2 + p(d_1 - \sigma)^2 \leq \varepsilon_1^2}} \left\{ -\log d_1 - \frac{1}{2d_1^2} (\|x - \hat{x}_i\|_2 + a)^2 - (p-1) \log \left( \sigma + \sqrt{\frac{\varepsilon_1^2 - a^2 - (d_1 - \sigma)^2}{p-1}} \right) \right\}.$$

Note that the gradient of the objective function is a non-Lipschitz function. Worse still, the gradient is even undefined on at the feasible point  $(d_1, a) = (\sigma, \varepsilon_1)$ . These properties induce numerical issues for the optimization algorithm. Therefore, we solve the following perturbed problem instead:

$$\min_{\substack{a \in \mathbb{R}_+, d_1 \in [\sigma, +\infty) \\ a^2 + p(d_1 - \sigma)^2 \leq \varepsilon_1^2}} \left\{ -\log d_1 - \frac{1}{2d_1^2} (\|x - \hat{\mu}\|_2 + a)^2 - (p-1) \log \left( \sigma + \sqrt{\frac{\zeta + \varepsilon_1^2 - a^2 - (d_1 - \sigma)^2}{p-1}} \right) \right\}, \quad (5)$$

for some small  $\zeta > 0$ . By Bonnans and Shapiro [2013, Proposition 4.4], the optimal value of problem (5) is continuous in  $\zeta$  and the optimal solution set is upper semi-continuous in  $\zeta$  as a set-valued mapping, see Bonnans and Shapiro [2013, Section 4.1].

We now derive a projected gradient descent algorithm with backtracking line search for solving problem (5). First, by letting

$$d_1 = u_2 + \sigma, \quad \text{and} \quad a = \sqrt{p}u_1,$$

we can equivalently transform problem (5) to the following one:

$$\min_{\substack{u_1, u_2 \geq 0 \\ u_1^2 + u_2^2 \leq (\varepsilon_1/\sqrt{p})^2}} F(u), \quad (6)$$

where the objective function is given by

$$F(u) = -\log(u_2 + \sigma) - \frac{1}{2(u_2 + \sigma)^2} (\|x - \hat{x}_i\|_2 + \sqrt{p}u_1)^2 - (p-1) \log \left( \sigma + \sqrt{\frac{\zeta + \varepsilon_1^2 - pu_1^2 - u_2^2}{p-1}} \right).$$

The upshot of problem (6) is that the feasible region is the intersection of the non-negative orthant with a circular disk of radius  $\varepsilon_1/\sqrt{p}$  centered at the origin. As we will see below, this enables easy computation of the projection and linear optimization oracle. Indeed, denoting by  $\mathcal{U} = \{u \in \mathbb{R}^2 : u_1, u_2 \geq 0, u_1^2 + u_2^2 \leq (\varepsilon_1/\sqrt{p})^2\}$  the feasible region of problem (6), the projection  $\text{Proj}_{\mathcal{U}}(u)$  can be computed in closed-form via

$$\text{Proj}_{\mathcal{U}}(u) = \begin{cases} u, & \text{if } u_1, u_2 \geq 0, u_1^2 + u_2^2 \leq (\varepsilon_1/\sqrt{p})^2, \\ \frac{(\varepsilon_1/\sqrt{p})}{\|u\|_2} u, & \text{if } u_1, u_2 \geq 0, u_1^2 + u_2^2 > (\varepsilon_1/\sqrt{p})^2, \\ (0, \frac{\varepsilon_1}{\sqrt{p}})^\top, & \text{if } u_1 < 0, u_2 > \frac{\varepsilon_1}{\sqrt{p}}, \\ (0, u_2)^\top, & \text{if } u_1 < 0, 0 \leq u_2 \leq \frac{\varepsilon_1}{\sqrt{p}}, \\ (\frac{\varepsilon_1}{\sqrt{p}}, 0)^\top, & \text{if } u_1 > \frac{\varepsilon_1}{\sqrt{p}}, u_2 < 0, \\ (u_1, 0)^\top, & \text{if } 0 \leq u_1 \leq \frac{\varepsilon_1}{\sqrt{p}}, u_2 < 0, \\ (0, 0)^\top, & \text{if } u_1, u_2 < 0. \end{cases}$$



A projected gradient descent algorithm can now be employed to solve problem (6).

## D RECOVERY OF THE ADVERSARIAL DISTRIBUTION

It is often instructive to recover and analyze the optimal distribution that maximizes the posterior probability odds ratio, or more directly, the likelihood ratio in (2). Equivalent, it suffices to characterize the distribution  $\mathbb{Q}_0^*$  that maximizes (3), and the distribution  $\mathbb{Q}_1^*$  that minimizes (4).

**Lemma D.1** (Likelihood maximizer). *For each  $i \in \mathcal{I}_0$ , let  $(a_i^*, d_{pi}^*)$  be the optimal solution of the following two-dimensional optimization problem*

$$\min_{\substack{a \in \mathbb{R}_+, d_p \in [\sigma, +\infty) \\ a^2 + (d_p - \sigma)^2 \leq \varepsilon_0^2}} \log d_p + \frac{(\|x - \hat{x}_i\|_2 - a)^2}{2d_p^2} + (p-1) \log \sigma.$$

*Then, the maximizer  $\mathbb{Q}_0^*$  of problem (3) is a Gaussian mixture with  $N_0$  components, and for  $i \in \mathcal{I}_0$ , the  $i$ -th components has mean*

$$\mu_i^* = \frac{a_i^*}{\|x - \hat{x}_i\|_2} x + \left(1 - \frac{a_i^*}{\|x - \hat{x}_i\|_2}\right) \hat{x}_i,$$

*and covariance matrix*

$$\Sigma_i^* = V_i^* \text{diag}(\sigma, \dots, \sigma, d_{pi}^*)^2 (V_i^*)^\top,$$

*where  $V_i^*$  is any orthogonal matrix with the  $p$ -th column given by  $\frac{x - \mu_i^*}{\|x - \mu_i^*\|_2}$ .*

*Proof of Lemma D.1.* The result follows directly by inspecting the proofs of Proposition 5.4, Lemma B.1 and Lemma B.2.  $\square$

**Lemma D.2** (Likelihood minimizer). *For each  $i \in \mathcal{I}_1$ , let  $(a_i^*, d_{1i}^*)$  be the optimal solution of the following two-dimensional optimization problem*

$$\min_{\substack{a \in \mathbb{R}_+, d_1 \in [\sigma, +\infty) \\ a^2 + p(d_1 - \sigma)^2 \leq \varepsilon_1^2}} \left\{ -\log d_1 - \frac{(\|x - \hat{x}_i\|_2 + a)^2}{2d_1^2} - (p-1) \log \left( \sigma + \sqrt{\frac{\varepsilon_1^2 - a^2 - (d_1 - \sigma)^2}{p-1}} \right) \right\}.$$

*Then, the minimizer  $\mathbb{Q}_1^*$  of problem (4) is a Gaussian mixture with  $N_1$  components, and for  $i \in \mathcal{I}_1$ , the  $i$ -th components has mean*

$$\mu_i^* = -\frac{a_i^*}{\|x - \hat{x}_i\|_2} x + \left(1 + \frac{a_i^*}{\|x - \hat{x}_i\|_2}\right) \hat{x}_i,$$

*and covariance matrix*

$$\Sigma_i^* = V_i^* \text{diag} \left( d_{1i}^*, \sigma + \sqrt{\frac{\varepsilon_1^2 - a_i^{*2} - (d_{1i}^* - \sigma)^2}{p-1}}, \dots, \sigma + \sqrt{\frac{\varepsilon_1^2 - a_i^{*2} - (d_{1i}^* - \sigma)^2}{p-1}} \right)^2 (V_i^*)^\top,$$

*where  $V_i^*$  is any orthogonal matrix with the 1st column given by  $\frac{x - \mu_i^*}{\|x - \mu_i^*\|_2}$ .*

*Proof of Lemma D.2.* The result follows directly by inspecting the proofs of Proposition 5.5, Lemma B.1 and Lemma B.2.  $\square$

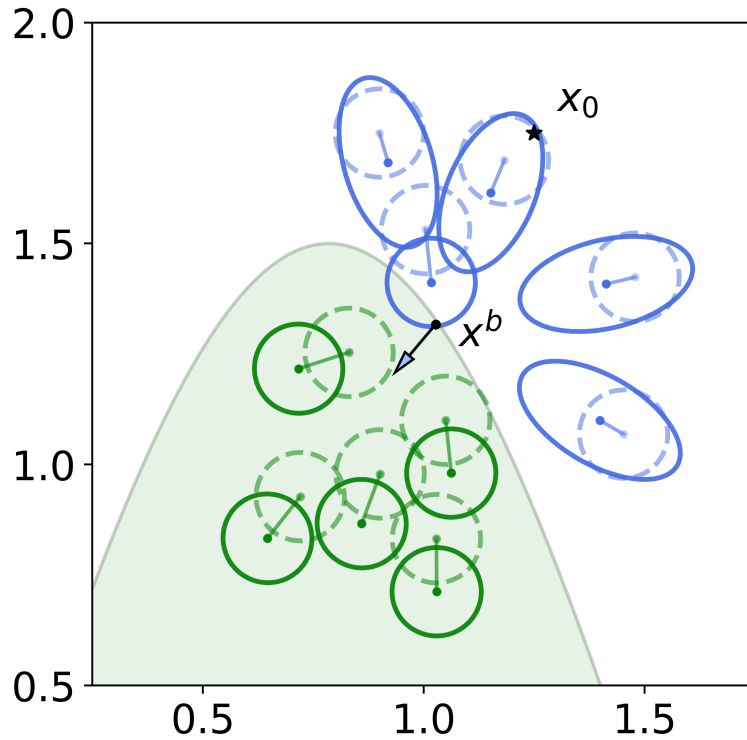


Figure 1: Visualization of the worst-case distributions on a toy dataset, color codes are similar to Figure 1. The dashed, opaque dots and circles represent the isotropic Gaussian around each data sample. The solid dots and circles represent the worst-case distributions corresponding to the boundary point  $x^b$ . For blue (unfavorably predicted) samples, the worst-case distribution is formed by perturbing the distribution towards  $x^b$  – which leads to maximizing the posterior probability of unfavorable prediction. For green (favorably predicted) samples, the worst-case distribution is formed by perturbing the distribution away from  $x^b$  – which leads to minimizing the posterior probability of favorable prediction. These worst-case distributions will maximize the posterior probability odds ratio.

## References

Amir Beck. *First-order Methods in Optimization*. SIAM, 2017.

J Frédéric Bonnans and Alexander Shapiro. *Perturbation Analysis of Optimization Problems*. Springer Science & Business Media, 2013.