

# Supplementary materials for “Optimizing Information-theoretical Generalization Bound via Anisotropic Noise in SGLD”

The supplementary materials are organized as follows. In Appendix A, we provide some basic lemmas which are used throughout the proofs in the rest of the materials. In Appendix B, we provide the proof of Lemma 2. In Appendix C, D, and E, we provide the detailed proofs of Lemmas and Theorems respectively in Section 3.2, Section 4, and Section 5. In Appendix F, we provide the detailed settings of the experiments in the main text together with an additional experiments to justify the result of Theorem 2.

## A Preliminaries

In this section, we provide some basic lemmas that will be used throughout the proof both from probability theory and from matrix analysis.

### A.1 Preparations in Probability Theory

The first lemma is a standard result characterizing the KL divergence between two Gaussian distributions.

**Lemma 5** (KL divergence between Gaussian distributions). *Let  $P_1$  and  $P_2$  are multivariate Gaussian distributions on  $\mathbb{R}^d$  with mean and covariance respectively  $\mu_1, \Sigma_1$  and  $\mu_2, \Sigma_2$ . Then the KL divergence between  $P_1$  and  $P_2$  are given as follows:*

$$\text{KL}(P_1||P_2) = \frac{1}{2} \left( \text{tr}(\Sigma_2^{-1}\Sigma_1) + (\mu_2 - \mu_1)^\top \Sigma_2^{-1}(\mu_2 - \mu_1) - d + \ln \left( \frac{\det \Sigma_2}{\det \Sigma_1} \right) \right).$$

We then provide a lemma which gives the expected difference between two uniform sampling variables.

**Lemma 6** (Two step sampling). *Suppose  $z$  is a discrete random variable with  $\mathbb{P}(z = z_i) = \frac{1}{N}$ ,  $\forall i = 1, 2, \dots, N$ , where the support set is  $\mathcal{Z} = \{z_1, \dots, z_N\} \subset \mathbb{R}^d$ . Suppose further  $\mathbf{U}$  is a random index set with size  $b$  and sampled uniformly without replacement from  $[N]$ . Suppose  $\mathbf{V}$  is another random index set independent of  $\mathbf{U}$  with size  $N - 1$  and sampled uniformly without replacement from  $[N]$ . Denote subset of  $\mathcal{Z}$  with index in  $\mathbf{U} \cap \mathbf{V}^c$  and  $\mathbf{V}$  respectively as  $\mathcal{Z}_{\mathbf{U} \cap \mathbf{V}^c} = \{z_i, i \in \mathbf{U} \cap \mathbf{V}^c\}$ ,  $\mathcal{Z}_{\mathbf{V}} = \{z_i, i \in \mathbf{V}\}$ , and the average of  $\mathcal{Z}_{\mathbf{U} \cap \mathbf{V}^c}$  and  $\mathcal{Z}_{\mathbf{V}}$  respectively as  $\bar{\mathcal{Z}}_{\mathbf{U} \cap \mathbf{V}^c}$  and  $\bar{\mathcal{Z}}_{\mathbf{V}}$ . Then the following equation holds:*

$$\mathbb{E}_{\mathbf{U}, \mathbf{V}} \left( \frac{(b - |\mathbf{U} \cap \mathbf{V}|)^2}{b^2} (\bar{\mathcal{Z}}_{\mathbf{V}} - \bar{\mathcal{Z}}_{\mathbf{U} \cap \mathbf{V}^c})(\bar{\mathcal{Z}}_{\mathbf{V}} - \bar{\mathcal{Z}}_{\mathbf{U} \cap \mathbf{V}^c})^\top \right) = \frac{1}{Nb} \left( \frac{N}{N-1} \right)^2 \text{Cov}(z).$$

*Proof.* We rewrite  $\mathbb{E}_{\mathbf{U}, \mathbf{V}} \left( \frac{(b - |\mathbf{U} \cap \mathbf{V}|)^2}{b^2} (\bar{\mathcal{Z}}_{\mathbf{V}} - \bar{\mathcal{Z}}_{\mathbf{U} \cap \mathbf{V}^c})(\bar{\mathcal{Z}}_{\mathbf{V}} - \bar{\mathcal{Z}}_{\mathbf{U} \cap \mathbf{V}^c})^\top \right)$  by taking conditional expectation with respect to  $|\mathbf{U} \cap \mathbf{V}|$  as follows:

$$\begin{aligned} & \mathbb{E}_{\mathbf{U}, \mathbf{V}} \left( \frac{(b - |\mathbf{U} \cap \mathbf{V}|)^2}{b^2} (\bar{\mathcal{Z}}_{\mathbf{V}} - \bar{\mathcal{Z}}_{\mathbf{U} \cap \mathbf{V}^c})(\bar{\mathcal{Z}}_{\mathbf{V}} - \bar{\mathcal{Z}}_{\mathbf{U} \cap \mathbf{V}^c})^\top \right) \\ &= \mathbb{E}_{|\mathbf{U} \cap \mathbf{V}^c|} \mathbb{E}_{\mathbf{U}, \mathbf{V}}^{|\mathbf{U} \cap \mathbf{V}^c|} \left( \frac{|\mathbf{U} \cap \mathbf{V}^c|^2}{b^2} (\bar{\mathcal{Z}}_{\mathbf{V}} - \bar{\mathcal{Z}}_{\mathbf{U} \cap \mathbf{V}^c})(\bar{\mathcal{Z}}_{\mathbf{V}} - \bar{\mathcal{Z}}_{\mathbf{U} \cap \mathbf{V}^c})^\top \right) \\ &= \mathbb{P}(|\mathbf{U} \cap \mathbf{V}^c| = 1) \mathbb{E}_{\mathbf{U}, \mathbf{V}}^{|\mathbf{U} \cap \mathbf{V}^c|=1} \left( \frac{1}{b^2} (\bar{\mathcal{Z}}_{\mathbf{V}} - \bar{\mathcal{Z}}_{\mathbf{V}^c})(\bar{\mathcal{Z}}_{\mathbf{V}} - \bar{\mathcal{Z}}_{\mathbf{V}^c})^\top \right) \\ &= \mathbb{P}(|\mathbf{U} \cap \mathbf{V}^c| = 1) \mathbb{E}_{\mathbf{V}} \left( \frac{1}{b^2} (\bar{\mathcal{Z}}_{\mathbf{V}} - \bar{\mathcal{Z}}_{\mathbf{V}^c})(\bar{\mathcal{Z}}_{\mathbf{V}} - \bar{\mathcal{Z}}_{\mathbf{V}^c})^\top \right) \\ &= \frac{1}{Nb} \left( \frac{N}{N-1} \right)^2 \text{Cov}(z). \end{aligned}$$

The proof is completed. □

We provide the following lemma for computing the KL divergence between two joint distributions.

**Lemma 7.** *Let  $\mathbf{X}$ ,  $\mathbf{Y}$ , and  $\mathbf{Z}$  be three random variables with  $\mathbf{X}$  and  $\mathbf{Y}$  having the same support set. Then the KL divergence between the joint distribution of  $(\mathbf{X}, \mathbf{Z})$  and  $(\mathbf{Y}, \mathbf{Z})$  can be decomposed into*

$$\text{KL}((\mathbf{X}, \mathbf{Z}) \| (\mathbf{Y}, \mathbf{Z})) = \mathbb{E}_{\mathbf{Z}} \text{KL}((\mathbf{X} | \mathbf{Z}) \| (\mathbf{Y} | \mathbf{Z})).$$

*Proof.* By the definition of KL divergence,

$$\begin{aligned} \text{KL}((\mathbf{X}, \mathbf{Z}) \| (\mathbf{Y}, \mathbf{Z})) &= \int \mathbb{P}(\mathbf{X}, \mathbf{Z}) \log \frac{\mathbb{P}(\mathbf{X}, \mathbf{Z})}{\mathbb{P}(\mathbf{Y}, \mathbf{Z})} \\ &= \int \mathbb{P}(\mathbf{Z}) \mathbb{P}^{\mathbf{Z}}(\mathbf{X}) \log \frac{\mathbb{P}^{\mathbf{Z}}(\mathbf{X}) \mathbb{P}(\mathbf{Z})}{\mathbb{P}^{\mathbf{Z}}(\mathbf{Y}) \mathbb{P}(\mathbf{Z})} \\ &= \int \mathbb{P}(\mathbf{Z}) \int \mathbb{P}^{\mathbf{Z}}(\mathbf{X}) \log \frac{\mathbb{P}^{\mathbf{Z}}(\mathbf{X})}{\mathbb{P}^{\mathbf{Z}}(\mathbf{Y})} \\ &= \mathbb{E}_{\mathbf{Z}} \text{KL}((\mathbf{X} | \mathbf{Z}) \| (\mathbf{Y} | \mathbf{Z})). \end{aligned}$$

The proof is completed.  $\square$

In the end of this section, we provide a proof of Lemma 1 using Lemma 7 for the completeness of this paper.

*Proof of Lemma 1.* By Lemma 7, we have

$$\begin{aligned} &\text{KL}(Q_{0:T} \| P_{0:T}) \\ &= \text{KL}(Q_{0:T} \| P_{0:T}) - \text{KL}(Q_{0:T} \| (Q_{0:T-1}, P_{T|[T-1]})) + \text{KL}(Q_{0:T} \| (Q_{0:T-1}, P_{T|[T-1]})) \\ &= \int Q_{0:T} \log \frac{Q_{0:T}}{P_{0:T}} - \int Q_{0:T} \log \frac{Q_{0:T}}{Q_{0:T-1} P_{T|[T-1]}} + \mathbb{E}_{Q_{T-1}} \text{KL}(Q_{T|[T-1]} \| P_{T|[T-1]}) \\ &= \int Q_{0:T} \log \frac{Q_{0:T-1}}{P_{0:T-1}} + \mathbb{E}_{Q_{T-1}} \text{KL}(Q_{T|[T-1]} \| P_{T|[T-1]}) \\ &= \int Q_{0:T-1} \log \frac{Q_{0:T-1}}{P_{0:T-1}} + \mathbb{E}_{Q_{T-1}} \text{KL}(Q_{T|[T-1]} \| P_{T|[T-1]}) \\ &= \text{KL}(Q_{0:T-1} \| P_{0:T-1}) + \mathbb{E}_{Q_{T-1}} \text{KL}(Q_{T|[T-1]} \| P_{T|[T-1]}). \end{aligned}$$

The proof is then completed by induction.  $\square$

**Remark 1.** *In this paper, we focus on the case where  $Q_{0:T}$  and  $P_{0:T}$  obeys the Markov Property, i.e., for any  $t$ ,*

$$Q_{t|[t-1]} = Q_{t|(t-1)}, P_{t|[t-1]} = P_{t|(t-1)}.$$

*Therefore, the result in Lemma 1 becomes*

$$\text{KL}(Q_{0:T} \| P_{0:T}) = \sum_{t=1}^T \mathbb{E}_{Q_{0:t-1}} [\text{KL}(Q_{t|(t-1)} \| P_{t|(t-1)})].$$

## A.2 Technical Lemmas in Matrix Analysis

We first provide a sufficient and necessary condition of that two symmetric matrices commute, and the proof can be found from any Linear Algebra Textbook (e.g. [32]).

**Lemma 8.** *Let  $\mathbf{A}$  and  $\mathbf{B}$  be two  $d \times d$  real symmetric matrices. Then,  $\mathbf{A}$  and  $\mathbf{B}$  commute (i.e.,  $\mathbf{AB} = \mathbf{BA}$ ), if and only if there exists an orthogonal matrix  $\mathbf{O}$  which can diagonalize  $\mathbf{A}$  and  $\mathbf{B}$  simultaneously, i.e., both  $\mathbf{O}^\top \mathbf{A} \mathbf{O}$  and  $\mathbf{O}^\top \mathbf{B} \mathbf{O}$  are diagonal.*

The next lemma is a key technique to obtain the optimal noise covariance of Theorem 2 and Theorem 3.

**Lemma 9.** Let  $\mathbf{B} \in \mathbb{R}^{d \times d}$  be a (fixed) positive definite matrix with eigenvalues  $(\beta_1, \dots, \beta_d)$ , where  $\beta_i \geq 0$ . Let  $\mathbf{G} \in \mathbb{R}^{d \times d}$  be a positive definite matrix variable with fixed trace  $\text{tr} \mathbf{G} = c$ , where  $c$  is a positive constant and  $c \leq \text{tr}(\mathbf{B})$ . Then the minimum of  $\text{tr}(\mathbf{G}^{-1} \mathbf{B}) + \ln \det(\mathbf{G})$  is achieved at  $\mathbf{G} = \mathbf{O}^\top \text{Diag}(\alpha_1, \dots, \alpha_d) \mathbf{O}$ , where

$$\alpha_i^* = \frac{\sqrt{1 - 4\lambda^* \beta_i} - 1}{-2\lambda^*},$$

$\mathbf{O}$  is any orthogonal matrix which diagonalize  $\mathbf{B}$  as

$$\mathbf{B} = \mathbf{O}^\top \text{Diag}(\beta_1, \dots, \beta_d) \mathbf{O},$$

and  $\lambda^* \leq 0$  is the unique solution of

$$\sum_{i=1}^d \frac{2\beta_i}{1 + \sqrt{1 - 4\lambda^* \beta_i}} = c. \quad (11)$$

**Remark 2.**  $f(\lambda) = \sum_{i=1}^d \frac{2\beta_i}{1 + \sqrt{1 - 4\lambda \beta_i}}$  is a monotonously increasing function with respect to  $\lambda$ , which guarantees the uniqueness of the solution of  $f(\lambda) = c$ .

Lemma 9 is proved via two steps: 1) we first prove for  $\mathbf{G}$  with fixed eigenvalues,  $\text{tr} \mathbf{G}^{-1} \mathbf{B} + \ln \det(\mathbf{G})$  is optimized if and only if  $\mathbf{G}$  and  $\mathbf{B}$  share the same eigenvectors; 2) we then calculate the eigenvalues of the optimal  $\mathbf{G}$  using the method of Lagrange multipliers. Theorem 3 can then be obtained by applying Lemma 9 and setting  $\mathbf{G} = \Sigma_t(\mathbf{S}, \mathbf{W})$  and  $\mathbf{B} = \sigma_t \mathbb{I} + \frac{\eta_t^2}{N b_t} \left( \frac{N}{N-1} \right)^2 \Sigma_{\mathbf{S}, \mathbf{W}}^{sd}$ . We first prove the eigenvectors of  $\mathbf{G}$  agree with those of  $\mathbf{B}$ .

**Lemma 10.** Let  $\mathbf{G} \in \mathbb{R}^{d \times d}$  be a positive definite matrix variable with fixed eigenvalues  $(\alpha_i)_{i=1}^d$ . Specifically, let  $\alpha_1 \geq \alpha_2 \geq \dots \geq \alpha_d > 0$  be all the eigenvalues of  $\mathbf{G}$ , and  $\mathbf{G}$  can be any element from the following set

$$\{\mathbf{Q}^\top \text{Diag}(\alpha_1, \dots, \alpha_d) \mathbf{Q} : \mathbf{Q} \text{ is orthogonal}\}.$$

Let  $\mathbf{B}$  be a fixed positive semi-definite matrix, with eigenvalues  $(\beta_i)_{i=1}^d$  satisfies  $\beta_1 \geq \beta_2 \geq \dots \geq \beta_d \geq 0$ . Then, the optimal (minimal) value of  $g(\mathbf{G}) = \text{tr}(\mathbf{G}^{-1} \mathbf{B})$  is achieved when

$$\mathbf{G}^* = \mathbf{O}^\top \text{Diag}(\alpha_1, \dots, \alpha_d) \mathbf{O},$$

where  $\mathbf{O}$  is any orthogonal matrix which diagonalizes  $\mathbf{B}$  as

$$\mathbf{B} = \mathbf{O}^\top \text{Diag}(\beta_1, \dots, \beta_d) \mathbf{O}$$

and the optimal value of  $g(\mathbf{G})$  is  $\sum_{i=1}^d \frac{\beta_i}{\alpha_i}$ .

*Proof.* Let  $\mathbf{G}^*$  be a optimal point of  $\text{tr}(\mathbf{G}^{-1} \mathbf{B})$ . We will then obtain the condition of  $\mathbf{G}^*$  by adding a disturbance. Specifically, let  $\mathbf{A}$  be an anti-symmetric matrix. Then,

$$(\mathbb{I} - \varepsilon \mathbf{A})(\mathbb{I} - \varepsilon \mathbf{A})^T = \mathbb{I} + \varepsilon^2 \mathbf{A} \mathbf{A}^\top.$$

As  $\varepsilon$  is small enough,  $\mathbb{I} + \varepsilon^2 \mathbf{A} \mathbf{A}^\top$  is inevitable, and positive definite. Therefore,  $(\mathbb{I} - \varepsilon \mathbf{A})(\mathbb{I} + \varepsilon^2 \mathbf{A} \mathbf{A}^\top)^{-\frac{1}{2}}$  is orthogonal. As  $\varepsilon \rightarrow 0$ ,

$$\lim_{\varepsilon \rightarrow 0} (\mathbb{I} - \varepsilon \mathbf{A})(\mathbb{I} + \varepsilon^2 \mathbf{A} \mathbf{A}^\top)^{-\frac{1}{2}} = \mathbb{I},$$

and

$$(\mathbb{I} - \varepsilon \mathbf{A})(\mathbb{I} + \varepsilon^2 \mathbf{A} \mathbf{A}^\top)^{-\frac{1}{2}} - \mathbb{I} = -\varepsilon \mathbf{A} + \mathbf{o}(\varepsilon).$$

Since  $\mathbf{G}^*$  is an optimal point of  $\text{tr}(\mathbf{G}^{-1} \mathbf{B})$ , we have

$$\begin{aligned} \text{tr}((\mathbf{G}^*)^{-1} \mathbf{B}) &\leq \text{tr} \left( (\mathbb{I} - \varepsilon \mathbf{A})(\mathbb{I} + \varepsilon^2 \mathbf{A} \mathbf{A}^\top)^{-\frac{1}{2}} (\mathbf{G}^*)^{-1} \left( (\mathbb{I} - \varepsilon \mathbf{A})(\mathbb{I} + \varepsilon^2 \mathbf{A} \mathbf{A}^\top)^{-\frac{1}{2}} \right)^\top \mathbf{B} \right) \\ &= \text{tr} \left( (\mathbb{I} - \varepsilon \mathbf{A})(\mathbb{I} + \varepsilon^2 \mathbf{A} \mathbf{A}^\top)^{-\frac{1}{2}} (\mathbf{G}^*)^{-1} (\mathbb{I} + \varepsilon^2 \mathbf{A} \mathbf{A}^\top)^{-\frac{1}{2}} (\mathbb{I} + \varepsilon \mathbf{A}) \mathbf{B} \right), \end{aligned}$$

which further leads to

$$-\varepsilon \operatorname{tr}(\mathbf{A}(\mathbf{G}^*)^{-1}\mathbf{B}) + \varepsilon \operatorname{tr}((\mathbf{G}^*)^{-1}\mathbf{A}\mathbf{B}) + o(\varepsilon) \geq 0.$$

By letting  $\varepsilon \rightarrow 0$ , we further have

$$-\operatorname{tr}(\mathbf{A}(\mathbf{G}^*)^{-1}\mathbf{B}) + \operatorname{tr}((\mathbf{G}^*)^{-1}\mathbf{A}\mathbf{B}) = 0,$$

which further leads to

$$\begin{aligned} 0 &= -\operatorname{tr}(\mathbf{A}(\mathbf{G}^*)^{-1}\mathbf{B}) + \operatorname{tr}((\mathbf{G}^*)^{-1}\mathbf{A}\mathbf{B}) \\ &= \operatorname{tr}(\mathbf{A}^\top(\mathbf{G}^*)^{-1}\mathbf{B}) + \operatorname{tr}(\mathbf{B}(\mathbf{G}^*)^{-1}\mathbf{A}) \\ &= 2 \operatorname{tr}(\mathbf{B}(\mathbf{G}^*)^{-1}\mathbf{A}). \end{aligned} \quad (12)$$

Since Eq.(12) holds for any anti-symmetry matrix  $\mathbf{A}$ , let  $\mathbf{A} = \mathbf{E}_{i,j} - \mathbf{E}_{j,i}$ , where  $i, j \in [d]$  and  $i \neq j$ . By Eq.(12), we have

$$(\mathbf{B}(\mathbf{G}^*)^{-1})_{i,j} = (\mathbf{B}(\mathbf{G}^*)^{-1})_{j,i},$$

which further leads to

$$\mathbf{B}(\mathbf{G}^*)^{-1} = (\mathbf{B}(\mathbf{G}^*)^{-1})^\top = (\mathbf{G}^*)^{-\top} \mathbf{B}^\top = (\mathbf{G}^*)^{-1} \mathbf{B}.$$

By simple rearranging, we have

$$\mathbf{G}^* \mathbf{B} = \mathbf{B} \mathbf{G}^*.$$

Therefore, by Lemma 8, we have that there exists an orthogonal matrix  $\mathbf{O}_0$ , such that both  $\mathbf{O}_0 \mathbf{G}^* \mathbf{O}_0^\top$  and  $\mathbf{O}_0 \mathbf{B} \mathbf{O}_0^\top$  are diagonal. By multiplying a permutation matrix, we further have there exists an orthogonal matrix  $\tilde{\mathbf{O}}$  such that  $\tilde{\mathbf{O}} \mathbf{G}^* \tilde{\mathbf{O}}^\top$  is diagonal, and

$$\tilde{\mathbf{O}} \mathbf{B} \tilde{\mathbf{O}}^\top = \operatorname{Diag}(\beta_1, \dots, \beta_d). \quad (13)$$

Since  $\tilde{\mathbf{O}} \mathbf{G}^* \tilde{\mathbf{O}}^\top$  is diagonal, there exists a permutation mapping  $\mathcal{T} : [d] \rightarrow [d]$ , such that

$$\tilde{\mathbf{O}} \mathbf{G}^* \tilde{\mathbf{O}}^\top = \operatorname{Diag}(\alpha_{\mathcal{T}(1)}, \dots, \alpha_{\mathcal{T}(d)}). \quad (14)$$

Denote the order of  $\beta_i$  ( $i = 1, 2, \dots, d$ ) as

$$\beta_1 = \dots = \beta_{s_1} > \beta_{s_1+1} = \dots = \beta_{s_1+s_2} > \dots > \beta_{\sum_{i=1}^{k-1} s_i+1} = \dots = \beta_{\sum_{i=1}^k s_i} > 0, \quad (15)$$

where  $\sum_{i=1}^k s_i = d$ , and we denote  $s_0 = 0$ . Since  $\mathbf{G}^*$  is the optimal point of  $\operatorname{tr}((\mathbf{G}^*)^{-1}\mathbf{B})$ , for any  $1 \leq i < j \leq d$  and  $\beta_i > \beta_j$ , we have  $\alpha_{\mathcal{T}(i)} > \alpha_{\mathcal{T}(j)}$ : otherwise, let

$$\mathbf{G}' = \tilde{\mathbf{O}}^\top \operatorname{Diag}(\alpha_{\mathcal{T}(1)}, \dots, \alpha_{\mathcal{T}(i-1)}, \alpha_{\mathcal{T}(j)}, \alpha_{\mathcal{T}(i+1)}, \dots, \alpha_{\mathcal{T}(j-1)}, \alpha_{\mathcal{T}(i)}, \alpha_{\mathcal{T}(j+1)}, \dots, \alpha_{\mathcal{T}(d)}) \tilde{\mathbf{O}},$$

we have

$$\operatorname{tr}((\mathbf{G}^*)^{-1}\mathbf{B}) > \operatorname{tr}((\mathbf{G}')^{-1}\mathbf{B}),$$

which contradicts that  $\mathbf{G}^*$  is optimal.

Therefore,  $\mathcal{T}(\sum_{i=1}^j s_i + 1), \dots, \mathcal{T}(\sum_{i=1}^{j+1} s_i)$  is then a permutation of  $\sum_{i=1}^j s_i + 1, \dots, \sum_{i=1}^{j+1} s_i$ , and there exists permutation matrix  $\mathbf{Q}$  such that

$$\mathbf{Q} = \operatorname{Diag}(\mathbf{Q}_1, \dots, \mathbf{Q}_k), \quad (16)$$

where  $\mathbf{Q}_i$  is a  $s_i \times s_i$  permutation sub-matrix, such that,

$$\mathbf{Q} \operatorname{Diag}(\alpha_{\mathcal{T}(1)}, \dots, \alpha_{\mathcal{T}(d)}) \mathbf{Q}^\top = \operatorname{Diag}(\alpha_1, \dots, \alpha_d). \quad (17)$$

Furthermore, by Eq.(16) and Eq.(15), we have

$$\mathbf{Q} \operatorname{Diag}(\beta_1, \dots, \beta_d) \mathbf{Q}^\top = \operatorname{Diag}(\beta_1, \dots, \beta_d). \quad (18)$$

Therefore, by Eqs.(13), (14), (17), and (18), we have

$$\begin{aligned}\mathbf{Q}\tilde{\mathbf{O}}\mathbf{B}\left(\mathbf{Q}\tilde{\mathbf{O}}\right)^{\top} &= \text{Diag}\left(\beta_1, \dots, \beta_d\right), \\ \mathbf{Q}\tilde{\mathbf{O}}\mathbf{G}^*\left(\mathbf{Q}\tilde{\mathbf{O}}\right)^{\top} &= \text{Diag}\left(\alpha_1, \dots, \alpha_d\right).\end{aligned}$$

Furthermore,

$$\begin{aligned}\text{tr}\left(\mathbf{G}^{-1}\mathbf{B}\right) &= \text{tr}\left(\left(\mathbf{Q}\tilde{\mathbf{O}}\right)^{\top}\text{Diag}\left(\alpha_1^{-1}, \dots, \alpha_d^{-1}\right)\left(\mathbf{Q}\tilde{\mathbf{O}}\right)\left(\mathbf{Q}\tilde{\mathbf{O}}\right)^{\top}\text{Diag}\left(\beta_1, \dots, \beta_d\right)\left(\mathbf{Q}\tilde{\mathbf{O}}\right)\right) \\ &= \sum_{i=1}^d \frac{\beta_i}{\alpha_i}.\end{aligned}$$

Therefore, the optimal value of  $\text{tr}(\mathbf{G}^{-1}\mathbf{B})$  is  $\sum_{i=1}^d \frac{\beta_i}{\alpha_i}$ , and the corresponding optimal point  $\mathbf{G}^*$  belongs to the following set

$$\mathcal{G} = \{\mathbf{O}^{\top}\text{Diag}(\alpha_1, \dots, \alpha_d)\mathbf{O} : \mathbf{B} = \mathbf{O}^{\top}\text{Diag}(\beta_1, \dots, \beta_d)\mathbf{O}\}.$$

On the other hand, it is easy to verify that for any element  $\mathbf{G} \in \mathcal{G}$ ,

$$\text{tr}(\mathbf{G}^{-1}\mathbf{B}) = \sum_{i=1}^d \frac{\beta_i}{\alpha_i}.$$

The proof is completed.  $\square$

Lemma 10 indicates that with eigenvalues fixed, the eigenvectors of  $\mathbf{G}$  should agree with those of  $\mathbf{B}$  by the order of eigenvalues. We then provide the following lemma to determine the optimal eigenvalues.

**Lemma 11.** *Let  $\beta_1, \beta_2, \dots, \beta_d$  be a series of fixed positive reals. Let  $\alpha_1, \alpha_2, \dots, \alpha_d \in \mathbb{R}^+$  be a series of real variables with constraint  $\sum_{i=1}^d \alpha_i = c$ , where  $c$  is a positive real constant which satisfies  $c \leq \sum_{i=1}^d \beta_i$ . Then the minimum of function*

$$f(\alpha_1, \dots, \alpha_d) = \sum_{i=1}^d \frac{\beta_i}{\alpha_i} + \sum_{i=1}^d \ln \alpha_i$$

is achieved at

$$\alpha_i^* = \frac{\sqrt{1 - 4\lambda^*\beta_i} - 1}{-2\lambda^*},$$

where  $\lambda^* \leq 0$  is the unique solution of

$$\sum_{i=1}^d \frac{2\beta_i}{1 + \sqrt{1 - 4\lambda^*\beta_i}} = c.$$

*Proof.* We find the minimum of  $f$  under the constraint that  $\alpha_1 + \dots + \alpha_d = c$  by the method of Lagrange Multiplier. Specifically, as for any  $i \in [d]$ ,  $\alpha_i \rightarrow 0^+$  or  $\alpha_i \rightarrow c^-$  will lead to  $f(\alpha_1, \dots, \alpha_d) \rightarrow \infty$ , we have that for any global optimal (minimal) point  $(\alpha_1^*, \dots, \alpha_d^*)$  of  $f$  under the constraint  $\alpha_1 + \dots + \alpha_d = c$ , we have that there exist a real  $\lambda^*$ , such that  $((\alpha_1^*, \dots, \alpha_d^*), \lambda^*)$  is a saddle point of  $\mathcal{L}((\alpha_1, \dots, \alpha_d), \lambda)$ , which is defined as

$$\mathcal{L}((\alpha_1, \dots, \alpha_d), \lambda) = f(\alpha_1, \dots, \alpha_d) + \lambda(c - \alpha_1 - \dots - \alpha_d).$$

By taking partial derivative of  $\mathcal{L}$  with respect to  $\alpha_i$ , we have

$$-\lambda^* = -\frac{1}{\alpha_i} + \frac{\beta_i}{\alpha_i^2} = \frac{\beta_i - \alpha_i}{\alpha_i^2}, \quad (19)$$

which further leads to

$$\sum_{i=1}^d \beta_i - c = \sum_{i=1}^d (\beta_i - \alpha_i) = -\lambda^* \left( \sum_{i=1}^d \alpha_i^2 \right).$$

Since  $\sum_{i=1}^d \beta_i \geq c$ , we have  $\lambda^* \leq 0$ . Therefore, for any  $i \in [d]$ , the quadratic equation  $\beta_i x^2 - x + \lambda^* = 0$  has only one positive solution  $\frac{1 + \sqrt{1 - 4\lambda^* \beta_i}}{2\beta_i}$ , and

$$\alpha_i^* = \frac{2\beta_i}{1 + \sqrt{1 - 4\lambda^* \beta_i}} = \frac{\sqrt{1 - 4\lambda^* \beta_i} - 1}{-2\lambda^*}.$$

On the other hand, by taking derivative of  $\mathcal{L}$  with respect to  $\lambda^*$ , we have

$$\sum_{i=1}^d \alpha_i^* = \sum_{i=1}^d \frac{2\beta_i}{1 + \sqrt{1 - 4\lambda^* \beta_i}} = c. \quad (20)$$

Since  $\sum_{i=1}^d \alpha_i^* = \sum_{i=1}^d \frac{2\beta_i}{1 + \sqrt{1 - 4\lambda^* \beta_i}}$  is a monotonously increasing function of  $\lambda^*$ , there is only one solution of  $\lambda^*$  of Eq.(20).

The proof is completed.  $\square$

The proof of Lemma 9 can then be obtained by combining Lemma 10 and Lemma 11 together.

*Proof of Lemma 9.* The original optimization problem can be written as

$$\min_{\text{tr}(\mathbf{G})=c} \text{tr}(\mathbf{G}^{-1}\mathbf{B}) + \ln(\det \mathbf{G}),$$

which can be further decomposed into

$$\begin{aligned} & \min_{\text{tr}(\mathbf{G})=c} \text{tr}(\mathbf{G}^{-1}\mathbf{B}) + \ln(\det \mathbf{G}) \\ &= \min_{\substack{\sum_{i=1}^d \alpha_i = c \\ \alpha_1 \geq \dots \geq \alpha_d > 0}} \min_{\mathbf{O} \in \mathcal{O}(d)} \left( \text{tr}(\mathbf{O}^\top \text{Diag}(\alpha_1^{-1}, \dots, \alpha_d^{-1}) \mathbf{O} \mathbf{B}) + \sum_{i=1}^d \ln \alpha_i \right) \\ &\stackrel{(*)}{=} \min_{\sum_{i=1}^d \alpha_i = c} \left( \sum_{i=1}^d \frac{\beta_i}{\alpha_i} + \sum_{i=1}^d \ln \alpha_i \right) \\ &\stackrel{(**)}{=} \sum_{i=1}^d \frac{1 + \sqrt{1 - 4\lambda^* \beta_i}}{2} + \sum_{i=1}^d \ln \frac{2\beta_i}{1 + \sqrt{1 - 4\lambda^* \beta_i}}, \end{aligned}$$

where Eq. (\*) is due to Lemma 10, Eq. (\*\*) is due to Lemma 11, and  $\lambda^* \leq 0$  is the unique solution of

$$\sum_{i=1}^d \frac{2\beta_i}{1 + \sqrt{1 - 4\lambda^* \beta_i}} = c.$$

Furthermore, the optimal point of  $\text{tr}(\mathbf{G}^{-1}\mathbf{B}) + \ln(\det \mathbf{G})$  can be calculated as

$$\begin{aligned} & \arg \min_{\text{tr}(\mathbf{G})=c} \text{tr}(\mathbf{G}^{-1}\mathbf{B}) + \ln(\det \mathbf{G}) \\ &= \left\{ \mathbf{O}^\top \text{Diag} \left( \frac{\sqrt{1 - 4\lambda^* \beta_1} - 1}{-2\lambda^*}, \dots, \frac{\sqrt{1 - 4\lambda^* \beta_d} - 1}{-2\lambda^*} \right) \mathbf{O} : \mathbf{B} = \mathbf{O}^\top \text{Diag}(\beta_1, \dots, \beta_d) \mathbf{O}, \right. \\ & \left. \lambda^* = \arg_{\lambda} \left( \sum_{i=1}^d \frac{2\beta_i}{1 + \sqrt{1 - 4\lambda \beta_i}} = c \right) \right\}. \end{aligned}$$

The proof is completed.  $\square$

## B Supplementary Materials of Section 3.1

*Proof of Lemma 2.* The  $\beta$ -smooth condition gives

$$\mathcal{R}_{\mathcal{S}}(\mathbf{W}_{t+1}) \leq \mathcal{R}_{\mathcal{S}}(\mathbf{W}_t) + \langle \nabla \mathcal{R}_{\mathcal{S}}(\mathbf{W}_t), \mathbf{W}_{t+1} - \mathbf{W}_t \rangle + \frac{\beta}{2} \|\mathbf{W}_{t+1} - \mathbf{W}_t\|^2. \quad (21)$$

Based on the update rule Eq.(3), we have

$$\mathbf{W}_{t+1} - \mathbf{W}_t = -\eta_{t+1} \nabla \mathcal{R}_{\mathcal{S}_{V_{t+1}}}(\mathbf{W}_t) + \varepsilon_{t+1}, \quad (22)$$

where  $\varepsilon_{t+1} \sim \mathcal{N}(0, \Sigma_{t+1}(\mathcal{S}, \mathbf{W}_t))$ .

Take expectation on Eq.(21) with respect to  $\mathbf{W}_{t+1} | \mathbf{W}_t$ , by  $\mathbb{E}^{\mathbf{W}_t}(\nabla \mathcal{R}_{\mathcal{S}_{V_{t+1}}}(\mathbf{W}_t)) = \nabla \mathcal{R}_{\mathcal{S}}(\mathbf{W}_t)$ ,

$$\mathbb{E}^{\mathbf{W}_t}[\mathcal{R}_{\mathcal{S}}(\mathbf{W}_{t+1})] \leq \mathcal{R}_{\mathcal{S}}(\mathbf{W}_t) - \eta_{t+1} \|\nabla \mathcal{R}_{\mathcal{S}}(\mathbf{W}_t)\|^2 + \frac{\beta}{2} \mathbb{E}^{\mathbf{W}_t} \|\mathbf{W}_{t+1} - \mathbf{W}_t\|^2. \quad (23)$$

Furthermore,

$$\begin{aligned} & \mathbb{E}^{\mathbf{W}_t} \|\mathbf{W}_{t+1} - \mathbf{W}_t\|^2 \\ &= \mathbb{E}^{\mathbf{W}_t} \|\eta_{t+1} \nabla \mathcal{R}_{\mathcal{S}_{V_{t+1}}}(\mathbf{W}_t) + \varepsilon_{t+1}\|^2 \\ &\stackrel{(*)}{=} \mathbb{E}^{\mathbf{W}_t} \|\eta_{t+1} \nabla \mathcal{R}_{\mathcal{S}_{V_t}}(\mathbf{W}_t)\|^2 + \mathbb{E}^{\mathbf{W}_t} \|\varepsilon_{t+1}\|^2 \\ &= \eta_{t+1}^2 \mathbb{E}^{\mathbf{W}_t} \|\nabla \mathcal{R}_{\mathcal{S}_{V_{t+1}}}(\mathbf{W}_t) - \nabla \mathcal{R}_{\mathcal{S}}(\mathbf{W}_t) + \nabla \mathcal{R}_{\mathcal{S}}(\mathbf{W}_t)\|^2 + \Sigma_{t+1}(\mathcal{S}, \mathbf{W}_t) \\ &= \eta_{t+1}^2 \mathbb{E}^{\mathbf{W}_t} \|\nabla \mathcal{R}_{\mathcal{S}_{V_{t+1}}}(\mathbf{W}_t) - \nabla \mathcal{R}_{\mathcal{S}}(\mathbf{W}_t)\|^2 + \eta_{t+1}^2 \|\nabla \mathcal{R}_{\mathcal{S}}(\mathbf{W}_t)\|^2 + \Sigma_{t+1}(\mathcal{S}, \mathbf{W}_t) \\ &= \frac{\eta_{t+1}^2}{N-1} \frac{N-b_{t+1}}{b_{t+1}} \Sigma_{\mathcal{S}, \mathbf{W}_t}^{sd} + \eta_{t+1}^2 \|\nabla \mathcal{R}_{\mathcal{S}}(\mathbf{W}_t)\|^2 + \Sigma_{t+1}(\mathcal{S}, \mathbf{W}_t), \end{aligned} \quad (24)$$

where Eq.(\*) is due to  $\varepsilon_{t+1}$  is independent of  $V_{t+1}$ , and  $\mathbb{E}^{\mathbf{W}_t} \varepsilon_{t+1} = 0$ .

Applying Eq.(24) back to Eq.(23) completes the proof.  $\square$

## C Supplementary Materials of Section 3.2

### C.1 Example to illustrate the difficulty to apply Proposition 1 to solve Problem 1

In this section, we show an example to demonstrate the difficulty for tackling **Problem 1** through Proposition 1. To start with, by the definition of state-dependent SGLD (Eq.(3)), covariance  $\Sigma_{[T]}$  is independent of  $\mathbf{J}$  and  $\mathbf{V}_{[T]}$ . Therefore, the square root separates the expectation with respect to  $\mathbf{V}_{[T]}$  and  $\mathbf{J}$  from the KL divergence term in the generalization bound

$$\mathbb{E}_{\mathcal{S}, \mathbf{V}_{[T]}, \mathbf{J}} \sqrt{\frac{(a_2 - a_1)^2}{2} \sum_{s=1}^T \mathbb{E}_{Q_{s-1}^{\mathcal{S}, \mathbf{V}_{[T]}}} \text{KL} \left( Q_{s|(s-1)}^{\mathcal{S}, \mathbf{V}_{[T]}} \parallel P_{s|(s-1)}^{\mathbf{J}, \mathcal{S}_{\mathbf{J}}, \mathbf{V}_{[T]}} \right)},$$

which makes the dependency of the bound on  $\Sigma_{[T]}$  even more complex. However, even though we change the optimization target into

$$\mathbb{E}_{\mathcal{S}} \sqrt{\frac{(a_2 - a_1)^2}{2} \mathbb{E}_{\mathbf{V}_{[T]}, \mathbf{J}} \sum_{s=1}^T \mathbb{E}_{Q_{s-1}^{\mathcal{S}, \mathbf{V}_{[T]}}} \text{KL} \left( Q_{s|(s-1)}^{\mathcal{S}, \mathbf{V}_{[T]}} \parallel P_{s|(s-1)}^{\mathbf{J}, \mathcal{S}_{\mathbf{J}}, \mathbf{V}_{[T]}} \right)}, \quad (25)$$

which is still a generalization bound by Jensen's Inequality, we demonstrate that the dependency on  $\Sigma_{[T]}$  is still too complex to tackle as follows.

To optimize Eq.(25) with respect to  $\Sigma_{[T]}(\mathcal{S}, \cdot)$  for fixed  $\mathcal{S}$ , we are actually seeking the optimal point of the following optimization problem:

$$\Sigma_{[T]}^*(\mathcal{S}, \cdot) = \arg \min_{\Sigma_{[T]}(\mathcal{S}, \cdot)} \sqrt{\mathbb{E}_{\mathbf{V}_{[T]}, \mathbf{J}} \sum_{s=1}^T \mathbb{E}_{Q_{s-1}^{\mathcal{S}, \mathbf{V}_{[T]}}} \text{KL} \left( Q_{s|(s-1)}^{\mathcal{S}, \mathbf{V}_{[T]}} \parallel P_{s|(s-1)}^{\mathbf{J}, \mathcal{S}_{\mathbf{J}}, \mathbf{V}_{[T]}} \right)}. \quad (26)$$

However, we will show it is technically hard to solve Eq. (26). As discussed in Section 3.2.1, for any fixed index  $i \in [T]$ , Eq.(26) depends on  $\Sigma_s(\mathbf{S}, \cdot)$  through both  $\mathbb{E}_{\mathbf{V}_{[T]}, \mathbf{J}} \mathbb{E}_{Q_{s-1}^{S, \mathbf{V}_{[T]}}}$   $\text{KL}(Q_{s|(s-1)}^{S, \mathbf{V}_{[T]}} \| P_{s|(s-1)}^{J, \mathbf{S}_J, \mathbf{V}_{[T]}})$  and  $\mathbb{E}_{\mathbf{V}_{[T]}, \mathbf{J}} \mathbb{E}_{Q_{i-1}^{S, \mathbf{V}_{[T]}}} \text{KL}(Q_{i|(i-1)}^{S, \mathbf{V}_{[T]}} \| P_{i|(i-1)}^{J, \mathbf{S}_J, \mathbf{V}_{[T]}})$  for  $\forall i > s$ . Specifically, we adopt the update rule for prior for all the steps and posterior for all steps  $t \neq s$  to be the isotropic SGLD in [25], i.e.,

$$\begin{aligned} \text{Posterior: } \mathbf{W}_t &= \mathbf{W}_{t-1} - \eta_t \nabla \mathcal{R}_{S_{V_t}}(\mathbf{W}_{t-1}) + \mathcal{N}(\mathbf{0}, \sigma_t \mathbb{I}) \\ \text{Prior: } \mathbf{W}_t &= \mathbf{W}_{t-1} - \eta_t \left( \frac{|\mathbf{V}_t \cap \mathbf{J}|}{|\mathbf{V}_t|} \nabla \mathcal{R}_{S_{V_t \cap J}}(\mathbf{W}_{t-1}) + \frac{|\mathbf{V}_t \cap \mathbf{J}^c|}{|\mathbf{V}_t|} \nabla \mathcal{R}_{S_J}(\mathbf{W}_{t-1}) \right) + \mathcal{N}(\mathbf{0}, \sigma_t \mathbb{I}), \end{aligned}$$

while we only optimize the noise covariance  $\Sigma_s(\mathbf{S}, \cdot)$  of step  $s$ :

$$\mathbf{W}_s = \mathbf{W}_{s-1} - \eta_s \nabla \mathcal{R}_{S_{V_s}}(\mathbf{W}_{s-1}) + \mathcal{N}(\mathbf{0}, \Sigma_s(\mathbf{S}, \mathbf{W}_{s-1})).$$

By simple calculation, for any step  $t \in [T]$ , given the same  $\mathbf{W}_{t-1}$ ,  $\mathbf{V}_t$ ,  $\mathbf{J}$ , and  $\mathbf{S}$ , the mean between the prior and posterior can be calculated as

$$\begin{aligned} &\mu^{S, \mathbf{V}_t, \mathbf{J}, \mathbf{W}_{t-1}} \\ &= -\eta_t \left( \frac{|\mathbf{V}_t \cap \mathbf{J}|}{|\mathbf{V}_t|} \nabla \mathcal{R}_{S_{V_t \cap J}}(\mathbf{W}_{t-1}) + \frac{|\mathbf{V}_t \cap \mathbf{J}^c|}{|\mathbf{V}_t|} \nabla \mathcal{R}_{S_J}(\mathbf{W}_{t-1}) \right) + \eta_t \nabla \mathcal{R}_{S_{V_t}}(\mathbf{W}_{t-1}) \\ &= \eta_t \frac{|\mathbf{V}_t \cap \mathbf{J}^c|}{|\mathbf{V}_t|} (\nabla \mathcal{R}_{S_{V_t \cap J^c}}(\mathbf{W}_{t-1}) - \nabla \mathcal{R}_{S_J}(\mathbf{W}_{t-1})). \end{aligned} \quad (27)$$

Therefore, by Lemma 5 and Lemma 6, the expected KL divergence  $\mathbb{E}_{\mathbf{V}_{[T]}, \mathbf{J}} \mathbb{E}_{Q_{i-1}^{S, \mathbf{V}_{[T]}}} \text{KL}(Q_{i|(i-1)}^{S, \mathbf{V}_{[T]}} \| P_{i|(i-1)}^{J, \mathbf{S}_J, \mathbf{V}_{[T]}})$  can be calculated as

$$\begin{aligned} &\mathbb{E}_{\mathbf{V}_{[T]}, \mathbf{J}} \mathbb{E}_{Q_{i-1}^{S, \mathbf{V}_{[T]}}} \text{KL} \left( Q_{i|(i-1)}^{S, \mathbf{V}_{[T]}} \parallel P_{i|(i-1)}^{J, \mathbf{S}_J, \mathbf{V}_{[T]}} \right) \\ &= \frac{1}{2} \mathbb{E}_{\mathbf{V}_{[T]}, \mathbf{J}} \mathbb{E}_{Q_{i-1}^{S, \mathbf{V}_{[T]}}} \left( \sigma_i^{-1} \mu^{S, \mathbf{V}_t, \mathbf{J}, \mathbf{W}_{i-1}} (\mu^{S, \mathbf{V}_t, \mathbf{J}, \mathbf{W}_{i-1}})^\top \right) \\ &= \frac{1}{2\sigma_i} \frac{1}{Nb_i} \left( \frac{N}{N-1} \right)^2 \mathbb{E}_{\mathbf{V}_{[i-1]}} \mathbb{E}_{Q_{i-1}^{S, \mathbf{V}_{[i-1]}}} \Sigma_{\mathbf{S}, \mathbf{W}_{i-1}}^{sd}. \end{aligned}$$

Therefore, the exact form of  $\mathbb{E}_{\mathbf{V}_{[T]}, \mathbf{J}} \mathbb{E}_{Q_{i-1}^{S, \mathbf{V}_{[T]}}} \text{KL}(Q_{i|(i-1)}^{S, \mathbf{V}_{[T]}} \| P_{i|(i-1)}^{J, \mathbf{S}_J, \mathbf{V}_{[T]}})$  requires taking expectation to  $\Sigma_{\mathbf{S}, \mathbf{W}_{i-1}}^{sd}$  with respect to Gaussian distribution with covariance  $\Sigma_s$ , and can be complex due to the complex structure of the model. Specifically, if  $i = s + 1$ , then  $\mathbb{E}_{\mathbf{V}_{[T]}, \mathbf{J}} \mathbb{E}_{Q_{i-1}^{S, \mathbf{V}_{[T]}}}$

$\text{KL}(Q_{i|(i-1)}^{S, \mathbf{V}_{[T]}} \| P_{i|(i-1)}^{J, \mathbf{S}_J, \mathbf{V}_{[T]}})$  can be further written as

$$\begin{aligned} &\mathbb{E}_{\mathbf{V}_{[T]}, \mathbf{J}} \mathbb{E}_{Q_s^{S, \mathbf{V}_{[T]}}} \text{KL} \left( Q_{s+1|s}^{S, \mathbf{V}_{[T]}} \parallel P_{s+1|s}^{J, \mathbf{S}_J, \mathbf{V}_{[T]}} \right) \\ &= \frac{1}{2} \frac{1}{Nb_{s+1}} \left( \frac{N}{N-1} \right)^2 \mathbb{E}_{\mathbf{V}_{[s]}} \mathbb{E}_{Q_{s-1}^{S, \mathbf{V}_{[s-1]}}} \mathbb{E}_{Q_{s|(s-1)}^{S, \mathbf{V}_s}} \Sigma_{\mathbf{S}, \mathbf{W}_s}^{sd}. \end{aligned}$$

Therefore, we need to optimize  $\mathbb{E}_{\mathbf{V}_s} \mathbb{E}_{Q_{s|(s-1)}^{S, \mathbf{V}_s}} \Sigma_{\mathbf{S}, \mathbf{W}_s}^{sd}$ , which can be further written as

$$\mathbb{E}_{\mathbf{V}_s} \mathbb{E}_{Q_{s|(s-1)}^{S, \mathbf{V}_s}} \Sigma_{\mathbf{S}, \mathbf{W}_s}^{sd} = \mathbb{E}_{\mathbf{V}_s} \mathbb{E}_{\mathcal{N}(-\eta_s \nabla \mathcal{R}_{S_{V_s}}(\mathbf{W}_{s-1}), \Sigma_s(\mathbf{S}, \mathbf{W}_{s-1}))} \Sigma_{\mathbf{S}, \mathbf{W}_s}^{sd}.$$

The explicit form of  $\mathbb{E}_{\mathbf{V}_s} \mathbb{E}_{\mathcal{N}(-\eta_s \nabla \mathcal{R}_{S_{V_s}}(\mathbf{W}_{s-1}), \Sigma_s(\mathbf{S}, \mathbf{W}_{s-1}))} \Sigma_{\mathbf{S}, \mathbf{W}_s}^{sd}$  can be obtained only when  $\Sigma_{\mathbf{S}, \mathbf{W}_s}^{sd}$  is some simple functions with respect to  $\mathbf{W}_s$  (e.g. quadratic functions), which makes the optimal of  $\mathbb{E}_{\mathbf{V}_s} \mathbb{E}_{\mathcal{N}(-\eta_s \nabla \mathcal{R}_{S_{V_s}}(\mathbf{W}_{s-1}), \Sigma_s(\mathbf{S}, \mathbf{W}_{s-1}))} \Sigma_{\mathbf{S}, \mathbf{W}_s}^{sd}$  complicated due to the complex structure of  $\mathcal{R}_S$  and  $\Sigma_{\mathbf{S}, \mathbf{W}}^{sd}$  in practical learning problems.



## C.2 Proof of Theorem 1

*Proof of Theorem 1.* For any two random measures  $P^{J, S_J, V_{[T]}}$ ,  $Q^{S, V_{[T]}}$ , by the Donsker-Varadhan variational formula [3], for any function  $g$  satisfying  $Q^{S, V_{[T]}}(\exp g) < \infty$ , we have

$$\text{KL}(P^{J, S_J, V_{[T]}} \| Q^{S, V_{[T]}}) \geq P^{J, S_J, V_{[T]}}(g) - Q^{S, V_{[T]}}(g) - \log Q^{S, V_{[T]}}(\exp(g - Q^{S, V_{[T]}}(g))).$$

Letting  $g(\mathbf{W}) = \lambda \left( \hat{\mathcal{R}}_{S_{J^c}}(\mathbf{W}) - \mathcal{R}_{\mathcal{D}}(\mathbf{W}) \right)$ , we further have

$$\begin{aligned} & \text{KL}(P^{J, S_J, V_{[T]}} \| Q^{S, V_{[T]}}) \\ & \geq \lambda \left( \mathcal{R}_{\mathcal{D}}(Q^{S, V_{[T]}}) - \hat{\mathcal{R}}_{S_{J^c}}(Q^{S, V_{[T]}}) - \left( \mathcal{R}_{\mathcal{D}}(P^{J, S_J, V_{[T]}}) - \hat{\mathcal{R}}_{S_{J^c}}(P^{J, S_J, V_{[T]}}) \right) \right) \\ & \quad - \log Q^{S, V_{[T]}} \left( \exp \left( \lambda \left( \hat{\mathcal{R}}_{S_{J^c}} - \mathcal{R}_{\mathcal{D}} - \left( \hat{\mathcal{R}}_{S_{J^c}}(Q^{S, V_{[T]}}) - \mathcal{R}_{\mathcal{D}}(Q^{S, V_{[T]}}) \right) \right) \right) \right). \end{aligned}$$

On the other hand, since  $\ell \in [a_1, a_2]$ ,  $\lambda \left( \hat{\mathcal{R}}_{S_{J^c}}(\mathbf{W}) - \mathcal{R}_{\mathcal{D}}(\mathbf{W}) \right)$  is  $\frac{\lambda(a_2 - a_1)}{2}$  subgaussian. Therefore,

$$\begin{aligned} & \left( \mathcal{R}_{\mathcal{D}}(Q^{S, V_{[T]}}) - \hat{\mathcal{R}}_{S_{J^c}}(Q^{S, V_{[T]}}) \right) - \left( \mathcal{R}_{\mathcal{D}}(P^{J, S_J, V_{[T]}}) - \hat{\mathcal{R}}_{S_{J^c}}(P^{J, S_J, V_{[T]}}) \right) \\ & \leq \inf_{\lambda > 0} \frac{\text{KL}(P^{J, S_J, V_{[T]}} \| Q^{S, V_{[T]}}) + \frac{1}{8} \lambda^2 (a_2 - a_1)^2}{\lambda}. \end{aligned}$$

Since  $P^{J, S_J, V_{[T]}}$  is independent of  $S_{J^c}$  then we have  $\mathbb{E}^{S_J, J, V_{[T]}} \left[ \mathcal{R}_{\mathcal{D}}(P^{J, S_J, V_{[T]}}) - \hat{\mathcal{R}}_{S_{J^c}}(P^{J, S_J, V_{[T]}}) \right] = 0$ . Hence, by averaging over  $S_{J^c}$  (equivalently, taking the conditional expectation conditional on  $(S_J, J, V_{[T]})$ ) we have, with probability one

$$\begin{aligned} & \mathbb{E}^{S_J, J, V_{[T]}} \left[ \mathcal{R}_{\mathcal{D}}(Q^{S, V_{[T]}}) - \hat{\mathcal{R}}_{S_{J^c}}(Q^{S, V_{[T]}}) \right] \\ & = \mathbb{E}^{S_J, J, V_{[T]}} \left[ \mathcal{R}_{\mathcal{D}}(Q^{S, V_{[T]}}) - \hat{\mathcal{R}}_{S_{J^c}}(Q^{S, V_{[T]}}) - \left( \mathcal{R}_{\mathcal{D}}(P^{J, S_J, V_{[T]}}) - \hat{\mathcal{R}}_{S_{J^c}}(P^{J, S_J, V_{[T]}}) \right) \right] \\ & \leq \mathbb{E}^{S_J, J, V_{[T]}} \left( \inf_{\lambda > 0} \frac{\text{KL}(P^{J, S_J, V_{[T]}} \| Q^{S, V_{[T]}}) + \frac{1}{8} \lambda^2 (a_2 - a_1)^2}{\lambda} \right) \end{aligned}$$

Finally, by taking the full expectation, since  $J \perp\!\!\!\perp Q^{S, V_{[T]}}$  we get:

$$\mathbb{E}_{S, V_{[T]}} \left[ \mathcal{R}_{\mathcal{D}}(Q^{S, V_{[T]}}) - \hat{\mathcal{R}}_S(Q^{S, V_{[T]}}) \right] \leq \mathbb{E}_{S, V_{[T]}, J} \left[ \inf_{\lambda > 0} \frac{\text{KL}(P^{J, S_J, V_{[T]}} \| Q^{S, V_{[T]}}) + \frac{1}{8} \lambda^2 (a_2 - a_1)^2}{\lambda} \right]$$

where the final  $\text{KL}(P^{J, S_J, V_{[T]}} \| Q^{S, V_{[T]}})$  on the right hand side is between two random measures, and hence is a random variable depending on  $(S, J, V_{[T]})$ ; and the expectation on the right hand side integrates over  $(S, J, V_{[T]})$ .

Since

$$\frac{\text{KL}(P^{J, S_J, V_{[T]}} \| Q^{S, V_{[T]}}) + \frac{1}{8} \lambda^2 (a_2 - a_1)^2}{\lambda} \geq \sqrt{\frac{1}{2} (a_2 - a_1)^2 \text{KL}(P^{J, S_J, V_{[T]}} \| Q^{S, V_{[T]}})},$$

the proof is completed.  $\square$

## D Supplementary of Section 4

In this section, we provide the proof of Theorem 2. Specifically, as mentioned in the main body, optimizing  $\text{Gen}_T$  with greedily selected prior involves three steps. (1). we first prove Lemma 3, which provides the optimal solution of noise covariance and prior for one single KL divergence term in the generalization bound  $\text{Gen}_T$ ; (2). as the optimal solution of noise covariance in Lemma 3 is independent of  $S_J$ ,  $V_{[T]}$ , and  $V_{[T]}$ , we are then able to obtain the greedy prior by Lemma 4; (3). applying the greedy prior back to  $\text{Gen}_T$ , we are finally able to derive Theorem 2.

We start by restating Lemma 12 and providing its proof.

**Lemma 12** (Lemma 3, restated). *For any  $s \in [T]$ ,  $\mathbf{J}$ ,  $\mathbf{S}_J$ , and  $\mathbf{V}_{[T]}$ , under Constraint 1,*

$$\min_{P_{s|(s-1)}^{\mathbf{J}, \mathbf{S}_J, \mathbf{V}_s}} \mathbb{E}_{\mathbf{S}_{J^c} \sim \mathcal{D}} \text{KL} \left( P_{s|(s-1)}^{\mathbf{J}, \mathbf{S}_J, \mathbf{V}_s} \parallel Q_{s|(s-1)}^{\mathbf{S}, \mathbf{V}_s} \right) \quad (28)$$

(1). *is independent of  $\Sigma_s$  when  $\mathbf{V}_s \cap \mathbf{J}^c = \emptyset$ , and (2). *is minimized at  $\Sigma_s(\mathbf{W}) = \lambda_s(\mathbf{W}) (\Sigma_{\mathbf{W}}^{\text{pop}})^{\frac{1}{2}}$ ,  $\forall \mathbf{W}$ , when  $\mathbf{V}_s \cap \mathbf{J}^c \neq \emptyset$ , where  $\lambda_s(\mathbf{W}) = c_s(\mathbf{W}) / \text{tr}((\Sigma_{\mathbf{W}}^{\text{pop}})^{\frac{1}{2}})$ .**

*Proof.* We first calculate  $\min_{P_{s|(s-1)}^{\mathbf{J}, \mathbf{S}_J, \mathbf{V}_s}} \mathbb{E}_{\mathbf{S}_{J^c} \sim \mathcal{D}} \text{KL} \left( P_{s|(s-1)}^{\mathbf{J}, \mathbf{S}_J, \mathbf{V}_s} \parallel Q_{s|(s-1)}^{\mathbf{S}, \mathbf{V}_s} \right)$  for any  $\Sigma_s$ . By applying the definition of the KL divergence, we have

$$\begin{aligned} & \arg \min_{P_{s|(s-1)}^{\mathbf{J}, \mathbf{S}_J, \mathbf{V}_s}} \mathbb{E}_{\mathbf{S}_{J^c} \sim \mathcal{D}} \text{KL} \left( P_{s|(s-1)}^{\mathbf{J}, \mathbf{S}_J, \mathbf{V}_s} \parallel Q_{s|(s-1)}^{\mathbf{S}, \mathbf{V}_{[T]}} \right) \\ &= \arg \min_{P_{s|(s-1)}^{\mathbf{J}, \mathbf{S}_J, \mathbf{V}_s}} \mathbb{E}_{\mathbf{S}_{J^c} \sim \mathcal{D}} \int P_{s|(s-1)}^{\mathbf{J}, \mathbf{S}_J, \mathbf{V}_s}(\mathbf{W}_s) \log \frac{P_{s|(s-1)}^{\mathbf{J}, \mathbf{S}_J, \mathbf{V}_s}(\mathbf{W}_s)}{Q_{s|(s-1)}^{\mathbf{S}, \mathbf{V}_{[T]}}(\mathbf{W}_s)} d\mathbf{W}_s \\ &\stackrel{(*)}{=} \arg \min_{P_{s|(s-1)}^{\mathbf{J}, \mathbf{S}_J, \mathbf{V}_s}} \int P_{s|(s-1)}^{\mathbf{J}, \mathbf{S}_J, \mathbf{V}_s}(\mathbf{W}_s) \log \frac{P_{s|(s-1)}^{\mathbf{J}, \mathbf{S}_J, \mathbf{V}_s}(\mathbf{W}_s)}{e^{\mathbb{E}_{\mathbf{S}_{J^c} \sim \mathcal{D}} \log Q_{s|(s-1)}^{\mathbf{S}, \mathbf{V}_{[T]}}(\mathbf{W}_s)}} d\mathbf{W}_s, \end{aligned} \quad (29)$$

where Eq. (\*) is due to the independence of  $P$  on  $\mathbf{S}_{J^c}$ .

Let

$$\tilde{Q}_{s|(s-1)}^{\mathbf{J}, \mathbf{S}_J, \mathbf{V}_{[T]}}(\mathbf{W}) = \frac{e^{\mathbb{E}_{\mathbf{S}_{J^c} \sim \mathcal{D}} \log Q_{s|(s-1)}^{\mathbf{S}, \mathbf{V}_{[T]}}(\mathbf{W})}}{\int e^{\mathbb{E}_{\mathbf{S}_{J^c} \sim \mathcal{D}} \log Q_{s|(s-1)}^{\mathbf{S}, \mathbf{V}_{[T]}}(\tilde{\mathbf{W}})} d\tilde{\mathbf{W}}}, \quad (30)$$

and  $\tilde{Q}_{s|(s-1)}^{\mathbf{J}, \mathbf{S}_J, \mathbf{V}_{[T]}}(\mathbf{W})$  is then a probability measure on  $\mathbb{R}^d$ . Applying Eq. (30) back to Eq. (29), we obtain

$$\begin{aligned} & \arg \min_{P_{s|(s-1)}^{\mathbf{J}, \mathbf{S}_J, \mathbf{V}_s}} \left( \int P_{s|(s-1)}^{\mathbf{J}, \mathbf{S}_J, \mathbf{V}_s}(\mathbf{W}_s) \log \frac{P_{s|(s-1)}^{\mathbf{J}, \mathbf{S}_J, \mathbf{V}_s}(\mathbf{W}_s)}{\tilde{Q}_{s|(s-1)}^{\mathbf{J}, \mathbf{S}_J, \mathbf{V}_{[T]}}(\mathbf{W}_s)} d\mathbf{W}_s \right. \\ & \quad \left. - \int P_{s|(s-1)}^{\mathbf{J}, \mathbf{S}_J, \mathbf{V}_s}(\mathbf{W}_s) \log \left( \int e^{\mathbb{E}_{\mathbf{S}_{J^c} \sim \mathcal{D}} \log Q_{s|(s-1)}^{\mathbf{S}, \mathbf{V}_{[T]}}(\tilde{\mathbf{W}})} d\tilde{\mathbf{W}} \right) d\mathbf{W}_s \right) \\ &= \arg \min_{P_{s|(s-1)}^{\mathbf{J}, \mathbf{S}_J, \mathbf{V}_s}} \left( \int P_{s|(s-1)}^{\mathbf{J}, \mathbf{S}_J, \mathbf{V}_s}(\mathbf{W}_s) \log \frac{P_{s|(s-1)}^{\mathbf{J}, \mathbf{S}_J, \mathbf{V}_s}(\mathbf{W}_s)}{\tilde{Q}_{s|(s-1)}^{\mathbf{J}, \mathbf{S}_J, \mathbf{V}_{[T]}}(\mathbf{W}_s)} d\mathbf{W}_s - \log \left( \int e^{\mathbb{E}_{\mathbf{S}_{J^c} \sim \mathcal{D}} \log Q_{s|(s-1)}^{\mathbf{S}, \mathbf{V}_{[T]}}(\tilde{\mathbf{W}})} d\tilde{\mathbf{W}} \right) \right) \\ &= \arg \min_{P_{s|(s-1)}^{\mathbf{J}, \mathbf{S}_J, \mathbf{V}_s}} \left( \int P_{s|(s-1)}^{\mathbf{J}, \mathbf{S}_J, \mathbf{V}_s}(\mathbf{W}_s) \log \frac{P_{s|(s-1)}^{\mathbf{J}, \mathbf{S}_J, \mathbf{V}_s}(\mathbf{W}_s)}{\tilde{Q}_{s|(s-1)}^{\mathbf{J}, \mathbf{S}_J, \mathbf{V}_{[T]}}(\mathbf{W}_s)} d\mathbf{W}_s \right) \\ &= \arg \min_{P_{s|(s-1)}^{\mathbf{J}, \mathbf{S}_J, \mathbf{V}_s}} \text{KL} \left( P \parallel \tilde{Q}_{s|(s-1)}^{\mathbf{J}, \mathbf{S}_J, \mathbf{V}_{[T]}} \right). \end{aligned} \quad (31)$$

The minimum of Eq.(31) is achieved if and only if  $P_{s|(s-1)}^{\mathbf{J}, \mathbf{S}_J, \mathbf{V}_s} = \tilde{Q}_{s|(s-1)}^{\mathbf{J}, \mathbf{S}_J, \mathbf{V}_{[T]}}$ , and we only need to calculate the exact form of  $\tilde{Q}_{s|(s-1)}^{\mathbf{J}, \mathbf{S}_J, \mathbf{V}_{[T]}}$ . Since  $\mathbf{W}_s | (\mathbf{W}_{s-1}, \mathbf{S}, \mathbf{V}_s) \sim \mathcal{N}(\mathbf{W}_{s-1} - \eta_s \nabla_{\mathbf{W}_{s-1}} \mathcal{R}_{\mathbf{S}_{V_s}}(\mathbf{W}_{s-1}), \Sigma_s(\mathbf{W}_{s-1}))$ , we have

$$\begin{aligned} & \exp \mathbb{E}_{\mathbf{S}_{J^c} \sim \mathcal{D}} \log Q_{s|(s-1)}^{\mathbf{S}, \mathbf{V}_{[T]}}(\mathbf{W}) \\ &= \exp \left( \mathbb{E}_{\mathbf{S}_{J^c} \sim \mathcal{D}} \left( -\frac{1}{2} (\mathbf{W} - \mathbf{W}_{s-1} + \eta_s \nabla_{\mathbf{W}_{s-1}} \mathcal{R}_{\mathbf{S}_{V_s}}(\mathbf{W}_{s-1}))^\top \Sigma_s(\mathbf{W}_{s-1})^{-1} (\mathbf{W} - \mathbf{W}_{s-1} \right. \right. \\ & \quad \left. \left. + \eta_s \nabla_{\mathbf{W}_{s-1}} \mathcal{R}_{\mathbf{S}_{V_s}}(\mathbf{W}_{s-1})) - \frac{d}{2} \log 2\pi - \frac{1}{2} \log \det(\Sigma_s(\mathbf{W}_{s-1})) \right) \right) \\ &= \exp \left( \mathbb{E}_{\mathbf{S}_{J^c} \sim \mathcal{D}} \left( -\frac{1}{2} (\mathbf{W} - \mathbf{W}_{s-1} + \eta_s \nabla_{\mathbf{W}_{s-1}} \mathcal{R}_{\mathbf{S}_{V_s}}(\mathbf{W}_{s-1}))^\top \Sigma_s(\mathbf{W}_{s-1})^{-1} (\mathbf{W} - \mathbf{W}_{s-1} \right. \right. \\ & \quad \left. \left. + \eta_s \nabla_{\mathbf{W}_{s-1}} \mathcal{R}_{\mathbf{S}_{V_s}}(\mathbf{W}_{s-1})) - \frac{d}{2} \log 2\pi - \frac{1}{2} \log \det(\Sigma_s(\mathbf{W}_{s-1})) \right) \right). \end{aligned} \quad (32)$$

On the other hand,

$$\begin{aligned}
& \mathbb{E}_{\mathcal{S}_{J^c} \sim \mathcal{D}} \left( -\frac{1}{2} (\mathbf{W} - \mathbf{W}_{s-1} + \eta_s \nabla_{\mathbf{W}_{s-1}} \mathcal{R}_{\mathcal{S}_{V_s}}(\mathbf{W}_{s-1}))^\top \boldsymbol{\Sigma}_s(\mathbf{W}_{s-1})^{-1} (\mathbf{W} - \mathbf{W}_{s-1} \right. \\
& \quad \left. + \eta_s \nabla_{\mathbf{W}_{s-1}} \mathcal{R}_{\mathcal{S}_{V_s}}(\mathbf{W}_{s-1})) \right) \\
&= -\frac{1}{2} \mathbb{E}_{\mathcal{S}_{J^c} \sim \mathcal{D}} \left( \mathbf{W} - \mathbf{W}_{s-1} + \eta_s \left( \frac{|\mathbf{V}_s \cap \mathbf{J}|}{|\mathbf{V}_s|} \nabla \mathcal{R}_{\mathcal{S}_{V_s \cap J}}(\mathbf{W}_{s-1}) + \frac{|\mathbf{V}_s \cap \mathbf{J}^c|}{|\mathbf{V}_s|} \nabla \mathcal{R}_{\mathcal{S}_{V_s \cap J^c}}(\mathbf{W}_{s-1}) \right) \right)^\top \\
& \quad \cdot \boldsymbol{\Sigma}_s(\mathbf{W}_{s-1})^{-1} \left( \mathbf{W} - \mathbf{W}_{s-1} + \eta_s \left( \frac{|\mathbf{V}_s \cap \mathbf{J}|}{|\mathbf{V}_s|} \nabla \mathcal{R}_{\mathcal{S}_{V_s \cap J}}(\mathbf{W}_{s-1}) + \frac{|\mathbf{V}_s \cap \mathbf{J}^c|}{|\mathbf{V}_s|} \nabla \mathcal{R}_{\mathcal{S}_{V_s \cap J^c}}(\mathbf{W}_{s-1}) \right) \right) \\
&= -\frac{1}{2} \left( \mathbf{W} - \mathbf{W}_{s-1} + \eta_s \left( \frac{|\mathbf{V}_s \cap \mathbf{J}|}{|\mathbf{V}_s|} \nabla \mathcal{R}_{\mathcal{S}_{V_s \cap J}}(\mathbf{W}_{s-1}) + \frac{|\mathbf{V}_s \cap \mathbf{J}^c|}{|\mathbf{V}_s|} \nabla \mathcal{R}_{\mathcal{D}}(\mathbf{W}_{s-1}) \right) \right)^\top \\
& \quad \cdot \boldsymbol{\Sigma}_s(\mathbf{W}_{s-1})^{-1} \left( \mathbf{W} - \mathbf{W}_{s-1} + \eta_s \left( \frac{|\mathbf{V}_s \cap \mathbf{J}|}{|\mathbf{V}_s|} \nabla \mathcal{R}_{\mathcal{S}_{V_s \cap J}}(\mathbf{W}_{s-1}) + \frac{|\mathbf{V}_s \cap \mathbf{J}^c|}{|\mathbf{V}_s|} \nabla \mathcal{R}_{\mathcal{D}}(\mathbf{W}_{s-1}) \right) \right) \\
& \quad - \frac{1}{2} \mathbb{E}_{\mathcal{S}_{J^c} \sim \mathcal{D}} \eta_s^2 \frac{|\mathbf{V}_s \cap \mathbf{J}^c|^2}{|\mathbf{V}_s|^2} (\nabla \mathcal{R}_{\mathcal{D}}(\mathbf{W}_{s-1}) - \nabla \mathcal{R}_{\mathcal{S}_{V_s \cap J^c}}(\mathbf{W}_{s-1}))^\top \boldsymbol{\Sigma}_s(\mathbf{W}_{s-1})^{-1} \\
& \quad \cdot (\nabla \mathcal{R}_{\mathcal{D}}(\mathbf{W}_{s-1}) - \nabla \mathcal{R}_{\mathcal{S}_{V_s \cap J^c}}(\mathbf{W}_{s-1})). \tag{33}
\end{aligned}$$

By combining Eq.(32) and Eq.(33), we further have

$$\begin{aligned}
& \exp \mathbb{E}_{\mathcal{S}_{J^c} \sim \mathcal{D}} \log Q_{s|(s-1)}^{\mathcal{S}, \mathbf{V}_{[T]}}(\mathbf{W}) \\
&= \frac{1}{(2\pi)^{-\frac{d}{2}} \det(\boldsymbol{\Sigma}_s(\mathbf{W}_{s-1}))^{\frac{1}{2}}} \exp \left( -\frac{1}{2} \left( \mathbf{W} - \mathbf{W}_{s-1} + \eta_s \left( \frac{|\mathbf{V}_s \cap \mathbf{J}|}{|\mathbf{V}_s|} \nabla \mathcal{R}_{\mathcal{S}_{V_s \cap J}}(\mathbf{W}_{s-1}) \right. \right. \right. \\
& \quad \left. \left. + \frac{|\mathbf{V}_s \cap \mathbf{J}^c|}{|\mathbf{V}_s|} \nabla \mathcal{R}_{\mathcal{D}}(\mathbf{W}_{s-1}) \right) \right)^\top \boldsymbol{\Sigma}_s(\mathbf{W}_{s-1})^{-1} \left( \mathbf{W} - \mathbf{W}_{s-1} + \eta_s \left( \frac{|\mathbf{V}_s \cap \mathbf{J}|}{|\mathbf{V}_s|} \nabla \mathcal{R}_{\mathcal{S}_{V_s \cap J}}(\mathbf{W}_{s-1}) \right. \right. \\
& \quad \left. \left. + \frac{|\mathbf{V}_s \cap \mathbf{J}^c|}{|\mathbf{V}_s|} \nabla \mathcal{R}_{\mathcal{D}}(\mathbf{W}_{s-1}) \right) \right) \exp \mathbb{E}_{\mathcal{S}_{J^c}} \left( -\frac{1}{2} \eta_s^2 \frac{|\mathbf{V}_s \cap \mathbf{J}^c|^2}{|\mathbf{V}_s|^2} (\nabla \mathcal{R}_{\mathcal{D}}(\mathbf{W}_{s-1}) - \nabla \mathcal{R}_{\mathcal{S}_{V_s \cap J^c}}(\mathbf{W}_{s-1}))^\top \right. \\
& \quad \left. \cdot \boldsymbol{\Sigma}_s(\mathbf{W}_{s-1})^{-1} (\nabla \mathcal{R}_{\mathcal{D}}(\mathbf{W}_{s-1}) - \nabla \mathcal{R}_{\mathcal{S}_{V_s \cap J^c}}(\mathbf{W}_{s-1})) \right). \tag{34}
\end{aligned}$$

Therefore, by taking integration with respect to  $\tilde{\mathbf{W}}$ , we have,

$$\begin{aligned}
& \int e^{\mathbb{E}_{\mathcal{S}_{J^c} \sim \mathcal{D}} \log Q_{s|(s-1)}^{\mathcal{S}, \mathbf{V}_{[T]}}(\tilde{\mathbf{W}})} d\tilde{\mathbf{W}} \\
&= \exp \mathbb{E}_{\mathcal{S}_{J^c}} \left( -\frac{1}{2} \eta_s^2 \frac{|\mathbf{V}_s \cap \mathbf{J}^c|^2}{|\mathbf{V}_s|^2} (\nabla \mathcal{R}_{\mathcal{D}}(\mathbf{W}_{s-1}) - \nabla \mathcal{R}_{\mathcal{S}_{V_s \cap J^c}}(\mathbf{W}_{s-1}))^\top \boldsymbol{\Sigma}_s(\mathbf{W}_{s-1})^{-1} \right. \\
& \quad \left. \cdot (\nabla \mathcal{R}_{\mathcal{D}}(\mathbf{W}_{s-1}) - \nabla \mathcal{R}_{\mathcal{S}_{V_s \cap J^c}}(\mathbf{W}_{s-1})) \right). \tag{35}
\end{aligned}$$

Therefore, by Eq.(30), Eq.(34), and Eq.(35), we have

$$\begin{aligned}
& \arg \min_{P^{\mathcal{J}, \mathcal{S}_{J^c}, \mathbf{V}_s}} \mathbb{E}_{\mathcal{S}_{J^c} \sim \mathcal{D}} \text{KL} \left( P_{s|(s-1)}^{\mathcal{J}, \mathcal{S}_{J^c}, \mathbf{V}_s} \parallel Q_{s|(s-1)}^{\mathcal{S}, \mathbf{V}_{[T]}} \right) = \tilde{Q}_{s|(s-1)}^{\mathcal{J}, \mathcal{S}_{J^c}, \mathbf{V}_{[T]}} \tag{36} \\
& \sim \mathcal{N} \left( \mathbf{W}_{s-1} - \eta_s \left( \frac{|\mathbf{V}_s \cap \mathbf{J}|}{|\mathbf{V}_s|} \nabla \mathcal{R}_{\mathcal{S}_{V_s \cap J}}(\mathbf{W}_{s-1}) + \frac{|\mathbf{V}_s \cap \mathbf{J}^c|}{|\mathbf{V}_s|} \nabla \mathcal{R}_{\mathcal{D}}(\mathbf{W}_{s-1}) \right), \boldsymbol{\Sigma}_s(\mathbf{W}_{s-1}) \right).
\end{aligned}$$

Applying Eq. (36) back to  $\mathbb{E}_{S_{J^c} \sim \mathcal{D}} \text{KL} \left( P_{s|(s-1)}^{J, S_J, \mathbf{V}_s} \parallel Q_{s|(s-1)}^{S, \mathbf{V}_{[T]}} \right)$ , we obtain

$$\begin{aligned}
& \min_{P_{s|(s-1)}^{J, S_J, \mathbf{V}_s}} \mathbb{E}_{S_{J^c} \sim \mathcal{D}} \text{KL} \left( P_{s|(s-1)}^{J, S_J, \mathbf{V}_s} \parallel Q_{s|(s-1)}^{S, \mathbf{V}_{[T]}} \right) = \mathbb{E}_{S_{J^c} \sim \mathcal{D}} \text{KL} \left( \tilde{Q}_{s|(s-1)}^{J, S_J, \mathbf{V}_{[T]}} \parallel Q_{s|(s-1)}^{S, \mathbf{V}_{[T]}} \right) \\
& = \int \tilde{Q}_{t|(t-1)}^{S_J, \mathbf{V}_{[s]}}(\mathbf{W}_s) \log \frac{\tilde{Q}_{t|(t-1)}^{S_J, \mathbf{V}_{[s]}}(\mathbf{W}_s)}{e^{\mathbb{E}_{S_{J^c} \sim \mathcal{D}} \log Q_{t|(t-1)}^{S, \mathbf{V}_{[s]}}(\mathbf{W}_s)}} d\mathbf{W}_s \\
& \stackrel{(\circ)}{=} - \int \tilde{Q}_{t|(t-1)}^{S_J, \mathbf{V}_{[s]}}(\mathbf{W}_s) \log \int e^{\mathbb{E}_{S_{J^c} \sim \mathcal{D}} \log Q_{t|(t-1)}^{S, \mathbf{V}_{[s]}}(\tilde{\mathbf{W}})} d\tilde{\mathbf{W}} d\mathbf{W}_s \\
& = - \log \int e^{\mathbb{E}_{S_{J^c} \sim \mathcal{D}} \log Q_{t|(t-1)}^{S, \mathbf{V}_{[s]}}(\tilde{\mathbf{W}})} d\tilde{\mathbf{W}} \\
& \stackrel{(\bullet)}{=} \frac{1}{2} \eta_t^2 \mathbb{E}_{S_{J^c} \sim \mathcal{D}} \frac{|\mathbf{V}_t \cap \mathbf{J}^c|^2}{|\mathbf{V}_t|^2} \left( \nabla \mathcal{R}_{\mathcal{D}}(\mathbf{W}_{t-1}) - \nabla \mathcal{R}_{S_{\mathbf{V}_t \cap \mathbf{J}^c}}(\mathbf{W}_{t-1}) \right)^\top \boldsymbol{\Sigma}_t(\mathbf{W}_{t-1})^{-1} \\
& \quad \cdot \left( \nabla \mathcal{R}_{\mathcal{D}}(\mathbf{W}_{t-1}) - \nabla \mathcal{R}_{S_{\mathbf{V}_t \cap \mathbf{J}^c}}(\mathbf{W}_{t-1}) \right) \\
& = \frac{1}{2} \eta_t^2 \mathbb{E}_{S_{J^c} \sim \mathcal{D}} \text{tr} \left( \boldsymbol{\Sigma}_t(\mathbf{W}_{t-1})^{-1} \frac{|\mathbf{V}_t \cap \mathbf{J}^c|^2}{|\mathbf{V}_t|^2} \left( \nabla \mathcal{R}_{\mathcal{D}}(\mathbf{W}_{t-1}) - \nabla \mathcal{R}_{S_{\mathbf{V}_t \cap \mathbf{J}^c}}(\mathbf{W}_{t-1}) \right)^\top \right. \\
& \quad \left. \cdot \left( \nabla \mathcal{R}_{\mathcal{D}}(\mathbf{W}_{t-1}) - \nabla \mathcal{R}_{S_{\mathbf{V}_t \cap \mathbf{J}^c}}(\mathbf{W}_{t-1}) \right) \right) \\
& \stackrel{(\diamond)}{=} \begin{cases} 0 & , \mathbf{V}_t \cap \mathbf{J}^c = \emptyset; \\ \frac{1}{2} \frac{\eta_t^2 N}{b_t(N-1)^2} \text{tr} \left( \boldsymbol{\Sigma}_t(\mathbf{W}_{t-1})^{-1} \boldsymbol{\Sigma}_{\mathbf{W}_{t-1}}^{pop} \right), & \mathbf{V}_t \cap \mathbf{J}^c \neq \emptyset. \end{cases} \tag{37}
\end{aligned}$$

where Eq. (◦) is due to the definition of  $\tilde{Q}_{t|(t-1)}^{S_J, \mathbf{V}_{[s]}}$  (Eq.(30)), Eq. (•) is due to Eq.(35) and Eq. (◊) is due to Lemma 6.

Therefore, when  $\mathbf{V}_t \cap \mathbf{J}^c = \emptyset$ , Eq.(28) is independent of  $\boldsymbol{\Sigma}_s$ . On the other hand, if  $\mathbf{V}_t \cap \mathbf{J}^c \neq \emptyset$ , we only need to solve

$$\boldsymbol{\Sigma}_s(\mathbf{W})^* = \arg \min_{\text{tr}(\boldsymbol{\Sigma}_s(\mathbf{W})) = c_s(\mathbf{W})} \text{tr} \left( \boldsymbol{\Sigma}_s(\mathbf{W})^{-1} \boldsymbol{\Sigma}_{\mathbf{W}}^{pop} \right), \text{ subject to Constraint 1.} \tag{38}$$

We complete the proof by solving Problem (38). Specifically, let the eigenvalues of  $\boldsymbol{\Sigma}_{\mathbf{W}}^{pop}$  be  $(\omega_i^{pop})_{i=1}^d$  (the value is by non-increasing order with respect to index) we first fix the eigenvalues of  $\boldsymbol{\Sigma}_s(\mathbf{W})$  to be  $\alpha_{[d]}$  with  $\alpha_i \geq 0$  (the value is by non-increasing order with respect to index),  $i \in [d]$ . Then, by Lemma 10, the minimum of  $\text{tr} \left( \boldsymbol{\Sigma}_s(\mathbf{W})^{-1} \boldsymbol{\Sigma}_{\mathbf{W}}^{pop} \right)$  is achieved when

$$\boldsymbol{\Sigma}_s(\mathbf{W}) \in \left\{ P^\top (\alpha_{[d]}) P : P \text{ is orthogonal and } \boldsymbol{\Sigma}_{\mathbf{W}}^{pop} = P^\top \left( \omega_{[d]}^{pop} \right) P \right\}, \tag{39}$$

and

$$\text{tr} \left( \boldsymbol{\Sigma}_s(\mathbf{W})^{-1} \boldsymbol{\Sigma}_{\mathbf{W}}^{pop} \right) = \sum_{i=1}^d \frac{\omega_i^{pop}}{\alpha_i}.$$

We then optimize  $\sum_{i=1}^d \frac{\omega_i^{pop}}{\alpha_i}$  under the constraint  $\sum_{i=1}^d \alpha_i = c_s(\mathbf{W}_{s-1})$ . By the Cauchy-Schwarz inequality,

$$c_s(\mathbf{W}_{s-1}) \left( \sum_{i=1}^d \frac{\omega_i^{pop}}{\alpha_i} \right) = \left( \sum_{i=1}^d \frac{\omega_i^{pop}}{\alpha_i} \right) \left( \sum_{i=1}^d \alpha_i \right) \stackrel{(*)}{\geq} \left( \sum_{i=1}^d \sqrt{\omega_i^{pop}} \right)^2, \tag{40}$$

where equality in inequality (\*) holds when  $\alpha_i^2/\omega_i^{pop}$  is invariant of  $i$ . By combining Eq.(39) and Eq.(40), the proof is completed.  $\square$

By Lemma 12, the optimal noise covariances  $\boldsymbol{\Sigma}_s$  of all KL divergence terms  $\mathbb{E}_{S_{J^c} \sim \mathcal{D}} \text{KL} \left( P_{s|(s-1)}^{J, S_J, \mathbf{V}_s} \parallel Q_{s|(s-1)}^{S, \mathbf{V}_s} \right)$  are the same regardless of  $\mathbf{V}_s$ ,  $\mathbf{J}$ , and  $S_J$ , which helps us to obtain Lemma 4.

*Proof of Lemma 4.* To begin with, denote the optimal noise covariance of first  $s$ -step in terms of the generalization bound  $\text{Gen}_s$  as  $\Sigma^s$  under Constraint 1, i.e.,

$$\Sigma_{[s]}^s \triangleq \arg \min_{\Sigma_{[s]}} \left( \min_P \text{Gen}_s(P, \Sigma_{[s]}) \right), \text{ subject to: Constraint 1,}$$

we also define  $Q^s$  accordingly as the posterior distribution with noise covariance  $\Sigma^s$ . Also, recall that  $P^s$  is the optimal prior in terms of the generalization bound  $\text{Gen}_s$  under Constraint 1, i.e.,

$$P^s = \arg \min_P \left( \min_{\Sigma_{[s]}} \text{Gen}_s(P, \Sigma_{[s]}) \right), \text{ subject to: Constraint 1.}$$

We would like to derive the form of  $\Sigma_s^s$  and  $P_{s|(s-1)}^s$ .

Specifically, we have

$$P_{s|(s-1)}^s = \arg \min_{P_{s|(s-1)}} \left( \text{Gen}_s(P, \Sigma_{[s]}^s) \right), \text{ subject to: } P_{t|(t-1)} = P_{t|(t-1)}^s (t < s),$$

and

$$\Sigma_s^s = \arg \min_{\Sigma_s} \left( \text{Gen}_s(P^s, \Sigma_{[s]}) \right), \text{ subject to: Constraint 1 and } \Sigma_t = \Sigma_{t|(t-1)}^s (t < s).$$

That is, to obtain the desired  $\Sigma_s^s$  and  $P_{s|(s-1)}^s$ , we only need to solve

$$\min_{\Sigma_s, P_{s|(s-1)}^s} \text{Gen}_s(P, \Sigma_{[s]}), \text{ subject to: } P_{t|(t-1)} = P_{t|(t-1)}^s (t < s) \text{ and } \Sigma_t = \Sigma_{t|(t-1)}^s (t < s).$$

On the other hand, with  $P_{t|(t-1)} = P_{t|(t-1)}^s (t < s)$  and  $\Sigma_t = \Sigma_{t|(t-1)}^s (t < s)$  and under Constraint 1, we have

$$\begin{aligned} & \min_{\Sigma_s, P_{s|(s-1)}^s} \text{Gen}_s(P, \Sigma_{[s]}) \\ &= \min_{\Sigma_s, P_{s|(s-1)}^s} \mathbb{E}_{\mathbf{S}_J, \mathbf{V}_{[s]}, \mathbf{J}} \sqrt{\frac{(a_2 - a_1)^2}{2} \mathbb{E}_{\mathbf{S}_{J^c}} \text{KL} \left( P^{\mathbf{J}, \mathbf{S}_J, \mathbf{V}_{[s]}} \parallel Q^{\mathbf{S}, \mathbf{V}_{[s]}} \right)} \\ &= \min_{\Sigma_s, P_{s|(s-1)}^s} \mathbb{E}_{\mathbf{S}_J, \mathbf{V}_{[s]}, \mathbf{J}} \sqrt{\frac{(a_2 - a_1)^2}{2} \mathbb{E}_{\mathbf{S}_{J^c}} \sum_{t=1}^s \mathbb{E}_{P_{t-1}^{\mathbf{J}, \mathbf{S}_J, \mathbf{V}_{[s]}}} \text{KL} \left( P_{t|(t-1)}^{\mathbf{J}, \mathbf{S}_J, \mathbf{V}_s} \parallel Q_{t|(t-1)}^{\mathbf{S}, \mathbf{V}_s} \right)} \\ &= \min_{\Sigma_s, P_{s|(s-1)}^s} \mathbb{E}_{\mathbf{S}_J, \mathbf{V}_{[s]}, \mathbf{J}} \left[ \sqrt{\frac{(a_2 - a_1)^2}{2} \mathbb{E}_{\mathbf{S}_{J^c}} \sum_{t=1}^s \mathbb{E}_{P_{t-1}^{\mathbf{J}, \mathbf{S}_J, \mathbf{V}_{[s]}}} \text{KL} \left( P_{t|(t-1)}^{\mathbf{J}, \mathbf{S}_J, \mathbf{V}_s} \parallel Q_{t|(t-1)}^{\mathbf{S}, \mathbf{V}_s} \right)} \right. \\ & \quad \left. + \frac{(a_2 - a_1)^2}{2} \mathbb{E}_{\mathbf{S}_{J^c}} \mathbb{E}_{P_{s-1}^{\mathbf{J}, \mathbf{S}_J, \mathbf{V}_{[s]}}} \text{KL} \left( P_{s|(s-1)}^{\mathbf{J}, \mathbf{S}_J, \mathbf{V}_s} \parallel Q_{s|(s-1)}^{\mathbf{S}, \mathbf{V}_s} \right) \right] \\ &\stackrel{(*)}{\geq} \mathbb{E}_{\mathbf{S}_J, \mathbf{V}_{[s]}, \mathbf{J}} \left[ \sqrt{\frac{(a_2 - a_1)^2}{2} \mathbb{E}_{\mathbf{S}_{J^c}} \sum_{t=1}^s \mathbb{E}_{P_{t-1}^{\mathbf{J}, \mathbf{S}_J, \mathbf{V}_{[s]}}} \text{KL} \left( P_{t|(t-1)}^{\mathbf{J}, \mathbf{S}_J, \mathbf{V}_s} \parallel Q_{t|(t-1)}^{\mathbf{S}, \mathbf{V}_s} \right)} \right. \\ & \quad \left. + \frac{(a_2 - a_1)^2}{2} \mathbb{E}_{\mathbf{S}_{J^c}} \mathbb{E}_{P_{s-1}^{\mathbf{J}, \mathbf{S}_J, \mathbf{V}_{[s]}}} \min_{\Sigma_s, P_{s|(s-1)}^s} \text{KL} \left( P_{s|(s-1)}^{\mathbf{J}, \mathbf{S}_J, \mathbf{V}_s} \parallel Q_{s|(s-1)}^{\mathbf{S}, \mathbf{V}_s} \right) \right]. \end{aligned}$$

By Lemma 12,  $\min_{\Sigma_s, P_{s|(s-1)}^s} \text{KL} \left( P_{s|(s-1)}^{\mathbf{J}, \mathbf{S}_J, \mathbf{V}_s} \parallel Q_{s|(s-1)}^{\mathbf{S}, \mathbf{V}_s} \right)$  is attained at  $\Sigma_s(\mathbf{W}) = \lambda_s(\mathbf{W}) (\Sigma_{\mathbf{W}}^{\text{pop}})^{\frac{1}{2}}$ , which is not dependent on  $\mathbf{J}, \mathbf{S}_J, \mathbf{V}_s$ , and

$$P_{s|(s-1)}^{\mathbf{J}, \mathbf{S}_J, \mathbf{V}_s} \sim \mathcal{N} \left( \mathbf{W}_{s-1} - \eta_s \left( \frac{|\mathbf{V}_s \cap \mathbf{J}|}{|\mathbf{V}_s|} \nabla \mathcal{R}_{\mathbf{S}_{\mathbf{V}_s \cap \mathbf{J}}}(\mathbf{W}_{s-1}) + \frac{|\mathbf{V}_s \cap \mathbf{J}^c|}{|\mathbf{V}_s|} \nabla \mathcal{R}_{\mathcal{D}}(\mathbf{W}_{s-1}) \right), \lambda_s(\mathbf{W}) (\Sigma_{\mathbf{W}}^{\text{pop}})^{\frac{1}{2}} \right).$$

Therefore, Inequality (\*) holds, and the proof is completed.  $\square$

By Lemma 4, we obtain the form of  $P^*$ , i.e.,

$$P_{s|(s-1)}^{*J, S_J, V_s} \sim \mathcal{N} \left( \mathbf{W}_{s-1} - \eta_s \left( \frac{|\mathbf{V}_s \cap \mathbf{J}|}{|\mathbf{V}_s|} \nabla \mathcal{R}_{S_{\mathbf{V}_s \cap \mathbf{J}}}(\mathbf{W}_{s-1}) + \frac{|\mathbf{V}_s \cap \mathbf{J}^c|}{|\mathbf{V}_s|} \nabla \mathcal{R}_{\mathcal{D}}(\mathbf{W}_{s-1}) \right), \lambda_s(\mathbf{W}) (\Sigma_{\mathbf{W}}^{pop})^{\frac{1}{2}} \right),$$

which allows us to further derive Theorem 2.

*Proof of Theorem 2.* By the definition of  $\text{Gen}_T$ , with prior the greedy prior and under Constraint 1, we have

$$\begin{aligned} & \min_{\Sigma_{[T]}} \text{Gen}_T(P^*, \Sigma_{[T]}) \\ &= \min_{\Sigma_{[T]}} \mathbb{E}_{S_J, V_{[T]}, J} \sqrt{\frac{(a_2 - a_1)^2}{2} \mathbb{E}_{S_{J^c}} \text{KL} \left( P_{S_J, V_{[T]}, J}^{*J, S_J, V_{[T]}} \parallel Q_{S, V_{[T]}} \right)} \\ &= \min_{\Sigma_{[T]}} \mathbb{E}_{S_J, V_{[T]}, J} \sqrt{\frac{(a_2 - a_1)^2}{2} \mathbb{E}_{S_{J^c}} \sum_{t=1}^T \mathbb{E}_{P_{t-1}^{*J, S_J, V_{[t-1]}}} \text{KL} \left( P_{t|(t-1)}^{*J, S_J, V_t} \parallel Q_{t|(t-1)}^{S, V_t} \right)} \\ &\stackrel{(\bullet)}{\geq} \mathbb{E}_{S_J, V_{[T]}, J} \sqrt{\frac{(a_2 - a_1)^2}{2} \sum_{t=1}^T \mathbb{E}_{P_{t-1}^{*J, S_J, V_{[t-1]}}} \min_{\Sigma_t} \mathbb{E}_{S_{J^c}} \text{KL} \left( P_{t|(t-1)}^{*J, S_J, V_t} \parallel Q_{t|(t-1)}^{S, V_t} \right)} \\ &\stackrel{(*)}{=} \mathbb{E}_{S_J, V_{[T]}, J} \sqrt{\frac{(a_2 - a_1)^2}{2} \sum_{t=1}^T \mathbb{E}_{P_{t-1}^{*J, S_J, V_{[t-1]}}} \min_{\Sigma_t, P_{t|(t-1)}^{*J, S_J, V_t}} \mathbb{E}_{S_{J^c}} \text{KL} \left( P_{t|(t-1)}^{*J, S_J, V_t} \parallel Q_{t|(t-1)}^{S, V_t} \right)}, \end{aligned}$$

where Eq. (\*) is due to that by the proof of Lemma 4,  $P_{s|(s-1)}^{*J, S_J, V_s}$  is the same as the prior minimizing  $\mathbb{E}_{S_{J^c} \sim \mathcal{D}} \text{KL} \left( P_{s|(s-1)}^{*J, S_J, V_s} \parallel Q_{s|(s-1)}^{S, V_s} \right)$  for any given  $J, S_J, V_s$ . Therefore,  $\min_{\Sigma_t} \mathbb{E}_{S_{J^c}} \text{KL} \left( P_{t|(t-1)}^{*J, S_J, V_t} \parallel Q_{t|(t-1)}^{S, V_t} \right)$  is attained when  $\Sigma_t(\mathbf{W}) = \lambda_s(\mathbf{W}) (\Sigma_{\mathbf{W}}^{pop})^{\frac{1}{2}}$ , which is independent of  $J, S_J, V_s$ , and Inequality (•) holds. Therefore,  $\min_{\Sigma_{[T]}} \text{Gen}_T(P^*, \Sigma_{[T]})$  is also attained at  $\Sigma_t(\mathbf{W}) = \lambda_s(\mathbf{W}) (\Sigma_{\mathbf{W}}^{pop})^{\frac{1}{2}}$ .

The proof is completed. □

## E Supplementary materials of Section 5

### E.1 Formal Description of the Prior in Section 5

In this section, we provide a detailed description of the update rule of the prior defined by Eq.(10).

---

**Algorithm 1:** Iteration of Prior

---

**Input:** Sample set  $\mathcal{S}$  with size  $N$ , initialization distribution  $\mathcal{W}_0$ , total step  $T$ , learning rate

$(\eta_t)_{t=1}^T$

**Output:**  $\mathbf{W}_{[T]}, \mathbf{J}$

- 1 Initialize  $\mathbf{W}_0$  according to  $\mathcal{W}_0$ ; initialize  $\mathbf{J}$  by uniformly sampling  $N - 1$  elements from  $[N]$  without replacement; set  $t = 0$
  - 2 **while**  $t < T$  **do**
  - 3     Uniformly sample index set  $\mathbf{V}_t \subset [N]$  such that  $|\mathbf{V}_t| = b_t$  without replacement and independent of  $\mathbf{J}$
  - 4     **if**  $\mathbf{V}_t \subset \mathbf{J}$  **then**
  - 5          $\mathbf{W}_t = \mathbf{W}_{t-1} - \eta_t \nabla \mathcal{R}_{S_{\mathbf{V}_t}}(\mathbf{W}_{t-1}) + \mathcal{N}(\mathbf{0}, \sigma_t \mathbb{I}_d)$
  - 6     **else**
  - 7          $\mathbf{W}_t = \mathbf{W}_{t-1} - \eta_t \frac{b_t - 1}{b_t} \nabla \mathcal{R}_{S_{\mathbf{V}_t \cap \mathbf{J}}}(\mathbf{W}_{t-1}) - \eta_t \frac{1}{b_t} \nabla \mathcal{R}_{S_{\mathbf{J}}}(\mathbf{W}_{t-1}) + \mathcal{N}(\mathbf{0}, \sigma_t \mathbb{I}_d)$
  - 8      $t = t + 1$
-

## E.2 Calculation of the Generalization Bound

To obtain the optimal noise covariance of **(P2)**, we first derive the explicit form of the generalization bound  $\widetilde{\text{Gen}}_T$  with the prior given by Eq. (10) as the following lemma:

**Lemma 13** (Calculate  $\widetilde{\text{Gen}}_T$ ). *Let Assumption 1 hold. Let the prior  $P$  is given by the update rule Eq. (10). Then, the generalization bound  $\widetilde{\text{Gen}}_T$  can be represented as*

$$\widetilde{\text{Gen}}_T = \mathbb{E}_{\mathbf{S}, \mathbf{J}} \sqrt{\frac{(a_2 - a_1)^2}{2} \sum_{t=1}^T \mathbb{E}_{\mathbf{S}_{J^c}, \mathbf{V}_{[t-1]}} \mathbb{E}_{P_{s-1}^{\mathbf{J}, \mathbf{S}_{\mathbf{J}}, \mathbf{V}_{[s-1]}}} A_t(\mathbf{S}, \mathbf{W}_{t-1}),}$$

where  $A_t(\mathbf{S}, \mathbf{W})$  is given as

$$\begin{aligned} A_t(\mathbf{S}, \mathbf{W}) &\triangleq \frac{1}{2} (\sigma_t(\mathbf{W}) \text{tr}(\boldsymbol{\Sigma}_t(\mathbf{S}, \mathbf{W})^{-1}) + \ln(\det \boldsymbol{\Sigma}_t(\mathbf{S}, \mathbf{W})) - d) \\ &\quad + \frac{\eta_t^2}{2Nb_t} \left( \frac{N}{N-1} \right)^2 \text{tr}(\boldsymbol{\Sigma}_t(\mathbf{S})^{-1} \boldsymbol{\Sigma}_{\mathbf{S}, \mathbf{W}_{t-1}}^{sd}) - \frac{1}{2} d \ln \sigma_t(\mathbf{W}). \end{aligned}$$

*Proof.* By the definition of  $\widetilde{\text{Gen}}_T$ , we have

$$\widetilde{\text{Gen}}_T = \mathbb{E}_{\mathbf{S}} \sqrt{\frac{(a_2 - a_1)^2}{2} \mathbb{E}_{\mathbf{V}_{[T]}, \mathbf{J}} \text{KL}(P^{\mathbf{J}, \mathbf{S}_{\mathbf{J}}, \mathbf{V}_{[T]}} \parallel Q^{\mathbf{S}, \mathbf{V}_{[T]}}),}$$

which by the decomposition of KL divergence (Lemma 1) further leads to

$$\begin{aligned} \widetilde{\text{Gen}}_T &= \mathbb{E}_{\mathbf{S}} \sqrt{\frac{(a_2 - a_1)^2}{2} \mathbb{E}_{\mathbf{V}_{[T]}, \mathbf{J}} \sum_{s=1}^T \mathbb{E}_{P_{s-1}^{\mathbf{J}, \mathbf{S}_{\mathbf{J}}, \mathbf{V}_{[T]}}} \text{KL}(P_{s|(s-1)}^{\mathbf{J}, \mathbf{S}_{\mathbf{J}}, \mathbf{V}_{[T]}} \parallel Q_{s|(s-1)}^{\mathbf{S}, \mathbf{V}_{[T]}})} \\ &= \mathbb{E}_{\mathbf{S}} \sqrt{\frac{(a_2 - a_1)^2}{2} \sum_{s=1}^T \mathbb{E}_{\mathbf{V}_{[s]}, \mathbf{J}} \mathbb{E}_{P_{s-1}^{\mathbf{J}, \mathbf{S}_{\mathbf{J}}, \mathbf{V}_{[s-1]}}} \text{KL}(P_{s|(s-1)}^{\mathbf{J}, \mathbf{S}_{\mathbf{J}}, \mathbf{V}_s} \parallel Q_{s|(s-1)}^{\mathbf{S}, \mathbf{V}_s)})} \\ &\stackrel{(*)}{=} \mathbb{E}_{\mathbf{S}} \sqrt{\frac{(a_2 - a_1)^2}{2} \sum_{s=1}^T \mathbb{E}_{\mathbf{V}_{[s-1]}} \mathbb{E}_{P_{s-1}^{\mathbf{J}, \mathbf{S}_{\mathbf{J}}, \mathbf{V}_{[s-1]}}} \mathbb{E}_{\mathbf{V}_s, \mathbf{J}} \text{KL}(P_{s|(s-1)}^{\mathbf{J}, \mathbf{S}_{\mathbf{J}}, \mathbf{V}_s} \parallel Q_{s|(s-1)}^{\mathbf{S}, \mathbf{V}_s)})} \\ &\stackrel{(**)}{=} \mathbb{E}_{\mathbf{S}} \sqrt{\frac{(a_2 - a_1)^2}{2} \sum_{s=1}^T \mathbb{E}_{\mathbf{V}_{[s-1]}, \mathbf{J}} \mathbb{E}_{P_{s-1}^{\mathbf{J}, \mathbf{S}_{\mathbf{J}}, \mathbf{V}_{[s-1]}}} \mathbb{E}_{\mathbf{V}_s, \mathbf{J}} \text{KL}(P_{s|(s-1)}^{\mathbf{J}, \mathbf{S}_{\mathbf{J}}, \mathbf{V}_s} \parallel Q_{s|(s-1)}^{\mathbf{S}, \mathbf{V}_s})}, \end{aligned}$$

where in Eq. (\*) we exchange the order between  $\mathbb{E}_{P_{s-1}^{\mathbf{J}, \mathbf{S}_{\mathbf{J}}, \mathbf{V}_{[s-1]}}}$  and  $\mathbb{E}_{\mathbf{V}_s, \mathbf{J}}$  due to Assumption 1, and Eq. (\*\*) is due to that  $\mathbb{E}_{P_{s-1}^{\mathbf{J}, \mathbf{S}_{\mathbf{J}}, \mathbf{V}_{[s-1]}}} \mathbb{E}_{\mathbf{V}_s, \mathbf{J}} \text{KL}(P_{s|(s-1)}^{\mathbf{J}, \mathbf{S}_{\mathbf{J}}, \mathbf{V}_s} \parallel Q_{s|(s-1)}^{\mathbf{S}, \mathbf{V}_s})$  is independent of  $\mathbf{J}$  by Assumption 1.

Therefore, we only need to prove  $\mathbb{E}_{\mathbf{V}_t, \mathbf{J}} \text{KL}(P_{t|(t-1)}^{\mathbf{J}, \mathbf{S}_{\mathbf{J}}, \mathbf{V}_t} \parallel Q_{t|(t-1)}^{\mathbf{S}, \mathbf{V}_t}) = A(t)$ , which can be obtained by

$$\begin{aligned}
& \mathbb{E}_{\mathbf{V}_t, \mathbf{J}} \text{KL} \left( P_{t|(t-1)}^{\mathbf{J}, \mathbf{S}_J, \mathbf{V}_t} \parallel Q_{t|(t-1)}^{\mathbf{S}, \mathbf{V}_t} \right) \\
& \stackrel{(\bullet)}{=} \frac{1}{2} \mathbb{E}_{\mathbf{V}_t, \mathbf{J}} \left( \left( \mu^{\mathbf{S}, \mathbf{V}_t, \mathbf{J}, \mathbf{W}_{t-1}} \right)^\top \Sigma_t(\mathbf{S}, \mathbf{W}_{t-1})^{-1} \mu^{\mathbf{S}, \mathbf{V}_t, \mathbf{J}, \mathbf{W}_{t-1}} + \ln \frac{\det \Sigma_t(\mathbf{S}, \mathbf{W}_{t-1})}{\sigma_t(\mathbf{W}_{t-1})^d} \right. \\
& \quad \left. + \text{tr} \left( \sigma_t(\mathbf{W}_{t-1}) \Sigma_t(\mathbf{S}, \mathbf{W}_{t-1})^{-1} \right) \right) - \frac{d}{2} \\
& = \frac{1}{2} \mathbb{E}_{\mathbf{V}_t, \mathbf{J}} \left( \text{tr} \left( \Sigma_t(\mathbf{S}, \mathbf{W}_{t-1})^{-1} \mu^{\mathbf{S}, \mathbf{V}_t, \mathbf{J}, \mathbf{W}_{t-1}} \left( \mu^{\mathbf{S}, \mathbf{V}_t, \mathbf{J}, \mathbf{W}_{t-1}} \right)^\top \right) + \ln \frac{\det \Sigma_t(\mathbf{S}, \mathbf{W}_{t-1})}{\sigma_t(\mathbf{W}_{t-1})^d} \right. \\
& \quad \left. + \text{tr} \left( \sigma_t(\mathbf{W}_{t-1}) \Sigma_t(\mathbf{S}, \mathbf{W}_{t-1})^{-1} \right) \right) - \frac{d}{2} \\
& = \frac{1}{2} \text{tr} \left( \Sigma_t(\mathbf{S}, \mathbf{W}_{t-1})^{-1} \mathbb{E}_{\mathbf{J}, \mathbf{V}_t} \mu^{\mathbf{S}, \mathbf{V}_t, \mathbf{J}, \mathbf{W}_{t-1}} \left( \mu^{\mathbf{S}, \mathbf{V}_t, \mathbf{J}, \mathbf{W}_{t-1}} \right)^\top \right) + \frac{1}{2} \ln \frac{\det \Sigma_t(\mathbf{S}, \mathbf{W}_{t-1})}{\sigma_t(\mathbf{W}_{t-1})^d} \\
& \quad + \frac{1}{2} \text{tr} \left( \sigma_t(\mathbf{W}_{t-1}) \Sigma_t(\mathbf{S}, \mathbf{W}_{t-1})^{-1} \right) - \frac{d}{2} \\
& \stackrel{(\circ)}{=} \frac{1}{2} \left( \sigma_t(\mathbf{W}_{t-1}) \text{tr} \left( \Sigma_t(\mathbf{S}, \mathbf{W}_{t-1})^{-1} \right) + \ln \left( \det \Sigma_t(\mathbf{S}, \mathbf{W}_{t-1}) \right) - d \right) - \frac{1}{2} d \ln \sigma_t(\mathbf{W}_{t-1}) \\
& \quad + \frac{\eta_t^2}{2N b_t} \left( \frac{N}{N-1} \right)^2 \text{tr} \left( \Sigma_t(\mathbf{S}, \mathbf{W}_{t-1})^{-1} \Sigma_{\mathbf{S}, \mathbf{W}_{t-1}}^{sd} \right),
\end{aligned}$$

where Eq. (•) is due to Lemma 5, where  $\mu^{\mathbf{S}, \mathbf{V}_t, \mathbf{J}, \mathbf{W}_{t-1}}$  is defined by Eq.(27), and Eq. (◦) is obtained by Lemma 6.

The proof is completed.  $\square$

By Lemma 13, for any  $t \in [T]$ ,  $\mathbf{S}$ , and  $\mathbf{W}_{t-1}$ ,  $\text{Gen}_{[T]}$  depend on  $\Sigma_t(\mathbf{W}, \text{Gen}_{[T]})$  only through  $A_t(\mathbf{S}, \mathbf{W})$ , and the solution of optimizing  $A_t$  with respect to  $\Sigma_t$  under Constraint 1 has already been given by Lemma 9. We then complete the proof of Theorem 3 in the next section by combining Lemma 13 and Lemma 9 together.

### E.3 Proof of Theorem 3

In this section, we first restate Theorem 3 with explicit form of  $\tilde{\omega}_i^{\mathbf{S}, \mathbf{W}}$  (omitted in the main text). We then provide the proof of the theorem by Lemma 13 and Lemma 9.

**Theorem 4.** *Let prior and posterior be defined as Eq.(10) and Eq.(3), respectively. Then, with Assumption 1, the solution of (P2) is given by*

$$\Sigma_t^*(\mathbf{S}, \mathbf{W}) = \mathbf{Q}_{\mathbf{S}, \mathbf{W}}^{sd} \text{Diag}(\tilde{\omega}_{t,1}^{\mathbf{S}, \mathbf{W}}, \dots, \tilde{\omega}_{t,d}^{\mathbf{S}, \mathbf{W}}) (\mathbf{Q}_{\mathbf{S}, \mathbf{W}}^{sd})^\top,$$

where

$$\tilde{\omega}_{t,i}^{\mathbf{S}, \mathbf{W}} = \frac{\sqrt{1 - 4\lambda^* \left( \mathbb{E}_{\mathbf{J}} \sigma_t(\mathbf{S}_J, \mathbf{W}) + \frac{\eta_i^2}{N b_t} \left( \frac{N}{N-1} \right)^2 \omega_i^{\mathbf{S}, \mathbf{W}} \right)} - 1}{-2\lambda^*},$$

$\lambda^*$  is determined by  $\sum_{i=1}^d \tilde{\omega}_{t,i}^{\mathbf{S}, \mathbf{W}} = c_t(\mathbf{S}, \mathbf{W})$ , and  $\mathbf{Q}_{\mathbf{S}, \mathbf{W}}^{sd}$  is the orthogonal matrix that diagonalizes  $\Sigma_{\mathbf{S}, \mathbf{W}}^{sd}$  as

$$\Sigma_{\mathbf{S}, \mathbf{W}}^{sd} = \mathbf{Q}_{\mathbf{S}, \mathbf{W}}^{sd} \text{Diag}(\omega_1^{\mathbf{S}, \mathbf{W}}, \dots, \omega_d^{\mathbf{S}, \mathbf{W}}) (\mathbf{Q}_{\mathbf{S}, \mathbf{W}}^{sd})^\top.$$



*Proof of Theorem 3.* By Lemma 13,  $\text{Gen}_T$  depends on  $\Sigma_t(\mathbf{S}, \mathbf{W})$  only through  $A_t(\mathbf{S}, \mathbf{W})$ , and we have

$$\begin{aligned}
& \Sigma_s^*(\mathbf{S}, \mathbf{W}) \\
&= \arg \min_{\text{Constraint 1}} A_s(\mathbf{S}, \mathbf{W}) \\
&= \arg \min_{\text{tr}(\Sigma)=c_s(\mathbf{S}, \mathbf{W})} \frac{1}{2} \left( \text{tr} \left( \Sigma^{-1} \left( \sigma_s(\mathbf{S}_J, \mathbf{W}_{s-1}) \mathbb{I} + \frac{\eta_s^2}{Nb_s} \left( \frac{N}{N-1} \right)^2 \Sigma_{\mathbf{S}, \mathbf{W}_{s-1}}^{sd} \right) \right) \right. \\
&\quad \left. - d \ln \sigma_s(\mathbf{S}_J, \mathbf{W}_{s-1}) - d + \ln(\det \Sigma) \right) \\
&= \arg \min_{\text{tr}(\Sigma)=c_s(\mathbf{S}, \mathbf{W})} \frac{1}{2} \left( \text{tr} \left( \Sigma^{-1} \left( \sigma_s(\mathbf{S}_J, \mathbf{W}_{s-1}) \mathbb{I} + \frac{\eta_s^2}{Nb_s} \left( \frac{N}{N-1} \right)^2 \Sigma_{\mathbf{S}, \mathbf{W}_{s-1}}^{sd} \right) \right) \right. \\
&\quad \left. + \ln(\det \Sigma) \right).
\end{aligned}$$

Applying Lemma 9 completes the proof. □

#### E.4 Smaller Condition Number

In this section, we demonstrate why the optimal noise of Theorem 3 has smaller condition number than  $\Sigma^{sd}$  as the following corollary.

**Corollary 1.** *The optimal noise covariance  $\Sigma^*$  given by Theorem 3 has smaller condition number than  $\Sigma^{sd}$ .*

*Proof.* We prove this claim following two steps.

Firstly, the noise covariance of the prior is isotropic, has condition number 1, and push the condition number of  $\sigma_t \mathbb{I} + \frac{\eta_t^2}{Nb_t} \left( \frac{N}{N-1} \right)^2 \Sigma_{\mathbf{S}, \mathbf{W}}^{sd}$  smaller than  $\Sigma_{\mathbf{S}, \mathbf{W}}^{sd}$ .

Secondly, the optimal solution  $G$  of Lemma 9 always has a smaller condition number than  $B$ , which implies that  $\Sigma_t^*(\mathbf{S}, \mathbf{W})$  has smaller condition number than  $B = \sigma_t \mathbb{I} + \frac{\eta_t^2}{Nb_t} \left( \frac{N}{N-1} \right)^2 \Sigma_{\mathbf{S}, \mathbf{W}}^{sd}$ . Hence the condition number of  $\Sigma_t^*(\mathbf{S}, \mathbf{W})$  is smaller than  $\Sigma_{\mathbf{S}, \mathbf{W}}^{sd}$ . □

## F Experiments

In this section, we introduce the settings of the experiments in Fig. (1) Fig.(2), Fig.(3), and Fig.(4). We further include an additional experiment comparing the generalization error between SGLD with square rooted empirical gradient covariance (SREC-SGLD) (the optimal noise covariance in Theorem 2) and SGLD with empirical gradient covariance (EC-SGLD) subject to Constraint 1.

### F.1 Experiment settings

For both Fig. (1), Fig.(2), Fig. (3), and Fig. (4), we adopt the same setting as the Fashion-MNIST experiment of [40, Section D.3] despite enlarging the training set. Specifically, we use the 4-layer convolutional neural network as our model to conduct multi-class classification on Fashion-MNIST [35]. Concretely, this convolutional neural network can be expressed in order as: convolutional layer with 10 channel and filter size  $5 \times 5$ , max-pool layer with kernel size 2 and stride 2, convolutional layer with 10 channel and filter size  $5 \times 5$ , max-pool with kernel size 2, two fully connected layer with width 50. Our training set consists of 10,000 examples uniformly sampled without replacement from the Fashion-MNIST dataset. Our training set is larger than that in [40] (which only contains 1200 samples), but is still one sixth of the whole Fashion-MNIST dataset due to the computational burden of gradient descent (without mini-batch) in the SGLD. The learning rates of all SGLD are set to 0.07, which is exactly the same as [40]. We also set the learning rate of SGD in Fig. (1) to 0.07 for fair comparison with SGLD.

**Empirical gradient covariance:** We use top 100 eigenvalues to approximate the empirical gradient covariance matrix. Specifically, we decompose the matrix  $\Sigma_{S,W}^{sd}$  into  $(Q_{S,W})^\top (\omega_{[d]}^{S,W}) Q_{S,W}$ , and use  $(Q_{S,W})^\top (\omega_{[100]}^{S,W}, \mathbf{0}_{d-100}) Q_{S,W}$  to approximate  $\Sigma_{S,W}^{sd}$ .

**Noise Scale:** In Fig. (1) and Fig.(2), the traces of all SGLDs are set to be  $\text{tr}(\Sigma_{S,W}^{sd})$ ; in Fig. (3), the traces are set to be  $\text{tr}((\Sigma_{S,W}^{sd})^{1/2})$  and  $5 \text{tr}((\Sigma_{S,W}^{sd})^{1/2})$ , respectively in (a) and (b); in Fig. (4), the traces are set to be  $\text{tr}((\Sigma_{S,W}^{sd})^{1/2})$ .

**Noise frequency:** Similar to [40], we re-estimate the noise structure of all SGLDs every 10 epochs to ease the computational burden.

## F.2 Comparison between EC-SGLD and SREC-SGLD

We further conduct an experiment to compare the generalization performance between Iso-SGLD, EC-SGLD and SREC-SGLD, with the traces of the covariance are all set to be  $5 \text{tr}((\Sigma_{S,W}^{sd})^{1/2})$ , and all other settings consistent with Appendix F.1. The generalization error along the iteration of SREC-SGLD, Iso-SGLD, and EC-SGLD is plotted as Fig. 5, where one can easily observe the generalization error of SREC-SGLD is the smallest, which supports Theorem 2.

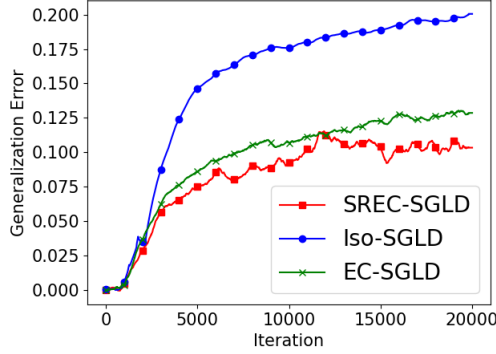


Figure 5: Comparison of generalization error for SGLDs with different noise structures. Traces of the covariances are all set to be  $5 \text{tr}((\Sigma_{S,W}^{sd})^{1/2})$ .