

825 A Appendix / supplemental material

Notations	
Symbol	Meaning
m	Number of neurons in hidden-layer of score network
C_{w_y, u_y}	Upper bound on $\ w_{y,i}\ _1, \ u_{y,i}\ _1$
826 F_T^2	Upper bound on $\mathbb{E}_{X_0} \mathbb{E}_{\xi_j} \left[\ \sigma(Wx(t) + Ue(t))\ _2^2 \right], 0 \leq t \leq T$
L_y	$\tilde{\mathcal{L}}_{reg}^{n_y}(\theta_y, \theta_{-y})$ is $\frac{L_y}{m^2}$ smooth w.r.t. θ_y
ϕ_y	Lipschitz constant of $\tilde{\mathcal{L}}_{mut}^{n_y}(\theta_y, \theta_{-y})$ w.r.t. θ_y
σ_y	$\tilde{\mathcal{L}}_{reg}^{n_y}(\theta_y, \theta_{-y})$ is σ_y Lipschitz in θ_y
B	Upper Bound of the Frobenius norm of A_y

827 B Sampling

828 Denote the backward time schedule as $\{t_j^{\leftarrow}\}_{0 \leq j \leq N}$ such that $0 = t_0^{\leftarrow} < t_1^{\leftarrow} < \dots, t_N^{\leftarrow} = T - \alpha$.
829 Lower case p_t represents the density of P_t . We consider the exponential integrator scheme for
830 simulating the backward SDE with

831 The generation algorithm can be expressed as a piecewise continuous-time SDE: for any $t \in$
832 $[t_j^{\leftarrow}, t_{j+1}^{\leftarrow})$.

$$d\bar{Y}_t = (\bar{Y}_t + 2s_{T-t_j^{\leftarrow}, \theta_y}(\bar{Y}_{t_j^{\leftarrow}}))dt + \sqrt{2}d\bar{W}_t \quad (23)$$

833 Denote $q_t := \text{Law}(\bar{Y}_t), \forall t \in [0, T - \delta]$.

834 **Theorem 2.** [3, Theorem 1] Let Assumption 1 hold. Then there exists a numerical constant $C_0 > 0$,
835 such that

$$D_{KL}(p_\alpha(\cdot|y) || q_{T-\alpha}(\cdot|y)) \leq C_0(E_S + E_D + E_F) \quad (24)$$

836 where $E_D \leq \kappa^2 N u_2^2 + \kappa T u_2^2$ is the discretization error due to the reverse SDE, $E_F \leq \exp(-2T)u_2^2$
837 is the error due to the convergence of the forward SDE and E_S is the score estimation error

$$E_S(\theta_y) = \sum_{j=0}^{N-1} \gamma_j \mathbb{E}_{x \sim p_{T-t_j^{\leftarrow}}} \left[\left\| \nabla \log p_{T-t_j^{\leftarrow}}(x|y) - s_{T-t_j^{\leftarrow}, \theta_y}(x) \right\|_2^2 \right] \quad (25)$$

838 where $\gamma_j := t_{j+1}^{\leftarrow} - t_j^{\leftarrow}, \forall j = 0, 1, \dots, N-1$ is the step-size of the generation algorithm.

839 When the training is done over the forward discretization given by $(t_{N-j} = T - t_j^{\leftarrow})_{j=0}^{N-1}$, we have

$$\begin{aligned} E_S &= \sum_{j=0}^{N-1} \frac{\bar{\sigma}_{t_{N-j}} \lambda(t_{N-j})}{\lambda(t_{N-j}) \bar{\sigma}_{t_{N-j}}} (t_{N-j} - t_{N-j-1}) \mathbb{E}_{X_0} \mathbb{E}_{X_{t_{N-j}} | X_0} \left\| \bar{\sigma}_{t_{N-j}} s_{t_{N-j}, \theta_y}(X_{t_{N-j}}) + \xi \right\|^2 \\ &\quad + \sum_{j=0}^{N-1} \frac{\bar{\sigma}_{t_{N-j}}}{\lambda(t_{N-j})} \lambda(t_{N-j}) (t_{N-j} - t_{N-j-1}) C_{t_{N-j}} \\ &\leq 2 \max_j \frac{\bar{\sigma}_{t_{N-j}}}{\lambda(t_{N-j})} \mathcal{L}^y(\theta_y) \end{aligned}$$

840 where

$$\begin{aligned} \mathcal{L}^y(\theta_y) &= \frac{1}{2} \sum_{j=0}^{N-1} \mathbb{E}_{X_0} \mathbb{E}_{X_{t_{N-j}} | X_0} \left[\lambda(t_{N-j}) (t_{N-j} - t_{N-j-1}) \right. \\ &\quad \left. \left\| \nabla_{x(t_{N-j})} \log p_{t_{N-j}}(x(t_{N-j}) | x_0) - s_{t_{N-j}, \theta_y}(x(t_{N-j})) \right\|_2^2 \right] \\ &\quad + \frac{1}{2} \sum_{j=0}^{N-1} \lambda(t_{N-j}) (t_{N-j} - t_{N-j-1}) C_{t_{N-j}}(y) \end{aligned} \quad (26)$$

841 **Theorem 3.** (Appendix B and [3, Theorem 1]) Let Assumption 1 hold. Then there exists a numerical
 842 constant $C_0 > 0$, such that

$$D_{KL}(p_\alpha(\cdot|y)||q_{T-\alpha}(\cdot|y)) \leq C_0(\mathcal{L}^y(\theta_y) + E_D + E_F) \quad (27)$$

843 where $E_D \leq \kappa^2 Nu_2^2 + \kappa T u_2^2$ is the discretization error due to the reverse SDE, $E_F \leq \exp(-2T)u_2^2$
 844 is the error due to the convergence of the forward SDE.

845 B.1 Decomposition of $\mathcal{L}^y(\theta_y)$

846 Let $\theta_y^* = \operatorname{argmin}_{\theta_y} \mathcal{L}^y(\theta_y)$. We further decompose $\mathcal{L}^y(\theta_y)$ as

$$\begin{aligned} \max_{y \in Y} \left(\mathcal{L}^y(\theta_y) - \mathcal{L}^y(\theta_y^*) \right) &\leq \max_{y \in Y} \left(\mathcal{L}^y(\theta_y) + \beta \mathcal{L}_{mut}^y(\theta_y, \theta_{-y}) - (\mathcal{L}^y(\theta_y^*) + \beta \mathcal{L}_{mut}^y(\theta_y^*, \theta_{-y})) \right) \\ &\quad + \beta (\mathcal{L}_{mut}^y(\theta_y^*, \theta_{-y}) - \mathcal{L}_{mut}^y(\theta_y, \theta_{-y})) \\ &\stackrel{(a)}{\leq} \max_{y \in Y} \left(\mathcal{L}^y(\theta_y) + \beta \mathcal{L}_{mut}^y(\theta_y, \theta_{-y}) - (\mathcal{L}^y(B(\theta_{-y})) \right. \\ &\quad \left. + \beta \mathcal{L}_{mut}^y(B(\theta_{-y}), \theta_{-y})) \right) + \beta \max_{y \in Y} \mathcal{L}_{mut}^y(B(\theta_{-y}), \theta_{-y}) \\ &\leq \max_{y \in Y} \left(\mathcal{L}_{reg}^y(\theta_y, \theta_{-y}) - \mathcal{L}_{reg}^y(B(\theta_{-y}), \theta_{-y}) \right) \\ &\quad + \beta \max_{y \in Y} \sup_{(\theta_y, \theta_{-y})} \mathcal{L}_{mut}^y(\theta_y, \theta_{-y}) \end{aligned}$$

847 where (a) follows from the fact that $\mathcal{L}_{reg}^y(B(\theta_{-y}), \theta_{-y}) \leq \mathcal{L}_{reg}^y(\theta_y, \theta_{-y})$. We further decompose
 848 this to obtain an upper bound on $\max_{y \in Y} \min_t \mathcal{L}^y(\theta_y^t)$

$$\begin{aligned} \max_{y \in Y} \mathcal{L}^y(\theta_y) &\leq \max_{y \in Y} \mathcal{L}^y(\theta_y^*) + \max_{y \in Y} \left(\mathcal{L}_{reg}^y(\theta_y, \theta_{-y}) - \mathcal{L}_{reg}^y(B(\theta_{-y}), \theta_{-y}) \right) \\ &\quad + \beta \max_{y \in Y} \sup_{(\theta_y, \theta_{-y})} \mathcal{L}_{mut}^y(\theta_y, \theta_{-y}) \\ &\stackrel{(a)}{\leq} \max_{y \in Y} \mathcal{L}^y(\theta_y^*) + \max_{y \in Y} | \mathcal{L}_{reg}^y(\theta_y, \theta_{-y}) - \bar{\mathcal{L}}_{reg}^{n_y}(\theta_y, \theta_{-y}) | \\ &\quad + \max_{y \in Y} | \bar{\mathcal{L}}_{reg}^{n_y}(\theta_y, \theta_{-y}) - \bar{\mathcal{L}}_{reg}^{n_y}(B(\theta_{-y}), \theta_{-y}) | \\ &\quad + \max_{y \in Y} | \mathcal{L}_{reg}^y(B(\theta_{-y}), \theta_{-y}) - \bar{\mathcal{L}}_{reg}^{n_y}(B(\theta_{-y}), \theta_{-y}) | \\ &\quad + \beta \max_{y \in Y} \sup_{(\theta_y, \theta_{-y})} \mathcal{L}_{mut}^y(\theta_y, \theta_{-y}) \\ &\stackrel{(b)}{\leq} \max_{y \in Y} \mathcal{L}^y(\theta_y^*) + 2 \max_{y \in Y} \sup_{(\theta_y, \theta_{-y})} | \mathcal{L}_{reg}^y(\theta_y, \theta_{-y}) - \bar{\mathcal{L}}_{reg}^{n_y}(\theta_y, \theta_{-y}) | \\ &\quad + \text{NE-gap}(\theta_y, \theta_{-y}) + \beta \max_{y \in Y} \sup_{(\theta_y, \theta_{-y})} \mathcal{L}_{mut}^y(\theta_y, \theta_{-y}) \\ &\implies \min_{\tau \in [T_{train}]} \max_{y \in Y} \mathcal{L}^y(\theta_y^\tau) \stackrel{(c)}{\leq} \max_{y \in Y} \mathcal{L}^y(\theta_y^*) + \\ &\quad + 2 \max_{y \in Y} \sup_{(\theta_y, \theta_{-y})} | \mathcal{L}_{reg}^y(\theta_y, \theta_{-y}) - \bar{\mathcal{L}}_{reg}^{n_y}(\theta_y, \theta_{-y}) | \\ &\quad + \min_{\tau \in [T_{train}]} \text{NE-gap}(\theta_y^\tau, \theta_{-y}^\tau) + \beta \max_{y \in Y} \sup_{(\theta_y, \theta_{-y})} \mathcal{L}_{mut}^y(\theta_y, \theta_{-y}) \end{aligned}$$

849 where (a) follows from adding and subtracting the empirical losses $\bar{\mathcal{L}}_{reg}^{n_y}(\theta_y, \theta_{-y})$ and
 850 $\bar{\mathcal{L}}_{reg}^{n_y}(B(\theta_{-y}), \theta_{-y})$ and using triangle inequality of the max norm, (b) follows from the gradi-

ent domination property for strongly convex functions, (c) follows from taking the minimum over the iterates of the algorithm.

B.2 Boundedness of Forward Dynamics

Lemma 1. *Consider the forward diffusion process with linear drift coefficients. For any $\delta > 0, \delta \ll 1$, w.p. (with probability) of at least $1 - \delta$, we have*

$$\|x(t)\|_\infty \leq C_T \left(\|x(0)\|_\infty + \sqrt{\log \frac{2}{\pi\delta^2}} \right) \quad (28)$$

where $C_T := \max_{t \in [0, T]} r(t), r(t)v(t)$.

Proof: The proof is similar to [15, Lemma 1] When the drift coefficient $f(\cdot, t) : \mathbb{R}^d \rightarrow \mathbb{R}^d$ is linear in x i.e. $f(x, t) = -f(t)x$, the transition kernel $p_{t|0}$ has a closed form

$$p_{t|0}(x(t)|x(0)) = \mathcal{N}(x(t); \mu(t)x(0), \bar{\sigma}^2(t)I_d) \quad (29)$$

where $\mu(t) := \exp(\int_0^t f(s)ds), \bar{\sigma}^2(t) := 2 \int_0^t \exp(2\mu_s - 2\mu_t) \sigma_s^2 ds$. Together we get,

$$x(t) = \mu(t)x(0) + \bar{\sigma}(t)z, z \sim \mathcal{N}(0, I_d) \quad (30)$$

For any $\epsilon \sim \mathcal{N}(0, 1), c > 1$, we have

$$\mathbb{P}\{\epsilon : |\epsilon| > c\} = 2 \int_c^\infty \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}} dx \leq \frac{1}{\sqrt{2\pi}} \int_c^\infty 2xe^{-\frac{x^2}{2}} dx = \frac{1}{\sqrt{2\pi}} \int_{c^2}^\infty e^{-\frac{x}{2}} dx = \sqrt{\frac{2}{\pi}} e^{-\frac{c^2}{2}} \quad (31)$$

Let $\delta = \sqrt{\frac{2}{\pi}} e^{-\frac{c^2}{2}}$, then

$$\mathbb{P}\{\epsilon : |\epsilon| \leq \sqrt{\log \frac{2}{\pi\delta^2}}\} \geq 1 - \delta \quad (32)$$

Hence, for any $\delta \in (0, 1)$ with $\delta \ll 1$, w.p. at least $1 - \delta$, we have

$$\|x(t)\|_\infty \leq C_T \left(\|x(0)\|_\infty + \sqrt{\log \frac{2}{\pi\delta^2}} \right) \quad (33)$$

where $C_T := \max_{t \in [0, T]} \{\mu(t), \bar{\sigma}(t)\}$. Let $C_{T, \delta} = C_T(K + \sqrt{\log \frac{2}{\pi\delta^2}})$

B.3 Boundedness of Loss function $\bar{\mathcal{L}}_{reg}^{n_y}(\theta_y, \theta_{-y})$

In this section, study some properties of the game defined by $\langle \mathcal{Y}, (\bar{\mathcal{L}}_{reg}^{n_y})_{y \in \mathcal{Y}}, (\Theta_y)_{y \in \mathcal{Y}} \rangle$. From Eq. 10, we have

$$\mathcal{L}_{conti, reg}^y(\theta_y, \theta_{-y}) = \mathcal{L}_{conti}^y(\theta_y) + \beta \mathcal{L}_{conti, mut}^y(\theta_y, \theta_{-y}) \quad (34)$$

where

$$\mathcal{L}_{conti, mut}^y(\theta_y, \theta_{-y}, \omega(\cdot)) = \frac{1}{2} \int_{t_0}^T \omega(t) \mathbb{E}_{x(t) \sim p_t} \mathbb{E}_{y' \sim Q} \left[\left\| s_{t, \theta_y}(x(t)) - s_{t, \theta_{y'}}(x(t)) \right\|_2^2 \right] dt$$

and

$$\mathcal{L}_{conti}^y(\theta_y, \theta_{-y}) = \frac{1}{2} \int_{t_0}^T \lambda(t) \mathbb{E}_{(x(t), y)} \left[\left\| \nabla_{x(t)} \log p_t(x(t)|y) - s_{t, \theta}(x(t), y) \right\|_2^2 \right] dt$$

Conditioning on X_0 and using law of iterated expectation, we can write [29, Appendix A], we get

$$\begin{aligned} \mathcal{L}_{conti, reg}^y(\theta_y, \theta_{-y}) &= \frac{1}{2} \int_{t_0}^T \mathbb{E}_{X_0} \mathbb{E}_{X_t|X_0, y} \left[\lambda(t) \left\| s_{t, \theta_y}(x(t)) - \nabla_{x(t)} \log p_t(x(t)|x_0) \right\|_2^2 \right. \\ &\quad \left. + \beta \omega(t) \mathbb{E}_{y' \sim Q} \left[\left\| s_{t, \theta_y}(x(t)) - s_{t, \theta_{y'}}(x(t)) \right\|_2^2 \right] \right] dt + \frac{1}{2} \int_{t_0}^T \left[\lambda(t) C_t(y) \right] dt \end{aligned}$$

870 where $C_t(y) = \mathbb{E}_{X_t} \|\nabla \log p_t(X_t|y)\|^2 - \mathbb{E}_{X_0} \mathbb{E}_{X_t|X_0} \|\nabla \log p_t(X_t|X_0, y)\|^2$ to learn the score
 871 $\nabla_{x(t)} \log p_t(x(t)|x_0, y)$.

872 Furthermore, we discretize the time points $0 = t_0 < t_1 < \dots < t_N = T$ to the objective function

$$\begin{aligned} \mathcal{L}_{reg}^y(\theta_y, \theta_{-y}) &= \mathcal{L}^y(\theta_y) + \beta \mathcal{L}_{mut}^y(\theta_y, \theta_{-y}) \\ &= \frac{1}{2} \sum_{j=1}^N \lambda(t_j)(t_j - t_{j-1}) \mathbb{E}_{X_0} \mathbb{E}_{X_{t_j}|X_0} \left[\left\| \nabla_{x(t_j)} \log p_t(x_i(t_j)|x_0) - s_{t_j, \theta_y}(x_i(t_j)) \right\|_2^2 \right] + \\ &+ \bar{C}(y) + \beta \frac{1}{2} \sum_{j=1}^N \omega(t_j)(t_j - t_{j-1}) \mathbb{E}_{X_0} \mathbb{E}_{X_{t_j}|X_0} \mathbb{E}_{y' \sim Q} \left[\left\| s_{t_j, \theta_y}(x_i(t_j)) - s_{t, \theta_{y'}}(x_i(t_j)) \right\|_2^2 \right] \end{aligned} \quad (35)$$

873 where $\bar{C}(y) = \frac{1}{2} \sum_{j=1}^N \lambda(t_j)(t_j - t_{j-1}) C_{t_j}(y)$ From [29, Appendix A], we have $X_t|X_0 \sim$
 874 $\mathcal{N}(e^{-\mu t} X_0, \bar{\sigma}_t^2 I)$ and its density function is

$$p_t(x|x_0) = (2\pi\bar{\sigma}_t^2)^{-\frac{d}{2}} \exp\left(-\frac{\|x - e^{-\mu t} x_0\|^2}{2\bar{\sigma}_t^2}\right)$$

875 Then,

$$\begin{aligned} \Delta &= \mathbb{E}_{X_0} \mathbb{E}_{X_t|X_0} \left\| s_{t_j, \theta_y}(x_i(t_j)) - \nabla_{x(t_j)} \log p_t(x(t_j)|x_0) \right\| \\ &= \mathbb{E}_{X_0} \mathbb{E}_{X_t|X_0} \left\| s_{t_j, \theta_y}(x_i(t_j)) - \nabla_x \left(-\frac{\|X_t - e^{-\mu t} X_0\|^2}{2\bar{\sigma}_t^2} \right) \right\|^2 \\ &= \mathbb{E}_{X_0} \mathbb{E}_{X_t|X_0} \left\| s_{t_j, \theta_y}(x_i(t_j)) + \frac{X_t - e^{\mu t} X_0}{\bar{\sigma}_t^2} \right\|^2 \\ &= \mathbb{E}_{X_0} \mathbb{E}_{\epsilon_t} \left\| s_{t_j, \theta_y}(x_i(t_j)) + \frac{\epsilon_t}{\bar{\sigma}_t^2} \right\|^2 \end{aligned}$$

876 Let $\xi = \frac{\epsilon_t}{\bar{\sigma}_t} \sim \mathcal{N}(0, I)$

$$\Delta = \frac{1}{\bar{\sigma}_t} \mathbb{E}_{X_0} \mathbb{E}_{\xi} \left\| \bar{\sigma}_t s_{t_j, \theta_y}(x_i(t_j)) + \xi \right\|^2 \quad (36)$$

877 Finally putting all of it together, we get the empirical loss function

$$\begin{aligned} \bar{\mathcal{L}}_{reg}^{n_y}(\theta_y, \theta_{-y}) &= \frac{1}{2n_y} \sum_{i=1}^{n_y} \sum_{j=1}^N \frac{\lambda(t_j)(t_j - t_{j-1})}{\bar{\sigma}_{t_j}} \left[\left\| \bar{\sigma}_{t_j} s_{t_j, \theta_y}(x_i(t_j)) + \xi_{ij} \right\|_2^2 \right. \\ &\quad \left. + \beta \omega(t_j)(t_j - t_{j-1}) \mathbb{E}_{y' \sim Q} \left[\left\| s_{t_j, \theta_y}(x_i(t_j)) - s_{t, \theta_{y'}}(x_i(t_j)) \right\|_2^2 \right] \right] \end{aligned} \quad (37)$$

878 We will show that the empirical loss function for the label $y \in \mathcal{Y}$, $\bar{\mathcal{L}}_{reg}^{n_y}(\theta_y, \theta_{-y})$ that is optimized is
 879 convex and smooth in θ_y with high probability.

880 **Lemma 2.** For $\delta > 0, \delta \ll 1$, wp. $1 - n_y N \delta$, the empirical loss function

$$\begin{aligned} \bar{\mathcal{L}}_{reg}^{n_y}(\theta_y, \theta_{-y}) &= \frac{1}{2n_y} \sum_{i=1}^{n_y} \sum_{j=1}^N \frac{\lambda(t_j)(t_j - t_{j-1})}{\bar{\sigma}_{t_j}} \left[\left\| \bar{\sigma}_{t_j} s_{t_j, \theta_y}(x_i(t_j)) + \xi_{ij} \right\|_2^2 \right. \\ &\quad \left. + \beta \omega(t_j)(t_j - t_{j-1}) \mathbb{E}_{y' \sim Q} \left[\left\| s_{t_j, \theta_y}(x_i(t_j)) - s_{t, \theta_{y'}}(x_i(t_j)) \right\|_2^2 \right] \right] \end{aligned} \quad (38)$$

881 is bounded i.e.

$$\bar{\mathcal{L}}_{reg}^{n_y}(\theta_y, \theta_{-y}) = \mathcal{O}\left(\sum_{j=1}^N \frac{\lambda(t_j)(t_j - t_{j-1})}{\bar{\sigma}_{t_j}} + \beta \omega(t_j)(t_j - t_{j-1})\right)$$

882 **Proof:** From Lemma 1, we have $\delta > 0, \delta \ll 1$

$$\mathbb{P}\{|\xi_{ij}| > \sqrt{\frac{2}{\pi\delta^2}}\} \leq \delta \quad (39)$$

883 Thus, w.p. $1 - n_y N \delta$, we have $|\xi_{ij}| \leq \sqrt{\frac{2}{\pi\delta^2}}$ and hence we have $\|x(t_j)\|_\infty \leq C_{t_N, \delta}, \forall i = 1, \dots, n_y$
 884 and $j = 1, \dots, N$

885 Thus, w.p. $1 - n_y N \delta$

$$\begin{aligned} \bar{\mathcal{L}}_{reg}^{n_y}(\theta_y, \theta_{-y}) &= \frac{1}{2n_y} \sum_{i=1}^{n_y} \sum_{j=1}^N \frac{\lambda(t_j)(t_j - t_{j-1})}{\bar{\sigma}_{t_j}} \left[\|\bar{\sigma}_{t_j} s_{t_j, \theta_y}(x_i(t_j)) + \xi_{ij}\|_2^2 \right. \\ &\quad \left. + \beta \omega(t_j)(t_j - t_{j-1}) \mathbb{E}_{y' \sim Q} \left[\left\| s_{t_j, \theta_y}(x_i(t_j)) - s_{t, \theta_{y'}}(x_i(t_j)) \right\|_2^2 \right] \right] \\ &\leq \frac{1}{n_y} \sum_{i=1}^{n_y} \sum_{j=1}^N \frac{\lambda(t_j)(t_j - t_{j-1})}{\bar{\sigma}_{t_j}} (\bar{\sigma}_{t_j}^2 \|s_{t_j, \theta_y}(x_i(t_j))\|_2^2 + \|\xi_{ij}\|_2^2) \\ &\quad + \beta \omega(t_j)(t_j - t_{j-1}) (\|s_{t_j, \theta_y}(x_i(t_j))\|_2^2 + \max_{y' \in \mathcal{Y}} \|s_{t_j, \theta_{y'}}(x_i(t_j))\|_2^2) \end{aligned}$$

886 For a bound on $\|s_{t_j, \theta_y}(x(t_j))\|_2$

$$\|s_{t_j, \theta_y}(x(t_j))\|_2 = \left\| \frac{1}{m} \sum_{i=1}^m a_{y,i} \sigma(w_{y,i}^T x(t_j) + u_{y,i}^T e(t_j)) \right\|_2 \quad (40)$$

$$\stackrel{(a)}{\leq} \frac{1}{m} \sum_{i=1}^m \|a_{y,i}\|_2 |\sigma(w_{y,i}^T x(t_j) + u_{y,i}^T e(t_j))| \quad (41)$$

$$\stackrel{(b)}{\leq} \frac{1}{m} \sum_{i=1}^m \|a_{y,i}\|_2 (\|w_{y,i}\|_1 \|x(t_j)\|_\infty + \|u_{y,i}\|_1 \|e(t_j)\|_\infty) \quad (42)$$

$$\leq \frac{1}{m} \sum_{i=1}^m \|a_{y,i}\|_2 (C_{t_N, \delta} \|w_{y,i}\|_1 + \max_j \|e(t_j)\|_\infty \|u_{y,i}\|_1) \quad (43)$$

$$\stackrel{(c)}{\leq} (C_{t_n, \delta} + C_{t_N, e}) C_{w_y, u_y} B \quad (44)$$

887 where (a) follows from triangle inequality for norms, (b) follows from the fact that the ReLU function
 888 satisfies $|\sigma(x)| \leq |x|$ and Holder inequality and (c) follows from the bounds on the embeddings and
 889 $x(t_j)$ with $\|w_{y,i}\|_1, \|u_{y,i}\|_1 \leq C_{w_y, u_y}, \forall i \in [m]$. Thus, for $\delta > 0, \delta \ll 1$, we have w.p. $1 - n_y N \delta$

$$\bar{\mathcal{L}}_{reg}^{n_y}(\theta_y, \theta_{-y}) \leq C_1 \sum_{j=1}^N \frac{\lambda(t_j)(t_j - t_{j-1})}{\bar{\sigma}_{t_j}} + \beta \omega(t_j)(t_j - t_{j-1}) \quad (45)$$

890 where $C_1 = (\bar{\sigma}_{t_N}^2 + 2)(C_{t_n, \delta} + C_{t_N, e})^2 C_{w_y, u_y}^2 B^2 + \frac{2}{\pi\delta^2}$. Since $\bar{\sigma}_{t_j}$ is non-decreasing in j , so
 891 $\max_j \bar{\sigma}_{t_j} = \bar{\sigma}_{t_N}$.

$$\begin{aligned}
& \|\nabla_{A_y} \bar{\mathcal{L}}_{reg}^{n_y}(\theta_y, \theta_{-y})\|_F^2 = \sum_{k=1}^d \|\nabla_{(A_y)_k} \bar{\mathcal{L}}_{reg}^{n_y}(\theta_y, \theta_{-y})\|^2 \\
& = \sum_{k=1}^d \left\| \frac{1}{n_y} \sum_{i=1}^{n_y} \sum_{j=1}^N \lambda(t_j)(t_j - t_{j-1})(\bar{\sigma}_{t_j} s_{t_j, \theta_y}(x(t_j)) + \xi_{ij})_k \sigma(W_y x(t_j) + U_y e(t_j)) \right. \\
& \quad \left. + \beta \omega(t_j)(t_j - t_{j-1}) \mathbb{E}_{y'}[(s_{t_j, \theta_y}(x(t_j)) - s_{t_j, \theta_{y'}}(x(t_j)))_k \sigma(W_y x(t_j) + U_y e(t_j))] \right\|^2 \\
& \leq 2 \sum_{k=1}^d \left\| \frac{1}{n_y} \sum_{i=1}^{n_y} \sum_{j=1}^N \lambda(t_j)(t_j - t_{j-1})(\bar{\sigma}_{t_j} s_{t_j, \theta_y}(x(t_j)) + \xi_{ij})_k \sigma(W_y x(t_j) + U_y e(t_j)) \right\|^2 \\
& \quad + 2\beta^2 \sum_{k=1}^d \left\| \frac{1}{n_y} \sum_{i=1}^{n_y} \sum_{j=1}^N \omega(t_j)(t_j - t_{j-1}) \mathbb{E}_{y'}[(s_{t_j, \theta_y}(x(t_j)) - s_{t_j, \theta_{y'}}(x(t_j)))_k \right. \\
& \quad \left. \sigma(W_y x(t_j) + U_y e(t_j))] \right\|^2 \\
& \leq 2 \frac{N}{n_y} \sum_{i=1}^{n_y} \sum_{j=1}^N \lambda(t_j)^2 (t_j - t_{j-1})^2 \sum_{k=1}^d \|\bar{\sigma}_{t_j} s_{t_j, \theta_y}(x(t_j)) + \xi_{ij}\|^2 \|\sigma(W_y x(t_j) + U_y e(t_j))\|^2 \\
& \quad + 2\beta^2 \frac{N}{n_y} \sum_{i=1}^{n_y} \sum_{j=1}^N \omega(t_j)^2 (t_j - t_{j-1})^2 \\
& \quad \sum_{k=1}^d \mathbb{E}_{y'}[\|s_{t_j, \theta_y}(x(t_j)) - s_{t_j, \theta_{y'}}(x(t_j))\|^2] \|\sigma(W_y x(t_j) + U_y e(t_j))\|^2 \\
& \leq 4Nd \|\sigma(W_y x(t_j) + U_y e(t_j))\|_2^2 \max_j \{\lambda(t_j)(t_j - t_{j-1}) \bar{\sigma}_{t_j}, \beta \omega(t_j)(t_j - t_{j-1})\} \bar{\mathcal{L}}_{reg}^{n_y}(\theta_y, \theta_{-y}) \\
& \leq 4Nd^2 (C_{t_N, \delta} + C_{t_N, e})^2 C_{w_y, u_y}^2 \max_j \{\lambda(t_j)(t_j - t_{j-1}) \bar{\sigma}_{t_j}, \beta \omega(t_j)(t_j - t_{j-1})\} \bar{\mathcal{L}}_{reg}^{n_y}(\theta_y, \theta_{-y})
\end{aligned}$$

893 Since $w.p.1 - n_y N \delta$ the empirical loss function $\bar{\mathcal{L}}_{reg}^{n_y}(\theta_y, \theta_{-y})$ is bounded, $\|\nabla_{A_y} \bar{\mathcal{L}}_{reg}^{n_y}(\theta_y, \theta_{-y})\|_F^2$
894 is bounded with the same probability.

895 This also shows that for fixed θ_{-y} , $(W_y, U_y)_{y \in \mathcal{Y}}$, $w.p.1 - n_y N \delta$, $\bar{\mathcal{L}}_{reg}^{n_y}(\theta_y, \theta_{-y})$ is a Lipschitz function
896 in θ_y with Lipschitz constant σ_y such that $\sigma_y^2 = 4C_1 N d^2 (C_{t_N, \delta} + C_{t_N, e})^2 C_{w_y, u_y}^2 \max_j \{\lambda(t_j)(t_j -$
897 $t_{j-1}) \bar{\sigma}_{t_j}, \beta \omega(t_j)(t_j - t_{j-1})\} \left(\sum_{j=1}^N \frac{\lambda(t_j)(t_j - t_{j-1})}{\bar{\sigma}_{t_j}} + \beta \omega(t_j)(t_j - t_{j-1}) \right)$

898 **B.5 Smoothness of Loss Function $\bar{\mathcal{L}}_{reg}^{n_y}(\theta_y, \theta_{-y})$**

899 **Lemma 3.** Let $(W_y, U_y)_{y \in \mathcal{Y}}, \theta_{-y}, \{t_j\}_{j=1}^N$ be fixed. Let $L_y = d(C_{t_N, \delta} +$
900 $C_{t_N, e})^2 C_{w_y, u_y}^2 \sum_{j=1}^N \left(\lambda(t_j)(t_j - t_{j-1}) \bar{\sigma}_{t_j} + \beta \omega(t_j)(t_j - t_{j-1}) \right)$. Then for $\delta > 0, \delta \ll 1$,
901 $w.p. 1 - n_y N \delta$, $\bar{\mathcal{L}}_{reg}^{n_y}(\theta_y, \theta_{-y})$ is $\frac{L_y}{m^2}$ smooth and convex in θ_y .

902 **Proof** We have,

$$\begin{aligned}
\bar{\mathcal{L}}_{reg}^{n_y}(\theta_y, \theta_{-y}) &= \frac{1}{2n_y} \sum_{i=1}^{n_y} \sum_{j=1}^N \frac{\lambda(t_j)(t_j - t_{j-1})}{\bar{\sigma}_{t_j}} \left[\|\bar{\sigma}_{t_j} s_{t_j, \theta_y}(x(t_j)) + \xi_{ij}\|_2^2 \right. \\
& \quad \left. + \beta \omega(t_j)(t_j - t_{j-1}) \mathbb{E}_{y' \sim Q} \left[\|s_{t_j, \theta_y}(x(t_j)) - s_{t_j, \theta_{y'}}(x(t_j))\|_2^2 \right] \right]
\end{aligned} \tag{46}$$

903 To show smoothness, we will show that the function $f(\theta_y) = \|\bar{\sigma}_{t_j} s_{t_j, \theta_y}(x(t_j)) + \xi_{ij}\|_2^2$ and
 904 $g(\theta_y) = \left\| s_{t_j, \theta_y}(x(t_j)) - s_{t_j, \theta_{y'}}(x(t_j)) \right\|_2^2$ are individually smooth. Once we prove this, it is easy
 905 to show $\bar{\mathcal{L}}_{reg}^{n_y}(\theta_y, \theta_{-y})$ is smooth as the linear combination of smooth functions is again smooth.
 906 To show smoothness, we need to show that $\|\nabla_{\theta_y}^2 f(\theta_y)\|$ and $\|\nabla_{\theta_y}^2 g(\theta_y)\|$ have a bounded norm.
 907 Recall that $s_{t, \theta_y}(x) = \frac{1}{m} A_y \sigma(W_y x(t) + U_y e(t))$. Let $h_1(x, t) := \sigma(W_y x + U_y e(t))$, $h_2(x, t) :=$
 908 $s_{t, \theta_{y'}}(x)$, $h_3(i, j) = \xi_{ij}$, we have

$$f(\theta_y) = \|\bar{\sigma}_{t_j} s_{t_j, \theta_y}(x(t_j)) + \xi_{ij}\|_2^2 \quad (47)$$

$$= \frac{\bar{\sigma}_{t_j}^2}{m^2} h_1^T(x(t_j), t_j) A_y^T A_y h_1(x(t_j), t_j) - 2\bar{\sigma}_{t_j} h_3^T(i, j) \left(\frac{A_y}{m} \right) h_1(x(t_j), t_j) \quad (48)$$

$$+ h_3^T(i, j) h_3(i, j) \quad (49)$$

$$\stackrel{a}{=} \frac{\bar{\sigma}_{t_j}^2}{m^2} \text{trace}(A_y^T A_y B_1) - \frac{2\bar{\sigma}_{t_j}}{m} \text{trace}(A_y B_3) + \text{constant} \quad (50)$$

$$\stackrel{b}{=} \frac{\bar{\sigma}_{t_j}^2}{m^2} \text{vec}(A_y)^T (B_1 \otimes I) \text{vec}(A_y) - \frac{2\bar{\sigma}_{t_j}}{m} \text{vec}(B_3^T)^T \text{vec}(A_y) + \text{constant} \quad (51)$$

909 where (a) follows from the identity $x^T A y = \text{trace}(B y x^T)$, (b) follows from the following identities

$$\text{trace}(A^T A B) = \text{trace}(A B A^T) = \text{vec}(A)^T (B \otimes I) \text{vec}(A)$$

$$\text{trace}(A B) = \text{vec}(A)^T \text{vec}(B^T)$$

910 and $B_3 = h_1(x(t_j), t_j) h_3^T(i, j)$.

911 Similarly, we have for $g(\theta_y)$

$$\begin{aligned} g(\theta_y) &= \left\| s_{t, \theta_y}(x(t_j)) - s_{t, \theta_{y'}}(x(t_j)) \right\|_2^2 \\ &= \frac{1}{m^2} h_1^T(x(t_j), t_j) A_y^T A_y h_1(x(t_j), t_j) - 2h_2^T(x(t_j), t_j) \left(\frac{A_y}{m} \right) h_1(x(t_j), t_j) \\ &\quad + h_2^T(x(t_j), t_j) h_2(x(t_j), t_j) \\ &\stackrel{a}{=} \frac{1}{m} \text{trace}(A_y^T A_y B_1) - \frac{2}{m} \text{trace}(A_y B_2) + \text{constant} \\ &\stackrel{b}{=} \frac{1}{m^2} \text{vec}(A_y)^T (B_1 \otimes I) \text{vec}(A_y) - \frac{2}{m} \text{vec}(B_2^T)^T \text{vec}(A) + \text{constant} \end{aligned}$$

912 where $B_1 := h_1(x(t_j), t_j) h_1^T(x(t_j), t_j)$ and $B_2 := h_1(x(t_j), t_j) h_2^T(x(t_j), t_j)$. Thus,

$$\frac{1}{\bar{\sigma}_{t_j}^2} \nabla_{\theta_y}^2 f(\theta_y) = \nabla_{\theta_y}^2 g(\theta_y) = \nabla_{\text{vec}(A_y)}^2 g(\theta_y) = \frac{2}{m^2} (B_1 \otimes I) \quad (52)$$

913 The eigenvalues of $(B_1 \otimes I)$ is the same as B_1 with multiplicity. Thus, to show smoothness, we
 914 need to bound the maximum eigenvalues of B_1 . For any $v \in \mathbb{R}^m$

$$0 \leq v^T B_1 v = (v^T h_1(x(t_j), t_j))^2 \leq d \|h_1(x(t_j), t_j)\|_\infty^2 v^T v \quad (53)$$

915 Now,

$$\|\sigma(W_y x(t_j) + U_y e(t_j))\|_\infty = \max_{i=1, \dots, m} \sigma(w_{y,i}^T x(t_j) + u_{y,i}^T e(t_j)) \quad (54)$$

$$\leq \max_{i=1, \dots, m} |w_{y,i}^T x(t_j) + u_{y,i}^T e(t_j)| \quad (55)$$

$$\leq \max_{i=1, \dots, m} \|w_{y,i}\|_1 \|x(t_j)\|_\infty + \|u_{y,i}\|_1 \|e(t_j)\|_\infty \quad (56)$$

$$\leq (C_{t_N, \delta} + C_{t_N, e}) C_{w_y, u_y} \quad (57)$$

916 Thus, we have for any $v \in \mathbb{R}^m$

$$0 \leq v^T B_1 v \leq d(C_{t_N, \delta} + C_{t_N, e})^2 C_{w_y, u_y}^2 v^T v \quad (58)$$

917 Since $w.p. 1 - n_y N \delta$ we have $\{\|x_{ij}\|_\infty \leq C_{t_N, \delta}\}_{i=1, j=1}^{n_y, N}$, we have with the same probability $f(\theta_y)$
 918 and $g(\theta_y)$ are smooth in θ_y for every $W_y, U_y, x(t_j), \theta_{-y}$.

919 Thus, $\bar{\mathcal{L}}_{reg}^{n_y}(\theta_y, \theta_{-y})$ is $L_y \frac{1}{m^2} = \frac{1}{m^2} d(C_{t_N, \delta} + C_{t_N, e})^2 C_{w_y, u_y}^2 \sum_{j=1}^N \left(\lambda(t_j)(t_j - t_{j-1}) \bar{\sigma}_{t_j} + \right.$
 920 $\left. \beta \omega(t_j)(t_j - t_{j-1}) \right)$ smooth.

921 B.6 Proof: First order convergence of the algorithm

922 **Proof** Our proof follows closely along the lines of [11]. Let $\bar{\mathcal{L}}_{reg}^{n_y}(\theta_y, \theta_{-y})$ be the empirical version
 923 of $\mathcal{L}_{reg}^y(\theta_y, \theta_{-y})$ with n_y samples. By L_y smoothness of $\bar{\mathcal{L}}_{reg}^{n_y}(\theta_y, \theta_{-y})$ we have, for any $y \in \mathcal{Y}$,

$$\begin{aligned} \bar{\mathcal{L}}_{reg}^{n_y}(\theta_y^{\tau+1}, \theta_{-y}^\tau) &\leq \bar{\mathcal{L}}_{reg}^{n_y}(\theta_y^\tau, \theta_{-y}^\tau) + \langle \nabla_{\theta_y} \bar{\mathcal{L}}_{reg}^{n_y}(\theta_y^\tau, \theta_{-y}^\tau), \theta_y^{\tau+1} - \theta_y^\tau \rangle \\ &\quad + \frac{L_y}{2} \|\theta_y^{\tau+1} - \theta_y^\tau\|^2 \end{aligned} \quad (59)$$

$$\begin{aligned} \implies \bar{\mathcal{L}}_{reg}^{n_y}(\theta_y^{\tau+1}, \theta_{-y}^\tau) &\leq \bar{\mathcal{L}}_{reg}^{n_y}(\theta_y^\tau, \theta_{-y}^\tau) - \eta_\tau \|\nabla_{\theta_y} \bar{\mathcal{L}}_{reg}^{n_y}(\theta_y^\tau, \theta_{-y}^\tau)\|^2 \\ &\quad + \frac{L_y}{2} \eta_\tau^2 \|\nabla_{\theta_y} \bar{\mathcal{L}}_{reg}^{n_y}(\theta_y^\tau, \theta_{-y}^\tau)\|^2 \end{aligned} \quad (60)$$

$$\begin{aligned} \implies \eta_\tau \|\nabla_{\theta_y} \bar{\mathcal{L}}_{reg}^{n_y}(\theta_y^\tau, \theta_{-y}^\tau)\|^2 &\leq \bar{\mathcal{L}}_{reg}^{n_y}(\theta_y^\tau, \theta_{-y}^\tau) - \bar{\mathcal{L}}_{reg}^{n_y}(\theta_y^{\tau+1}, \theta_{-y}^\tau) \\ &\quad + \frac{L_y}{2} \eta_\tau^2 \|\nabla_{\theta_y} \bar{\mathcal{L}}_{reg}^{n_y}(\theta_y^\tau, \theta_{-y}^\tau)\|^2 \end{aligned} \quad (61)$$

$$\begin{aligned} \implies \sum_{\tau=1}^{T_{train}} \eta_\tau \|\nabla_{\theta_y} \bar{\mathcal{L}}_{reg}^{n_y}(\theta_y^\tau, \theta_{-y}^\tau)\|^2 &\leq \bar{\mathcal{L}}_{reg}^{n_y}(\theta_y^1) - \bar{\mathcal{L}}_{reg}^{n_y}(\theta_y^{\tau+1}) + \beta \sum_{\tau=1}^{T_{train}} \psi(\theta_y^{\tau+1}, \theta_y^\tau, \theta_{-y}^\tau) \\ &\quad + \sum_{\tau=1}^{T_{train}} \frac{L_y}{2} \eta_\tau^2 \|\nabla_{\theta_y} \bar{\mathcal{L}}_{reg}^{n_y}(\theta_y^\tau, \theta_{-y}^\tau)\|^2 \end{aligned} \quad (62)$$

$$\begin{aligned} \implies \sum_{\tau=1}^{T_{train}} \eta_\tau \|\nabla_{\theta_y} \bar{\mathcal{L}}_{reg}^{n_y}(\theta_y^\tau, \theta_{-y}^\tau)\|^2 &\leq \bar{\mathcal{L}}_{reg}^{n_y}(\theta_y^1) - \bar{\mathcal{L}}_{reg}^{n_y}(\theta_y^{\tau+1}) + \sum_{\tau=1}^{T_{train}} \frac{L_y}{2} \eta_\tau^2 \sigma_y^2 \\ &\quad + \beta \sum_{\tau=1}^{T_{train}} \psi(\theta_y^{\tau+1}, \theta_y^\tau, \theta_{-y}^\tau) \end{aligned} \quad (63)$$

$$\begin{aligned} \implies \min_{\tau \in [T_{train}]} \|\nabla_{\theta_y} \bar{\mathcal{L}}_{reg}^{n_y}(\theta_y^\tau, \theta_{-y}^\tau)\|^2 &\leq \frac{\bar{\mathcal{L}}_{reg}^{n_y}(\theta_y^1) - \bar{\mathcal{L}}_{reg}^{n_y}(\theta_y^*) + \frac{L_y}{2} \sigma_y^2 \sum_{\tau=1}^{T_{train}} \eta_\tau^2}{\sum_{\tau=1}^{T_{train}} \eta_\tau} \\ &\quad + \beta \frac{\sum_{t=1}^{T_{train}} \psi(\theta_y^{\tau+1}, \theta_y^\tau, \theta_{-y}^\tau)}{\sum_{\tau=1}^T \eta_\tau} \end{aligned} \quad (64)$$

924 where $\psi(\theta_y^{\tau+1}, \theta_y^\tau, \theta_{-y}^\tau) = \bar{\mathcal{L}}_{mut}^{n_y}(\theta_y^\tau, \theta_{-y}^\tau) - \bar{\mathcal{L}}_{mut}^{n_y}(\theta_y^{\tau+1}, \theta_{-y}^\tau)$.

925 B.6.1 Analyzing the Bias Term

926 **Lemma 4.** Suppose $\theta_{-y}, (W_y, U_y)_{y \in \mathcal{Y}}$ are fixed. Let $\phi_y = d^{1.5} N (C_{t_N, \delta} +$
 927 $C_{t_N, e})^2 C_{w_y, u_y}^2 B \max_j \omega(t_j)(t_j - t_{j-1})$. Then for $\delta > 0, \delta \ll 1, w.p. 1 - n_y N \delta$, we have

$$\bar{\mathcal{L}}_{mut}^{n_y}(\theta_y, \theta_{-y}) = \frac{1}{2n_y} \sum_{i=1}^{n_y} \sum_{j=1}^N \omega(t_j)(t_j - t_{j-1}) \mathbb{E}_{y' \sim Q} \left[\left\| s_{t_j, \theta_y}(x(t_j)) - s_{t, \theta_{y'}}(x(t_j)) \right\|_2^2 \right] \quad (65)$$

928 is ϕ_y Lipschitz in θ_y .

Proof:

$$\begin{aligned}
& \|\nabla_{A_y} \bar{\mathcal{L}}_{reg,mut}^{n_y}(\theta_y, \theta_{-y})\|_F^2 = \sum_{k=1}^d \|\nabla_{(A_y)_k} \bar{\mathcal{L}}_{mut}^{n_y}(\theta_y, \theta_{-y})\|^2 \\
& = \sum_{k=1}^d \left\| \frac{1}{n_y} \sum_{i=1}^{n_y} \sum_{j=1}^N \omega(t_j) (t_j - t_{j-1}) \mathbb{E}_{y'}[(s_{t_j, \theta_y}(x(t_j)) - s_{t_j, \theta_{y'}}(x(t_j)))_k \sigma(W_y x(t_j) + U_y e(t_j))] \right\|^2 \\
& \leq \frac{N}{n_y} \sum_{i=1}^{n_y} \sum_{j=1}^N \omega(t_j)^2 (t_j - t_{j-1})^2 \\
& \quad \sum_{k=1}^d \mathbb{E}_{y'}[\|s_{t_j, \theta_y}(x(t_j)) - s_{t_j, \theta_{y'}}(x(t_j))\|^2] \|\sigma(W_y x(t_j) + U_y e(t_j))\|^2 \\
& \leq \|\sigma(W_y x(t_j) + U_y e(t_j))\|^2 N d \max_j \omega(t_j) (t_j - t_{j-1}) \bar{\mathcal{L}}_{mut}^{n_y}(\theta_y, \theta_{-y}) \\
& \leq 4d^2 N (C_{t_N, \delta} + C_{t_N, e})^2 C_{w_y, u_y}^2 \max_j \omega(t_j) (t_j - t_{j-1}) \bar{\mathcal{L}}_{mut}^{n_y}(\theta_y, \theta_{-y}) \\
& \leq d^3 N (C_{t_N, \delta} + C_{t_N, e})^4 C_{w_y, u_y}^4 B^2 \max_j \omega(t_j) (t_j - t_{j-1}) \sum_{j=1}^N \omega(t_j) (t_j - t_{j-1})
\end{aligned}$$

929 Since $\bar{\mathcal{L}}_{mut}^{n_y}(\theta_y, \theta_{-y}) \leq \bar{\mathcal{L}}_{reg}^{n_y}(\theta_y, \theta_{-y})$ and $w.p.1 - n_y N \delta$, $\bar{\mathcal{L}}_{reg}^{n_y}(\theta_y, \theta_{-y})$ is bounded. Thus,
930 $\|\nabla_{A_y} \bar{\mathcal{L}}_{mut}^{n_y}(\theta_y, \theta_{-y})\|_F^2$ is bounded and hence $\bar{\mathcal{L}}_{mut}^{n_y}(\theta_y, \theta_{-y})$ is Lipschitz in θ_y with $\phi_y =$
931 $d^{1.5} N (C_{t_N, \delta} + C_{t_N, e})^2 C_{w_y, u_y}^2 B \max_j \omega(t_j) (t_j - t_{j-1})$ Here,

$$\psi(\theta_y^{\tau+1}, \theta_y^\tau, \theta_{-y}^\tau) = \left| \bar{\mathcal{L}}_{mut}^{n_y}(\theta_y^\tau, \theta_{-y}^\tau) - \bar{\mathcal{L}}_{mut}^{n_y}(\theta_y^{\tau+1}, \theta_{-y}^\tau) \right| \quad (66)$$

$$\leq \phi_y \|\theta_y^\tau - \theta_y^{\tau+1}\| \quad (67)$$

$$\leq \phi_y \eta_t \|\nabla_{A_y} \bar{\mathcal{L}}_{reg}^{n_y}(\theta_y^\tau, \theta_{-y}^\tau)\| \quad (68)$$

$$\leq \phi_y \eta_\tau \sigma_y \quad (69)$$

932 By taking $\eta_\tau \leq \frac{m^2}{\max_{y \in \mathcal{Y}} L_y \sqrt{T_{train}}}$, $\forall y \in \mathcal{Y}$

$$\begin{aligned}
\max_{y \in \mathcal{Y}} \min_{\tau \in [T_{train}]} \|\nabla_{\theta_y} \bar{\mathcal{L}}_{reg}^{n_y}(\theta_y^\tau, \theta_{-y}^\tau)\|^2 &= \max_{y \in \mathcal{Y}} \mathcal{O} \left(\frac{2(\bar{\mathcal{L}}^{n_y}(\theta_y^0) - \bar{\mathcal{L}}^{n_y}(\theta_y^*))}{\max_{y \in \mathcal{Y}} L_y \sqrt{T_{train}}} + \frac{\sigma_y^2}{\sqrt{T_{train}}} + \beta \phi_y \sigma_y \right) \\
&= \mathcal{O} \left(\frac{m^2}{\sqrt{T_{train}}} + \beta \right)
\end{aligned} \quad (70)$$

$$= \mathcal{O} \left(\frac{m^2}{\sqrt{T_{train}}} + \beta \right) \quad (71)$$

933 For (θ_y, θ_{-y}) , we have

$$\text{NE-gap}(\theta_y, \theta_{-y}) = \max_{y \in \mathcal{Y}} |\bar{\mathcal{L}}_{reg}^{n_y}(\theta_y, \theta_{-y}) - \bar{\mathcal{L}}_{reg}^{n_y}(B(\theta_{-y}), \theta_{-y})| \quad (72)$$

$$\leq \max_{y \in \mathcal{Y}} \|\nabla_{\theta_y} \bar{\mathcal{L}}_{reg}^{n_y}(\theta_y, \theta_{-y})\|^2 \|\theta_y - B(\theta_{-y})\|_2^2 \quad (73)$$

934 Since the strategy space for θ_y is bounded in norm. We have

$$\text{NE-gap}(\theta_y, \theta_{-y}) \lesssim \max_{y \in \mathcal{Y}} \|\nabla_{\theta_y} \bar{\mathcal{L}}_{reg}^{n_y}(\theta_y, \theta_{-y})\|^2 \quad (74)$$

$$\Rightarrow \min_{\tau \in [T_{train}]} \text{NE-gap}(\theta_y^\tau, \theta_{-y}^\tau) = \mathcal{O} \left(\frac{m^2}{\sqrt{T_{train}}} + \beta \right) \quad (75)$$

935 **B.7 Monte Carlo Error of the Finite Neural Network**

936 Observe that

$$\mathcal{L}^y(\theta_y) = \frac{1}{2} \sum_{j=1}^N \lambda(t_j)(t_j - t_{j-1}) \mathbb{E}_{X_0} \mathbb{E}_{X_{t_j} | X_0} \left[\left\| \nabla_{x(t_j)} \log p_t(x_i(t_j) | x_0) - s_{t_j, \theta_y}(x_i(t_j)) \right\|_2^2 \right] \quad (76)$$

$$+ \frac{1}{2} \sum_{j=1}^N \lambda(t_j)(t_j - t_{j-1}) C_{t_j} \quad (77)$$

$$= \frac{1}{2} \sum_{j=1}^N \lambda(t_j)(t_j - t_{j-1}) \mathbb{E}_{x(t_j) \sim p_{t_j}} \left[\left\| s_{t_j, \theta_y}(x(t_j)) - \nabla_x \log p_{t_j}(x(t_j)) \right\|_2^2 \right] \quad (78)$$

937 For each $y \in \mathcal{Y}$, $\mathcal{L}^y(\theta_y^*)$ is the optimal loss function for the unregularized version under the current
 938 hypothesis class. Let $\mathcal{L}^y(\bar{\theta}_y^*)$ be the optimal unregularized loss function under the continuous version
 939 of the random feature model. Then,

$$\mathcal{L}^y(\theta_y^*) = \frac{1}{2} \sum_{j=1}^N \lambda(t_j)(t_j - t_{j-1}) \mathbb{E}_{x(t_j) \sim p_{t_j}} \left[\left\| s_{t_j, \theta_y^*}(x(t_j)) - \nabla_x \log p_{t_j}(x(t_j)) \right\|_2^2 \right] \quad (79)$$

$$\leq 2 \left(\frac{1}{2} \sum_{j=1}^N \lambda(t_j)(t_j - t_{j-1}) \mathbb{E}_{x(t_j) \sim p_{t_j}} \left[\left\| \bar{s}_{t_j, \bar{\theta}_y^*}(x(t_j)) - \nabla_x \log p_{t_j}(x(t_j)) \right\|_2^2 \right] \right) \quad (80)$$

$$+ \frac{1}{2} \sum_{j=1}^N \lambda(t_j)(t_j - t_{j-1}) \mathbb{E}_{x(t_j) \sim p_{t_j}} \left[\left\| \bar{s}_{t_j, \bar{\theta}_y^*}(x(t_j)) - s_{t_j, \theta_y^*}(x(t_j)) \right\|_2^2 \right] \quad (81)$$

$$\leq 2\mathcal{L}^y(\bar{\theta}_y^*) + Err_{MC}(\theta_y^*, \bar{\theta}_y^*, \{t_j\}_{j=1}^N, \{\lambda(t_j)\}_{j=1}^N) \quad (82)$$

940 **Proposition 2. Monte Carlo estimates.** Define the Monte Carlo error

$$Err_{MC}(\theta, \bar{\theta}, \{t_j\}_{j=1}^N, \{\lambda(t_j)\}_{j=1}^N) := \sum_{j=1}^N \lambda(t_j)(t_j - t_{j-1}) \mathbb{E}_{x(t_j) \sim p_{t_j}} \left[\left\| \bar{s}_{t_j, \bar{\theta}}(x(t_j)) - s_{t_j, \theta}(x(t_j)) \right\|_2^2 \right] \quad (83)$$

941 Suppose that $\|X(0)\|_\infty \leq K$ and the trainable parameter a and embedding functions $W, U, e(\cdot)$ are
 942 both bounded. Then, given any $\bar{\theta}$, for any $\delta > 0, \delta \ll 1$, with probability of at least $1 - 2N\delta$, there
 943 exists θ such that

$$Err_{MC}(\theta, \bar{\theta}, \{t_j\}_{j=1}^N, \{\lambda(t_j)\}_{j=1}^N) \leq \frac{2C_{w,u}^2 B^2 (C_{t_N, \delta} + C_{t_N, e})^2 d^2}{m} \log\left(\frac{2}{\delta}\right) \sum_{j=1}^N \lambda(t_j)(t_j - t_{j-1}) \quad (84)$$

944 **Proof.** The proof closely along the line of [15]. Fix any $\bar{\theta}$. For notational convenience, we will drop
 945 y from θ_y and $\bar{\theta}_y$. For $k = 1, 2, \dots, d$, define

$$Z_{t,k}(W, U) := \left\| s_{t, \theta, k}(x) - \bar{s}_{t, \bar{\theta}, k}(x) \right\|_{L^2(p_t)} = \mathbb{E}_{x \sim p_t}^{1/2} \left[|s_{t, \theta, k}(x) - \bar{s}_{t, \bar{\theta}, k}(x)|^2 \right] \quad (85)$$

$$= \mathbb{E}_{x \sim p_t} \left[\left| \frac{1}{m} \sum_{i=1}^m a_{i,k} \sigma(w_i^T x + u_i^T e(t)) - \mathbb{E}_{(w,u)} [a_k(w, u) \sigma(w^T x + u^T e(t))] \right|^2 \right] \quad (86)$$

946 Then, we have

$$\mathbb{E}_{x \sim p_t} \left[\left\| s_{t, \theta_y}(x) - \bar{s}_{t, \bar{\theta}_y}(x) \right\|_2^2 \right] = \sum_{k=1}^d \mathbb{E}_{x \sim p_t} \left[|s_{t, \theta_{y,k}}(x) - \bar{s}_{t, \bar{\theta}_{y,k}}(x)|^2 \right] \quad (87)$$

$$= \sum_{k=1}^d Z_{t,k}^2(W, U) \quad (88)$$

$$\leq \sum_{k=1}^d \left(|Z_{t,k}(W, U) - \mathbb{E}_{W,U}[Z_{t,k}]| + |\mathbb{E}_{W,U}[Z_{t,k}(W, U)]| \right)^2 \quad (89)$$

$$\stackrel{(a)}{\leq} 2 \sum_{k=1}^d \left(|Z_{t,k}(W, U) - \mathbb{E}_{W,U}[Z_{t,k}(W, U)]|^2 \right) \quad (90)$$

$$+ \mathbb{E}_{W,U}[Z_{t,k}^2(W, U)] \quad (91)$$

947 where (a) follows from the fact that $(a+b)^2 \leq 2(a^2+b^2)$ and Jensen's Inequality $\mathbb{E}^2[Z_{t,k}(W, U)] \leq$
 948 $\mathbb{E}_{W,U}[Z_{t,k}^2(W, U)]$. According to Lemma 1. for any $\delta > 0$, $\delta \ll 1$, w.p. atleast $1 - \delta$, we have

$$\|x(t)\|_\infty \leq C_{t_N, \delta} \quad (92)$$

949 If (\tilde{W}, \tilde{U}) is different from (W, U) at only one component indexed by i , we have w.p. $1 - \delta$

$$|Z_{t,k}(W, U) - Z_{t,k}(\tilde{W}, \tilde{U})| \quad (93)$$

$$= \left| \left\| s_{t, \theta, k}(x) - \bar{s}_{t, \bar{\theta}, k}(x) \right\|_{L^2(p_t)} - \left\| s_{t, \tilde{\theta}, k}(x) - \bar{s}_{t, \bar{\theta}, k}(x) \right\|_{L^2(p_t)} \right| \quad (94)$$

$$\stackrel{(a)}{\leq} \left\| s_{t, \tilde{\theta}, k}(x) - s_{t, \theta, k}(x) \right\|_{L^2(p_t)} \quad (95)$$

$$= \frac{1}{m} \left\| a_{i,k} \sigma(w_i^T x + u_i^T e(t)) - \tilde{a}_{i,k} \sigma(\tilde{w}_i^T x + \tilde{u}_i^T e(t)) \right\|_{L^2(p_t)} \quad (96)$$

$$\stackrel{(b)}{\leq} \frac{1}{m} \left(|a_{i,k}| \left\| \sigma(w_i^T x + u_i^T e(t)) \right\|_{L^2(p_t)} + |\tilde{a}_{i,k}| \left\| \sigma(\tilde{w}_i^T x + \tilde{u}_i^T e(t)) \right\|_{L^2(p_t)} \right) \quad (97)$$

$$\stackrel{(c)}{\leq} \frac{1}{m} \left(|a_{i,k}| \left\| w_i^T x + u_i^T e(t) \right\|_{L^2(p_t)} + |\tilde{a}_{i,k}| \left\| (\tilde{w}_i^T x + \tilde{u}_i^T e(t)) \right\|_{L^2(p_t)} \right) \quad (98)$$

$$\stackrel{(d)}{\leq} \frac{1}{m} \left(|a_{i,k}| (\|w_i\|_1 C_{t_N, \delta} + \|u_i\|_1 \|e(t)\|_\infty) + |\tilde{a}_{i,k}| (\|\tilde{w}_i\|_1 C_{t_N, \delta} + \|\tilde{u}_i\|_1 \|e(t)\|_\infty) \right) \quad (99)$$

$$\stackrel{(e)}{\leq} \frac{2}{m} B C_{w,u} (C_{t_N, \delta} + C_{t_N, e}) \quad (100)$$

950 where (a) and (b) follows from triangle inequality $||a| - |b|| \leq \|a - b\|$ and $\|a - b\| \leq \|a\| + \|b\|$,
 951 (c) follows from the fact that $|\sigma(y)| \leq |y|$, (d) follows from Lemma 1 and Holder Inequality, (e)
 952 follows from the bounds on $\|w_i\|_1, \|u_i\|_1, x, |a_{i,k}|, e(t_j)$.

953 Thus, w.p. $1 - \delta$, $Z_{t,k}(W, U)$ has bounded increment property. Using McDiarmid's inequality,
 954 w.p. $1 - 2\delta$, we have

$$|Z_{t,k}(W, U) - \mathbb{E}_{W,U}[Z_{t,k}(W, U)]| \leq \frac{B}{m} C_{w,u} (C_{t_N, \delta} + C_{t_N, e}) \sqrt{d \log \left(\frac{2}{\delta} \right)} \quad (101)$$

955 Now we compute

$$\begin{aligned}
& \mathbb{E}_{W,U}[Z_{t,k}^2(W,U)] \\
&= \mathbb{E}_{W,U} \left[\mathbb{E}_{x \sim p_t} [|s_{t,\theta,k}(x) - \bar{s}_{t,\bar{\theta},k}(x)|^2] \right] \\
&\stackrel{(b)}{=} \mathbb{E}_{x \sim p_t} \left[\mathbb{E}_{W,U} [|s_{t,\theta,k}(x) - \bar{s}_{t,\bar{\theta},k}(x)|^2] \right] \\
&= \frac{1}{m^2} \mathbb{E}_{x \sim p_t} \left[\mathbb{E}_{W,U} \left[\left| \sum_{i=1}^m (a_{i,k} \sigma(w_i^T x + u_i^T e(t)) - \mathbb{E}_{w,u} [a_k(w,u) \sigma(w^T x + u^T e(t))]) \right|^2 \right] \right] \\
&\quad + \frac{1}{m^2} \mathbb{E}_{x \sim p_t} \left[\mathbb{E}_{W,U} \left[\sum_{i \neq j} (a_{i,k} \sigma(w_i^T x + u_i^T e(t)) - \mathbb{E}_{w,u} [a_k(w,u) \sigma(w^T x + u^T e(t))]) \right. \right. \\
&\quad \times (a_{j,k} \sigma(w_j^T x + u_j^T e(t)) - \mathbb{E}_{w,u} [a_k(w,u) \sigma(w^T x + u^T e(t))]) \left. \left. \right] \right] \\
&\stackrel{(c)}{=} \frac{1}{m^2} \mathbb{E}_{x \sim p_t} \left[\mathbb{E}_{W,U} \left[\sum_{i=1}^m (a_{i,k} \sigma(w_i^T x + u_i^T e(t)) - \mathbb{E}_{w,u} [a_k(w,u) \sigma(w^T x + u^T e(t))])^2 \right] \right] \\
&\quad + \frac{1}{m^2} \mathbb{E}_{x \sim p_t} \left[\sum_{i \neq j} \mathbb{E}_{w_i, u_i} \left[(a_{i,k} \sigma(w_i^T x + u_i^T e(t)) - \mathbb{E}_{w,u} [a_k(w,u) \sigma(w^T x + u^T e(t))]) \right. \right. \\
&\quad \times \mathbb{E}_{w_j, u_j} \left[(a_{j,k} \sigma(w_j^T x + u_j^T e(t)) - \mathbb{E}_{w,u} [a_k(w,u) \sigma(w^T x + u^T e(t))]) \right] \left. \left. \right] \right] \\
&\stackrel{(d)}{=} \frac{1}{m^2} \mathbb{E}_{x \sim p_t} \left[\sum_{i=1}^m \mathbb{E}_{W,U} \left[(a_{i,k} \sigma(w_i^T x + u_i^T e(t)) - \mathbb{E}_{w,u} [a_k(w,u) \sigma(w^T x + u^T e(t))])^2 \right] \right] \\
&\stackrel{(e)}{\leq} \frac{1}{m} \mathbb{E}_{x \sim p_t} \left[\mathbb{E}_{w,u} \left[(a_k(w,u) \sigma(w^T x + u^T e(t)))^2 \right] \right] \\
&\stackrel{(f)}{\leq} \frac{1}{m} \mathbb{E}_{x \sim p_t} \left[\mathbb{E}_{w,u} \left[(|a_{y,k}(w,u)| (\|w\|_1 C_{t_N,\delta} + \|u\|_1 \|e(t)\|_\infty))^2 \right] \right] \\
&\leq \frac{1}{m} (C_{t_N,\delta} + C_{t_N,e})^2 C_{w,u}^2 B^2
\end{aligned}$$

956 where (b) is due to Fubini's theorem, (c) is due to independence of sampling (w_i, u_i) and (w_j, u_j) ,
957 (d) is due to $a_{j,k} \sigma(w_j^T x + u_j^T e(t))$ being an unbiased estimator of the continuous version of score
958 network, (e) follows from $\text{Var}(X) \leq \mathbb{E}[X^2]$, (f) follows from $|\sigma(y)| \leq |y|$ and Holder's inequality.
959 Thus. $w.p.1 - 2\delta$,

$$\mathbb{E}_{x \sim p_t} \left[\left\| s_{t,\theta_y}(x) - \bar{s}_{t,\bar{\theta}_y}(x) \right\|_2^2 \right] \leq \frac{2C_{w,u}^2 B^2 (C_{t_N,\delta} + C_{t_N,e})^2 d^2}{m} \log\left(\frac{2}{\delta}\right) \quad (102)$$

960 Finally, we have $w.p.1 - 2N\delta$

$$\text{Err}_{MC}(\theta, \bar{\theta}, \{t_j\}_{j=1}^N, \{\lambda(t_j)\}_{j=1}^N) \leq \frac{2C_{w,u}^2 B^2 (C_{t_N,\delta} + C_{t_N,e})^2 d^2}{m} \log\left(\frac{2}{\delta}\right) \sum_{j=1}^N \lambda(t_j) (t_j - t_{j-1}) \quad (103)$$

961 B.8 Radamacher Complexity

962 In this section, we will bound the term related to the generalization bound

$$\sup_{(\theta_y, \theta_{-y})} |\mathcal{L}_{reg}^y(\theta_y, \theta_{-y}) - \bar{\mathcal{L}}_{reg}^{n_y}(\theta_y, \theta_{-y})| \quad (104)$$

963 The Rademacher complexity of a real valued function class \mathcal{F} is defined as:

$$\mathcal{R}_n(\mathcal{F}) := \mathbb{E}_{x_1, \dots, x_n} \mathbb{E}_{\sigma_1, \dots, \sigma_n} \left[\sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \sigma_i f(x_i) \right], \quad (105)$$

964 The variables $\sigma_1, \dots, \sigma_m$ are iid Bernoulli random variables that take values $\{+1, -1\}$ with equal
 965 probability and are independent of x_1, \dots, x_m . However, for our random feature model, we have a
 966 vector valued function class

$$\hat{\mathcal{F}}_{W,U} := \left\{ f(x) = \frac{A}{m} \Phi(x, W, U) = \frac{1}{m} \sum_{k=1}^m \alpha_k \phi(x, w_k, u_k) \mid \|A\|_F \leq B \right\} \quad (106)$$

967 **Theorem 4.** [17, Theorem 3] Let X be nontrivial, symmetric and subgaussian. Then there exists a
 968 constant $C < \infty$, depending only on the distribution of X , such that for any countable set S and
 969 functions $\psi_i : S \rightarrow \mathbb{R}, \phi_i : S \rightarrow l_2, 1 \leq i \leq n$ satisfying

$$\forall s, s' \in S, \psi_i(s) - \psi_i(s') \leq \|\phi_i(s) - \phi_i(s')\| \quad (107)$$

970 we have

$$\mathbb{E} \sup_{s \in S} \sum_i \epsilon_i \psi_i(s) \leq C \mathbb{E} \sup_{s \in S} \sum_{i,k} X_{ik} \phi_i(s)_k \quad (108)$$

971 where the X_{ik} are independent copies of X for $1 \leq i \leq n$ and $1 \leq k \leq \infty$ and $\phi_i(s)_k$ is the k -th
 972 coordinate of $\phi_i(s)$. If X is a Rademacher variable we may choose $C = \sqrt{2}$, if X is a standard
 973 normal $C = \sqrt{\frac{\pi}{2}}$.

974 **Corollary 2.** [17, Corollary 4] Let \mathcal{X} be any set, $(x_1, \dots, x_n) \in \mathcal{X}^n$, let F be a class of functions
 975 $f : \mathcal{X} \rightarrow l_2$ and let $h_i : l_2 \rightarrow \mathbb{R}$ have Lipschitz norm L . Then

$$\mathbb{E} \sup_{f \in F} \sum_i \epsilon_i h_i(f(x_i)) \leq \sqrt{2} L \mathbb{E} \sup_{f \in F} \sum_{i,k} \epsilon_{ik} f_k(x_i) \quad (109)$$

976 where ϵ_{ik} is an independent doubly indexed Rademacher sequence and $f_k(x_i)$ is the k -th component
 977 of $f(x_i)$.

978 **Lemma 5.** [17] Consider the function class $\mathcal{F} = \{x \rightarrow \frac{A}{m} \phi(x, W, U) : A \in \mathcal{B}(H, \mathbb{R}), \|A\|_F \leq B\}$.
 979 Then the empirical Rademacher complexity of F is

$$\hat{\text{Rad}}_n(\mathcal{F}) = \mathbb{E} \sup_{f \in F} \sum_{i,k} \epsilon_{ik} f_k(x_i) \leq \frac{B}{\sqrt{m}} \sqrt{\sum_i \|\phi(x_i, W, U)\|^2} \quad (110)$$

980 Moreover, if $\mathbb{E}_x \|\phi(x, W, U)\|^2 \leq C^2$, the Rademacher Complexity of \mathcal{F} is

$$\mathcal{R}_n(\mathcal{F}) \leq \frac{BC}{\sqrt{mn}} \quad (111)$$

Proof:

$$\hat{\text{Rad}}_n(\mathcal{F}) = \mathbb{E} \sup_{f \in F} \sum_{i,k} \epsilon_{ik} f_k(x_i) = \frac{1}{m} \mathbb{E} \sup_{\|A\|_F \leq B} \sum_k \langle a_k, \sum_i \epsilon_{ik} x_i \rangle \quad (112)$$

$$= \frac{1}{m} \mathbb{E} \sup_{\|A\|_F \leq B} \text{tr}(D^* A) \leq B \mathbb{E} \|D^*\|_* \quad (113)$$

981 where $D \in \mathcal{B}(H, \mathbb{R}^K)$ is the random transformation

$$v \rightarrow \left(\langle v, \sum_i \epsilon_{i1} x_i \rangle, \dots, \langle v, \sum_i \epsilon_{iK} x_i \rangle \right) \quad (114)$$

982 Thus,

$$\mathbb{E} \|D^*\|_* = \mathbb{E} \sqrt{\sum_m \left\| \sum_i \epsilon_{ik} \phi(x_i, W, U) \right\|^2} \leq \sqrt{m \sum_i \|\phi(x_i, W, U)\|^2} \quad (115)$$

983 Thus,

$$\mathcal{R}_n(\mathcal{F}) = \mathbb{E}_{x_1, \dots, x_n} \frac{1}{n} \hat{Rad}_n(\mathcal{F}) \leq \frac{B}{\sqrt{mn}} \mathbb{E}_{x_1, \dots, x_n} \sqrt{\sum_i \|\phi(x_i, W, U)\|^2} \quad (116)$$

$$\leq \frac{B}{n\sqrt{m}} \sqrt{\sum_i \mathbb{E}_{x_1, \dots, x_n} \|\phi(x_i, W, U)\|^2} \quad (117)$$

$$\leq \frac{BC}{\sqrt{mn}} \quad (118)$$

984 Suppose $0 < t_1 < \dots < t_N = T$ are the chosen points of discretization for training, we have from
985 the forward process

$$X(t) = e^{-t} X(0) + \sqrt{1 - e^{-2t}} Z, Z \sim N(0, 1) \quad (119)$$

$$\implies \mathbb{E}_Z[X^2(t)] = e^{-2t} x^2(0) + \frac{1 - e^{-2t}}{2} \quad (120)$$

$$\implies \mathbb{E}_{X(0)} \mathbb{E}_Z[X^2(t_j)] \leq K^2 + \frac{1 - e^{-2T}}{2}, \forall 0 < t_j < T \quad (121)$$

986 Using the above bounds along with bounded support of embedding matrices W, U and embedding
987 function $e(t)$ and Assumption 1, it is easy to show that

$$\mathbb{E}_{X_0} \mathbb{E}_{\xi_j} \left[\|\sigma(Wx(t) + Ue(t))\|_2^2 \right] \leq F_T^2, \forall 0 < t \leq T \quad (122)$$

988 for some constant F_T^2 and $x(t) = e^{-t} x(0) + \sqrt{1 - e^{-2t}} \xi_j, \xi_j \sim \mathcal{N}(0, I)$

989 **Lemma 6.** *The term*

$$\mathcal{L}_{mut}^y(\theta_y, \theta_{-y}) = \sum_{j=1}^N \omega(t_j)(t_j - t_{j-1}) \mathbb{E}_{X_0} \mathbb{E}_{X_{t_j}|X_0} \mathbb{E}_{y' \sim Q} \left[\left\| s_{t_j, \theta_y}(x_i(t_j)) - s_{t, \theta_{y'}}(x_i(t_j)) \right\|_2^2 \right] \quad (123)$$

990 is $\mathcal{O}\left(F_T B \sum_{j=1}^N \omega(t_j)(t_j - t_{j-1})\right)$

991 **Proof:** Using the fact of bounded support of embedding matrices W, U and embedding function
992 $e(t)$, bounded strategy space and Assumption 1 and eq 122, we get the desired bounded.

993 **Lemma 7.** *Suppose $L_{C_1} = \bar{\sigma}_{t_j}^2 B F_T + \sqrt{d} \sqrt{\log \frac{2}{\pi \delta^2}}$. Then, with probability $1 - \delta$, the function*
994 $h : \mathcal{A} \subset \mathbb{R}^d \rightarrow \mathbb{R}$

$$h(x) = \|\bar{\sigma}_{t_j} x + \xi_{ij}\|^2 \quad (124)$$

995 is Lipschitz in x , where $\mathcal{A} = \{x \in \mathbb{R}^d : \|x\|_2 \leq F_T B\}$.

996 **Proof.** It is sufficient to show the norm of the gradient of $h(x)$ is bounded for $x \in \mathcal{A}$. With
997 probability $1 - \delta$,

$$\|\nabla_x h(x)\|_2 = \bar{\sigma}_{t_j} \|\bar{\sigma}_{t_j} x + \xi_{ij}\|_2 \leq \bar{\sigma}_{t_j}^2 F_T B + \sqrt{d} \sqrt{\log \frac{2}{\pi \delta^2}} \quad (125)$$

$$(126)$$

998 **Lemma 8.** *Suppose $L_{C_2} = 2F_T B |\mathcal{Y}|$. Define $g : \mathcal{A}^{\mathcal{Y}} \subset \mathbb{R}^{d|\mathcal{Y}|} \rightarrow \mathbb{R}$ where*

$$g(x_1, x_2, \dots, x_{|\mathcal{Y}|}) = \mathbb{E}_{y'} [\|x_i - x_{y'}\|^2], y' \in \{1, 2, \dots, |\mathcal{Y}|\} - i \quad (127)$$

999 is Lipschitz in x , where $\mathcal{A} = \{x \in \mathbb{R}^d : \|x\|_2 \leq F_T B\}$.

Proof:

$$\nabla_{x_i} g(x_1, x_2, \dots, x_{|\mathcal{Y}|}) = 2\mathbb{E}_{y'}[(x_i - x_{y'})] \quad (128)$$

$$\nabla_{x_j} g(x_1, x_2, \dots, x_{|\mathcal{Y}|}) = 2p(x_j)(x_j - x_i), j \neq i \quad (129)$$

$$\|\nabla_x g(x)\| \leq \|\nabla_{x_i} g(x_1, x_2, \dots, x_{|\mathcal{Y}|})\| + \sum_{j \neq i} \|\nabla_{x_j} g(x_1, x_2, \dots, x_{|\mathcal{Y}|})\| \quad (130)$$

$$\leq 2\mathbb{E}_{y'}[\|x_i - x_{y'}\|] + 2 \sum_{k \neq i} \|x_k - x_i\| \leq 2F_T B |\mathcal{Y}| \quad (131)$$

1000 We know

$$\mathcal{L}^y(\theta_y) = \frac{1}{2} \sum_{j=1}^N \frac{\lambda(t_j)(t_j - t_{j-1})}{\bar{\sigma}_{t_j}} \mathbb{E}_{X_0} \mathbb{E}_{\xi_j} \left[\|\bar{\sigma}_{t_j} s_{t_j, \theta_y}(x(t_j)) + \xi_j\|_2^2 \right] \quad (132)$$

$$+ \frac{1}{2} \sum_{j=1}^N \lambda(t_j)(t_j - t_{j-1}) C_{t_j}(y) \quad (133)$$

1001 where $C_t(y) = \mathbb{E}_{X_t} \|\nabla \log p_t(\cdot|y)\|^2 - \mathbb{E}_{X_0} \mathbb{E}_{X_t|X_0} \|\nabla \log p_t(x_t|x_0, y)\|^2$. Let $\bar{C}(y) =$
 1002 $\frac{1}{2} \sum_{j=1}^N \lambda(t_j)(t_j - t_{j-1}) C_{t_j}(y)$

1003 **Lemma 9.** With probability $1 - Nn_y\delta$, an upper bound for the generalization gap i.e.

$$\sup_{(\theta_y, \theta_{-y})} |\mathcal{L}_{reg}^y(\theta_y, \theta_{-y}) - \bar{\mathcal{L}}_{reg}^{n_y}(\theta_y, \theta_{-y})| \quad (134)$$

1004 is

$$\frac{2\sqrt{2}BF_T}{\sqrt{mn_y}} L_{C_1} \sum_{j=1}^N \frac{\lambda(t_j)(t_j - t_{j-1})}{\bar{\sigma}_{t_j}} + \frac{2\sqrt{2}BF_T|\mathcal{Y}|^2}{\sqrt{mn_y}} L_{C_2} \sum_{j=1}^N \omega(t_j)(t_j - t_{j-1}) + \bar{C} \quad (135)$$

1005 where $L_{C_1} = \bar{\sigma}_{t_j}^2 BF_T + \sqrt{d} \sqrt{\log \frac{2}{\pi\delta^2}}, L_{C_2} = 2F_T B |\mathcal{Y}|, \bar{C} = \max_{y \in \mathcal{Y}} |\bar{C}(y)|$

1006 **Proof.** Observe that, we can rewrite Eq. 134 using triangle inequality as

$$\begin{aligned} \sup_{(\theta_y, \theta_{-y})} |\mathcal{L}_{reg}^y(\theta_y, \theta_{-y}) - \bar{\mathcal{L}}_{reg}^{n_y}(\theta_y, \theta_{-y})| &\leq \sup_{\theta_y} |\mathcal{L}^y(\theta_y) - \bar{\mathcal{L}}^{n_y}(\theta_y)| \\ &\quad + \beta \sup_{(\theta_y, \theta_{-y})} |\mathcal{L}_{mut}^y(\theta_y, \theta_{-y}) - \bar{\mathcal{L}}_{mut}^{n_y}(\theta_y, \theta_{-y})| \end{aligned} \quad (136)$$

1007 Further decomposing them, we get

$$\sup_{\theta_y} |\mathcal{L}^y(\theta_y) - \bar{\mathcal{L}}^{n_y}(\theta_y)| \leq \sum_{j=1}^N \frac{\lambda(t_j)(t_j - t_{j-1})}{\bar{\sigma}_{t_j}} \sup_{\theta_y} |\mathcal{L}^y(\theta_y)(j) - \bar{\mathcal{L}}^{n_y}(\theta_y)(j)| + \bar{C} \quad (137)$$

1008 where $|\mathcal{L}^y(\theta_y)(j) - \bar{\mathcal{L}}^{n_y}(\theta_y)(j)| = \left| \frac{1}{2n_y} \sum_{i=1}^{n_y} \left[\|\bar{\sigma}_{t_j} s_{t_j, \theta_y}(x_i(t_j)) + \xi_{ij}\|_2^2 - \right. \right.$
 1009 $\left. \mathbb{E}_{X_0} \mathbb{E}_{\xi_j} \left[\|\bar{\sigma}_{t_j} s_{t_j, \theta_y}(e^{-t_j} X_0 + \sqrt{1 - e^{-2t_j}} \xi_j) + \xi_j\|_2^2 \right] \right|$ and

$$\begin{aligned} \sup_{(\theta_y, \theta_{-y})} |\mathcal{L}_{mut}^y(\theta_y, \theta_{-y}) - \bar{\mathcal{L}}_{mut}^{n_y}(\theta_y, \theta_{-y})| &\leq \\ \sum_{j=1}^N \omega(t_j)(t_j - t_{j-1}) \sup_{(\theta_y, \theta_{-y})} |\mathcal{L}_{mut}^y(\theta_y, \theta_{-y})(j) - \bar{\mathcal{L}}_{mut}^{n_y}(\theta_y, \theta_{-y})(j)| \end{aligned} \quad (138)$$

1010 where

$$\begin{aligned}
& |\mathcal{L}_{mut}^y(\theta_y, \theta_{-y})(j) - \bar{\mathcal{L}}_{mut}^{n_y}(\theta_y, \theta_{-y})(j)| = \left| \frac{1}{2n_y} \sum_{i=1}^{n_y} \left[\mathbb{E}_{y' \sim Q} \left[\left\| s_{t_j, \theta_y}(x_i(t_j)) - s_{t_j, \theta_{y'}}(x_i(t_j)) \right\|_2^2 \right] \right. \right. \\
& \quad \left. \left. - \mathbb{E}_{X_0} \mathbb{E}_{\xi_j} \mathbb{E}_{y' \sim Q} \left[\left\| s_{t_j, \theta_y}(e^{-t_j} X_0 + \sqrt{1 - e^{-2t_j}} \xi_j) - s_{t_j, \theta_{y'}}(e^{-t_j} X_0 + \sqrt{1 - e^{-2t_j}} \xi_j) \right\|_2^2 \right] \right] \right| \quad (139)
\end{aligned}$$

1011 Finally using Corollary 2, Lemmas 5 7,8, [17, Section 4.1] we have

$$\sup_{\theta_y} |\mathcal{L}^y(\theta_y) - \bar{\mathcal{L}}^{n_y}(\theta_y)| \leq \frac{2\sqrt{2}BF_T}{\sqrt{mn_y}} L_{C_1} \sum_{j=1}^N \lambda(t_j)(t_j - t_{j-1}) + \bar{\mathcal{C}} \quad (140)$$

1012 and

$$\sup_{(\theta_y, \theta_{-y})} |\mathcal{L}_{mut}^y(\theta_y, \theta_{-y}) - \bar{\mathcal{L}}_{mut}^{n_y}(\theta_y, \theta_{-y})| \leq \frac{2\sqrt{2}BF_T|\mathcal{Y}|^2}{\sqrt{mn_y}} L_{C_2} \sum_{j=1}^N \omega(t_j)(t_j - t_{j-1}) \quad (141)$$

1013 C Numerical Experiments

1014 **Computing resources.** The numerical experiments were conducted on a MacBook Air (2023) and
1015 Gilbreth. Gilbreth has heterogeneous hardware comprising of Nvidia V100, A100, A10, and A30
1016 GPUs in separate sub-clusters. All the nodes are connected by 100 Gbps Infiniband interconnects.
1017 We used sub-cluster B with 16 nodes, 24 cores per node, 192 GB memory per node, 3 A30 (24 GB)
1018 per node. For more information follow this link.

1019 The width of network $m = 16$, learning rate $\eta_\tau = 10^{-4}, \forall \tau, T_{train} = 5000$ is fixed for Adam
1020 optimizer. We set $\lambda(t) = \bar{\sigma}_t, \omega(t) = e^t$, total number of training samples is 50.

1021 **Case one** We perform more empirical experiments on $d = 1$, imbalance ratio $r = 2.5, \beta = 0.01$.
1022 We compute the KL-divergence between the ground truth distribution and the learned model using the
1023 procedure in [15]. $P(x|y = 1) \sim \mathcal{N}(-\mu, \sigma^2)$ and class 2 is $P(x|y = 2) \sim \mathcal{N}(\mu, \sigma^2)$. We observe
1024 Fig. 2 the worst case KL divergence for the mutual learning case is lower than the vanilla when we
1025 change the distance between mean and the variance of each class label. The performance of head
1026 class doesn't worsened for small μ . However, the head class performance suffers for mutual learning
1027 case when the distance between the mean increases. This might be because when the support of class
1028 distribution are farther apart mutual learning is not advantageous as transfer of knowledge between
1029 the class is not useful.

1030 **Case two** We now consider a case with two classes with imbalance ratio $r = 2.5, \beta = 0.01$. Class 1
1031 itself is a uniform mixture of two Gaussian i.e $P(x|y = 1) \sim \frac{1}{2}\mathcal{N}(-4, 3) + \frac{1}{2}\mathcal{N}(4, 3)$ and class 2 is
1032 $P(x|y = 2) \sim \mathcal{N}(0, 2)$ as in Fig. 3. We observe the Mutual Learning objective with our formulation
1033 have lower KL-divergence for both the classes compared to the vanilla diffusion models trained on
1034 each class. In this case, mutual learning allows useful transfer of knowledge between the classes
1035 increasing the performance for both. We hypothesize that under some notion of similarity between
1036 various class distributions, mutual learning is advantageous in improving the performance of all
1037 classes.

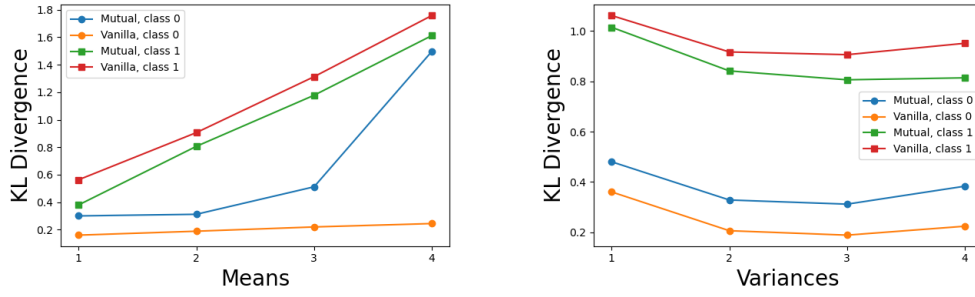


Figure 2: Case one: (Left) The first plot shows the KL-divergence for each class with and without mutual learning objective as μ is varied. (Right) shows the KL-divergence for each class with and without mutual learning objective as σ is varied ($\mu = 2$ fixed).

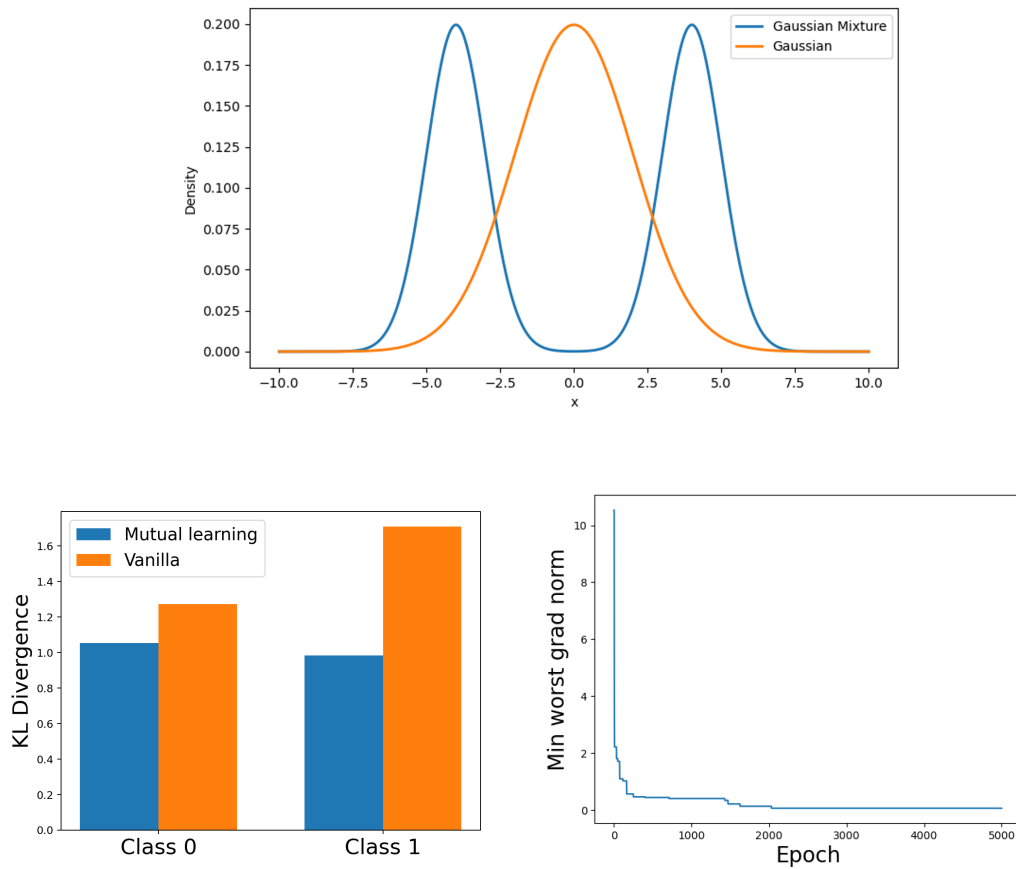


Figure 3: Case two: (Top) The first plot shows class 1 as a gaussian mixture with class 2 as Gaussian. (Bottom Left) Shows the KL-divergence for each class with and without mutual learning objective. (Bottom Right) Shows $\min_{\tau} \max_{y \in \mathcal{Y}} \|\nabla \tilde{\mathcal{L}}_{reg}^{n_y}(\theta_y^{\tau}, \theta_{-y}^{\tau})\|$ decreasing with training epoch.