

## APPENDIX

## A RELATED WORKS

**Semi-supervised learning:** There have been many existing results discussing about various methods of SSL. The book by [Chapelle et al. \(2006\)](#) presented a comprehensive overview of the SSL methods both theoretically and practically. [Chawla & Karakoulas \(2005\)](#) presented an empirical study of various SSL techniques on a variety of datasets and investigated sample-selection bias when the labelled and unlabelled data are from different distributions. [Zhu \(2008\)](#) classified various SSL methods into six main classes: generative models, low-density separation methods, graph-based methods, self-training and co-training. Pseudo-labelling is a technique among the self-training and co-training ([Zhu & Goldberg, 2009](#)). In self-training, the model is initially trained by the limited number of labelled data and generate pseudo-labels to the unlabelled data. Subsequently, the model is retrained with the pseudo-labelled data and repeats the process iteratively. It is a simple and effective SSL method without restrictions on the data samples ([Triguero et al., 2015](#)). A variety of works have also shown the benefits of utilizing the unlabelled data. [Singh et al. \(2008\)](#) developed a finite sample analysis that characterized how the unlabelled data improves the excess risk compared to the supervised learning, with respect to the number of unlabelled data and the margin between different classes. [Li et al. \(2019\)](#) studied multi-class classification with unlabelled data and provided a sharper generalization error bound using the notion of Rademacher complexity that yields a faster convergence rate. [Carmon et al. \(2019\)](#) proved that using unlabelled data can help to achieve high robust accuracy as well as high standard accuracy at the same time. [Dupre et al. \(2019\)](#) considered iteratively pseudo-labelling the whole unlabelled dataset with a confidence threshold and showed that the accuracy converges relatively quickly. [Oymak & Gulcu \(2021\)](#), in which part of our analysis hinges on, studied SSL under the binary Gaussian mixture model setup and characterized the correlation between the learned and the optimal estimators concerning the margin and the regularization factor. However, these works do not investigate how the unlabelled data affects the generalization error over the iterations.

**Generalization error bounds:** The traditional way of analyzing generalization error includes the Vapnik-Chervonenkis or VC dimension ([Vapnik, 2000](#)) and the Rademacher complexity ([Boucheron et al., 2005](#)). Recently, [Russo & Zou \(2016\)](#) proposed using mutual information between the estimated output of an algorithm and the actual realized value of the estimates to analyze and bound the bias in data analysis, which can be regarded equivalent to the generalization error. This new approach is simpler and can handle a wider range of loss functions compared to the abovementioned methods and other methods like differential privacy, total-variation information and so on. It also paves a new way to improving generalization capability of learning algorithms from an information-theoretic aspects. Following [Russo & Zou \(2016\)](#), [Xu & Raginsky \(2017\)](#) derived upper bounds on generalization error of learning algorithms with mutual information between the input dataset and the output hypothesis, which formalizes the intuition that less information that a learning algorithm can extract from training dataset leads to less overfitting. Later [Pensia et al. \(2018\)](#) derived generalization error bounds for noisy and iterative algorithms and the key contribution is to bound the mutual information between input data and output hypothesis. [Negrea et al. \(2019\)](#) improved mutual information bounds for Stochastic Gradient Langevin Dynamics (SGLD) via data-dependent estimates compared to distribution-dependent bounds.

However, one major shortcoming of the aforementioned mutual information bounds is that the bounds go to infinity for (deterministic) learning algorithms without noise, e.g., Stochastic Gradient Descent (SGD). Some other works have tried to overcome this problem. [Lopez & Jog \(2018\)](#) derived upper bounds on the generalization error using the Wasserstein distance involving the distributions of input data and output hypothesis, which are shown to be tighter under some natural cases. [Esposito et al. \(2021\)](#) derived generalization error bounds via Rényi-,  $f$ -divergences and maximal leakage. [Steinke & Zakyntinou \(2020\)](#) proposed using Conditional Mutual Information (CMI) to bound the generalization error, which can still preserve the chain rule property. [Bu et al. \(2020\)](#) provided a tightened upper bound based on the *individual* mutual information (IMI) between the *individual* data sample and the output. [Wu et al. \(2020\)](#) extended [Bu et al. \(2020\)](#)'s result to the transfer learning problems and characterized the upper bound based on IMI and KL-divergence. In a similar manner, [Jose & Simeone \(2020\)](#) provided a tightened bound on transfer generalization error based on the Jensen-Shannon divergence.

## B PROOF OF THEOREM 1

Before the proof, let us define some notation. The *cumulant generating function* (CGF) of a random variable  $L \in \mathbb{R}$  is  $\Lambda_L(\lambda) := \log \mathbb{E}_L[e^{\lambda(L - \mathbb{E}[L])}]$  for all  $\lambda \in \mathbb{R}$ . Note that  $\Lambda_L(0) = \Lambda'_L(0) = 0$  and  $\Lambda_L(\lambda)$  is convex. Then for any  $L \sim \text{subG}(R)$ , it means  $\Lambda_L(\lambda) \leq \frac{R^2 \lambda^2}{2}$ , for all  $\lambda \in \mathbb{R}$ .

For any convex function  $\psi : [0, b) \mapsto \mathbb{R}$ , its *Legendre dual*  $\psi^*$  is defined as  $\psi^*(x) := \sup_{\lambda \in [0, b)} \lambda x - \psi(\lambda)$  for all  $x \in [0, \infty)$ . According to Boucheron et al. (2013, Lemma 2.4), when  $\psi(0) = \psi'(0) = 0$ ,  $\psi^*(x)$  is a nonnegative convex and nondecreasing function on  $[0, \infty)$ . Moreover, for every  $y \geq 0$ , its generalized inverse function  $\psi^{*-1}(y) := \inf\{x \geq 0 : \psi^*(x) \geq y\}$  is concave and can be rewritten as  $\psi^{*-1}(y) = \inf_{\lambda \in [0, b)} \frac{y + \psi(\lambda)}{\lambda}$ .

We first introduce the following theorem that is applicable to more general loss functions.

**Theorem 4.** For any  $\tilde{\theta}_t \in \Theta$ , let  $\psi_-(\lambda, \tilde{\theta}_t)$  and  $\psi_+(\lambda, \tilde{\theta}_t)$  be convex functions of  $\lambda$  and  $\psi_+(0, \tilde{\theta}_t) = \psi'_+(0, \tilde{\theta}_t) = \psi_-(0, \tilde{\theta}_t) = \psi'_-(0, \tilde{\theta}_t) = 0$ . Assume that  $\Lambda_{l(\tilde{\theta}_t, \tilde{Z})}(\lambda, \tilde{\theta}_t) \leq \psi_+(\lambda, \tilde{\theta}_t)$  for all  $\lambda \in [0, b_+)$  and  $\Lambda_{l(\tilde{\theta}_t, \tilde{Z})}(\lambda, \tilde{\theta}_t) \leq \psi_-(\lambda, \tilde{\theta}_t)$  for  $\lambda \in (b_-, 0]$  under distribution  $P_{\tilde{Z}|\theta^{(t-1)}} = P_Z$ , where  $0 < b_+ \leq \infty$  and  $-\infty \leq b_- < 0$ . Let  $\psi_+(\lambda) = \sup_{\tilde{\theta}_t} \psi_+(\lambda, \tilde{\theta}_t)$  and  $\psi_-(\lambda) = \sup_{\tilde{\theta}_t} \psi_-(\lambda, \tilde{\theta}_t)$ . We have

$$\begin{aligned}
& \text{gen}_t(P_Z, P_X, \{P_{\theta_k|S_1, S_u}\}_{k=0}^t, \{f_{\theta_k}\}_{k=0}^{t-1}) \\
& \leq \frac{w}{n} \sum_{i=1}^n \mathbb{E}_{\theta^{(t-1)}} [\psi_-^{*-1}(I_{\theta^{(t-1)}}(\theta_t; Z_i))] \\
& \quad + \frac{1-w}{m} \sum_{i=(t-1)m+1}^{tm} \mathbb{E}_{\theta^{(t-1)}} [\psi_-^{*-1}(I_{\theta^{(t-1)}}(\theta_t; X'_i, \hat{Y}'_i) + D_{\theta^{(t-1)}}(P_{X'_i, \hat{Y}'_i} \| P_Z))] , \quad (29) \\
& - \text{gen}_t(P_Z, P_X, \{P_{\theta_k|S_1, S_u}\}_{k=0}^t, \{f_{\theta_k}\}_{k=0}^{t-1}) \\
& \leq \frac{w}{n} \sum_{i=1}^n \mathbb{E}_{\theta^{(t-1)}} [\psi_+^{*-1}(I_{\theta^{(t-1)}}(\theta_t; Z_i))] \\
& \quad + \frac{1-w}{m} \sum_{i=(t-1)m+1}^{tm} \mathbb{E}_{\theta^{(t-1)}} [\psi_+^{*-1}(I_{\theta^{(t-1)}}(\theta_t; X'_i, \hat{Y}'_i) + D_{\theta^{(t-1)}}(P_{X'_i, \hat{Y}'_i} \| P_Z))] , \quad (30)
\end{aligned}$$

where  $P_{X'_i, \hat{Y}'_i|\theta^{(t-1)}}(x, y|\hat{\theta}^{(t-1)}) = P_X(x)\mathbb{1}\{y = f_{\hat{\theta}_{t-1}}(x)\}$  for any  $x \in \mathcal{X}$ ,  $y \in \mathcal{Y}$  and  $\hat{\theta}^{(t-1)} \in \Theta^{t-1}$ , and  $P_{Z|\theta^{(t-1)}} = P_Z$ .

*Proof.* Consider the Donsker–Varadhan variational representation of KL-divergence between any two distributions  $P$  and  $Q$  on  $\mathcal{X}$ :

$$D(P\|Q) = \sup_{g \in \mathcal{G}} \{\mathbb{E}_{X \sim P}[g(X)] - \log \mathbb{E}_{X \sim Q}[e^{g(X)}]\} \quad (31)$$

where the supremum is taken over the set of measurable functions in  $\mathcal{G} = \{g : \mathcal{X} \mapsto \mathbb{R} : \mathbb{E}_{X \sim Q}[e^{g(X)}] < \infty\}$ .

Recall that  $\tilde{\theta}_t$  and  $\tilde{Z}$  are independent copies of  $\theta_t$  and  $Z$  respectively, such that  $P_{\tilde{\theta}_t, \tilde{Z}} = Q_{\theta_t} \otimes P_Z$ ,  $P_{\tilde{\theta}_t, \tilde{Z}|\theta^{(t-1)}} = P_{\theta_t|\theta^{(t-1)}} \otimes P_Z$ . For any iterative SSL algorithm, by applying the law of total expectation, the generalization error can be rewritten as

$$\begin{aligned}
& \text{gen}_t(P_Z, P_X, \{P_{\theta_k|S_1, S_u}\}_{k=0}^t, \{f_{\theta_k}\}_{k=0}^{t-1}) \\
& = w \left( \mathbb{E}_{\theta_t} [\mathbb{E}_Z[l(\theta_t, Z)]] - \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{\theta_t, Z_i} [l(\theta_t, Z_i)] \right) \\
& \quad + (1-w) \left( \mathbb{E}_{\theta_t} [\mathbb{E}_Z[l(\theta_t, Z)]] - \frac{1}{m} \sum_{i=(t-1)m+1}^{tm} \mathbb{E}_{\theta_t, X'_i, \hat{Y}'_i} [l(\theta_t, (X'_i, \hat{Y}'_i))] \right) \quad (32)
\end{aligned}$$

$$\begin{aligned}
&= \frac{w}{n} \sum_{i=1}^n \left( \mathbb{E}_{\tilde{\theta}_t, \tilde{Z}} [l(\tilde{\theta}_t, \tilde{Z})] - \mathbb{E}_{\theta_t, Z_i} [l(\theta_t, Z_i)] \right) \\
&\quad + \frac{1-w}{m} \sum_{i=(t-1)m+1}^{tm} \left( \mathbb{E}_{\tilde{\theta}_t, \tilde{Z}} [l(\tilde{\theta}_t, \tilde{Z})] - \mathbb{E}_{\theta_t, X'_i, \hat{Y}'_i} [l(\theta_t, (X'_i, \hat{Y}'_i))] \right) \quad (33)
\end{aligned}$$

$$\begin{aligned}
&= \frac{w}{n} \sum_{i=1}^n \mathbb{E}_{\theta^{(t-1)}} \left[ \mathbb{E}_{\tilde{\theta}_t, \tilde{Z}} [l(\tilde{\theta}_t, \tilde{Z}) | \theta^{(t-1)}] - \mathbb{E}_{\theta_t, Z_i} [l(\theta_t, Z_i) | \theta^{(t-1)}] \right] \\
&\quad + \frac{1-w}{m} \sum_{i=(t-1)m+1}^{tm} \mathbb{E}_{\theta^{(t-1)}} \left[ \mathbb{E}_{\tilde{\theta}_t, \tilde{Z}} [l(\tilde{\theta}_t, \tilde{Z}) | \theta^{(t-1)}] - \mathbb{E}_{\theta_t, X'_i, \hat{Y}'_i} [l(\theta_t, (X'_i, \hat{Y}'_i)) | \theta^{(t-1)}] \right]. \quad (34)
\end{aligned}$$

Note that  $\psi_+(\lambda) = \sup_{\tilde{\theta}_t} \psi_+(\lambda, \tilde{\theta}_t)$  and  $\psi_-(\lambda) = \sup_{\tilde{\theta}_t} \psi_-(\lambda, \tilde{\theta}_t)$  are convex, and so their Legendre duals  $\psi_-^*$ ,  $\psi_+^*$ , and the corresponding inverses are well-defined.

Let  $\tilde{l}(\theta, z) = l(\theta, z) - \mathbb{E}_Z[l(\theta, Z)]$ . We have the fact that  $\mathbb{E}_{\tilde{Z}}[\tilde{l}(\tilde{\theta}_t, \tilde{Z})] = 0$  for any  $\tilde{\theta}_t$ . Again, by the Donsker–Varadhan variational representation of the KL-divergence, for any fixed  $\theta^{(t-1)}$  and any  $\lambda \in [0, b_+)$ , we have

$$\begin{aligned}
I_{\theta^{(t-1)}}(\theta_t; Z) &= D(P_{\theta_t, Z | \theta^{(t-1)}} \| P_{\theta_t | \theta^{(t-1)}} \otimes P_Z) \\
&\geq \mathbb{E}_{\theta_t, Z} [\lambda \tilde{l}(\theta_t, Z) | \theta^{(t-1)}] - \log \mathbb{E}_{\tilde{\theta}_t, \tilde{Z}} [e^{\lambda \tilde{l}(\tilde{\theta}_t, \tilde{Z})} | \theta^{(t-1)}] \quad (35)
\end{aligned}$$

$$= \mathbb{E}_{\theta_t, Z} [\lambda \tilde{l}(\theta_t, Z) | \theta^{(t-1)}] - \log \mathbb{E}_{\tilde{\theta}_t | \theta^{(t-1)}} \mathbb{E}_{\tilde{Z}} [e^{\lambda \tilde{l}(\tilde{\theta}_t, \tilde{Z})}] \quad (36)$$

$$= \mathbb{E}_{\theta_t, Z} [\lambda \tilde{l}(\theta_t, Z) | \theta^{(t-1)}] - \log \mathbb{E}_{\tilde{\theta}_t | \theta^{(t-1)}} [\exp(\Lambda_{l(\tilde{\theta}_t, \tilde{Z})}(\lambda, \tilde{\theta}_t))] \quad (37)$$

$$\geq \lambda \mathbb{E}_{\theta_t, Z} [l(\theta_t, Z) - \mathbb{E}_Z[l(\theta_t, Z)] | \theta^{(t-1)}] - \log \mathbb{E}_{\tilde{\theta}_t | \theta^{(t-1)}} [\exp(\psi_+(\lambda, \tilde{\theta}_t))] \quad (38)$$

$$\geq \lambda \mathbb{E}_{\theta_t, Z} [l(\theta_t, Z) - \mathbb{E}_Z[l(\theta_t, Z)] | \theta^{(t-1)}] - \psi_+(\lambda) \quad (39)$$

$$= \lambda (\mathbb{E}_{\theta_t, Z} [l(\theta_t, Z) | \theta^{(t-1)}] - \mathbb{E}_{\tilde{\theta}_t, \tilde{Z}} [l(\tilde{\theta}_t, \tilde{Z}) | \theta^{(t-1)}]) - \psi_+(\lambda). \quad (40)$$

where (37) follows from the definition of  $\Lambda_{l(\tilde{\theta}_t, \tilde{Z})}(\lambda, \tilde{\theta}_t)$  in (52), (38) follows from the assumption that  $\Lambda_{l(\tilde{\theta}_t, \tilde{Z})}(\lambda, \tilde{\theta}_t) \leq \psi_+(\lambda, \tilde{\theta}_t)$  for all  $\lambda \in [0, b_+)$  and (39) follows since  $\psi_+(\lambda) = \sup_{\tilde{\theta}_t} \psi_+(\lambda, \tilde{\theta}_t)$ . Thus, we have

$$\begin{aligned}
&\mathbb{E}_{\theta_t, Z} [l(\theta_t, Z) | \theta^{(t-1)}] - \mathbb{E}_{\tilde{\theta}_t, \tilde{Z}} [l(\tilde{\theta}_t, \tilde{Z}) | \theta^{(t-1)}] \\
&\leq \inf_{\lambda \in [0, b_+)} \frac{I_{\theta^{(t-1)}}(\theta_t; Z) + \psi_+(\lambda)}{\lambda} \quad (41)
\end{aligned}$$

$$= \psi_+^{*-1}(I_{\theta^{(t-1)}}(\theta_t; Z)). \quad (42)$$

Similarly, for  $\lambda \in (b_-, 0]$ ,

$$\begin{aligned}
&\mathbb{E}_{\tilde{\theta}_t, \tilde{Z}} [l(\tilde{\theta}_t, \tilde{Z}) | \theta^{(t-1)}] - \mathbb{E}_{\theta_t, Z} [l(\theta_t, Z) | \theta^{(t-1)}] \\
&\leq \inf_{\lambda \in [0, -b_-)} \frac{I_{\theta^{(t-1)}}(\theta_t; Z) + \psi_-(\lambda)}{\lambda} \quad (43)
\end{aligned}$$

$$= \psi_-^{*-1}(I_{\theta^{(t-1)}}(\theta_t; Z)). \quad (44)$$

By applying the same techniques, for any pair of pseudo-labelled random variables  $(X', \hat{Y}')$  used at iteration  $t$  and any  $\lambda \in [0, b_+)$ , we have

$$\begin{aligned}
&I_{\theta^{(t-1)}}(\theta_t; X', \hat{Y}') + D_{\theta^{(t-1)}}(P_{X', \hat{Y}'} \| P_Z) \\
&= D_{\theta^{(t-1)}}(P_{\theta_t, X', \hat{Y}'} \| P_{\theta_t} \otimes P_{X', \hat{Y}'}) + D_{\theta^{(t-1)}}(P_{\theta_t} \otimes P_{X', \hat{Y}'} \| P_{\theta_t} \otimes P_Z) \quad (45)
\end{aligned}$$

$$\begin{aligned}
&\geq \mathbb{E}_{\theta_t, X', \hat{Y}'} [\lambda l(\theta_t, (X', \hat{Y}')) | \theta^{(t-1)}] - \log \mathbb{E}_{\theta_t} [\mathbb{E}_{X', \hat{Y}'} [e^{\lambda l(\theta_t, (X', \hat{Y}'))} | \theta^{(t-1)}] | \theta^{(t-1)}] \\
&\quad + \mathbb{E}_{\theta_t} [\mathbb{E}_{X', \hat{Y}'} [\lambda l(\theta_t, (X', \hat{Y}')) | \theta^{(t-1)}] | \theta^{(t-1)}] - \log \mathbb{E}_{\theta_t} [\mathbb{E}_Z [e^{\lambda l(\theta_t, Z)}] | \theta^{(t-1)}] \quad (46)
\end{aligned}$$

$$\geq \mathbb{E}_{\theta_t, X', \hat{Y}'} [\lambda l(\theta_t, (X', \hat{Y}')) | \theta^{(t-1)}] - \log \mathbb{E}_{\theta_t} [\mathbb{E}_Z [e^{\lambda l(\theta_t, Z)}] | \theta^{(t-1)}] \quad (47)$$

$$= \lambda \left( \mathbb{E}_{\theta_t, X', \hat{Y}'} [\lambda l(\theta_t, (X', \hat{Y}')) | \theta^{(t-1)}] - \mathbb{E}_{\theta_t} [\mathbb{E}_Z [l(\theta_t, Z)] | \theta^{(t-1)}] \right) - \log \mathbb{E}_{\tilde{\theta}_t | \theta^{(t-1)}} [\exp(\Lambda_{l(\tilde{\theta}_t, \tilde{Z})}(\lambda, \tilde{\theta}_t))] \quad (48)$$

$$\geq \lambda \left( \mathbb{E}_{\theta_t, X', \hat{Y}'} [\lambda l(\theta_t, (X', \hat{Y}'))] - \mathbb{E}_{\theta_t} [\mathbb{E}_Z [l(\theta_t, Z)]] \right) - \psi_+(\lambda), \quad (49)$$

where (47) follows from the Jensen's inequality. Thus, we get

$$\begin{aligned} & \mathbb{E}_{\theta_t, X'_i, \hat{Y}'_i} [l(\theta_t, (X', \hat{Y}')) | \theta^{(t-1)}] - \mathbb{E}_{\tilde{\theta}_t, \tilde{Z}} [l(\tilde{\theta}_t, \tilde{Z}) | \theta^{(t-1)}] \\ & \leq \psi_+^{*-1}(I_{\theta^{(t-1)}}(\theta_t; X', \hat{Y}') + D_{\theta^{(t-1)}}(P_{X', \hat{Y}'} \| P_Z)) \end{aligned} \quad (50)$$

and

$$\begin{aligned} & \mathbb{E}_{\tilde{\theta}_t, \tilde{Z}} [l(\tilde{\theta}_t, \tilde{Z}) | \theta^{(t-1)}] - \mathbb{E}_{\theta_t, X'_i, \hat{Y}'_i} [l(\theta_t, (X', \hat{Y}')) | \theta^{(t-1)}] \\ & \leq \psi_-^{*-1}(I_{\theta^{(t-1)}}(\theta_t; X', \hat{Y}') + D_{\theta^{(t-1)}}(P_{X', \hat{Y}'} \| P_Z)). \end{aligned} \quad (51)$$

The proof is completed by applying inequalities (42), (44), (50) and (51) to the expansion of  $\text{gen}_t$  in (34).  $\square$

Let  $\tilde{\theta}_t$  and  $\tilde{Z}$  be independent copies of  $\theta_t$  and  $Z$  respectively, such that  $P_{\tilde{\theta}_t, \tilde{Z}} = Q_{\theta_t} \otimes P_Z$ , where  $Q_{\theta_t}$  is the marginal distribution of  $\theta_t$ . For any fixed  $\tilde{\theta}_t \in \Theta$ , let the cumulant generating function (CGF) of  $l(\tilde{\theta}_t, \tilde{Z})$  be

$$\Lambda_{l(\tilde{\theta}_t, \tilde{Z})}(\lambda, \tilde{\theta}_t) := \log \mathbb{E}_{\tilde{Z}} [e^{\lambda(l(\tilde{\theta}_t, \tilde{Z}) - \mathbb{E}_{\tilde{Z}}[l(\tilde{\theta}_t, \tilde{Z})])}]. \quad (52)$$

When the loss function  $l(\theta, Z) \sim \text{subG}(R)$  under  $Z \sim P_Z$  for any  $\theta \in \Theta$ , we have  $\Lambda_{l(\tilde{\theta}_t, \tilde{Z})}(\lambda, \tilde{\theta}_t) \leq \frac{R^2 \lambda^2}{2}$  for all  $\lambda \in \mathbb{R}$ . Then we can let  $\psi_-(\lambda, \tilde{\theta}_t) = \psi_+(\lambda, \tilde{\theta}_t) = \frac{R^2 \lambda^2}{2}$  for all  $\tilde{\theta}_t \in \Theta$ . Hence,  $\psi_+(\lambda) = \psi_-(\lambda) = \sup_{\tilde{\theta}_t \in \Theta} \frac{R^2 \lambda^2}{2} = \frac{R^2 \lambda^2}{2}$  and  $\psi_+^{*-1}(y) = \psi_-^{*-1}(y) = \sqrt{2R^2 y}$  for any  $y \geq 0$ . Finally, Theorem 1 can then be directly obtained from Theorem 4.

## C PROOF OF THEOREM 2

Theorem 2 can be proved iteratively by applying Theorem 1. For simplicity, in the following proofs, we abbreviate  $\text{gen}_t(P_Z, P_X, \{P_{\theta_k | S_1, S_u}\}_{k=0}^t, \{f_{\theta_k}\}_{k=0}^{t-1})$  as  $\text{gen}_t$ .

1. **Initial round**  $t = 0$ : Since  $Y_i \mathbf{X}_i \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mu, \sigma^2 \mathbf{I}_d)$ , we have  $\theta_0 \sim \mathcal{N}(\mu, \frac{\sigma^2}{n} \mathbf{I}_d)$  and for some constant  $c \in \mathbb{R}_+$ ,

$$\Pr(\theta_0 \in \Theta_{\mu, c}) = \Pr(\|\theta_0 - \mu\|_\infty \leq c) = \left(1 - 2\Phi\left(\frac{-\sqrt{nc}}{\sigma}\right)\right)^d =: 1 - \delta_{\sqrt{nc}, d}. \quad (53)$$

By choosing  $c$  large enough,  $\delta_{\sqrt{nc}, d}$  can be made arbitrarily small. Consider  $\tilde{\theta}_0$  and  $(\tilde{\mathbf{X}}, \tilde{Y})$  as independent copies of  $\theta_0 \sim Q_{\theta_0}$  and  $(\mathbf{X}, Y) \sim P_{\mathbf{X}, Y} = P_Y \otimes \mathcal{N}(Y\mu, \sigma^2 \mathbf{I}_d)$ , respectively, such that  $P_{\tilde{\theta}_0, \tilde{\mathbf{X}}, \tilde{Y}} = Q_{\theta_0} \otimes P_{\mathbf{X}, Y}$ . Then the probability that  $l(\theta_0, (\mathbf{X}, Y)) \sim \text{subG}((c_2 - c_1)/2)$  under  $(\mathbf{X}, Y) \sim P_{\mathbf{X}, Y}$  is given as follows

$$\Pr\left(\Lambda_{l(\tilde{\theta}_0, (\tilde{\mathbf{X}}, \tilde{Y}))}(\lambda, \tilde{\theta}_0) \leq \frac{\lambda^2(c_2 - c_1)^2}{8}\right) \quad (54)$$

$$\geq \Pr\left(\Lambda_{l(\tilde{\theta}_0, (\tilde{\mathbf{X}}, \tilde{Y}))}(\lambda, \tilde{\theta}_0) \leq \frac{\lambda^2(c_2 - c_1)^2}{8} \text{ and } \tilde{\theta}_0 \in \Theta_{\mu, c}\right) \quad (55)$$

$$= \Pr(\tilde{\theta}_0 \in \Theta_{\mu, c}) \Pr\left(\Lambda_{l(\tilde{\theta}_0, (\tilde{\mathbf{X}}, \tilde{Y}))}(\lambda, \tilde{\theta}_0) \leq \frac{\lambda^2(c_2 - c_1)^2}{8} \middle| \tilde{\theta}_0 \in \Theta_{\mu, c}\right) \quad (56)$$

$$= (1 - \delta_{\sqrt{nc}, d})(1 - \delta_{r, d}), \quad (57)$$

where the last equality follows from (14) and (53).

Fix some  $d \in \mathbb{N}$ ,  $\epsilon > 0$  and  $\delta \in (0, 1)$ . There exists  $n_0(d, \delta) \in \mathbb{N}$ ,  $r_0(d, \delta) \in \mathbb{R}_+$  such that for all  $n > n_0$ ,  $r > r_0$ ,  $\delta_{\sqrt{nc}, d} < \frac{\delta}{3}$ ,  $\delta_{r, d} < \frac{\delta}{3}$ , and then with probability at least  $1 - \delta$ , the absolute generalization error can be upper bounded as follows

$$|\text{gen}_0| \leq \frac{1}{n} \sum_{i=1}^n \sqrt{\frac{(c_2 - c_1)^2}{2}} I(\theta_0; \mathbf{X}_i, Y_i). \quad (58)$$

Then mutual information can be calculated as follows

$$I(\theta_0; \mathbf{X}_i, Y_i) = h(\theta_0) - h(\theta_0 | \mathbf{X}_i, Y_i) \quad (59)$$

$$= h\left(\frac{1}{n} \sum_{j=1}^n Y_j \mathbf{X}_j\right) - h\left(\frac{1}{n} \sum_{j=1}^n Y_j \mathbf{X}_j \middle| \mathbf{X}_i, Y_i\right) \quad (60)$$

$$= \frac{d}{2} \log \left( \frac{2\pi e \sigma^2}{n} \right) - h\left(\frac{1}{n} \sum_{j \in [n], j \neq i} Y_j \mathbf{X}_j\right) \quad (61)$$

$$= \frac{d}{2} \log \left( \frac{2\pi e \sigma^2}{n} \right) - \frac{d}{2} \log \left( \frac{2\pi e (n-1) \sigma^2}{n^2} \right) \quad (62)$$

$$= \frac{d}{2} \log \frac{n}{n-1}. \quad (63)$$

Thus we can get (23).

2. **Pseudo-label using  $\theta_0$ :** For any  $i \in [1 : m]$  and  $X'_i \in S_u$ , the pseudo-label is given by

$$\hat{Y}'_i = \text{sgn}(\theta_0^\top \mathbf{X}'_i). \quad (64)$$

Given any pair of  $(\xi_0, \mu^\perp)$ ,  $\theta_0$  is fixed and  $\{\hat{Y}'_i\}_{i \in [1:m]}$  are conditionally i.i.d. from  $P_{\hat{Y}'_i | \xi_0, \mu^\perp} \in \mathcal{P}(\mathcal{Y})$ . Recall the pseudo-labelled dataset is defined as  $\hat{S}_{u,1} = \{(\mathbf{X}'_i, \hat{Y}'_i)\}_{i=1}^m$ .

Since  $\theta_0 \sim \mathcal{N}(\mu, \frac{\sigma^2}{n} \mathbf{I}_d)$ , inspired by Oymak & Gulcu (2021), we can decompose it as follows:

$$\theta_0 = \mu + \frac{\sigma}{\sqrt{n}} \xi \quad (65)$$

$$= \mu + \frac{\sigma}{\sqrt{n}} (\xi_0 \mu + \mu^\perp) \quad (66)$$

$$= \left(1 + \frac{\sigma}{\sqrt{n}} \xi_0\right) \mu + \frac{\sigma}{\sqrt{n}} \mu^\perp, \quad (67)$$

where  $\xi \sim \mathcal{N}(0, \mathbf{I}_d)$ ,  $\xi_0 \sim \mathcal{N}(0, 1)$ ,  $\mu^\perp \perp \mu$ ,  $\mu^\perp \sim \mathcal{N}(0, \mathbf{I}_d - \mu \mu^\top)$  and  $\mu^\perp$  is independent of  $\xi_0$ .

Recall the correlation between  $\theta_0$  and  $\mu$  given in (19), the decomposition of  $\bar{\theta}_0$  in (20) and  $\alpha, \beta$ . Since  $\text{sgn}(\theta_0^\top \mathbf{X}'_i) = \text{sgn}(\bar{\theta}_0^\top \mathbf{X}'_i)$ , in the following we can analyze the normalized parameter  $\bar{\theta}_0$  instead.

Given any  $(\xi_0, \mu^\perp)$ ,  $\alpha$  is fixed, and for any  $i \in \mathbb{N}$ , let us define a Gaussian noise vector  $\mathbf{g}_i \sim \mathcal{N}(0, \mathbf{I}_d)$  and decompose it as follows

$$\mathbf{g}_i = g_{0,i} \mu + g_i \mathbf{v} + \mathbf{g}_i^\perp, \quad (68)$$

where  $g_{0,i}, g_i \sim \mathcal{N}(0, 1)$ ,  $\mathbf{g}_i^\perp \sim \mathcal{N}(0, \mathbf{I}_d - \mu \mu^\top - \mathbf{v} \mathbf{v}^\top)$ ,  $\mathbf{g}_i^\perp \perp \mu$ ,  $\mathbf{g}_i^\perp \perp \mathbf{v}$  and  $g_{0,i}, g_i, \mathbf{g}_i^\perp$  are mutually independent.

For any sample  $\mathbf{X}'_i \sim \mathcal{N}(\mu, \sigma^2 \mathbf{I}_d)$ , we can decompose it as

$$\mathbf{X}'_i = \mu + \sigma \mathbf{g}_i = \mu + \sigma (g_{0,i} \mu + g_i \mathbf{v} + \mathbf{g}_i^\perp). \quad (69)$$

Then we have

$$\bar{\theta}_0^\top \mathbf{X}'_i = (\alpha \mu + \beta \mathbf{v})^\top (\mu + \sigma \mathbf{g}_i) \quad (70)$$

$$= \alpha + \sigma (\alpha \mu + \beta \mathbf{v})^\top (g_{0,i} \mu + g_i \mathbf{v} + \mathbf{g}_i^\perp) \quad (71)$$

$$= \alpha + \sigma (\alpha g_{0,i} + \beta g_i) \quad (72)$$

$$=: \alpha + \sigma h_i. \quad (73)$$

Note that  $h_i \sim \mathcal{N}(0, 1)$  for any  $\alpha \in [-1, 1]$ . Similarly, for any sample  $\mathbf{X}'_i \sim \mathcal{N}(-\mu, \sigma^2 \mathbf{I}_d)$ , we have

$$\mathbf{X}'_i = -\mu + \sigma \mathbf{g}_i \quad (74)$$

and

$$\bar{\theta}_0^\top \mathbf{X}'_i = -\alpha + \sigma h_i. \quad (75)$$

Denote the true label of  $\mathbf{X}'_i$  as  $Y'_i$  and  $P_{Y'_i} = P_Y \sim \text{unif}(\{-1, +1\})$ . The probability that the pseudo-label  $\hat{Y}'_i$  is equal to 1 is given by

$$\Pr(\hat{Y}'_i = 1) = \Pr(\bar{\theta}_0^\top \mathbf{X}'_i > 0) \quad (76)$$

$$= \frac{1}{2} \Pr(\bar{\theta}_0^\top \mathbf{X}'_i > 0 | Y'_i = 1) + \frac{1}{2} \Pr(\bar{\theta}_0^\top \mathbf{X}'_i > 0 | Y'_i = -1) \quad (77)$$

$$= \frac{1}{2} \mathbb{E}_\alpha [\Pr(\alpha + \sigma h_i > 0)] + \frac{1}{2} \mathbb{E}_\alpha [\Pr(-\alpha + \sigma h_i > 0)] \quad (78)$$

$$= \frac{1}{2} \mathbb{E}_\alpha \left[ \mathbb{Q}\left(\frac{-\alpha}{\sigma}\right) \right] + \frac{1}{2} \mathbb{E}_\alpha \left[ \mathbb{Q}\left(\frac{\alpha}{\sigma}\right) \right] = \frac{1}{2}. \quad (79)$$

We also have  $\Pr(\hat{Y}'_i = -1) = 1 - \Pr(\hat{Y}'_i = 1) = 1/2$ , and so  $P_{\hat{Y}'_i} = P_Y$ .

3. **Iteration  $t = 1$ :** Recall (18) and the new model parameter learned from the pseudo-labelled dataset  $\hat{\mathcal{S}}_{u,1}$  is given by

$$\theta_1 = \frac{1}{m} \sum_{i=1}^m \hat{Y}'_i \mathbf{X}'_i = \frac{1}{m} \sum_{i=1}^m \text{sgn}(\bar{\theta}_0^\top \mathbf{X}'_i) \mathbf{X}'_i = \frac{1}{m} \sum_{i=1}^m \text{sgn}(\bar{\theta}_0^\top \mathbf{X}'_i) \mathbf{X}'_i. \quad (80)$$

- (a) First let us calculate the conditional expectation of  $\theta_1$  given  $\theta_0$ .

Given any  $(\xi_0, \mu^\perp)$ , for any  $j \in [1 : m]$ , let  $\mu_1^{\xi_0, \mu^\perp} := \mathbb{E}[\text{sgn}(\bar{\theta}_0^\top \mathbf{X}'_j) \mathbf{X}'_j | \xi_0, \mu^\perp]$  and  $\mathbb{P}_{\xi_0, \mu^\perp}$  denotes the probability measure under the parameters  $(\xi_0, \mu^\perp)$ .

The expectation  $\mu_1^{\xi_0, \mu^\perp}$  can be calculated as follows:

$$\mu_1^{\xi_0, \mu^\perp} = \mathbb{E}[\text{sgn}(\bar{\theta}_0^\top \mathbf{X}'_j) \mathbf{X}'_j | \xi_0, \mu^\perp] \quad (81)$$

$$= \mathbb{E}_{Y'_j} [\mathbb{E}[\text{sgn}(\bar{\theta}_0^\top \mathbf{X}'_j) \mathbf{X}'_j | \xi_0, \mu^\perp, Y'_j]] \quad (82)$$

$$= \frac{1}{2} \mathbb{E}[\text{sgn}(\bar{\theta}_0^\top \mathbf{X}'_j) \mathbf{X}'_j | \xi_0, \mu^\perp, Y'_j = -1] + \frac{1}{2} \mathbb{E}[\text{sgn}(\bar{\theta}_0^\top \mathbf{X}'_j) \mathbf{X}'_j | \xi_0, \mu^\perp, Y'_j = 1]. \quad (83)$$

Different from (68), here we decompose the Gaussian random vector  $\mathbf{g}_j \sim \mathcal{N}(0, \mathbf{I}_d)$  in another way

$$\mathbf{g}_j = \tilde{g}_j \bar{\theta}_0 + \tilde{\mathbf{g}}_j^\perp, \quad (84)$$

where  $\tilde{g}_j \sim \mathcal{N}(0, 1)$ ,  $\tilde{\mathbf{g}}_j^\perp \sim \mathcal{N}(0, \mathbf{I}_d - \bar{\theta}_0 \bar{\theta}_0^\top)$ ,  $\tilde{g}_j$  and  $\tilde{\mathbf{g}}_j^\perp$  are mutually independent and  $\tilde{\mathbf{g}}_j^\perp \perp \bar{\theta}_0$ .

Then we decompose  $\mathbf{X}'_j$  and  $\bar{\theta}_0 \mathbf{X}'_j$  as

$$\mathbf{X}'_j = Y'_j \mu + \sigma \tilde{g}_j \bar{\theta}_0 + \sigma \tilde{\mathbf{g}}_j^\perp, \text{ and} \quad (85)$$

$$\bar{\theta}_0^\top \mathbf{X}'_j = Y'_j \alpha + \sigma \tilde{g}_j. \quad (86)$$

Then we have

$$\begin{aligned} & \mathbb{E}[\text{sgn}(\bar{\theta}_0^\top \mathbf{X}'_j) \mathbf{X}'_j | \xi_0, \mu^\perp, Y'_j = -1] \\ &= \mathbb{E}[\text{sgn}(-\alpha + \sigma \tilde{g}_j) (-\mu + \sigma \tilde{g}_j \bar{\theta}_0 + \sigma \tilde{\mathbf{g}}_j^\perp) | \xi_0, \mu^\perp] \end{aligned} \quad (87)$$

$$\begin{aligned} &= -\mathbb{E}[\text{sgn}(-\alpha + \sigma \tilde{g}_j) | \xi_0, \mu^\perp] \mu + \sigma \mathbb{E}[\text{sgn}(-\alpha + \sigma \tilde{g}_j) \tilde{g}_j | \xi_0, \mu^\perp] \bar{\theta}_0 \\ &\quad + \sigma \mathbb{E}[\text{sgn}(-\alpha + \sigma \tilde{g}_j) \tilde{\mathbf{g}}_j^\perp | \xi_0, \mu^\perp] \end{aligned} \quad (88)$$

$$= -\mathbb{E}[\text{sgn}(-\alpha + \sigma \tilde{g}_j) | \xi_0, \mu^\perp] \mu + \sigma \mathbb{E}[\text{sgn}(-\alpha + \sigma \tilde{g}_j) \tilde{g}_j | \xi_0, \mu^\perp] \bar{\theta}_0, \quad (89)$$

where (89) follows since  $\tilde{\mathbf{g}}^\perp$  is independent of  $\tilde{g}_j$  and  $\mathbb{E}[\tilde{\mathbf{g}}^\perp] = 0$ . For the first term in (89), recall  $\tilde{g}_j \sim \mathcal{N}(0, 1)$  and we have

$$-\mathbb{E}[\text{sgn}(-\alpha + \sigma \tilde{g}_j) | \xi_0, \boldsymbol{\mu}^\perp] \boldsymbol{\mu} = \left(1 - 2\text{Q}\left(\frac{\alpha}{\sigma}\right)\right) \boldsymbol{\mu}. \quad (90)$$

For the second term in (89), we have

$$\begin{aligned} & \mathbb{E}[\text{sgn}(-\alpha + \sigma \tilde{g}_j) \tilde{g}_j | \xi_0, \boldsymbol{\mu}^\perp] \bar{\boldsymbol{\theta}}_0 \\ &= \left( - \int_{-\infty}^{\frac{\alpha}{\sigma}} \frac{1}{\sqrt{2\pi}} e^{-\frac{g^2}{2}} g \, dg + \int_{\frac{\alpha}{\sigma}}^{\infty} \frac{1}{\sqrt{2\pi}} e^{-\frac{g^2}{2}} g \, dg \right) \bar{\boldsymbol{\theta}}_0 \end{aligned} \quad (91)$$

$$= \frac{2}{\sqrt{2\pi}} \exp\left(-\frac{\alpha^2}{2\sigma^2}\right) \bar{\boldsymbol{\theta}}_0. \quad (92)$$

By combining (90) and (92), we have

$$\mathbb{E}[\text{sgn}(\bar{\boldsymbol{\theta}}_0^\top \mathbf{X}'_j) \mathbf{X}'_j | \xi_0, \boldsymbol{\mu}^\perp, Y'_j = -1] = \left(1 - 2\text{Q}\left(\frac{\alpha}{\sigma}\right)\right) \boldsymbol{\mu} + \frac{2\sigma}{\sqrt{2\pi}} \exp\left(-\frac{\alpha^2}{2\sigma^2}\right) \bar{\boldsymbol{\theta}}_0, \quad (93)$$

and similarly,

$$\mathbb{E}[\text{sgn}(\bar{\boldsymbol{\theta}}_0^\top \mathbf{X}'_j) \mathbf{X}'_j | \xi_0, \boldsymbol{\mu}^\perp, Y'_j = 1] = \left(2\text{Q}\left(\frac{-\alpha}{\sigma}\right) - 1\right) \boldsymbol{\mu} + \frac{2\sigma}{\sqrt{2\pi}} \exp\left(-\frac{\alpha^2}{2\sigma^2}\right) \bar{\boldsymbol{\theta}}_0. \quad (94)$$

Thus, recall  $\bar{\boldsymbol{\theta}}_0 = \alpha \boldsymbol{\mu} + \beta \mathbf{v}$  and  $\boldsymbol{\mu}_1^{\xi_0, \boldsymbol{\mu}^\perp}$  is given by

$$\begin{aligned} \boldsymbol{\mu}_1^{\xi_0, \boldsymbol{\mu}^\perp} &= \mathbb{E}[\text{sgn}(\bar{\boldsymbol{\theta}}_0^\top \mathbf{X}'_j) \mathbf{X}'_j | \xi_0, \boldsymbol{\mu}^\perp] \\ &= \left(1 - 2\text{Q}\left(\frac{\alpha}{\sigma}\right)\right) \boldsymbol{\mu} + \frac{2\sigma}{\sqrt{2\pi}} \exp\left(-\frac{\alpha^2}{2\sigma^2}\right) \bar{\boldsymbol{\theta}}_0 \end{aligned} \quad (95)$$

$$= \left(1 - 2\text{Q}\left(\frac{\alpha}{\sigma}\right) + \frac{2\sigma\alpha}{\sqrt{2\pi}} \exp\left(-\frac{\alpha^2}{2\sigma^2}\right)\right) \boldsymbol{\mu} + \frac{2\sigma\beta}{\sqrt{2\pi}} \exp\left(-\frac{\alpha^2}{2\sigma^2}\right) \mathbf{v}. \quad (96)$$

The  $l_\infty$  norm between  $\boldsymbol{\mu}_1^{\xi_0, \boldsymbol{\mu}^\perp}$  and  $\boldsymbol{\mu}$  can be upper bounded by

$$\begin{aligned} & \|\boldsymbol{\mu}_1^{\xi_0, \boldsymbol{\mu}^\perp} - \boldsymbol{\mu}\|_\infty \\ & \leq \sqrt{\left(-2\text{Q}\left(\frac{\alpha}{\sigma}\right) + \frac{2\sigma\alpha}{\sqrt{2\pi}} \exp\left(-\frac{\alpha^2}{2\sigma^2}\right)\right)^2 + \frac{2\sigma^2\beta^2}{\pi} \exp\left(-\frac{2\alpha^2}{2\sigma^2}\right)} \end{aligned} \quad (97)$$

$$< \sqrt{\left(2\Phi\left(\frac{1}{\sigma}\right) + \frac{2\sigma}{\sqrt{2\pi}}\right)^2 + \frac{2\sigma^2}{\pi}} =: \tilde{c}_1, \quad (98)$$

where  $\tilde{c}_1$  is a constant only dependent on  $\sigma$ .

- (b) Next, we need to calculate the probability that  $l(\boldsymbol{\theta}_1, (\mathbf{X}, Y)) \sim \text{subG}((c_2 - c_1/2))$  under  $(\mathbf{X}, Y) \sim P_{\mathbf{X}, Y}$ .

Let  $\mathbf{V}_i = \text{sgn}(\bar{\boldsymbol{\theta}}_0^\top \mathbf{X}'_i) \mathbf{X}'_i - \boldsymbol{\mu}_1^{\xi_0, \boldsymbol{\mu}^\perp}$ . For any  $k \in [1 : d]$ , let  $V_{i,k}$ ,  $\theta_{1,k}$ ,  $\mu_{1,k}$  denote the  $k$ -th components of  $\mathbf{V}_i$ ,  $\boldsymbol{\theta}_1$  and  $\boldsymbol{\mu}_1^{\xi_0, \boldsymbol{\mu}^\perp}$ , respectively. Recall the decompositions  $\mathbf{X}'_i = Y'_i \boldsymbol{\mu} + \sigma \tilde{g}_i \bar{\boldsymbol{\theta}}_0 + \sigma \tilde{\mathbf{g}}_i^\perp$  in (85) and  $\bar{\boldsymbol{\theta}}_0 \mathbf{X}'_i = Y'_i \alpha + \sigma \tilde{g}_i$  in (86). Suppose the basis of  $\mathbb{R}^d$  is denoted by  $B = \{\mathbf{v}_1, \dots, \mathbf{v}_d\}$  and let  $\mathbf{v}_1 = \bar{\boldsymbol{\theta}}_0$ . Then we have

$$\tilde{\mathbf{g}}_i^\perp = \tilde{g}_{i,2}^\perp \mathbf{v}_2 + \dots + \tilde{g}_{i,d}^\perp \mathbf{v}_d, \quad (99)$$

where  $\tilde{g}_{i,k}^\perp \sim \mathcal{N}(0, 1)$  for any  $k \in [2 : d]$  and  $\{\tilde{g}_{i,k}\}_{k=2}^d$  are mutually independent. We also let  $\boldsymbol{\mu} = (\mu_{0,1}, \dots, \mu_{0,d})$ .

Given any  $(\xi_0, \boldsymbol{\mu}^\perp)$ , the moment generating function (MGF) of  $V_{i,1}$  is given as follows: for any  $s_1 > 0$ ,

$$\begin{aligned} & \mathbb{E}_{V_{i,1}}[e^{s_1 V_{i,1}}] \\ &= Q\left(\frac{-\alpha}{\sigma}\right) \mathbb{E}_{\tilde{g}_i}\left[e^{s_1(\mu_{0,1}-\mu_{1,1}+\sigma\tilde{g}_i)} \middle| \tilde{g}_i > \frac{-\alpha}{\sigma}\right] + Q\left(\frac{\alpha}{\sigma}\right) \mathbb{E}_{\tilde{g}_i}\left[e^{s_1(-\mu_{0,1}-\mu_{1,1}+\sigma\tilde{g}_i)} \middle| \tilde{g}_i > \frac{\alpha}{\sigma}\right] \end{aligned} \quad (100)$$

$$= e^{s_1(\mu_{0,1}-\mu_{1,1})} e^{\frac{\sigma^2 s_1^2}{2}} \Phi\left(\frac{\alpha}{\sigma} + \sigma s_1\right) + e^{s_1(-\mu_{0,1}-\mu_{1,1})} e^{\frac{\sigma^2 s_1^2}{2}} \Phi\left(\frac{-\alpha}{\sigma} + \sigma s_1\right). \quad (101)$$

The final equality follows from the fact that the MGF of a zero-mean univariate Gaussian truncated to  $(a, b)$  is  $e^{\sigma^2 s^2/2} \left[ \frac{\Phi(b-\sigma s) - \Phi(a-\sigma s)}{\Phi(b) - \Phi(a)} \right]$ . The second derivative of  $\log \mathbb{E}_{V_{i,1}}[e^{s_1 V_{i,1}}]$  is given as

$$\tilde{R}_1(s_1) := \frac{d^2 \log \mathbb{E}_{V_{i,1}}[e^{s_1 V_{i,1}}]}{ds_1^2} \quad (102)$$

$$\leq \sigma^2 + \frac{\text{const.}}{\left(\Phi\left(\frac{\alpha}{\sigma} + \sigma s_1\right) e^{s_k \mu_{0,k}} + \Phi\left(\frac{-\alpha}{\sigma} + \sigma s_1\right) e^{-s_k \mu_{0,k}}\right)^2} < \infty. \quad (103)$$

For  $k \in [2 : d]$  and any  $s_k > 0$ , the MGF of  $V_{i,k}$  is given as

$$\mathbb{E}_{V_{i,k}}[e^{s_k V_{i,k}}] = \mathbb{E}_{\sigma \tilde{g}_{i,k}^\perp, Y'_i}\left[e^{s_k(Y'_i \mu_{0,k} - \mu_{1,k} + \sigma \tilde{g}_{i,k}^\perp)}\right] \quad (104)$$

$$= Q\left(\frac{-\alpha}{\sigma}\right) e^{s_k(\mu_{0,k}-\mu_{1,k})} e^{\frac{\sigma^2 s_k^2}{2}} + Q\left(\frac{\alpha}{\sigma}\right) e^{s_k(-\mu_{0,k}-\mu_{1,k})} e^{\frac{\sigma^2 s_k^2}{2}}, \quad (105)$$

and the second derivative of  $\log \mathbb{E}_{V_{i,k}}[e^{s_k V_{i,k}}]$  is given by

$$\tilde{R}_k(s_k) := \frac{d^2 \log \mathbb{E}_{V_{i,k}}[e^{s_k V_{i,k}}]}{ds_k^2} = \sigma^2 + \frac{4\mu_{0,k}^2 Q\left(\frac{-\alpha}{\sigma}\right) Q\left(\frac{\alpha}{\sigma}\right)}{\left(Q\left(\frac{-\alpha}{\sigma}\right) e^{s_k \mu_{0,k}} + Q\left(\frac{\alpha}{\sigma}\right) e^{-s_k \mu_{0,k}}\right)^2}. \quad (106)$$

Fix  $k \in [1 : d]$ . According to Taylor's theorem, we have

$$\log \mathbb{E}_{V_{i,k}}[e^{s_k V_{i,k}}] = \frac{\tilde{R}_k(\xi_{L,k})}{2} s_k^2, \quad (107)$$

for some  $\xi_{L,k} \in (0, s_k)$  and  $\tilde{R}_k(\xi_{L,k}) < \infty$ . Then the Cramér transform of  $\log \mathbb{E}_{V_{i,k}}[e^{s_k V_{i,k}}]$  can be lower bounded as follows: for any  $\varepsilon > 0$ ,

$$\sup_{s_k > 0} \left( s_k \varepsilon - \log \mathbb{E}_{V_{i,k}}[e^{s_k V_{i,k}}] \right) \geq \sup_{s_k > 0} \left( s_k \varepsilon - \frac{\tilde{R}_k(\xi_{L,k}) s_k^2}{2} \right) = \frac{\varepsilon^2}{2\tilde{R}_k(\xi_{L,k})}. \quad (108)$$

Let  $\tilde{R}^* = \max_{\xi_0, \boldsymbol{\mu}^\perp} \min_{k \in [1:d]} \tilde{R}_k(\xi_{L,k})$ , which is a finite constant only dependent on  $\sigma$ . Since  $\{\text{sgn}(\tilde{\boldsymbol{\theta}}_0^\top \mathbf{X}'_i) \mathbf{X}'_i\}_{i=1}^m$  are i.i.d. random variables conditioned on  $(\xi_0, \boldsymbol{\mu}^\perp)$ , by applying Chernoff-Cramér inequality, we have for all  $\varepsilon > 0$

$$\begin{aligned} & \mathbb{P}_{\xi_0, \boldsymbol{\mu}^\perp} \left( \|\boldsymbol{\theta}_1 - \boldsymbol{\mu}_1^{\xi_0, \boldsymbol{\mu}^\perp}\|_\infty > \varepsilon \right) \\ &= \mathbb{P}_{\xi_0, \boldsymbol{\mu}^\perp} \left( \max_{k \in [1:d]} |\theta_{1,k} - \mu_{1,k}| > \varepsilon \right) \end{aligned} \quad (109)$$

$$\leq \sum_{k=1}^d \mathbb{P}_{\xi_0, \boldsymbol{\mu}^\perp} \left( |\theta_{1,k} - \mu_{1,k}| > \varepsilon \right) \quad (110)$$

$$= \sum_{k=1}^d \mathbb{P}_{\xi_0, \boldsymbol{\mu}^\perp} \left( \left| \frac{1}{m} \sum_{i=1}^m V_{i,k} \right| > \varepsilon \right) \quad (111)$$

$$\leq \sum_{k=1}^d 2 \exp \left( -m \sup_{s_k > 0} \left( s_k \varepsilon - \log \mathbb{E}_{V_{i,k}}[e^{s_k V_{i,k}}] \right) \right) \quad (112)$$

$$\leq 2d \exp \left( -\frac{m \varepsilon^2}{2\tilde{R}^*} \right) \quad (113)$$

$$=: \delta_{m, \varepsilon, d}, \quad (114)$$



where  $\delta_{m,\varepsilon,d} \xrightarrow{\text{a.s.}} 0$  as  $m \rightarrow \infty$  and does not depend on  $\xi_0, \mu^\perp$ .

Choose some  $c \in (\tilde{c}_1, \infty)$  ( $\tilde{c}_1$  defined in (98)). We have

$$\mathbb{P}_{\xi_0, \mu^\perp}(\theta_1 \in \Theta_{\mu,c}) \geq \mathbb{P}_{\xi_0, \mu^\perp}(\|\theta_1 - \mu_1^{\xi_0, \mu^\perp}\|_\infty \leq c - \tilde{c}_1) \geq 1 - \delta_{m,c-\tilde{c}_1,d}. \quad (115)$$

Consider  $\tilde{\theta}_1$  as an independent copy of  $\theta_1$  and independent of  $(\tilde{\mathbf{X}}, \tilde{Y})$ . Then the probability that  $l(\theta_1, (\mathbf{X}, Y)) \sim \text{subG}((c_2 - c_1)/2)$  under  $(\mathbf{X}, Y) \sim P_{\mathbf{X},Y}$  is given as follows

$$\mathbb{P}_{\xi_0, \mu^\perp} \left( \Lambda_{l(\tilde{\theta}_1, (\tilde{\mathbf{X}}, \tilde{Y}))}(\lambda, \tilde{\theta}_1) \leq \frac{\lambda^2(c_2 - c_1)^2}{8} \right) \quad (116)$$

$$\geq \mathbb{P}_{\xi_0, \mu^\perp}(\tilde{\theta}_1 \in \Theta_{\mu,c}) \mathbb{P}_{\xi_0, \mu^\perp} \left( \Lambda_{l(\tilde{\theta}_1, (\tilde{\mathbf{X}}, \tilde{Y}))}(\lambda, \tilde{\theta}_1) \leq \frac{\lambda^2(c_2 - c_1)^2}{8} \middle| \tilde{\theta}_1 \in \Theta_{\mu,c} \right) \quad (117)$$

$$= (1 - \delta_{m,c,d})(1 - \delta_{r,d}). \quad (118)$$

Thus, for some  $c \in (\tilde{c}_1, \infty)$ , with probability at least  $(1 - \delta_{m,c-\tilde{c}_1,d})(1 - \delta_{r,d})$ , the absolute generalization error can be upper bounded as follows:

$$|\text{gen}_1| = |\mathbb{E}[L_{P_{\mathbf{Z}}}(\theta_1) - L_{\hat{S}_{u,1}}(\theta_1)]| \quad (119)$$

$$= \left| \frac{1}{m} \sum_{i=1}^m \mathbb{E}_{\xi_0, \mu^\perp} \left[ \mathbb{E} \left[ l(\theta_1, (\mathbf{X}, Y)) - l(\theta_1, (\mathbf{X}'_i, \hat{Y}'_i)) \middle| \xi_0, \mu^\perp \right] \right] \right| \quad (120)$$

$$\leq \frac{1}{m} \sum_{i=1}^m \mathbb{E}_{\xi_0, \mu^\perp} \left[ \sqrt{\frac{(c_2 - c_1)^2}{2} \left( I_{\xi_0, \mu^\perp}(\theta_1, (\mathbf{X}'_i, \hat{Y}'_i)) + D_{\xi_0, \mu^\perp}(P_{\mathbf{X}'_i, \hat{Y}'_i} \| P_{\mathbf{X},Y}) \right)} \right], \quad (121)$$

where  $P_{\theta_1, (\mathbf{X}, Y) | \xi_0, \mu^\perp} = Q_{\theta_1 | \xi_0, \mu^\perp} \otimes P_{\mathbf{X},Y}$  and  $Q_{\theta_1 | \xi_0, \mu^\perp}$  denotes the marginal distribution of  $\theta_1$  under parameters  $(\xi_0, \mu^\perp)$ .

In the following, we derive the closed form expressions of the mutual information and KL-divergence in (121). For any  $j \in [1 : m]$ :

- **Calculate  $I_{\xi_0, \mu^\perp}(\theta_1; \mathbf{X}'_j, \hat{Y}'_j)$ :** For arbitrary random variables  $X$  and  $U$ , we define the *disintegrated conditional differential entropy* of  $X$  given  $U$  as

$$h_U(X) := h(P_{X|U}). \quad (122)$$

Conditioned on a certain pair of  $(\xi_0, \mu^\perp)$ , the mutual information between  $\theta_1$  and  $(\mathbf{X}'_j, \hat{Y}'_j)$  is given by

$$\begin{aligned} I_{\xi_0, \mu^\perp}(\theta_1; \mathbf{X}'_j, \hat{Y}'_j) &= h_{\xi_0, \mu^\perp} \left( \frac{1}{m} \sum_{i=1}^m \text{sgn}(\theta_0^\top \mathbf{X}'_i) \mathbf{X}'_i \right) - h_{\xi_0, \mu^\perp} \left( \frac{1}{m} \sum_{j=1}^m \hat{Y}'_j \mathbf{X}'_j \middle| \mathbf{X}'_j, \hat{Y}'_j \right) \end{aligned} \quad (123)$$

$$= h_{\xi_0, \mu^\perp} \left( \frac{1}{m} \sum_{i=1}^m \text{sgn}(\theta_0^\top \mathbf{X}'_i) \mathbf{X}'_i \right) - h_{\xi_0, \mu^\perp} \left( \frac{1}{m} \sum_{i \in [m], i \neq j} \text{sgn}(\theta_0^\top \mathbf{X}'_i) \mathbf{X}'_i \right) \quad (124)$$

$$\begin{aligned} &= h_{\xi_0, \mu^\perp} \left( \frac{1}{m} \sum_{i=1}^m \text{sgn}(\theta_0^\top \mathbf{X}'_i) \mathbf{X}'_i \right) \\ &\quad - h_{\xi_0, \mu^\perp} \left( \frac{1}{m-1} \sum_{i \in [m], i \neq j} \text{sgn}(\theta_0^\top \mathbf{X}'_i) \mathbf{X}'_i \right) - d \log \frac{m-1}{m}. \end{aligned} \quad (125)$$

As  $m \rightarrow \infty$ ,  $I_{\xi_0, \mu^\perp}(\theta_1; \mathbf{X}'_j, \hat{Y}'_j) \rightarrow 0$  almost surely and hence, in probability. Thus, for any  $\epsilon > 0$ , and there exists  $m_0(\epsilon, d, \delta) \in \mathbb{N}$  such that for all  $m > m_0$ ,

$$\mathbb{P}_{\xi_0, \mu^\perp}(I_{\xi_0, \mu^\perp}(\theta_1; \mathbf{X}'_j, \hat{Y}'_j) > \epsilon) \leq \delta. \quad (126)$$

- **Calculate  $D_{\xi_0, \mu^\perp}(P_{\mathbf{X}'_j, \hat{Y}'_j} \| P_{\mathbf{X},Y})$ :** First of all, since  $P_{\hat{Y}'_j} = P_Y$  (cf. (79)) regardless of the values of  $(\xi_0, \mu^\perp)$ , the *disintegrated conditional KL-divergence* can be rewritten as

$$\begin{aligned} &D_{\xi_0, \mu^\perp}(P_{\mathbf{X}'_j, \hat{Y}'_j} \| P_{\mathbf{X},Y}) \\ &= P_{\hat{Y}'_j}(-1) D_{\xi_0, \mu^\perp}(P_{\mathbf{X}'_j | \hat{Y}'_j = -1} \| P_{\mathbf{X} | Y = -1}) + P_{\hat{Y}'_j}(1) D_{\xi_0, \mu^\perp}(P_{\mathbf{X}'_j | \hat{Y}'_j = 1} \| P_{\mathbf{X} | Y = 1}). \end{aligned} \quad (127)$$

Recall the decomposition of a Gaussian vector  $\tilde{\mathbf{g}}_j \sim \mathcal{N}(0, \mathbf{I}_d)$  in (84). Note that  $\text{rank}(\text{Cov}(\tilde{\mathbf{g}}_j^\perp)) = \text{rank}(\mathbf{I}_d - \tilde{\boldsymbol{\theta}}\tilde{\boldsymbol{\theta}}^\top) = d - 1$ .

For any pair of labelled data sample  $(\mathbf{X}, Y)$ , from (85), we similarly decompose  $\mathbf{X}$  as  $\mathbf{X} = Y\boldsymbol{\mu} + \sigma(\tilde{g}\tilde{\boldsymbol{\theta}}_0 + \tilde{\mathbf{g}}^\perp)$ , where  $\tilde{g} \sim \mathcal{N}(0, 1)$  and  $\tilde{\mathbf{g}}^\perp \sim \mathcal{N}(0, \mathbf{I}_d - \tilde{\boldsymbol{\theta}}_0\tilde{\boldsymbol{\theta}}_0^\top)$ . Let  $p_{\tilde{g}}$  and  $p_{\tilde{\mathbf{g}}^\perp}$  denote the probability density functions of  $\tilde{g}$  and  $\tilde{\mathbf{g}}^\perp$ , respectively. For any  $\mathbf{x} = \boldsymbol{\mu} + \sigma(u\tilde{\boldsymbol{\theta}}_0 + \mathbf{u}^\perp) \in \mathbb{R}^d$ , the joint probability distribution at  $(\mathbf{X}, Y) = (\mathbf{x}, 1)$  is given by

$$\begin{aligned} P_{\mathbf{X}, Y}(\mathbf{x}, 1) &= P_Y(1)p_{\boldsymbol{\mu}}(\mathbf{x}|1) \\ &= \frac{P_Y(1)}{\sqrt{(2\pi)^d \sigma^d}} \exp\left(-\frac{1}{2\sigma^2}(\mathbf{x} - y\boldsymbol{\mu})^\top(\mathbf{x} - y\boldsymbol{\mu})\right) \end{aligned} \quad (128)$$

$$= \frac{P_Y(1)}{\sqrt{(2\pi)^d \sigma^d}} \exp\left(-\frac{1}{2\sigma^2}(\sigma u\tilde{\boldsymbol{\theta}} + \sigma\mathbf{u}^\perp)^\top(\sigma u\tilde{\boldsymbol{\theta}} + \sigma\mathbf{u}^\perp)\right) \quad (129)$$

$$= \frac{P_Y(y)}{\sqrt{(2\pi)^d \sigma^d}} \exp\left(-\frac{u^2}{2}\right) \exp\left(-\frac{(\mathbf{u}^\perp)^\top \mathbf{u}^\perp}{2}\right) \quad (130)$$

$$= P_Y(1)p_{\tilde{g}}(u)p_{\tilde{\mathbf{g}}^\perp}(\mathbf{u}^\perp). \quad (131)$$

Similarly, for any  $\mathbf{x} = -\boldsymbol{\mu} + \sigma(u\tilde{\boldsymbol{\theta}}_0 + \mathbf{u}^\perp) \in \mathbb{R}^d$ , the joint probability density evaluated at  $(\mathbf{X}, Y) = (\mathbf{x}, -1)$  is given by

$$P_{\mathbf{X}, Y}(\mathbf{x}, -1) = P_Y(-1)p_{\boldsymbol{\mu}}(\mathbf{x}|-1) = P_Y(-1)p_{\tilde{g}}(u)p_{\tilde{\mathbf{g}}^\perp}(\mathbf{u}^\perp). \quad (132)$$

Second, we have  $P_{\mathbf{X}'_j|\hat{Y}'_j} = \sum_{y \in \{-1, +1\}} P_{\mathbf{X}'_j|\hat{Y}'_j, Y'_j=y} P_{Y'_j=y|\hat{Y}'_j}$ . The conditional probability distribution  $P_{Y'_j|\hat{Y}'_j}$  can be calculated as follows

$$P_{Y'_j|\hat{Y}'_j} = \frac{P_{\hat{Y}'_j|Y'_j} P_{Y'_j}}{P_{\hat{Y}'_j}} = P_{\hat{Y}'_j|Y'_j}, \quad (133)$$

where the last equality follows since  $P_{Y'_j}(-1) = P_{Y'_j}(1) = P_{\hat{Y}'_j}(-1) = P_{\hat{Y}'_j}(1) = 1/2$ . Since  $\hat{Y}'_j = \text{sgn}(Y'_j\alpha + \sigma\tilde{g}_j)$  (cf. (86)), we have

$$P_{\hat{Y}'_j|Y'_j}(-1|-1) = \Pr(Y'_j\alpha + \sigma\tilde{g}_j < 0 | Y'_j = -1) = \mathbb{Q}\left(\frac{-\alpha}{\sigma}\right), \quad (134)$$

and similarly,

$$P_{\hat{Y}'_j|Y'_j}(1|-1) = \mathbb{Q}\left(\frac{\alpha}{\sigma}\right), \quad P_{\hat{Y}'_j|Y'_j}(-1|1) = \mathbb{Q}\left(\frac{\alpha}{\sigma}\right), \quad P_{\hat{Y}'_j|Y'_j}(1|1) = \mathbb{Q}\left(\frac{-\alpha}{\sigma}\right). \quad (135)$$

Thus, we conclude that

$$P_{Y'_j|\hat{Y}'_j}(y'_j|\hat{y}'_j) = \begin{cases} \mathbb{Q}\left(\frac{-\alpha}{\sigma}\right) & y'_j = \hat{y}'_j \\ \mathbb{Q}\left(\frac{\alpha}{\sigma}\right) & y'_j \neq \hat{y}'_j. \end{cases} \quad (136)$$

To calculate the conditional probability distribution  $P_{\mathbf{X}'_j|\hat{Y}'_j, Y'_j}$ , recall the decomposition of  $\mathbf{X}'_j$  and  $\tilde{\boldsymbol{\theta}}_0^\top \mathbf{X}'_j$  in (85) and (86). Since the event  $\{\hat{Y}'_j = -1, Y'_j = -1\}$  is equivalent to  $\{\tilde{g}_j < \alpha/\sigma\}$  and  $\tilde{g}_j \sim \mathcal{N}(0, 1)$ , the conditional density of  $\tilde{g}_j$  given  $\hat{Y}'_j = -1, Y'_j = -1$  is given by

$$p_{\tilde{g}_j|\hat{Y}'_j, Y'_j}(u|-1, -1) = p_{\tilde{g}_j|\tilde{g}_j \leq \alpha/\sigma}(u) = \frac{\mathbb{1}\{u \leq \alpha/\sigma\} p_{\tilde{g}_j}(u)}{\Phi(\alpha/\sigma)}, \quad \forall u \in \mathbb{R}. \quad (137)$$

Similarly, for any  $u \in \mathbb{R}$

$$p_{\tilde{g}_j|\hat{Y}'_j, Y'_j}(u|-1, 1) = p_{\tilde{g}_j|\tilde{g}_j \leq -\alpha/\sigma}(u) = \frac{\mathbb{1}\{u \leq -\alpha/\sigma\} f_{\tilde{g}_j}(u)}{\Phi(-\alpha/\sigma)}, \quad (138)$$

$$p_{\tilde{g}_j|\hat{Y}'_j, Y'_j}(u|1, -1) = p_{\tilde{g}_j|\tilde{g}_j > \alpha/\sigma}(u) = \frac{\mathbb{1}\{u > \alpha/\sigma\} f_{\tilde{g}_j}(u)}{\mathbb{Q}(\alpha/\sigma)}, \quad (139)$$

$$p_{\tilde{g}_j|\hat{Y}'_j, Y'_j}(u|1, 1) = p_{\tilde{g}_j|\tilde{g}_j > -\alpha/\sigma}(u) = \frac{\mathbb{1}\{u > -\alpha/\sigma\} p_{\tilde{g}_j}(u)}{\mathbb{Q}(-\alpha/\sigma)}. \quad (140)$$

For any  $\mathbf{x} = \boldsymbol{\mu} + \sigma(u\bar{\boldsymbol{\theta}}_0 + \mathbf{u}^\perp) \in \mathbb{R}^d$ , given  $\hat{Y}'_j = 1, Y'_j = 1$ , the conditional probability distribution at  $\mathbf{X}'_j = \mathbf{x}$  is given by

$$P_{\mathbf{X}'_j|\hat{Y}'_j,Y'_j}(\mathbf{x}|1,1) = P_{\boldsymbol{\mu}+\sigma\tilde{g}_j\bar{\boldsymbol{\theta}}_0+\sigma\tilde{\mathbf{g}}_j^\perp|\hat{Y}'_j,Y'_j}(\boldsymbol{\mu} + \sigma(u\bar{\boldsymbol{\theta}}_0 + \mathbf{u}^\perp)|1,1) \quad (141)$$

$$= P_{\sigma\tilde{g}_j\bar{\boldsymbol{\theta}}_0+\sigma\tilde{\mathbf{g}}_j^\perp|\hat{Y}'_j,Y'_j}(\sigma(u\bar{\boldsymbol{\theta}}_0 + \mathbf{u}^\perp)|1,1) \quad (142)$$

$$= p_{\tilde{g}_j|\hat{Y}'_j,Y'_j}(u|1,1)p_{\tilde{\mathbf{g}}_j^\perp}(\mathbf{u}^\perp), \quad (143)$$

where (143) follows since  $\tilde{g}_j$  and  $\tilde{\mathbf{g}}_j^\perp$  are mutually independent and  $\bar{\boldsymbol{\theta}}_0 \perp \tilde{\mathbf{g}}_j^\perp$ . Since we can decompose  $2\boldsymbol{\mu}/\sigma$  as

$$\frac{2\boldsymbol{\mu}}{\sigma} = \frac{2\alpha\bar{\boldsymbol{\theta}}_0 + 2\beta^2\boldsymbol{\mu} - 2\alpha\beta\mathbf{v}}{\sigma} = \frac{2\alpha}{\sigma}\bar{\boldsymbol{\theta}}_0 + \bar{\boldsymbol{\theta}}_0^\perp, \quad (144)$$

given  $\hat{Y}'_j = 1, Y'_j = -1$ , the conditional probability distribution at  $\mathbf{X}'_j = \mathbf{x}$  is given by

$$P_{\mathbf{X}'_j|\hat{Y}'_j,Y'_j}(\mathbf{x}|1,-1) = P_{-\boldsymbol{\mu}+\sigma\tilde{g}_j\bar{\boldsymbol{\theta}}_0+\sigma\tilde{\mathbf{g}}_j^\perp|\hat{Y}'_j,Y'_j}(\boldsymbol{\mu} + \sigma(u\bar{\boldsymbol{\theta}}_0 + \mathbf{u}^\perp)|1,-1) \quad (145)$$

$$= P_{\sigma\tilde{g}_j\bar{\boldsymbol{\theta}}_0+\sigma\tilde{\mathbf{g}}_j^\perp|\hat{Y}'_j,Y'_j}\left(\sigma\left(\frac{2\boldsymbol{\mu}}{\sigma} + u\bar{\boldsymbol{\theta}}_0 + \mathbf{u}^\perp\right)\middle|1,-1\right) \quad (146)$$

$$= p_{\tilde{g}_j|\hat{Y}'_j,Y'_j}\left(u + \frac{2\alpha}{\sigma}\middle|1,-1\right)p_{\tilde{\mathbf{g}}_j^\perp}(\mathbf{u}^\perp + \bar{\boldsymbol{\theta}}_0^\perp). \quad (147)$$

Similarly, for any  $\mathbf{x} = -\boldsymbol{\mu} + \sigma(u\bar{\boldsymbol{\theta}}_0 + \mathbf{u}^\perp) \in \mathbb{R}^d$ , given  $\hat{Y}'_j = -1, Y'_j = 1$ , the conditional distribution at  $\mathbf{X}'_j = \mathbf{x}$  is given by

$$P_{\mathbf{X}'_j|\hat{Y}'_j,Y'_j}(\mathbf{x}|-1,1) = P_{\boldsymbol{\mu}+\sigma\tilde{g}_j\bar{\boldsymbol{\theta}}_0+\sigma\tilde{\mathbf{g}}_j^\perp|\hat{Y}'_j,Y'_j}(-\boldsymbol{\mu} + \sigma(u\bar{\boldsymbol{\theta}}_0 + \mathbf{u}^\perp)|-1,1) \quad (148)$$

$$= p_{\tilde{g}_j|\hat{Y}'_j,Y'_j}\left(u - \frac{2\alpha}{\sigma}\middle|-1,1\right)p_{\tilde{\mathbf{g}}_j^\perp}(\mathbf{u}^\perp - \bar{\boldsymbol{\theta}}_0^\perp); \quad (149)$$

and given  $\hat{Y}'_j = -1, Y'_j = -1$ ,

$$P_{\mathbf{X}'_j|\hat{Y}'_j,Y'_j}(\mathbf{x}|-1,-1) = P_{-\boldsymbol{\mu}+\sigma\tilde{g}_j\bar{\boldsymbol{\theta}}_0+\sigma\tilde{\mathbf{g}}_j^\perp|\hat{Y}'_j,Y'_j}(-\boldsymbol{\mu} + \sigma(u\bar{\boldsymbol{\theta}}_0 + \mathbf{u}^\perp)|-1,-1) \quad (150)$$

$$= p_{\tilde{g}_j|\hat{Y}'_j,Y'_j}(u|-1,-1)p_{\tilde{\mathbf{g}}_j^\perp}(\mathbf{u}^\perp). \quad (151)$$

Furthermore, for any  $\mathbf{x} = -\boldsymbol{\mu} + \sigma(u\bar{\boldsymbol{\theta}}_0 + \mathbf{u}^\perp) \in \mathbb{R}^d$ , we have

$$P_{\mathbf{X}'_j|\hat{Y}'_j=-1}(\mathbf{x}) = \sum_{y \in \{-1, +1\}} P_{\mathbf{X}'_j|\hat{Y}'_j=-1, Y'_j=y}(\mathbf{x})P_{Y'_j|\hat{Y}'_j=-1}(y) \quad (152)$$

$$= P_{Y'_j|\hat{Y}'_j=-1}(1)p_{\tilde{g}_j|\hat{Y}'_j,Y'_j}\left(u - \frac{2\alpha}{\sigma}\middle|-1,1\right)p_{\tilde{\mathbf{g}}_j^\perp}(\mathbf{u}^\perp - \bar{\boldsymbol{\theta}}_0^\perp) \\ + P_{Y'_j|\hat{Y}'_j=-1}(-1)p_{\tilde{g}_j|\hat{Y}'_j,Y'_j}(u|-1,-1)p_{\tilde{\mathbf{g}}_j^\perp}(\mathbf{u}^\perp) \quad (153)$$

$$= \mathbb{1}\left\{u \leq \frac{\alpha}{\sigma}\right\}p_{\tilde{g}_j}\left(u - \frac{2\alpha}{\sigma}\right)p_{\tilde{\mathbf{g}}_j^\perp}(\mathbf{u}^\perp - \bar{\boldsymbol{\theta}}_0^\perp) + \mathbb{1}\left\{u \leq \frac{\alpha}{\sigma}\right\}p_{\tilde{g}_j}(u)p_{\tilde{\mathbf{g}}_j^\perp}(\mathbf{u}^\perp); \quad (154)$$

for any  $\mathbf{x} = \boldsymbol{\mu} + \sigma(u\bar{\boldsymbol{\theta}}_0 + \mathbf{u}^\perp) \in \mathbb{R}^d$ , we have

$$P_{\mathbf{X}'_j|\hat{Y}'_j=1}(\mathbf{x}) = \sum_{y \in \{-1, +1\}} P_{\mathbf{X}'_j|\hat{Y}'_j=1, Y'_j=y}(\mathbf{x})P_{Y'_j|\hat{Y}'_j=1}(y) \quad (155)$$

$$= \mathbb{1}\left\{u > -\frac{\alpha}{\sigma}\right\}p_{\tilde{g}_j}\left(u + \frac{2\alpha}{\sigma}\right)p_{\tilde{\mathbf{g}}_j^\perp}(\mathbf{u}^\perp + \bar{\boldsymbol{\theta}}_0^\perp) + \mathbb{1}\left\{u > -\frac{\alpha}{\sigma}\right\}p_{\tilde{g}_j}(u)p_{\tilde{\mathbf{g}}_j^\perp}(\mathbf{u}^\perp). \quad (156)$$

Define the set  $\mathcal{U}_0^\perp(\xi_0, \boldsymbol{\mu}^\perp) := \{\mathbf{u}^\perp \in \mathbb{R}^d : \mathbf{u}^\perp \perp \boldsymbol{\theta}_0\}$ . We also use  $\mathcal{U}_0^\perp$  to represent  $\mathcal{U}_0^\perp(\xi_0, \boldsymbol{\mu}^\perp)$ , if there is no risk of confusion. Recall (21) and note that  $\int_{\mathcal{U}_0^\perp} p_{\tilde{\mathbf{g}}_j^\perp}(\mathbf{u}^\perp) d\mathbf{u}^\perp = 1$ .

Finally, the KL-divergence is given by

$$\begin{aligned} D_{\xi_0, \mu^\perp}(P_{\mathbf{X}'_j|\hat{Y}'_j=-1} \| P_{\mathbf{X}|Y=-1}) \\ = \int_{\mathcal{U}_0^\perp} \int_{-\infty}^{\frac{\alpha}{\sigma}} \left( p_{\tilde{g}_j} \left( u - \frac{2\alpha}{\sigma} \right) p_{\tilde{\mathbf{g}}_j^\perp}(\mathbf{u}^\perp - \bar{\boldsymbol{\theta}}_0^\perp) + p_{\tilde{g}_j}(u) p_{\tilde{\mathbf{g}}_j^\perp}(\mathbf{u}^\perp) \right) \\ \times \log \left( 1 + \frac{p_{\tilde{g}_j} \left( u - \frac{2\alpha}{\sigma} \right) p_{\tilde{\mathbf{g}}_j^\perp}(\mathbf{u}^\perp - \bar{\boldsymbol{\theta}}_0^\perp)}{p_{\tilde{g}_j}(u) p_{\tilde{\mathbf{g}}_j^\perp}(\mathbf{u}^\perp)} \right) du d\mathbf{u}^\perp \end{aligned} \quad (157)$$

$$= G_\sigma(\alpha, \xi_0, \mu^\perp) \quad (158)$$

and

$$\begin{aligned} D_{\xi_0, \mu^\perp}(P_{\mathbf{X}'_j|\hat{Y}'_j=1} \| P_{\mathbf{X}|Y=1}) \\ = \int_{\mathcal{U}_0^\perp} \int_{-\frac{\alpha}{\sigma}}^{+\infty} \left( p_{\tilde{g}_j} \left( u + \frac{2\alpha}{\sigma} \right) p_{\tilde{\mathbf{g}}_j^\perp}(\mathbf{u}^\perp + \bar{\boldsymbol{\theta}}_0^\perp) + p_{\tilde{g}_j}(u) p_{\tilde{\mathbf{g}}_j^\perp}(\mathbf{u}^\perp) \right) \\ \times \log \left( 1 + \frac{p_{\tilde{g}_j} \left( u + \frac{2\alpha}{\sigma} \right) p_{\tilde{\mathbf{g}}_j^\perp}(\mathbf{u}^\perp + \bar{\boldsymbol{\theta}}_0^\perp)}{p_{\tilde{g}_j}(u) p_{\tilde{\mathbf{g}}_j^\perp}(\mathbf{u}^\perp)} \right) du d\mathbf{u}^\perp \end{aligned} \quad (159)$$

$$= G_\sigma(\alpha, \xi_0, \mu^\perp), \quad (160)$$

where (160) follows from since  $p_{\tilde{g}_j}$  and  $p_{\tilde{\mathbf{g}}_j^\perp}$  are zero-mean Gaussian distributions. Then from (127), we have

$$D_{\xi_0, \mu^\perp}(P_{\mathbf{X}'_j, \hat{Y}'_j} \| P_{\mathbf{X}, Y}) = G_\sigma(\alpha, \xi_0, \mu^\perp). \quad (161)$$

Thus, by combining the aforementioned results, we get the closed-form expression of the upper bound for  $|\text{gen}_1|$ . Indeed, if we fix some  $d \in \mathbb{N}$ ,  $\epsilon > 0$  and  $\delta \in (0, 1)$ , there exists  $n_0(d, \delta) \in \mathbb{N}$ ,  $m_0(\epsilon, d, \delta) \in \mathbb{N}$ ,  $c_0(d, \delta) \in (\bar{c}_1, \infty)$ ,  $r_0(d, \delta) \in \mathbb{R}_+$  such that for all  $n > n_0$ ,  $m > m_0$ ,  $c > c_0$ ,  $r > r_0$ ,  $\delta_{m, c - \bar{c}_1, d} < \frac{\delta}{3}$ ,  $\delta_{r, d} < \frac{\delta}{3}$ , and with probability at least  $1 - \delta$ ,

$$|\text{gen}_1| \leq \sqrt{\frac{(c_2 - c_1)^2}{2}} \mathbb{E}_{\xi_0, \mu^\perp} \left[ \sqrt{G_\sigma(\alpha(\xi_0, \mu^\perp), \xi_0, \mu^\perp) + \epsilon} \right]. \quad (162)$$

**4. Pseudo-label using  $\boldsymbol{\theta}_1$ :** Let  $\bar{\boldsymbol{\theta}}_1 := \boldsymbol{\theta}_1 / \|\boldsymbol{\theta}_1\|_2$ . For any  $i \in [m+1 : 2m]$ , the pseudo-labels are given by

$$\hat{Y}'_i = \text{sgn}(\boldsymbol{\theta}_1^\top \mathbf{X}'_i) = \text{sgn}(\bar{\boldsymbol{\theta}}_1^\top \mathbf{X}'_i). \quad (163)$$

It can be seen that the pseudo-labels  $\{\hat{Y}'_i\}_{i=m+1}^{2m}$  are conditionally i.i.d. given  $\boldsymbol{\theta}_1$  and let us denote the conditional distribution under fixed  $\boldsymbol{\theta}_1$  as  $P_{\hat{Y}'_i|\boldsymbol{\theta}_1} \in \mathcal{P}(\mathcal{Y})$ . The pseudo-labelled dataset is denoted as  $\hat{S}_{u,2} = \{(\mathbf{X}'_i, \hat{Y}'_i)\}_{i=m+1}^{2m}$ .

For any fixed  $\bar{\boldsymbol{\theta}}_1 \in \Theta$ , we can decompose it as  $\bar{\boldsymbol{\theta}}_1 = \alpha'_1 \boldsymbol{\mu} + \beta'_1 \mathbf{v}$ , where  $\alpha'_1 \in [-1, 1]$  and  $\beta'_1 = \sqrt{1 - (\alpha'_1)^2}$ . Recall the decomposition of  $\mathbf{X}'_i$  and  $\bar{\boldsymbol{\theta}}_0^\top \mathbf{X}'_i$  in (69) and (73). Similarly, we have

$$\bar{\boldsymbol{\theta}}_1^\top \mathbf{X}'_i =: Y'_i \alpha'_1 + \sigma h_i^1, \quad (164)$$

where  $h_i^1 \sim \mathcal{N}(0, 1)$ . Note that  $P_{\hat{Y}'_i|\boldsymbol{\theta}_1, \xi_0, \mu^\perp} = P_{\hat{Y}'_i|\boldsymbol{\theta}_1}$  and then the conditional probability  $P_{\hat{Y}'_i|\boldsymbol{\theta}_1, \xi_0, \mu^\perp}$  can be given by

$$P_{\hat{Y}'_i|\boldsymbol{\theta}_1, \xi_0, \mu^\perp}(1) = P_{\hat{Y}'_i|\boldsymbol{\theta}_1}(1) = \mathbb{P}_{\boldsymbol{\theta}_1}(\bar{\boldsymbol{\theta}}_1^\top \mathbf{X}'_i > 0) \quad (165)$$

$$= \frac{1}{2} \mathbb{P}_{\boldsymbol{\theta}_1}(\alpha'_1 + \sigma h_i^1 > 0) + \frac{1}{2} \mathbb{P}_{\boldsymbol{\theta}_1}(\alpha'_1 + \sigma h_i^1 \leq 0) = \frac{1}{2}, \quad (166)$$

and  $P_{\hat{Y}'_i|\boldsymbol{\theta}_1, \xi_0, \mu^\perp}(-1) = 1/2$ , where  $\mathbb{P}_{\boldsymbol{\theta}_1}$  denotes the probability measure under parameter  $\boldsymbol{\theta}_1$ .

5. **Iteration  $t = 2$ :** Recall (18) and the new model parameter learned from the pseudo-labelled dataset  $\hat{\mathcal{S}}_{u,2}$  is given by

$$\theta_2 = \frac{1}{m} \sum_{i=m+1}^{2m} \hat{Y}'_i \mathbf{X}'_i = \frac{1}{m} \sum_{i=m+1}^{2m} \text{sgn}(\bar{\theta}_1^\top \mathbf{X}'_i) \mathbf{X}'_i, \quad (167)$$

where  $\{\text{sgn}(\bar{\theta}_1^\top \mathbf{X}'_i) \mathbf{X}'_i\}_{i=m+1}^{2m}$  are conditionally i.i.d. random variables given  $\theta_1, \xi_0, \mu^\perp$ .

Given any  $(\theta_1, \xi_0, \mu^\perp)$ , for any  $j \in [m+1 : 2m]$ , let  $\mu_2^{\theta_1, \xi_0, \mu^\perp} := \mathbb{E}[\text{sgn}(\bar{\theta}_1^\top \mathbf{X}'_j) \mathbf{X}'_j | \theta_1, \xi_0, \mu^\perp]$  and  $\mathbb{P}_{\theta_1, \xi_0, \mu^\perp}$  denotes the probability measure under the parameters  $\theta_1, \xi_0, \mu^\perp$ . Following the similar steps that derive (114), for any  $\varepsilon > 0$ , we have

$$\mathbb{P}_{\theta_1, \xi_0, \mu^\perp} (\|\theta_2 - \mu_2^{\theta_1, \xi_0, \mu^\perp}\|_\infty > \varepsilon) \leq \delta_{m, \varepsilon, d}. \quad (168)$$

From (98), no matter what  $\theta_1$  is, we always have  $\|\mu_2^{\theta_1, \xi_0, \mu^\perp} - \mu^{\xi_0, \mu^\perp}\| \leq \tilde{c}_1$ . Then, for some  $c \in (\tilde{c}_1, \infty)$ ,

$$\mathbb{P}_{\theta_1, \xi_0, \mu^\perp} (\theta_2 \in \Theta_{\mu, c}) \geq 1 - \delta_{m, c - \tilde{c}_1, d}. \quad (169)$$

With probability at least  $(1 - \delta_{m, c - \tilde{c}_1, d})(1 - \delta_{r, d})$ , the absolute generalization error can be upper bounded as follows:

$$|\text{gen}_2| = |\mathbb{E}[L_{P_Z}(\theta_2) - L_{\hat{\mathcal{S}}_{u,2}}(\theta_2)]| \quad (170)$$

$$= \left| \frac{1}{m} \sum_{i=m+1}^{2m} \mathbb{E}_{\theta_1, \xi_0, \mu^\perp} \left[ \mathbb{E} \left[ l(\theta_2, (\mathbf{X}, Y)) - l(\theta_2, (\mathbf{X}'_i, \hat{Y}'_i)) | \theta_1, \xi_0, \mu^\perp \right] \right] \right| \quad (171)$$

$$\leq \sqrt{\frac{(c_2 - c_1)^2}{2}} \times \frac{1}{m} \sum_{i=m+1}^{2m} \mathbb{E}_{\theta_1, \xi_0, \mu^\perp} \left[ \sqrt{I_{\theta_1, \xi_0, \mu^\perp}(\theta_2; (\mathbf{X}'_i, \hat{Y}'_i)) + D_{\theta_1, \xi_0, \mu^\perp}(P_{\mathbf{X}'_i, \hat{Y}'_i} \| P_{\mathbf{X}, Y})} \right], \quad (172)$$

where  $P_{\theta_2, \mathbf{X}, Y | \theta_1, \xi_0, \mu^\perp} = P_{\theta_2 | \theta_1, \xi_0, \mu^\perp} \otimes P_{\mathbf{X}, Y}$ .

Similar to (126), for any  $\epsilon > 0$  and  $\delta \in (0, 1)$ , there exists  $m_1(\epsilon, d, \delta)$  such that for all  $m > m_1$ ,

$$\mathbb{P}_{\theta_1, \xi_0, \mu^\perp} (I_{\theta_1, \xi_0, \mu^\perp}(\theta_2; (\mathbf{X}'_i, \hat{Y}'_i)) > \epsilon) \leq \delta. \quad (173)$$

Recall (166) that  $P_{\hat{Y}'_i | \theta_1, \xi_0, \mu^\perp} \sim \text{unif}(\{-1, +1\})$ . For any fixed  $(\theta_1, \xi_0, \mu^\perp)$ , let  $\bar{\theta}_1$  be decomposed as  $\bar{\theta}_1 = \alpha'_1(\xi_0, \mu^\perp) \mu + \beta'_1(\xi_0, \mu^\perp) v$ , where  $\alpha'_1(\xi_0, \mu^\perp) \in [-1, 1]$  and  $\beta'_1(\xi_0, \mu^\perp) = \sqrt{1 - (\alpha'_1(\xi_0, \mu^\perp))^2}$ .

By following the similar steps in the first iteration, the disintegrated conditional KL-divergence between pseudo-labelled distribution and true distribution is given by

$$\begin{aligned} & D_{\theta_1, \xi_0, \mu^\perp} (P_{\mathbf{X}'_i, \hat{Y}'_i} \| P_{\mathbf{X}, Y}) \\ &= \frac{1}{2} D_{\theta_1, \xi_0, \mu^\perp} (P_{\mathbf{X}'_i | \hat{Y}'_i = -1} \| P_{\mathbf{X} | Y = -1}) + \frac{1}{2} D_{\theta_1, \xi_0, \mu^\perp} (P_{\mathbf{X}'_i | \hat{Y}'_i = 1} \| P_{\mathbf{X} | Y = 1}) \end{aligned} \quad (174)$$

$$= G_\sigma(\alpha'_1(\xi_0, \mu^\perp), \xi_0, \mu^\perp), \quad (175)$$

Given any pair of  $(\xi_0, \mu^\perp)$ , recall the decomposition of  $\mu_1^{\xi_0, \mu^\perp}$  in (96). Then the correlation between  $\mu_1^{\xi_0, \mu^\perp}$  and  $\mu$  is given by

$$\rho(\mu_1^{\xi_0, \mu^\perp}, \mu) = \frac{1 - 2Q(\frac{\alpha}{\sigma}) + \frac{2\sigma\alpha}{\sqrt{2\pi}} \exp(-\frac{\alpha^2}{2\sigma^2})}{\sqrt{(1 - 2Q(\frac{\alpha}{\sigma}) + \frac{2\sigma\alpha}{\sqrt{2\pi}} \exp(-\frac{\alpha^2}{2\sigma^2}))^2 + \frac{2\sigma^2(1-\alpha^2)}{\pi} \exp(-\frac{\alpha^2}{\sigma^2})}} \quad (176)$$

$$= F_\sigma(\alpha(\xi_0, \mu^\perp)). \quad (177)$$

By the strong law of large numbers, we have  $\alpha'_1(\xi_0, \boldsymbol{\mu}^\perp) \xrightarrow{\text{a.s.}} F_\sigma(\alpha(\xi_0, \boldsymbol{\mu}^\perp))$  as  $m \rightarrow \infty$ . Then for any  $\epsilon > 0$  and  $\delta \in (0, 1)$ , there exists  $m_2(\epsilon, d, \delta)$  such that for all  $m > m_2$ ,

$$\mathbb{P}_{\boldsymbol{\theta}_1, \xi_0, \boldsymbol{\mu}^\perp} \left( \left| G_\sigma(\alpha'_1(\xi_0, \boldsymbol{\mu}^\perp), \xi_0, \boldsymbol{\mu}^\perp) - G_\sigma(F_\sigma(\alpha(\xi_0, \boldsymbol{\mu}^\perp)), \xi_0, \boldsymbol{\mu}^\perp) \right| > \epsilon \right) \leq \delta. \quad (178)$$

Therefore, fix some  $d \in \mathbb{N}$ ,  $\epsilon > 0$  and  $\delta \in (0, 1)$ . There exists  $n_0(d, \delta) \in \mathbb{N}$ ,  $m_3(\epsilon, d, \delta) \in \mathbb{N}$ ,  $c_0(d, \delta) \in (\tilde{c}_1, \infty)$ ,  $r_0(d, \delta) \in \mathbb{R}_+$  such that for all  $n > n_0, m > m_3, c > c_0, r > r_0$ ,  $\delta_{m, c-\tilde{c}_1, d} < \frac{\delta}{3}$ ,  $\delta_{r, d} < \frac{\delta}{3}$ , and then with probability at least  $1 - \delta$ , the absolute generalization error at  $t = 2$  can be upper bounded as follows:

$$|\text{gen}_2| \leq \sqrt{\frac{(c_2 - c_1)^2}{2}} \mathbb{E}_{\xi_0, \boldsymbol{\mu}^\perp} \left[ \sqrt{G_\sigma(F_\sigma(\alpha(\xi_0, \boldsymbol{\mu}^\perp)), \xi_0, \boldsymbol{\mu}^\perp)} + \epsilon \right]. \quad (179)$$

**6. Any iteration  $t \in [3 : \tau]$ :** By similarly repeating the calculation in iteration  $t = 2$ , we obtain the upper bound for  $|\text{gen}_t|$  in (24).

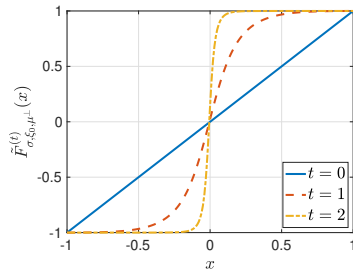
**Remark 2** (Extra remarks about Theorem 2). *In the other extreme case, when  $\alpha = \rho(\boldsymbol{\theta}_0, \boldsymbol{\mu}) = -1$  and  $\bar{\boldsymbol{\theta}}_0 = -\boldsymbol{\mu}$ , the error probability  $\Pr(Y'_j \neq Y_j) = 1 - Q(1/\sigma) > \frac{1}{2}$  (for all  $\sigma > 0$ ) and  $D_{\xi_0, \boldsymbol{\mu}^\perp}(P_{\mathbf{X}'_j, \hat{Y}'_j} \| P_{\mathbf{X}, Y}) < \infty$ , so in this other extreme (flipped) scenario, we have more mistakes than correct pseudo-labels. The reason why  $D_{\alpha, \boldsymbol{\mu}^\perp}(P_{\mathbf{X}'_j, \hat{Y}'_j} \| P_{\mathbf{X}, Y})$  is finite is that when  $P_{\mathbf{X}, Y}(\mathbf{x}, y)$  is small, it means that  $\mathbf{x}$  is far from both  $-\boldsymbol{\mu}$  and  $\boldsymbol{\mu}$ , and then  $P_{\mathbf{X}}(\mathbf{x})$  is also small. Thus,  $P_{\mathbf{X}'_j, \hat{Y}'_j}(\mathbf{x}, y) = P_{\hat{Y}'_j | \mathbf{X}'_j}(y | \mathbf{x}) P_{\mathbf{X}}(\mathbf{x})$  is also small.*

## D REUSING $S_1$ IN EACH ITERATION

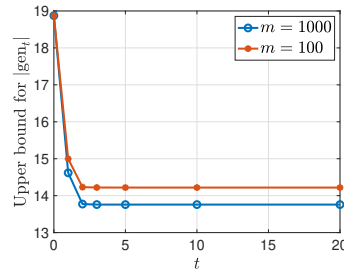
If the labelled data  $S_1$  are reused in each iteration and  $w = \frac{n}{n+m}$  (cf. (5)), for each  $t \in [1 : \tau]$ , the learned model parameter is given by

$$\boldsymbol{\theta}_t = \frac{n}{n+m} \boldsymbol{\theta}_{t-1} + \frac{1}{n+m} \sum_{i=(t-1)m+1}^{tm} \hat{Y}'_i \mathbf{X}'_i \quad (180)$$

$$= \frac{n}{n+m} \boldsymbol{\theta}_{t-1} + \frac{1}{n+m} \sum_{i=(t-1)m+1}^{tm} \text{sgn}(\bar{\boldsymbol{\theta}}_{t-1}^\top \mathbf{X}'_i) \mathbf{X}'_i. \quad (181)$$



**Figure 10:**  $\tilde{F}_{\sigma, \xi_0, \boldsymbol{\mu}^\perp}^{(t)}(x)$  versus  $x$  under  $t \in \{0, 1, 2\}$  when  $\sigma = 0.5$ ,  $\xi_0 = 0$ ,  $\|\boldsymbol{\mu}^\perp\|_2 = 1$ ,  $n = 10$ ,  $m = 1000$ .



**Figure 11:** Upper bound for  $|\text{gen}_t|$  versus  $t$  for  $m = 100$  and  $m = 1000$ , when  $n = 10$ ,  $\sigma = 0.6$ ,  $d = 2$ ,  $\boldsymbol{\mu} = (1, 0)$ .

Recall the definition of the function  $\tilde{F}_{\sigma, \xi_0, \boldsymbol{\mu}^\perp}$  in (27). Let the  $t$ -th iterate of  $\tilde{F}_{\sigma, \xi_0, \boldsymbol{\mu}^\perp}$  be denoted as  $\tilde{F}_{\sigma, \xi_0, \boldsymbol{\mu}^\perp}^{(t)}$  with initial condition  $\tilde{F}_{\sigma, \xi_0, \boldsymbol{\mu}^\perp}^{(0)}(x) = x$ . As shown in Figure 10, we can see that for any fixed  $(\sigma, \xi_0, \boldsymbol{\mu}^\perp)$ ,  $\tilde{F}_{\sigma, \xi_0, \boldsymbol{\mu}^\perp}^{(t)}$  has a similar behaviour as  $F_\sigma^{(t)}$  as  $t$  increases, which implies that the upper bound in (28) in Corollary 3 also decreases as  $t$  increases. As a result,  $\tilde{F}_{\sigma, \xi_0, \boldsymbol{\mu}^\perp}^{(t)}$  represents the improvement of the model parameter  $\boldsymbol{\theta}_t$  over the iterations.

As shown in Figure 11, under the same setup as Figure 6(c), when the labelled data  $S_1$  are reused in each iteration, the upper bound for  $|\text{gen}_t|$  is also a decreasing function of  $t$ . When  $m = 1000$ , the upper bound is almost the same as that one in Figure 6(c), which means that for large enough  $m/n$ , reusing the labelled data does not necessarily help to improve the generalization performance. Moreover, when  $m = 100$ , the upper bound is higher than that for  $m = 1000$ , which coincides with the intuition that increasing the number of unlabelled data helps to reduce the generalization error.

## E PROOF OF COROLLARY 3

Following the similar steps in Appendix C we first derive the upper bound for  $|\text{gen}_1|$  as follows.

At  $t = 1$ , from (65) and (96), the expectation  $\mu_1^{\xi_0, \mu^\perp} = \mathbb{E}[\theta_1 | \xi_0, \mu^\perp]$  is rewritten as

$$\begin{aligned} \mu_1^{\xi_0, \mu^\perp} &= \frac{n}{n+m} \theta_0 + \frac{1}{n+m} \sum_{i=1}^m \mathbb{E}[\text{sgn}(\bar{\theta}_0^\top \mathbf{X}'_i) \mathbf{X}'_i | \xi_0, \mu^\perp] \\ &= \frac{n}{n+m} \left( \left( 1 + \frac{\sigma}{\sqrt{n}} \xi_0 \right) \mu + \frac{\sigma}{\sqrt{n}} \mu^\perp \right) \\ &\quad + \frac{m}{n+m} \left( \left( 1 - 2Q\left(\frac{\alpha}{\sigma}\right) + \frac{2\sigma\alpha}{\sqrt{2\pi}} \exp\left(-\frac{\alpha^2}{2\sigma^2}\right) \right) \mu + \frac{2\sigma\beta}{\sqrt{2\pi}} \exp\left(-\frac{\alpha^2}{2\sigma^2}\right) \mathbf{v} \right) \end{aligned} \quad (182)$$

$$\begin{aligned} &= \left( 1 + \frac{\sqrt{n}\sigma\xi_0}{n+m} + \frac{m}{n+m} \left( -2Q\left(\frac{\alpha}{\sigma}\right) + \frac{2\sigma\alpha}{\sqrt{2\pi}} \exp\left(-\frac{\alpha^2}{2\sigma^2}\right) \right) \right) \mu \\ &\quad + \left( \frac{\sqrt{n}\sigma\|\mu^\perp\|_2}{n+m} + \frac{m}{n+m} \frac{2\sigma\beta}{\sqrt{2\pi}} \exp\left(-\frac{\alpha^2}{2\sigma^2}\right) \right) \mathbf{v}. \end{aligned} \quad (183)$$

Then the correlation between  $\mu_1^{\xi_0, \mu^\perp}$  and  $\mu$  is given by

$$\rho(\mu_1^{\xi_0, \mu^\perp}, \mu) = \tilde{F}_{\sigma, \xi_0, \mu^\perp}(\alpha). \quad (184)$$

Let  $\theta'_1 = \frac{1}{m} \sum_{i=1}^m \text{sgn}(\bar{\theta}_0^\top \mathbf{X}'_i) \mathbf{X}'_i$ . For some  $c \in (\tilde{c}_1, \infty)$ , from (53) and (114),

$$\begin{aligned} &\Pr\left(\|\theta_1 - \mu\|_\infty > c\right) \\ &\leq \Pr\left(\frac{n}{n+m} \|\theta_0 - \mu\|_\infty + \frac{m}{n+m} \|\theta'_1 - \mu\|_\infty > c\right) \end{aligned} \quad (185)$$

$$\leq \Pr\left(\|\theta_0 - \mu\|_\infty > c\right) + \Pr\left(\|\theta'_1 - \mu\|_\infty > c\right) \quad (186)$$

$$\leq \delta_{\sqrt{n}c, d} + \delta_{m, c-\tilde{c}_1, d} \quad (187)$$

Thus, from Theorem 1 for some  $c \in (\tilde{c}_1, \infty)$ , with probability at least  $(1 - \delta_{\sqrt{n}c, d} - \delta_{m, c-\tilde{c}_1, d})(1 - \delta_{r, d})$ , the absolute generalization error can be upper bounded as follows:

$$\begin{aligned} &|\text{gen}_1| \\ &= \left| \frac{1}{n+m} \sum_{i=1}^n \mathbb{E}[l(\theta_1, (\mathbf{X}, Y)) - l(\theta_1, (\mathbf{X}_i, Y_i))] \right. \\ &\quad \left. + \frac{1}{n+m} \sum_{i=1}^m \mathbb{E}_{\xi_0, \mu^\perp} \left[ \mathbb{E}[l(\theta_1, (\mathbf{X}, Y)) - l(\theta_1, (\mathbf{X}'_i, \hat{Y}'_i)) | \xi_0, \mu^\perp] \right] \right| \end{aligned} \quad (188)$$

$$\begin{aligned} &\leq \frac{1}{n+m} \sum_{i=1}^n \sqrt{\frac{(c_2 - c_1)^2}{2} I(\theta_1; (\mathbf{X}_i, Y_i))} \\ &\quad + \frac{1}{n+m} \sum_{i=1}^m \mathbb{E}_{\xi_0, \mu^\perp} \left[ \sqrt{\frac{(c_2 - c_1)^2}{2} \left( I_{\xi_0, \mu^\perp}(\theta_1; (\mathbf{X}'_i, \hat{Y}'_i)) + D_{\xi_0, \mu^\perp}(P_{\mathbf{X}'_i, \hat{Y}'_i} \| P_{\mathbf{X}, Y}) \right)} \right] \end{aligned} \quad (189)$$

$$\begin{aligned}
&\leq \frac{1}{n+m} \sum_{i=1}^n \sqrt{\frac{(c_2 - c_1)^2}{2}} I(\theta_0; (\mathbf{X}_i, Y_i)) \\
&\quad + \frac{1}{n+m} \sum_{i=1}^m \mathbb{E}_{\xi_0, \mu^\perp} \left[ \sqrt{\frac{(c_2 - c_1)^2}{2}} \left( I_{\xi_0, \mu^\perp}(\theta_1; (\mathbf{X}'_i, \hat{Y}'_i)) + D_{\xi_0, \mu^\perp}(P_{\mathbf{X}'_i, \hat{Y}'_i} \| P_{\mathbf{X}, Y}) \right) \right] \\
&\hspace{15em} (190)
\end{aligned}$$

$$\begin{aligned}
&= \frac{1}{n+m} \sum_{i=1}^n \sqrt{\frac{(c_2 - c_1)^2 d}{4}} \log \frac{n}{n-1} \\
&\quad + \frac{1}{n+m} \sum_{i=1}^m \mathbb{E}_{\xi_0, \mu^\perp} \left[ \sqrt{\frac{(c_2 - c_1)^2}{2}} \left( I_{\xi_0, \mu^\perp}(\theta_1; (\mathbf{X}'_i, \hat{Y}'_i)) + D_{\xi_0, \mu^\perp}(P_{\mathbf{X}'_i, \hat{Y}'_i} \| P_{\mathbf{X}, Y}) \right) \right], \\
&\hspace{15em} (191)
\end{aligned}$$

where  $P_{\theta_1, (\mathbf{X}, Y) | \xi_0, \mu^\perp} = Q_{\theta_1 | \xi_0, \mu^\perp} \otimes P_{\mathbf{X}, Y}$  and  $Q_{\theta_1 | \xi_0, \mu^\perp}$  denotes the marginal distribution of  $\theta_1$  under parameters  $(\xi_0, \mu^\perp)$ , and (190) follows since  $(\mathbf{X}_i, Y_i) - \theta_0 - \theta_1$  forms a Markov chain.

In (191), the KL-divergence is already given in (161) and the disintegrated conditional mutual information can be calculated as follows. Since we have

$$I_{\xi_0, \mu^\perp}(\theta_1; (\mathbf{X}'_i, \hat{Y}'_i)) = I_{\xi_0, \mu^\perp} \left( \frac{1}{m} \sum_{i=1}^m \text{sgn}(\bar{\theta}_{i-1}^\top \mathbf{X}'_i) \mathbf{X}'_i; (\mathbf{X}'_i, \hat{Y}'_i) \right), \quad (192)$$

from (126), for any  $\epsilon > 0$  and any  $\delta \in (0, 1)$ , there exists  $m'_1(\epsilon, \delta, d) \in \mathbb{N}$  such that for all  $m > m'_1(\epsilon, \delta, d)$ ,

$$\mathbb{P}_{\xi_0, \mu^\perp} \left( I_{\xi_0, \mu^\perp}(\theta_1; \mathbf{X}'_i, \hat{Y}'_i) > \epsilon \right) \leq 1 - \delta. \quad (193)$$

Therefore, fix  $d \in \mathbb{N}$ , any  $\epsilon > 0$  and any  $\delta \in (0, 1)$ , and there exists  $n_0(d, \delta) \in \mathbb{N}$ ,  $m_4(\epsilon, d, \delta) \in \mathbb{N}$ ,  $c_0(d, \delta) \in (\tilde{c}_1, \infty)$ ,  $r_0(d, \delta) \in \mathbb{R}_+$  such that for all  $n > n_0, m > m_3, c > c_0, r > r_0$ ,  $\delta_{\sqrt{nc}, d} < \frac{\delta}{6}$ ,  $\delta_{m, c - \tilde{c}_1, d} < \frac{\delta}{6}$ ,  $\delta_{r, d} < \frac{\delta}{3}$ , and with probability at least  $1 - \delta$ , the absolute generalization error  $|\text{gen}_1|$  can be upper bounded as follows:

$$\begin{aligned}
|\text{gen}_1| &\leq w \sqrt{\frac{(c_2 - c_1)^2 d}{4}} \log \frac{n}{n-1} \\
&\quad + (1-w) \sqrt{\frac{(c_2 - c_1)^2}{2}} \mathbb{E}_{\xi_0, \mu^\perp} \left[ \sqrt{G_\sigma(\alpha(\xi_0, \mu^\perp), \xi_0, \mu^\perp) + \epsilon} \right]. \\
&\hspace{15em} (194)
\end{aligned}$$

For  $t \geq 2$ , the only difference from the derivation in Appendix C is the correlation function  $\tilde{F}_{\sigma, \xi_0, \mu^\perp}(\cdot)$  (compared to (177)). Thus by replacing  $F_\sigma(\cdot)$  with  $\tilde{F}_{\sigma, \xi_0, \mu^\perp}(\cdot)$ , we obtain the upper bound in (28), completing the proof of Corollary 3.

## F ADDITIONAL EXPERIMENTS

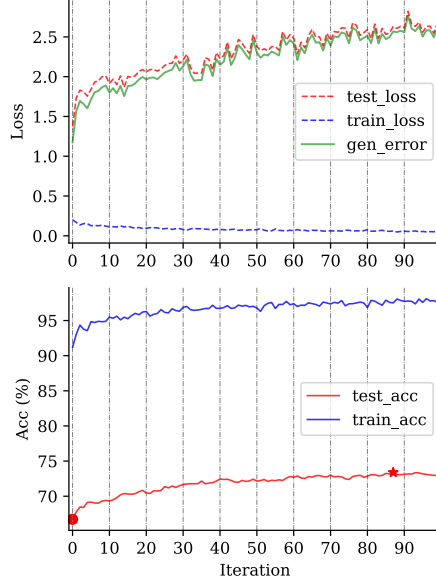
**Table 1:** The  $l_2$  distances between the RGB-mean and RGB-variance of different pairs of classes from the CIFAR10 dataset.

Classes	RGB-mean $l_2$ distance	RGB-variance $l_2$ distance	Difficulty
horse-ship	0.0180	3.90e-05	Easy
automobile-truck	0.0038	7.06e-05	Moderate
cat-dog	0.0007	4.95e-05	Challenging

In Table 1, we display the RGB means and variances of the test data in six classes taken from the CIFAR10 dataset. We observe that the RGB variances of each pair are almost 0 (and small compared to the RGB-mean  $l_2$  distances), and thus, the RGB-mean  $l_2$  distance is indicative of the difficulty of



the classification task. Indeed, a smaller RGB-mean  $l_2$  distance implies a higher overlap of the two classes and consequently, greater difficulty in distinguishing them. Therefore, the “cat-dog” pair, which is more difficult to disambiguate compared to the “horse-ship” and “automobile-truck” pairs, is analogous to the bGMM with large variance (i.e. large overlap between the  $\{\pm 1\}$  classes).



**Figure 12:** Binary classification on “cat” and “dog” from the CIFAR10 dataset.

Under the same experimental settings as in Section 5, we perform another classification experiment on the “cat-dog” pair (from the CIFAR10 dataset). As shown in Figure 12, the test accuracy at the initial point when only labelled data are used is about 65% and the test loss is about 1.4; these are much worse than the performances of the classification tasks as shown in Figures 7 and 8, which means that the two classes are more challenging to classify.

It can be observed from Figure 12 that although the training loss decreases and the test and training accuracies increase as the iteration count increases (which are expected), the test loss and the generalization error both increase. The fact that both the test loss and test accuracy appear to increase with  $t$  is, in fact, not contradictory. To intuitively explain this, in binary classification using the softmax (hence, logistic) function to predict the output classes, the learned probability of a data example belonging to its true class is  $p \in [0, 1]$  and if  $p \in (1/2, 1]$ , the classification is correct. In other words, the accuracy is 100%. However, when  $p$  (i.e., the classification confidence) decreases towards  $(1/2)^+$ , the corresponding decision margin  $2p - 1$  (Cao et al., 2019) also decreases and the test loss  $-\log p$  increases commensurately. Thus, when the decision margin is small, even though the test accuracy may increase as the iteration counter  $t$  increases, the test loss (representing our lack of confidence) may also increase at the same time.

In summary, for the classification task involving the “cat” and “dog” classes, our above observations correspond to that for the bGMM in Figure 6(d), namely that the unlabelled data does not help to improve the generalization error when the classification task is challenging and the initialization with the labelled data  $S_1$  does not already result in a relatively high accuracy.