
How does semi-supervised learning with pseudo-labelers work? A case study

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Semi-supervised learning is a popular machine learning paradigm that utilizes
2 a large amount of unlabeled data as well as a small amount of labeled data to
3 facilitate learning tasks. While semi-supervised learning has achieved great success
4 in training neural networks, its theoretical understanding remains largely open. In
5 this paper, we aim to theoretically understand a semi-supervised learning approach
6 based on pre-training and linear probing. We prove that, under a certain data
7 generation model and two-layer convolutional neural network, the semi-supervised
8 learning approach can achieve nearly zero test loss, while a neural network directly
9 trained by supervised learning on the same amount of labeled data can only achieve
10 constant test loss. Through this case study, we demonstrate a separation between
11 semi-supervised learning and supervised learning in terms of test loss provided the
12 same amount of labeled data.

13 1 Introduction

14 *Semi-supervised learning* (Scudder, 1965; Fralick, 1967; Agrawala, 1970), which leverages both a
15 small amount of labeled data and a large amount of unlabeled data to improve learning performance,
16 is one of the most widely used approaches. It has been shown to achieve promising performance
17 for a wide variety of tasks, including image classification (Rasmus et al., 2015; Springenberg, 2015;
18 Laine and Aila, 2016), image generation (Kingma et al., 2014; Odena, 2016; Salimans et al., 2016),
19 domain adaptation (Saito et al., 2017; Shu et al., 2018; Lee et al., 2019), and word embedding (Turian
20 et al., 2010; Peters et al., 2017). One of the popular semi-supervised learning approaches is *pseudo-*
21 *labeling* (Lee et al., 2013; Xie et al., 2020; Pham et al., 2021b; Rizve et al., 2021), which generates
22 pseudo-labels of unlabeled data for pre-training. This approach has been remarkably successful in
23 improving performance on many tasks. In this paper, we attempt to theoretically explain the success
24 of semi-supervised learning with pseudo-labelers in training neural networks. The contributions of
25 our work are summarized as follows.

- 26 • We theoretically show that with the help of pseudo-labelers, CNN can learn the feature representa-
27 tion during the pre-training stage. Moreover, the learned feature is highly correlated with the true
28 labels of the data, even though the true labels are not used during the pre-training stage.
- 29 • Based on our analysis of the pre-training process, we further show that when linear-probing the
30 pre-trained model in the downstream task, the final classifier can achieve near-zero test loss and
31 test error. Notably, these guarantees of small test loss and error only require a very small number
32 of labeled training data.
- 33 • As a comparison, we show that standard supervised learning cannot learn a good classifier under
34 the same setting. Specifically, we show that, even when the training process converges to a global
35 minimum of the training loss, the learned two-layer CNN can only achieve constant level test
36 loss. This, together with the aforementioned results for semi-supervised learning, demonstrates the
37 advantage of semi-supervised learning over standard supervised learning.

Submitted to 37th Conference on Neural Information Processing Systems (NeurIPS 2023). Do not distribute.

38 2 Problem Setup and Preliminaries

39 In this section, we will introduce our data model, the convolutional neural network, and the details of
 40 the training algorithms considered in this paper. Inspired by recent work (Allen-Zhu and Li, 2020b;
 41 Zou et al., 2021; Shen et al., 2022; Cao et al., 2022), we consider a data model where each data input
 42 \mathbf{x} consists of two patches $\mathbf{x}^{(1)}$ and $\mathbf{x}^{(2)}$, where each patch has d dimensions. We focus on the binary
 43 classification task and present our data distribution \mathcal{D} as follows.

44 **Data distribution.** Each data point (\mathbf{x}, y) with $\mathbf{x} = [\mathbf{x}^{(1)\top}, \mathbf{x}^{(2)\top}]^\top \in \mathbb{R}^{2d}$ and $y \in \{-1, +1\}$ is
 45 generated as follows: the label y is generated as a Rademacher random variable; one of $\mathbf{x}^{(1)}, \mathbf{x}^{(2)}$ is
 46 given by the feature vector $y \cdot \mathbf{v}$, the other is given by a noise vector $\boldsymbol{\xi}$ that is generated from a $2d$ -
 47 dimensional Gaussian distribution $\mathcal{N}(\mathbf{0}, \sigma_p^2(\mathbf{I} - \mathbf{v}\mathbf{v}^\top / \|\mathbf{v}\|_2^2))$. We denote by \mathcal{D} the joint distribution
 48 of (\mathbf{x}, y) , and denote by $\mathcal{D}_{\mathbf{x}}$ the marginal distribution of \mathbf{x} .

49 2.1 Supervised Learning Models

50 For supervised learning, we consider a two-layer CNN whose filters are applied to the patches $\mathbf{x}^{(1)}$
 51 and $\mathbf{x}^{(2)}$ respectively and parameters in the second layers are set to be ± 1 . Then the CNN can be
 52 written as $f_{\mathbf{W}}(\mathbf{x}) = f_{\mathbf{W}}^{+1}(\mathbf{x}) - f_{\mathbf{W}}^{-1}(\mathbf{x})$ where $f_{\mathbf{W}}(\mathbf{x})^{+1}, f_{\mathbf{W}}(\mathbf{x})^{-1}$ are formulated as

$$f_{\mathbf{W}}^{+1}(\mathbf{x}) = \sum_{j=1}^m \left[\sigma(\langle \mathbf{w}_j, \mathbf{x}^{(1)} \rangle) + \sigma(\langle \mathbf{w}_j, \mathbf{x}^{(2)} \rangle) \right], f_{\mathbf{W}}^{-1}(\mathbf{x}) = \sum_{j=m+1}^{2m} \left[\sigma(\langle \mathbf{w}_j, \mathbf{x}^{(1)} \rangle) + \sigma(\langle \mathbf{w}_j, \mathbf{x}^{(2)} \rangle) \right]. \quad (2.1)$$

53 Here σ is activation function $\text{ReLU}^q(\cdot) = [\cdot]_+^q (q > 2)$, m is the width of the network, $\mathbf{w}_j \in \mathbb{R}^d$
 54 denotes the j -th filter, and \mathbf{W} is the collection of all filters $\{\mathbf{w}_j\}_{j=1}^{2m}$. Given labeled training dataset
 55 $S' = \{(\mathbf{x}'_i, y'_i)\}_{i=1}^{n_1}$, we train the CNN model by minimizing the empirical cross-entropy loss

$$L_{S'}(\mathbf{W}) = \frac{1}{n_1} \sum_{i=1}^{n_1} L_i(\mathbf{W}),$$

56 where $L_i(\mathbf{W}) = \ell(y'_i \cdot f_{\mathbf{W}}(\mathbf{x}'_i))$ with $\ell(z) = \log(1 + \exp(-z))$ denotes the individual loss for the
 57 training example (\mathbf{x}_i, y_i) . We minimize the empirical function $L_{S'}(\mathbf{W})$ with gradient descent as
 58 follows

$$\mathbf{w}_j^{(t+1)} = \mathbf{w}_j^{(t)} - \eta \cdot \nabla_{\mathbf{w}_j} L_{S'}(\mathbf{W}^{(t)}), \quad \mathbf{w}_j^{(0)} \sim \mathcal{N}(\mathbf{0}, \sigma_0^2 \mathbf{I}), \quad j \in [2m],$$

59 where $\eta > 0$ is the learning rate and σ_0 defines the scale of random initialization.

60 2.2 Semi-supervised Learning Models

61 For semi-supervised pre-training, we assume that we have access to K pseudo-labelers $\{f_k^w\}_{k=1}^K$.
 62 The accuracy of k -th pseudo-labeler is $p_k \in (1/2, 1)$. Then we use K pseudo-labelers to generate K
 63 pseudo-labeled dataset $\{S_k\}_{k=1}^K$, where $S_k := \{(\mathbf{x}_i, \hat{y}_{k,i}) \mid \hat{y}_{k,i} = f_k^w(\mathbf{x}_i)\}_{i=1}^{n_u}$. Next we solve K pre-
 64 training tasks with two-layer CNN models $\{f_{\mathbf{W}_k}\}_{k=1}^K$ defined in (2.1) using $\{S_k\}_{k=1}^K$ respectively.
 65 We consider learning the model parameter \mathbf{W}_k by optimizing the empirical loss of both pseudo-
 66 labeled dataset S_k and labeled dataset $S' = \{(\mathbf{x}'_i, y'_i)\}_{i=1}^{n_1}$ with weight decay regularization

$$L_{S_k \cup S'}(\mathbf{W}_k) = \frac{1}{n_u + n_1} \left(\sum_{i=1}^{n_u} L_i(\mathbf{W}_k) + \sum_{i'=1}^{n_1} L_{i'}(\mathbf{W}_k) \right) + \frac{\lambda}{2} \|\mathbf{W}_k\|_F^2,$$

67 where $\lambda \geq 0$ is the regularization parameter, $L_i(\mathbf{W}_k) = \ell(\hat{y}_{k,i} \cdot f_{\mathbf{W}_k}(\mathbf{x}_i))$ denotes the individual loss
 68 for the pseudo-labeled data $L_{i'}(\mathbf{W}_k) = \ell(y'_{i'} \cdot f_{\mathbf{W}_k}(\mathbf{x}'_{i'}))$ denotes the individual loss for the labeled
 69 data $(\mathbf{x}'_{i'}, y'_{i'})$. We also use gradient descent to minimize the regularized loss function $L_{S_k \cup S'}(\mathbf{W}_k)$
 70 starting from $\mathbf{w}_{k,j}^{(0)} \sim \mathcal{N}(\mathbf{0}, \sigma_0^2 \mathbf{I}_d)$.

71 **Downstream Task: Linear Model.** The semi-supervised pre-training gives us K CNN models
 72 with parameters $\{\mathbf{W}_k^*\}_{k=1}^K$. Based on them, for the downstream task, we consider a linear model

$$g_{\mathbf{a}}(\mathbf{x}) = \sum_{k=1}^K a_k f_{\mathbf{W}_k^*}(\mathbf{x}),$$

73 where $a_k \in \mathbb{R}$ denotes the trainable weight for the k -th pre-trained model. Then, given $\{f_{\mathbf{W}_k}\}_{k=1}^K$
 74 and labeled training data $S' = \{(\mathbf{x}'_i, y'_i)\}_{i=1}^n$, we consider learning the downstream linear model
 75 parameter \mathbf{a} by optimizing the following empirical loss

$$L_{S'}(\mathbf{a}) = \frac{1}{n} \sum_{i=1}^n \ell(y'_i \cdot g_{\mathbf{a}}(\mathbf{x}'_i)).$$

76 We initialize \mathbf{a} as an all-zero vector and optimize the empirical loss by gradient descent with learning
 77 rate η , i.e.,

$$\mathbf{a}^{(t+1)} = \mathbf{a}^{(t)} - \eta \cdot \nabla_{\mathbf{a}} L_{S'}(\mathbf{a}^{(t)}), \quad \mathbf{a}^{(0)} = \mathbf{0}.$$

78 3 Main Results

79 In this section, we start with a condition that is required by our analysis.

80 **Condition 3.1.** The strength of the signal is $\|\mathbf{v}\|_2^2 = \Theta(d)$, the noise variance is $\sigma_p = \Theta(d^\epsilon)$,
 81 where $0 < \epsilon < 1/8$ is a small constant, and the width of the network satisfies $m = \text{polylog}(d)$.
 82 We also assume that the size of the unlabeled dataset $n_u = \Omega(d^{4\epsilon})$, and labeled data $n_l = \Theta(1)$.
 83 For both supervised learning and semi-supervised learning settings, we initialize the weight with
 84 $\sigma_0 = \Theta(d^{-3/4})$. For semi-supervised learning, we require $\lambda = o(d^{3/4})$ and assume that there exists
 85 a constant C such that for all pseudo-labelers, their test accuracy $p_k > 1/2 + C$.

86 Next, we present the main theoretical results in this paper.

87 **Theorem 3.2** (Semi-supervised Learning: Pre-training). Let $k \in [K]$ and consider the semi-
 88 supervised pre-training of $f_{\mathbf{W}_k}(\mathbf{x})$. For any test data point (\mathbf{x}, y) , denote $\hat{y} = f_k^w(\mathbf{x})$. Then
 89 under Condition 3.1, after $T_0 = \tilde{\Theta}(d^{q/4-3/2}\eta^{-1})$ training iterations with learning rate $\eta = O(d^{-1.1})$,
 90 the trained neural network $f_{\mathbf{W}_k^{(T_0)}}(\mathbf{x})$ can achieve nearly 0 test error on the distribution \mathcal{D} .

91 Theorem 3.2 characterizes the prediction power of the feature representation learned in the pre-trained
 92 models using unlabeled data. For any test data point (\mathbf{x}, y) , the sign of y can be predicted based on
 93 $f_{\mathbf{W}_k^{(T_0)}}(\mathbf{x})$ with high probability.

94 **Theorem 3.3** (Semi-supervised Learning: Downstream). Let $\{f_{\mathbf{W}_k^{(T_0^k)}}\}_{k=1}^d$ be the neural networks
 95 trained according to the K pre-training tasks, and consider the learning of the downstream task based
 96 in $\{f_{\mathbf{W}_k^{(T_0^k)}}\}_{k=1}^d$. Under Condition 3.1, after $T' = \Theta(d^{0.1}/\eta)$ iterations with learning rate $\eta = \Theta(1)$,
 97 with probability $1 - o(1)$, the obtained $\mathbf{a}^{(T')}$ satisfies:

- 98 • Training error is 0: $\frac{1}{n} \sum_{i=1}^n \mathbb{1}[y_i \cdot g_{\mathbf{a}^{(T')}}(\mathbf{x}_i) \leq 0] = 0$.
- 99 • Test error and loss are nearly 0: $\mathbb{P}_{(\mathbf{x}, y) \sim \mathcal{D}}[y \cdot g_{\mathbf{a}^{(T')}}(\mathbf{x}) \leq 0] = o(1), L_{\mathcal{D}}(\mathbf{a}^{(T')}) = o(1)$.

100 Theorem 3.3 shows that the feature representation learned based on the semi-supervised pre-training
 101 can ensure small training and test errors for the supervised downstream task. Notably, this result
 102 holds even though we assume that there are only a constant number of labeled data. This shows
 103 that semi-supervised learning can significantly reduce the need for a large labeled training dataset.
 104 For comparison, we also have the following guarantees on the performance of standard supervised
 105 learning of CNNs.

106 **Theorem 3.4** (Supervised Learning). Under supervised learning setting, after gradient descent for
 107 $T = \tilde{\Theta}(d^{(1/4-\epsilon)q-3/2}\eta^{-1})$ iterations with learning rate $\eta = O(d^{-1-2\epsilon})$, then there exists $t \leq T$
 108 such that with probability $1 - o(1)$ the CNNs defined in (2.1) with parameter $\mathbf{W}^{(t)}$ satisfies:

- 109 • Training loss is nearly zero: $L_{S'}(\mathbf{W}^{(t)}) = o(1)$.
- 110 • Test loss is high: $L_{\mathcal{D}}(\mathbf{W}^{(t)}) = \Theta(1)$.

111 Theorem 3.4 shows that although standard supervised learning can train a CNN model with nearly
 112 zero training loss, the obtained CNN model generalizes poorly to test data. Comparing Theorem 3.4
 113 with Theorem 3.3 shows that the generalization of semi-supervised learning and supervised learning
 114 are largely different. The reason behind this difference is that the pre-training, with a relatively large
 115 number of unlabeled training data, helps learn a feature representation that captures the feature in

Table 1: Training error and loss, test error and loss for semi-supervised and supervised learning.

	Semi-supervised		Supervised
	Pre-train	Downstream	
Training error	0.1753 ± 0.0259	0	0
Test error	0	0	0.4982 ± 0.0208
Training loss	0.4155 ± 0.0418	0.0150 ± 0.0022	$(6.473 \pm 5.031) \times 10^{-7}$
Test loss	0.2200 ± 0.0886	0.0182 ± 0.0021	0.6931 ± 0.0005

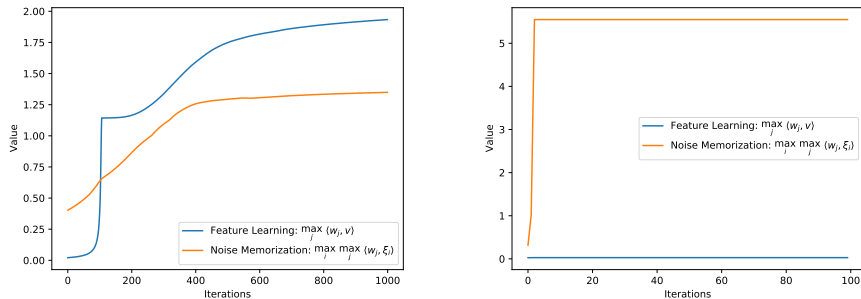


Figure 1: Visualization of the feature learning and noise memorization in the training process. (Left: Semi-supervised, Right: Supervised)

116 our data model, while direct application of supervised learning can only memorize the noises in the
 117 training dataset, which is independent of the labels of the data.

118 4 Experiments

119 In this section, we perform numerical experiments on synthetic datasets, generated according to the
 120 data distribution in Section 2, to verify our main theoretical results. The detailed experiment setting
 121 can be seen from Appendix B.

122 For semi-supervised learning, we first use a plain classifier to generate n_u pseudo-labels for unlabeled
 123 samples in order to help semi-supervised learning. After that, for pre-training, we use these pseudo-
 124 labeled samples and n_l labeled samples together to train a CNN. After 200 iterations, we can obtain a
 125 CNN model with a training error close to the error of pseudo-labeler and zero test error, according
 126 to Table 1. For the downstream task, we use n_l labeled samples to train a linear probe. After 100
 127 iterations, we can obtain a final model with low training and test loss as well as 100% training
 128 accuracy and test accuracy. For supervised learning, we directly use n_l labeled data to train the same
 129 CNN model. After 200 iterations, we obtain a CNN with 0 training error and small training loss,
 130 about 0.5 test error, and high test loss, which indicates supervised learning will give a model that
 131 behaves badly and even no better than a random guess.

132 Moreover, we also calculate the inner products representing feature learning and noise memorization
 133 respectively, to verify our key lemmas. The results are reported in Figure 1. It can be seen from Figure
 134 1 that under semi-supervised learning setting the algorithm will the feature learning will dominate the
 135 noise memorization though the noise patch has a larger norm than the signal patch, while under the
 136 supervised learning setting, the algorithm will entirely forget the feature but fit noise.

137 5 Conclusion

138 In this paper, we study semi-supervised learning with pseudo-labelers and provide a theoretical
 139 understanding of the success of semi-supervised learning. We show the advantage of semi-supervised
 140 learning over supervised learning through a case study. By considering a simple data model and two-
 141 layer CNN, we present a comprehensive analysis of the training procedure from a beyond-NTK
 142 feature learning perspective. We prove that the final classifier of a semi-supervised learning scenario
 143 can achieve near-zero test loss and error with only a small number of labeled training data, while its
 144 supervised-learned counterpart fails to achieve the same performance with the same data complexity.

145 **References**

- 146 AGRAWALA, A. (1970). Learning with a probabilistic teacher. *IEEE Transactions on Information*
147 *Theory* **16** 373–379.
- 148 ALLEN-ZHU, Z. and LI, Y. (2020a). Feature purification: How adversarial training performs robust
149 deep learning. *arXiv preprint arXiv:2005.10190* .
- 150 ALLEN-ZHU, Z. and LI, Y. (2020b). Towards understanding ensemble, knowledge distillation and
151 self-distillation in deep learning. *arXiv preprint arXiv:2012.09816* .
- 152 BALCAN, M.-F. and BLUM, A. (2010). A discriminative model for semi-supervised learning.
153 *Journal of the ACM (JACM)* **57** 1–46.
- 154 BELKIN, M., NIYOGI, P. and SINDHWANI, V. (2006). Manifold regularization: A geometric
155 framework for learning from labeled and unlabeled examples. *Journal of machine learning*
156 *research* **7**.
- 157 BENNETT, K. and DEMIRIZ, A. (1998). Semi-supervised support vector machines. *Advances in*
158 *Neural Information processing systems* **11**.
- 159 BLUM, A. and MITCHELL, T. (1998). Combining labeled and unlabeled data with co-training. In
160 *Proceedings of the eleventh annual conference on Computational learning theory*.
- 161 BROCK, A., DONAHUE, J. and SIMONYAN, K. (2018). Large scale gan training for high fidelity
162 natural image synthesis. *arXiv preprint arXiv:1809.11096* .
- 163 CAO, Y., CHEN, Z., BELKIN, M. and GU, Q. (2022). Benign overfitting in two-layer convolutional
164 neural networks. *arXiv preprint arXiv:2202.06526* .
- 165 CARON, M., MISRA, I., MAIRAL, J., GOYAL, P., BOJANOWSKI, P. and JOULIN, A. (2020).
166 Unsupervised learning of visual features by contrasting cluster assignments. *Advances in Neural*
167 *Information Processing Systems* **33** 9912–9924.
- 168 CASTELLI, V. and COVER, T. (1996). The relative value of labeled and unlabeled samples in pattern
169 recognition with an unknown mixing parameter. *IEEE Transactions on Information Theory* **42**
170 2102–2117.
- 171 CASTELLI, V. and COVER, T. M. (1995). On the exponential value of labeled samples. *Pattern*
172 *Recognition Letters* **16** 105–111.
- 173 CHAPELLE, O., SCHOLKOPF, B. and ZIEN, A. (2010). *Semi-supervised learning*. The MIT Press.
- 174 CHEN, T., KORNBLITH, S., NOROUZI, M. and HINTON, G. (2020). A simple framework for
175 contrastive learning of visual representations. In *International conference on machine learning*.
176 PMLR.
- 177 DARNSTÄDT, M. (2015). *An investigation on the power of unlabeled data*. Ph.D. thesis, Bochum,
178 Ruhr-Universität Bochum, Diss., 2015.
- 179 FRALICK, S. (1967). Learning to recognize patterns without a teacher. *IEEE Transactions on*
180 *Information Theory* **13** 57–64.
- 181 FREI, S., CHATTERJI, N. S. and BARTLETT, P. L. (2022a). Benign overfitting without linearity:
182 Neural network classifiers trained by gradient descent for noisy linear data. *arXiv preprint*
183 *arXiv:2202.05928* .
- 184 FREI, S., ZOU, D., CHEN, Z. and GU, Q. (2022b). Self-training converts weak learners to strong
185 learners in mixture models. In *International Conference on Artificial Intelligence and Statistics*.
186 PMLR.
- 187 GIDARIS, S., SINGH, P. and KOMODAKIS, N. (2018). Unsupervised representation learning by
188 predicting image rotations. *arXiv preprint arXiv:1803.07728* .

- 189 GLOBERSON, A., LIVNI, R. and SHALEV-SHWARTZ, S. (2017). Effective semisupervised learning
190 on manifolds. In *Conference on Learning Theory*. PMLR.
- 191 GOODFELLOW, I., POUGET-ABADIE, J., MIRZA, M., XU, B., WARDE-FARLEY, D., OZAIR,
192 S., COURVILLE, A. and BENGIO, Y. (2014). Generative adversarial nets. *Advances in neural
193 information processing systems* **27**.
- 194 HAOCHEN, J. Z., WEI, C., GAIDON, A. and MA, T. (2021). Provable guarantees for self-supervised
195 deep learning with spectral contrastive loss. *Advances in Neural Information Processing Systems*
196 **34**.
- 197 HE, K., FAN, H., WU, Y., XIE, S. and GIRSHICK, R. (2020). Momentum contrast for unsupervised
198 visual representation learning. In *Proceedings of the IEEE/CVF conference on computer vision
199 and pattern recognition*.
- 200 JING, L., YANG, X., LIU, J. and TIAN, Y. (2018). Self-supervised spatiotemporal feature learning
201 via video rotation prediction. *arXiv preprint arXiv:1811.11387* .
- 202 KARRAS, T., LAINE, S. and AILA, T. (2019). A style-based generator architecture for generative
203 adversarial networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern
204 Recognition*.
- 205 KINGMA, D. P., MOHAMED, S., JIMENEZ REZENDE, D. and WELLING, M. (2014). Semi-
206 supervised learning with deep generative models. *Advances in neural information processing
207 systems* **27**.
- 208 LAINE, S. and AILA, T. (2016). Temporal ensembling for semi-supervised learning. *arXiv preprint
209 arXiv:1610.02242* .
- 210 LEE, D.-H. ET AL. (2013). Pseudo-label: The simple and efficient semi-supervised learning method
211 for deep neural networks. In *Workshop on challenges in representation learning, ICML*, vol. 3.
- 212 LEE, H.-Y., HUANG, J.-B., SINGH, M. and YANG, M.-H. (2017). Unsupervised representation
213 learning by sorting sequences. In *Proceedings of the IEEE International Conference on Computer
214 Vision*.
- 215 LEE, J. D., LEI, Q., SAUNSHI, N. and ZHUO, J. (2020). Predicting what you already know helps:
216 Provable self-supervised learning. *arXiv preprint arXiv:2008.01064* .
- 217 LEE, S., KIM, D., KIM, N. and JEONG, S.-G. (2019). Drop to adapt: Learning discriminative
218 features for unsupervised domain adaptation. In *Proceedings of the IEEE/CVF International
219 Conference on Computer Vision*.
- 220 LI, Y., PAN, Q., WANG, S., PENG, H., YANG, T. and CAMBRIA, E. (2019). Disentangled variational
221 auto-encoder for semi-supervised learning. *Information Sciences* **482** 73–85.
- 222 MILLER, D. J. and UYAR, H. (1996). A mixture of experts classifier with learning based on both
223 labelled and unlabelled data. In *Advances in Neural Information Processing Systems* (M. Mozer,
224 M. Jordan and T. Petsche, eds.), vol. 9. MIT Press.
- 225 MISRA, I., ZITNICK, C. L. and HEBERT, M. (2016). Shuffle and learn: unsupervised learning using
226 temporal order verification. In *European Conference on Computer Vision*. Springer.
- 227 MITROVIC, J., MCWILLIAMS, B., WALKER, J., BUESING, L. and BLUNDELL, C. (2020). Repre-
228 sentation learning via invariant causal mechanisms. *arXiv preprint arXiv:2010.07922* .
- 229 NIGAM, K., MCCALLUM, A. K., THRUN, S. and MITCHELL, T. (2000). Text classification from
230 labeled and unlabeled documents using em. *Machine learning* **39** 103–134.
- 231 NIYOGI, P. (2013). Manifold regularization and semi-supervised learning: Some theoretical analyses.
232 *Journal of Machine Learning Research* **14**.

- 233 NOROOZI, M. and FAVARO, P. (2016). Unsupervised learning of visual representations by solving
234 jigsaw puzzles. In *European conference on computer vision*. Springer.
- 235 ODENA, A. (2016). Semi-supervised learning with generative adversarial networks. *arXiv preprint*
236 *arXiv:1606.01583* .
- 237 OYMAK, S. and GULCU, T. C. (2021). Statistical and algorithmic insights for semi-supervised
238 learning with self-training. In *International Conference on Artificial Intelligence and Statistics*
239 *(AISTATS)*.
- 240 PATHAK, D., KRAHENBUHL, P., DONAHUE, J., DARRELL, T. and EFROS, A. A. (2016). Context
241 encoders: Feature learning by inpainting. In *Proceedings of the IEEE conference on computer*
242 *vision and pattern recognition*.
- 243 PETERS, M. E., AMMAR, W., BHAGAVATULA, C. and POWER, R. (2017). Semi-supervised
244 sequence tagging with bidirectional language models. *arXiv preprint arXiv:1705.00108* .
- 245 PHAM, H., DAI, Z., XIE, Q. and LE, Q. V. (2021a). Meta pseudo labels. In *Proceedings of the*
246 *IEEE/CVF Conference on Computer Vision and Pattern Recognition*.
- 247 PHAM, H., DAI, Z., XIE, Q., LUONG, M.-T. and LE, Q. V. (2021b). Meta pseudo labels. In *IEEE*
248 *Conference on Computer Vision and Pattern Recognition (CVPR)*.
- 249 RASMUS, A., BERGLUND, M., HONKALA, M., VALPOLA, H. and RAIKO, T. (2015). Semi-
250 supervised learning with ladder networks. *Advances in neural information processing systems*
251 **28**.
- 252 RIZVE, M. N., DUARTE, K., RAWAT, Y. S. and SHAH, M. (2021). In defense of pseudo-labeling: An
253 uncertainty-aware pseudo-label selection framework for semi-supervised learning. In *International*
254 *Conference on Learning Representations (ICLR)*.
- 255 SAITO, K., USHIKU, Y. and HARADA, T. (2017). Asymmetric tri-training for unsupervised domain
256 adaptation. In *International Conference on Machine Learning*. PMLR.
- 257 SAJJADI, M., JAVANMARDI, M. and TASDIZEN, T. (2016). Regularization with stochastic trans-
258 formations and perturbations for deep semi-supervised learning. *Advances in neural information*
259 *processing systems* **29**.
- 260 SALIMANS, T., GOODFELLOW, I., ZAREMBA, W., CHEUNG, V., RADFORD, A. and CHEN, X.
261 (2016). Improved techniques for training gans. *Advances in neural information processing systems*
262 **29**.
- 263 SAUNSHI, N., ASH, J., GOEL, S., MISRA, D., ZHANG, C., ARORA, S., KAKADE, S. and
264 KRISHNAMURTHY, A. (2022). Understanding contrastive learning requires incorporating inductive
265 biases. *arXiv preprint arXiv:2202.14037* .
- 266 SAUNSHI, N., PLEVRAKIS, O., ARORA, S., KHODAK, M. and KHANDEPARKAR, H. (2019). A the-
267 oretical analysis of contrastive unsupervised representation learning. In *International Conference*
268 *on Machine Learning*. PMLR.
- 269 SCUDDER, H. (1965). Probability of error of some adaptive pattern-recognition machines. *IEEE*
270 *Transactions on Information Theory* **11** 363–371.
- 271 SHEN, R., BUBECK, S. and GUNASEKAR, S. (2022). Data augmentation as feature manipulation.
272 In *International Conference on Machine Learning*. PMLR.
- 273 SHU, R., BUI, H. H., NARUI, H. and ERMON, S. (2018). A dirt-t approach to unsupervised domain
274 adaptation. *arXiv preprint arXiv:1802.08735* .
- 275 SINGH, A., NOWAK, R. and ZHU, J. (2008). Unlabeled data: Now it helps, now it doesn't. In
276 *Advances in Neural Information Processing Systems* (D. Koller, D. Schuurmans, Y. Bengio and
277 L. Bottou, eds.), vol. 21. Curran Associates, Inc.

- 278 SPRINGENBERG, J. T. (2015). Unsupervised and semi-supervised learning with categorical genera-
279 tive adversarial networks. *arXiv preprint arXiv:1511.06390* .
- 280 TARVAINEN, A. and VALPOLA, H. (2017). Mean teachers are better role models: Weight-averaged
281 consistency targets improve semi-supervised deep learning results. *Advances in neural information*
282 *processing systems* **30**.
- 283 TIAN, Y., YU, L., CHEN, X. and GANGULI, S. (2020). Understanding self-supervised learning with
284 dual deep networks. *arXiv preprint arXiv:2010.00578* .
- 285 TOSH, C., KRISHNAMURTHY, A. and HSU, D. (2021a). Contrastive estimation reveals topic
286 posterior information to linear models. *Journal of Machine Learning Research* **22** 1–31.
- 287 TOSH, C., KRISHNAMURTHY, A. and HSU, D. (2021b). Contrastive estimation reveals topic
288 posterior information to linear models. *Journal of Machine Learning Research* **22** 1–31.
- 289 TSAI, Y.-H. H., WU, Y., SALAKHUTDINOV, R. and MORENCY, L.-P. (2020). Demystifying
290 self-supervised learning: An information-theoretical framework. *arXiv preprint arXiv:2006.05576*
291 .
- 292 TULYAKOV, S., LIU, M.-Y., YANG, X. and KAUTZ, J. (2018). Mocogan: Decomposing motion
293 and content for video generation. In *Proceedings of the IEEE conference on computer vision and*
294 *pattern recognition*.
- 295 TURIAN, J., RATINOV, L. and BENGIO, Y. (2010). Word representations: a simple and general
296 method for semi-supervised learning. In *Proceedings of the 48th annual meeting of the association*
297 *for computational linguistics*.
- 298 VONDRICK, C., PIRSIAVASH, H. and TORRALBA, A. (2016). Generating videos with scene
299 dynamics. *Advances in neural information processing systems* **29** 613–621.
- 300 WANG, T. and ISOLA, P. (2020). Understanding contrastive representation learning through align-
301 ment and uniformity on the hypersphere. In *International Conference on Machine Learning*.
302 PMLR.
- 303 WEI, C., SHEN, K., CHEN, Y. and MA, T. (2020). Theoretical analysis of self-training with deep
304 networks on unlabeled data. *arXiv preprint arXiv:2010.03622* .
- 305 WEI, D., LIM, J. J., ZISSERMAN, A. and FREEMAN, W. T. (2018). Learning and using the arrow
306 of time. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*.
- 307 WEN, Z. and LI, Y. (2021). Toward understanding the feature learning process of self-supervised
308 contrastive learning. In *International Conference on Machine Learning*. PMLR.
- 309 XIE, Q., LUONG, M.-T., HOVY, E. and LE, Q. V. (2020). Self-training with noisy student improves
310 imagenet classification. In *Proceedings of the IEEE/CVF conference on computer vision and*
311 *pattern recognition*.
- 312 XU, Z., JIN, R., ZHU, J., KING, I. and LYU, M. (2007). Efficient convex relaxation for transductive
313 support vector machine. *Advances in neural information processing systems* **20**.
- 314 XU, Z., JIN, R., ZHU, J., KING, I., LYU, M. and YANG, Z. (2009). Adaptive regularization for
315 transductive support vector machine. *Advances in Neural Information Processing Systems* **22**.
- 316 ZHAI, X., OLIVER, A., KOLESNIKOV, A. and BEYER, L. (2019). S4l: Self-supervised semi-
317 supervised learning. In *Proceedings of the IEEE/CVF International Conference on Computer*
318 *Vision*.
- 319 ZHANG, R., ISOLA, P. and EFROS, A. A. (2016). Colorful image colorization. In *European*
320 *conference on computer vision*. Springer.
- 321 ZHOU, D., BOUSQUET, O., LAL, T., WESTON, J. and SCHÖLKOPF, B. (2003). Learning with local
322 and global consistency. *Advances in neural information processing systems* **16**.

- 323 ZHU, X., GHAMRANI, Z. and LAFFERTY, J. D. (2003). Semi-supervised learning using gaussian
324 fields and harmonic functions. In *Proceedings of the 20th International conference on Machine*
325 *learning (ICML-03)*.
- 326 ZHU, X. and GOLDBERG, A. B. (2009). Introduction to semi-supervised learning. *Synthesis lectures*
327 *on artificial intelligence and machine learning* **3** 1–130.
- 328 ZOU, D., CAO, Y., LI, Y. and GU, Q. (2021). Understanding the generalization of adam in learning
329 neural networks with proper regularization. *arXiv preprint arXiv:2108.11371* .

330 A Related Work

331 **Semi-supervised learning methods in practice.** Since the invention of semi-supervised learning
332 in Scudder (1965); Fralick (1967); Agrawala (1970), a wide range of semi-supervised learning
333 approaches have been proposed, including generative models (Miller and Uyar, 1996; Nigam et al.,
334 2000), semi-supervised support vector machines (Bennett and Demiriz, 1998; Xu et al., 2007, 2009),
335 graph-based methods (Zhu et al., 2003; Belkin et al., 2006; Zhou et al., 2003), and co-training (Blum
336 and Mitchell, 1998), etc. For a comprehensive review of classical semi-supervised learning methods,
337 please refer to Chapelle et al. (2010); Zhu and Goldberg (2009). In the past years, a number of
338 deep semi-supervised learning approaches have been proposed, such as generative methods (Odena,
339 2016; Li et al., 2019), consistency regularization methods (Sajjadi et al., 2016; Laine and Aila, 2016;
340 Rasmus et al., 2015; Tarvainen and Valpola, 2017) and pseudo-labeling methods (Lee et al., 2013;
341 Zhai et al., 2019; Xie et al., 2020; Pham et al., 2021a). In this work, we will focus on pseudo-labeling
342 methods.

343 **Theory of semi-supervised learning.** To understand semi-supervised learning, Castelli and Cover
344 (1995, 1996) studied the relative value of labeled data over unlabeled data under a parametric
345 assumption on the marginal distribution of input features. Later, a series of works proved that
346 semi-supervised learning can possess better sample complexity or generalization performance than
347 supervised learning under certain assumptions on the marginal distribution (Niyogi, 2013; Globerson
348 et al., 2017) or the ratio of labeled and unlabeled samples (Singh et al., 2008; Darnstädt, 2015), while
349 Balcan and Blum (2010) provided a unified PAC framework able to analyze both sample-complexity
350 and algorithmic issues. Oymak and Gulcu (2021); Frei et al. (2022b) considered semi-supervised
351 learning with pseudo-labelers by learning a linear classifier for mixture models and convergence to
352 Bayes-optimal predictor.

353 **Self-supervised learning in practice.** A closely related learning paradigm to semi-supervised
354 learning is called self-supervised learning, which creates human-designed supervised learning prob-
355 lems to leverage natural structures and learn representations from unlabeled data. Representative
356 self-supervised learning approaches include *contrastive learning* and *pretext-based self-supervised*
357 *learning*. Contrastive learning (Caron et al., 2020; He et al., 2020; Chen et al., 2020) aims to group
358 similar examples closer and dissimilar examples far from each other by utilizing a similarity metric,
359 while pretext-based self-supervised tries to learn a good representation from *pretext tasks* generated
360 from the unlabeled data to facilitate *downstream learning tasks*. In practice, various pretext tasks
361 have been proposed, which include (1) generation-based ones such as colorizing grayscale images
362 (Zhang et al., 2016), image inpainting (Pathak et al., 2016), image and video generation with GAN
363 (Goodfellow et al., 2014; Brock et al., 2018; Karras et al., 2019; Vondrick et al., 2016; Tulyakov et al.,
364 2018); and (2) context-based ones such as image jigsaw puzzle (Noroozi and Favaro, 2016), geometric
365 transformation (Gidaris et al., 2018; Jing et al., 2018), frame order verification and recognition (Lee
366 et al., 2017; Misra et al., 2016; Wei et al., 2018). The semi-supervised learning approach with
367 pseudo-labelers studied in this paper is related to pretext-based self-supervised learning because the
368 unlabeled data with pseudo-labels can be seen as a particular pretext task.

369 **Theory of self-supervised learning.** In order to understand self-supervised learning, there is a line
370 of work towards understanding *contrastive learning* (Saunshi et al., 2019; Tsai et al., 2020; Mitrovic
371 et al., 2020; Tian et al., 2020; Wang and Isola, 2020; Tosh et al., 2021a,b; HaoChen et al., 2021;
372 Wen and Li, 2021; Saunshi et al., 2022), which is one of the most used self-supervised learning
373 approaches based on data augmentation. Unlike contrastive learning, the theoretical understanding
374 of pretext-based self-supervised learning is still rather limited. The only notable works are Lee
375 et al. (2020) and Wei et al. (2020). Lee et al. (2020) proved generalization guarantees for self-
376 supervised algorithms using empirical risk minimization on the pretext task under certain conditional
377 independence assumptions. Wei et al. (2020) proved that under an “expansion” assumption, the
378 minimizer of the population loss based on self-training and input-consistency regularization will
379 achieve high prediction accuracy. Since semi-supervised learning with pseudo-labelers can be seen
380 as a special case of pretext-based self-supervised learning (the pretext task is generated by the
381 pseudo-labelers), we believe the case study in the current paper and its theoretical understanding can
382 shed light on pretext-based self-supervised learning as well.

383 **Feature learning by neural networks.** Our work is also closely related to several recent works that
384 study how neural networks learn the features. Allen-Zhu and Li (2020a) showed that adversarial

385 training purifies the learned features by removing certain “dense mixtures” in the hidden layer weights
 386 of the network. Allen-Zhu and Li (2020b) studied how ensemble and knowledge distillation work in
 387 deep learning when the data have “multi-view” features. Zou et al. (2021) studied an aspect of feature
 388 learning by Adam and GD and showed that GD can learn the sparse features while Adam may fail even
 389 with proper regularization. Notably, there are two concurrent works studying the benign overfitting
 390 phenomenon in learning neural networks: Frei et al. (2022a) established theoretical guarantees for
 391 benign overfitting of two-layer fully connected neural networks with zero training error and test error
 392 close to the Bayes-optimal error, while Cao et al. (2022) studied the benign overfitting phenomenon
 393 in training a two-layer convolutional neural network (CNN), achieving arbitrarily small training and
 394 test loss. Our work studies a different aspect of feature learning afforded by semi-supervised learning
 395 versus supervised learning: given a small amount of labeled data, semi-supervised learning can learn
 396 the features with the help of pseudo-labelers, while supervised learning fails to learn the features and
 397 tends to overfit the noise in the training data.

398 **Comparison with related work.** A recent line of work (Oymak and Gulcu, 2021; Frei et al., 2022b)
 399 studies the semi-supervised learning methods with pseudo-labelers. Our results are different from
 400 theirs in several aspects: (i) we are considering learning with CNNs rather than a linear model, so
 401 the problem is highly non-convex with various local minima, which makes the optimization analysis
 402 more challenging; (ii) the Bayesian optimal predictor is no longer unique for CNNs. Therefore, we
 403 measure the quality of the learned features via downstream task instead of making a comparison with
 404 the Bayesian optimal predictor; (iii) They can only deal with the case where the teacher network
 405 (pseudo-labeler) is the same as the student network (Frei et al., 2022b) or the case where the teacher
 406 network (pseudo-labeler) is at least as complex as the student network (Oymak and Gulcu, 2021).
 407 However, our teacher network (pseudo-labeler) is not specified and can be any structure, such as a
 408 linear network. Therefore we can handle the case where the student network is more complex than
 409 the teacher network, one of the most natural settings for semi-supervised learning with pseudo-labeler
 410 (Xie et al., 2020).

411 B Experiment Setting

412 In particular, we set the problem dimension $d = 10000$, labeled training sample size $n_l = 20$ (10
 413 positive samples and 10 negative samples), pseudo-labeled training sample size $n_u = 20000$ (10000
 414 positive samples and 10000 negative samples), feature vector \mathbf{v} sampled from distribution $\mathcal{N}(\mathbf{0}, \mathbf{I})$
 415 and noise vector sampled from distribution $\mathcal{N}(\mathbf{0}, \sigma_p^2 \mathbf{I})$ where $\sigma_p = 10d^{0.01}$.

416 For semi-supervised learning tasks, we have a linear pseudo-labeler with test error 0.196 ± 0.044 .
 417 Then, we use this classifier to generate pseudo-labels for $n_u = 20000$ unlabeled samples in order
 418 to help semi-supervised learning. After that, for pre-training, we use these pseudo-labeled samples
 419 and n_l labeled samples together to train a CNN with network width $m = 20$, activation function
 420 $\sigma(z) = [z]_+^3$, regularization parameter $\lambda = 0.1$ and learning rate $\eta = 1 \times 10^{-4}$. Besides, we initialize
 421 CNN parameters from $\mathcal{N}(0, \sigma_0^2)$, where $\sigma_0 = 0.1 \times d^{-3/4}$. After 200 iterations, we can obtain a
 422 CNN model with a training error close to the error of pseudo-labeler and zero test error, according
 423 to Table 1. For a downstream task, we use n_l labeled samples to train a linear probe. By applying
 424 learning rate $\eta = 0.1$ and after $T = 100$ iterations, we can obtain a final model with low training and
 425 test loss as well as 100% training accuracy and test accuracy.

426 For supervised learning task, we directly use n_l labeled data to train a CNN with network width
 427 $m = 20$, activation function $\sigma(z) = [z]_+^3$, learning rate $\eta = 1 \times 10^{-4}$. After 200 iterations, we
 428 obtain a CNN with 0 training error and small training loss, about 0.5 test error, and high test loss,
 429 which indicates supervised learning will give a model that behaves badly and even no better than a
 430 random guess.

431 C Proof for Semi-supervised Learning Setting

432 We consider learning K functions $f_{\mathbf{W}_k}(\mathbf{x}), k \in [K]$ based on the pre-training. Since the learning
 433 process of these K functions can be analyzed in exactly the same way, here we only focus on the
 434 learning of one of these functions. For simplicity of notation, we drop the subscript k in the following
 435 proof for Sections C.2, C.3, C.4, C.5, C.6, C.7 and C.8. We start with a condition that is required by
 436 our analysis.

437 **Condition C.1.** The strength of the signal is $\|\mathbf{v}\|_2^2 = \Theta(d)$, the noise variance is $\sigma_p = \Theta(d^\epsilon)$,
438 where $0 < \epsilon < 1/8$ is a small constant, and the width of the network satisfies $m = \text{polylog}(d)$.
439 We also assume that the size of the unlabeled dataset $n_u = \Omega(d^{4\epsilon})$, and labeled data $n_l = \tilde{\Theta}(1)$.
440 For both supervise learning and semi-supervised learning settings, we initialize the weight with
441 $\sigma_0 = \Theta(d^{-3/4})$. For semi-supervised learning, we require $\lambda = o(d^{3/4})$ and assume that there exists
442 a constant C such that for all pseudo-labelers, their test accuracy $p_k > 1/2 + C$.

443 Since we generate the noise patch from the Gaussian distribution, the strength of the noise patch is
444 $\|\boldsymbol{\xi}\|_2^2 \approx d^{1+\epsilon}$ by standard concentration inequalities, which is larger than the strength of the signal
445 patch $\|\mathbf{v}\|_2^2 = \Theta(d)$. Therefore, Condition 3.1 defines a setting with large noises. The condition of
446 $d \gg n_u \gg n_l$ further ensures that learning is in a sufficiently over-parameterized setting. Here we
447 only require the neural network width m to be polylogarithmic in the dimension d and require the
448 pseudo-labelers to perform better than a random guess.

449 **Notation.** We use lower case letters, lower case bold face letters, and upper case bold face letters to
450 denote scalars, vectors, and matrices respectively. For a scalar x , we use $[x]_+$ to denote $\max\{x, 0\}$.
451 For a vector $\mathbf{v} = (v_1, \dots, v_d)^\top$, we denote by $\|\mathbf{v}\|_2 := (\sum_{i=1}^d v_i^2)^{\frac{1}{2}}$ its ℓ_2 norm, and use $\text{supp}(\mathbf{v}) :=$
452 $\{j : v_j \neq 0\}$ to denote its support. For two sequences $\{a_k\}$ and $\{b_k\}$, we denote $a_k = O(b_k)$ if
453 $|a_k| \leq C|b_k|$ for some absolute constant C , denote $a_k = \Omega(b_k)$ if $b_k = O(a_k)$, and denote
454 $a_k = \Theta(b_k)$ if $|a_k| \leq C|b_k|$ and $a_k = \Omega(b_k)$. We also denote $a_k = o(b_k)$ if $\lim |a_k/b_k| = 0$. Finally,
455 we use $\tilde{\Theta}(\cdot)$, $\tilde{O}(\cdot)$ and $\tilde{\Omega}(\cdot)$ to omit logarithmic terms in the notations.

456 C.1 Proof Sketch

457 In this section, we present the proof sketch for the semi-supervised learning setting.

458 **Semi-supervised Pre-training.** We consider learning K functions $f_{\mathbf{W}_k}(\mathbf{x})$, $k \in [K]$ based on the
459 pre-training. Since the learning process of these K functions can be analyzed in exactly the same
460 way, here we only focus on the learning of one of these functions. For simplicity of notation, we drop
461 the subscript k in the following proof sketch.

462 Our study of the pre-training focuses on two aspects of the training process: *feature learning* and
463 *noise memorization*. Specifically, we aim to monitor how the filters in the CNN model learn the
464 feature vector \mathbf{v} and the noise vectors $\boldsymbol{\xi}_i$'s. Therefore, we introduce the following notations.

$$\begin{aligned}
\widehat{\Lambda}_1^{(t)} &:= \max_{1 \leq j \leq m} \langle \mathbf{w}_j^{(t)}, \mathbf{v} \rangle, \bar{\Lambda}_1^{(t)} := \max_{1 \leq j \leq m} -\langle \mathbf{w}_j^{(t)}, \mathbf{v} \rangle, \\
\widehat{\Lambda}_{-1}^{(t)} &:= \max_{m+1 \leq j \leq 2m} -\langle \mathbf{w}_j^{(t)}, \mathbf{v} \rangle, \bar{\Lambda}_{-1}^{(t)} := \max_{m+1 \leq j \leq 2m} \langle \mathbf{w}_j^{(t)}, \mathbf{v} \rangle, \\
\Gamma_i^{(t)} &:= \max_{1 \leq j \leq 2m} \langle \mathbf{w}_j^{(t)}, \boldsymbol{\xi}_i \rangle, \Gamma_i^{\prime(t)} := \max_{1 \leq j \leq 2m} \langle \mathbf{w}_j^{(t)}, \boldsymbol{\xi}_i' \rangle, \Gamma^{(t)} = \max \left\{ \max_{i \in [n_u]} \Gamma_i^{(t)}, \max_{i \in [n_l]} \Gamma_i^{\prime(t)} \right\}.
\end{aligned} \tag{C.1}$$

465 Based on the above definitions for $r \in \{\pm 1\}$, a larger $\widehat{\Lambda}_r^{(t)}$ implies better feature learning along the
466 positive feature direction \mathbf{v} , while a larger $\bar{\Lambda}_r^{(t)}$ implies better feature learning along the negative
467 feature direction $-\mathbf{v}$. Moreover, a larger $\Gamma^{(t)}$ implies a higher level of noise memorization.

468 Based on the update rule of gradient descent, for the inner products $\langle \mathbf{w}_j^{(t)}, \mathbf{v} \rangle$ and $\langle \mathbf{w}_j^{(t)}, \boldsymbol{\xi}_l \rangle$, for
469 $j \in [2m]$, $l \in [n_u]$, we can obtain iterative equations in (C.2).

470 With the help of the iterative equations and definitions in (C.1), we can further show the following
471 lemma.

472 **Lemma C.2.** Assume we use both unlabeled data with pseudo-labels generated by the pseudo-labeler
473 and labeled data for the training of our CNN model. Then for $r \in \{\pm 1\}$, let T_r be the first iteration
474 that $r\widehat{\Lambda}_r^{(t)}$ reaches $\Theta(1/m)$, then for $t \in [0, T_r]$, we have

$$\begin{aligned}
\widehat{\Lambda}_r^{(t+1)} &\geq (1 - \eta\lambda) \cdot \widehat{\Lambda}_r^{(t)} + \eta \cdot C \cdot \Theta(d) \cdot (\widehat{\Lambda}_r^{(t)})^{q-1}, r \in \{\pm 1\}, \\
\bar{\Lambda}_r^{(t+1)} &\leq (1 - \eta\lambda) \cdot \bar{\Lambda}_r^{(t)}, r \in \{\pm 1\}, \\
\Gamma^{(t+1)} &\leq (1 - \eta\lambda) \cdot \Gamma^{(t)} + \eta \cdot \tilde{\Theta}(d^{1-2\epsilon}) \cdot (\Gamma^{(t)})^{q-1},
\end{aligned}$$

475 where C is defined in Condition 3.1.

476 **Lemma C.3.** Assume we use only labeled data for the training of our CNN model. Then for $i \in [n_1]$,
 477 let T'_i be the first iteration that $\Gamma_i^{(t)}$ reaches $\Theta(1/m)$, then we have

$$\begin{aligned}\widehat{\Lambda}_r^{(t+1)} &\leq (1 - \eta\lambda) \cdot \widehat{\Lambda}_r^{(t)} + \eta \cdot \Theta(d) \cdot ((\widehat{\Lambda}_r^{(t)})^{q-1} + (\bar{\Lambda}_r^{(t)})^{q-1}), r \in \{\pm 1\}, \\ \bar{\Lambda}_r^{(t+1)} &\leq (1 - \eta\lambda) \cdot \bar{\Lambda}_r^{(t)}, r \in \{\pm 1\}, \\ \Gamma_i^{(t+1)} &\geq (1 - \eta\lambda) \cdot \Gamma_i^{(t)} + \eta \cdot \tilde{\Theta}(d^{1+2\epsilon}) \cdot (\Gamma_i^{(t)})^{q-1}, i \in [n_1], \text{ for } t \in [0, T'_i].\end{aligned}$$

478 Based on the results in Lemma C.2, we can observe that if both pseudo-labeled and labeled data are
 479 used for training, the CNN will learn the positive direction of the feature vector \mathbf{v} , while barely tending
 480 to fit the negative direction of the feature vector or memorize the noise. And if only labeled data are
 481 used, the CNN will fit noise faster than a feature, which can be seen from Lemma C.3. Leveraging
 482 Lemmas C.2 and C.3, we can obtain the following Lemmas C.4 and C.5, which characterize the
 483 magnitude of feature learning and noise memorization.

484 **Lemma C.4.** If both pseudo-labeled and labeled data are used to train CNN, for $r \in \{\pm 1\}$, let T_r be
 485 the first iteration that $\widehat{\Lambda}_r^{(t)}$ reaches $\Theta(1/m)$ respectively. Let $T_0 = \max_{r \in \{\pm 1\}} \{T_r\}$. Then, it holds
 486 that $\widehat{\Lambda}_r^{(T_0)} = \tilde{\Theta}(1)$, $\bar{\Lambda}_r^{(T_0)} = \tilde{O}(d^{-\frac{1}{4}})$ and $\Gamma^{(T_0)} = \tilde{O}(d^{-\frac{1}{4} + \epsilon})$ for all $t \in [0, T_0]$.

487 **Lemma C.5.** If only labeled data are used to train CNN, for $i \in [n_1]$, let T'_i be the first iteration that
 488 $\Gamma_i^{(t)}$ reaches $\Theta(1/m)$. Let $T'_0 = \max_{i \in [n_1]} T'_i$. Then, it holds that $\widehat{\Lambda}_r = \tilde{O}(d^{-\frac{1}{4}})$, $\bar{\Lambda}_r = \tilde{O}(d^{-\frac{1}{4}})$ for
 489 $r \in \{\pm 1\}$ and $\Gamma_i^{(T'_0)} = \tilde{\Theta}(1)$ for $i \in [n_1]$.

490 The above results indicate the deviation between the two settings. The reason is that assume we
 491 consider a sequence $\{x_t\}$ with iterative equation $x_{t+1} = x_t + \eta \cdot C_t x_t^{q-1}$. If we only use labeled
 492 data, as shown in Lemma C.3, $\Gamma_i^{(t)}$ has $C_t = \tilde{\Theta}(d^{1+2\epsilon})$ while $\widehat{\Lambda}_r^{(t)}$ has $C_t = \Theta(d)$, therefore $\Gamma_i^{(t)}$
 493 increases faster than $\widehat{\Lambda}_r^{(t)}$. In contrast, if we use both labeled data and pseudo-labeled data, C_t will be
 494 $\tilde{\Theta}(d^{1-2\epsilon})$ for $\Gamma_i^{(t)}$ and $\Theta(d)$ for $\widehat{\Lambda}_r^{(t)}$, leading to a slower increasing speed of $\Gamma_i^{(t)}$.

495 **Downstream task.** After the pre-training, we have obtained K CNN classifiers $\{f_{\mathbf{W}_k^{(T_0^k)}}\}_{k=1}^K$. Now
 496 we train the second-layer parameters \mathbf{a} with the training data whose true labels are available. The
 497 following lemma shows that the l_1 -norm of \mathbf{a} will increase with a logarithmic order.

498 **Lemma C.6.** For any learning rate $\eta = \Theta(1)$, we have $\|\mathbf{a}^{(t)}\|_1 = \log(t)/\tilde{\Theta}(1)$. For any labeled data
 499 $(\mathbf{x}'_i, y'_i) \in S'$, we have with high probability that $y'_i \cdot f_{\mathbf{W}^{(t)}}(\mathbf{x}'_i) = \|\mathbf{a}^{(t)}\|_1 \cdot \tilde{\Theta}(1)$. For any newly
 500 generated data $(\mathbf{x}, y) \sim \mathcal{D}$, we also have with high probability that $y \cdot f_{\mathbf{W}^{(t)}}(\mathbf{x}) = \|\mathbf{a}^{(t)}\|_1 \cdot \tilde{\Theta}(1)$.

501 With the help of the above lemma and note that training error and test error are related to $y \cdot f_{\mathbf{W}^{(T_0)}}(\mathbf{x})$
 502 and test loss is related to $\|\mathbf{a}^{(T_0)}\|_1$, we can prove that after $T = \Theta(d^{0.1}/\eta)$ iterations with learning
 503 rate $\eta = \Theta(1)$, the model can achieve nearly zero training error, test error, training loss and test loss.

504 C.2 Gradient Calculation

505 **Lemma C.7** (Gradient Calculation). The gradient of loss function $L_S(\mathbf{W})$ with respect to weight
 506 parameters \mathbf{w}_j is

$$\begin{aligned}\nabla_{\mathbf{w}_j} L_{SUS'}(\mathbf{W}) &= -\frac{q}{n_1 + n_u} \left(\sum_{i=1}^{n_u} c_i \widehat{y}_i ([\langle \mathbf{w}_j, y_i \cdot \mathbf{v} \rangle]_+^{q-1} \cdot y_i \cdot \mathbf{v} + [\langle \mathbf{w}_j, \boldsymbol{\xi}_i \rangle]_+^{q-1} \cdot \boldsymbol{\xi}_i) \right. \\ &\quad \left. + \sum_{i=1}^{n_1} b_i y'_i ([\langle \mathbf{w}_j, y'_i \cdot \mathbf{v} \rangle]_+^{q-1} \cdot y'_i \cdot \mathbf{v} + [\langle \mathbf{w}_j, \boldsymbol{\xi}'_i \rangle]_+^{q-1} \cdot \boldsymbol{\xi}'_i) \right) + \lambda \cdot \mathbf{w}_j,\end{aligned}$$

507 for $1 \leq j \leq m$; and

$$\nabla_{\mathbf{w}_j} L_{SUS'}(\mathbf{W}) = \frac{q}{n_1 + n_u} \left(\sum_{i=1}^{n_u} c_i \widehat{y}_i ([\langle \mathbf{w}_j, y_i \cdot \mathbf{v} \rangle]_+^{q-1} \cdot y_i \cdot \mathbf{v} + [\langle \mathbf{w}_j, \boldsymbol{\xi}_i \rangle]_+^{q-1} \cdot \boldsymbol{\xi}_i) \right)$$

$$+ \sum_{i=1}^{n_1} b_i y'_i \left([\langle \mathbf{w}_j, y'_i \cdot \mathbf{v} \rangle]_+^{q-1} \cdot y'_i \cdot \mathbf{v} + [\langle \mathbf{w}_j, \boldsymbol{\xi}'_i \rangle]_+^{q-1} \cdot \boldsymbol{\xi}'_i \right) + \lambda \cdot \mathbf{w}_j,$$

508 for $m+1 \leq j \leq 2m$, where $-\ell'(\widehat{y}_i \cdot f_{\mathbf{W}}(\mathbf{x}_i)) = \exp[-\widehat{y}_i \cdot f_{\mathbf{W}}(\mathbf{x}_i)] / (1 + \exp[-\widehat{y}_i \cdot f_{\mathbf{W}}(\mathbf{x}_i)])$ is
 509 denoted by c_i and $-\ell'(y'_i \cdot f_{\mathbf{W}}(\mathbf{x}'_i)) = \exp[-y'_i \cdot f_{\mathbf{W}}(\mathbf{x}'_i)] / (1 + \exp[-y'_i \cdot f_{\mathbf{W}}(\mathbf{x}'_i)])$ is denoted by b_i .

510 *Proof of Lemma C.7.* When $1 \leq j \leq m$,

$$\begin{aligned} \nabla_{\mathbf{w}_j} \ell(\widehat{y}_i \cdot f_{\mathbf{W}}(\mathbf{x}_i)) &= \ell'(\widehat{y}_i \cdot f_{\mathbf{W}}(\mathbf{x}_i)) \cdot \widehat{y}_i \cdot \nabla_{\mathbf{w}_j} f_{\mathbf{W}}(\mathbf{x}_i) \\ &= -c_i \cdot \widehat{y}_i \cdot \nabla_{\mathbf{w}_j} f_{\mathbf{W}}(\mathbf{x}_i) \\ &= -c_i \widehat{y}_i \cdot (\sigma'(\langle \mathbf{w}_j, y_i \cdot \mathbf{v} \rangle) \cdot y_i \cdot \mathbf{v} + \sigma'(\langle \mathbf{w}_j, \boldsymbol{\xi}_i \rangle) \cdot \boldsymbol{\xi}_i) \\ &= -q c_i \widehat{y}_i ([\langle \mathbf{w}_j, y_i \cdot \mathbf{v} \rangle]_+^{q-1} \cdot y_i \cdot \mathbf{v} + [\langle \mathbf{w}_j, \boldsymbol{\xi}_i \rangle]_+^{q-1} \cdot \boldsymbol{\xi}_i) \end{aligned}$$

511

$$\begin{aligned} \nabla_{\mathbf{w}_j} \ell(y'_i \cdot f_{\mathbf{W}}(\mathbf{x}'_i)) &= \ell'(y'_i \cdot f_{\mathbf{W}}(\mathbf{x}'_i)) \cdot y'_i \cdot \nabla_{\mathbf{w}_j} f_{\mathbf{W}}(\mathbf{x}'_i) \\ &= -b_i \cdot y'_i \cdot \nabla_{\mathbf{w}_j} f_{\mathbf{W}}(\mathbf{x}'_i) \\ &= -b_i y'_i \cdot (\sigma'(\langle \mathbf{w}_j, y'_i \cdot \mathbf{v} \rangle) \cdot y'_i \cdot \mathbf{v} + \sigma'(\langle \mathbf{w}_j, \boldsymbol{\xi}'_i \rangle) \cdot \boldsymbol{\xi}'_i) \\ &= -q b_i y'_i ([\langle \mathbf{w}_j, y'_i \cdot \mathbf{v} \rangle]_+^{q-1} \cdot y'_i \cdot \mathbf{v} + [\langle \mathbf{w}_j, \boldsymbol{\xi}'_i \rangle]_+^{q-1} \cdot \boldsymbol{\xi}'_i) \end{aligned}$$

512 and when $m+1 \leq j \leq 2m$,

$$\begin{aligned} \nabla_{\mathbf{w}_j} \ell(\widehat{y}_i \cdot f_{\mathbf{W}}(\mathbf{x}_i)) &= q c_i \widehat{y}_i ([\langle \mathbf{w}_j, y_i \cdot \mathbf{v} \rangle]_+^{q-1} \cdot y_i \cdot \mathbf{v} + [\langle \mathbf{w}_j, \boldsymbol{\xi}_i \rangle]_+^{q-1} \cdot \boldsymbol{\xi}_i) \\ \nabla_{\mathbf{w}_j} \ell(y'_i \cdot f_{\mathbf{W}}(\mathbf{x}'_i)) &= q b_i y'_i ([\langle \mathbf{w}_j, y'_i \cdot \mathbf{v} \rangle]_+^{q-1} \cdot y'_i \cdot \mathbf{v} + [\langle \mathbf{w}_j, \boldsymbol{\xi}'_i \rangle]_+^{q-1} \cdot \boldsymbol{\xi}'_i) \end{aligned}$$

513 Note that $\nabla_{\mathbf{w}_j} L_{S \cup S'}(\mathbf{W}) = (\sum_{i=1}^{n_u} \nabla_{\mathbf{w}_j} \ell(\widehat{y}_i \cdot f_{\mathbf{W}}(\mathbf{x}_i)) + \sum_{i=1}^{n_1} \nabla_{\mathbf{w}_j} \ell(y'_i \cdot f_{\mathbf{W}}(\mathbf{x}'_i))) / (n_1 + n_u) +$
 514 $\lambda \cdot \mathbf{w}_j$, we have proved the lemma. \square

515 C.3 Inner Product Update Rule Calculation

516 When the model is trained by gradient descent, the update rule can be formulated by

$$\mathbf{w}_j^{(t+1)} = \mathbf{w}_j^{(t)} - \eta \cdot \nabla_{\mathbf{w}_j} L_S(\mathbf{W}^{(t)}), \quad j \in [2m]. \quad (\text{C.2})$$

517 We study the performance of entire training process from two perspective: feature learning and noise
 518 memorization. Mathematically, we will focus on two quantities: $\langle \mathbf{w}_j^{(t)}, \mathbf{v} \rangle$ and $\langle \mathbf{w}_j^{(t)}, \boldsymbol{\xi}_i \rangle$. And then
 519 we have following lemma for the inner product update rule.

520 **Lemma C.8** (Inner Product Update Rule). The feature learning and noise memorization performance
 521 of gradient descent can be formulated by

$$\begin{aligned} \langle \mathbf{w}_j^{(t+1)}, \mathbf{v} \rangle &= (1 - \eta\lambda) \cdot \langle \mathbf{w}_j^{(t)}, \mathbf{v} \rangle + \frac{q\eta u_j}{n_1 + n_u} \left(\sum_{i=1}^{n_u} y_i \widehat{y}_i c_i^{(t)} [\langle \mathbf{w}_j^{(t)}, y_i \cdot \mathbf{v} \rangle]_+^{q-1} \|\mathbf{v}\|_2^2 \right. \\ &\quad \left. + \sum_{i=1}^{n_1} b_i^{(t)} [\langle \mathbf{w}_j^{(t)}, y'_i \cdot \mathbf{v} \rangle]_+^{q-1} \|\mathbf{v}\|_2^2 \right), \end{aligned}$$

522

$$\begin{aligned} \langle \mathbf{w}_j^{(t+1)}, \boldsymbol{\xi}_i \rangle &= (1 - \eta\lambda) \cdot \langle \mathbf{w}_j^{(t)}, \boldsymbol{\xi}_i \rangle + \frac{q\eta u_j}{n_1 + n_u} \left(\sum_{i=1}^{n_u} \widehat{y}_i c_i^{(t)} [\langle \mathbf{w}_j^{(t)}, \boldsymbol{\xi}_i \rangle]_+^{q-1} \langle \boldsymbol{\xi}_i, \boldsymbol{\xi}_i \rangle \right. \\ &\quad \left. + \sum_{i=1}^{n_1} y'_i b_i^{(t)} [\langle \mathbf{w}_j^{(t)}, \boldsymbol{\xi}'_i \rangle]_+^{q-1} \langle \boldsymbol{\xi}'_i, \boldsymbol{\xi}_i \rangle \right), \end{aligned}$$

523

$$\langle \mathbf{w}_j^{(t+1)}, \boldsymbol{\xi}'_i \rangle = (1 - \eta\lambda) \cdot \langle \mathbf{w}_j^{(t)}, \boldsymbol{\xi}'_i \rangle + \frac{q\eta u_j}{n_1 + n_u} \left(\sum_{i=1}^{n_u} \widehat{y}_i c_i^{(t)} [\langle \mathbf{w}_j^{(t)}, \boldsymbol{\xi}_i \rangle]_+^{q-1} \langle \boldsymbol{\xi}_i, \boldsymbol{\xi}'_i \rangle \right)$$

$$+ \sum_{i=1}^{n_1} y'_i b_i^{(t)} [\langle \mathbf{w}_j^{(t)}, \boldsymbol{\xi}'_i \rangle_+^{q-1} \langle \boldsymbol{\xi}'_i, \boldsymbol{\xi}'_l \rangle],$$

524 where $j \in [2m]$, $l \in [n_u]$ and $u_j := \mathbb{1}_{[1 \leq j \leq m]} - \mathbb{1}_{[m+1 \leq j \leq 2m]}$.

525 *Proof of Lemma C.8.* According to Lemma C.7 and gradient descent update rule (C.2), we have

$$\begin{aligned} \mathbf{w}_j^{(t+1)} &= (1 - \eta\lambda) \cdot \mathbf{w}_j^{(t)} + \frac{q\eta u_j}{n_1 + n_u} \cdot \left(\sum_{i=1}^{n_u} c_i \widehat{y}_i ([\langle \mathbf{w}_j, y_i \cdot \mathbf{v} \rangle_+^{q-1} \cdot y_i \cdot \mathbf{v} + [\langle \mathbf{w}_j, \boldsymbol{\xi}_i \rangle_+^{q-1} \cdot \boldsymbol{\xi}_i] \right. \\ &\quad \left. + \sum_{i=1}^{n_1} b_i y'_i ([\langle \mathbf{w}_j, y'_i \cdot \mathbf{v} \rangle_+^{q-1} \cdot y'_i \cdot \mathbf{v} + [\langle \mathbf{w}_j, \boldsymbol{\xi}'_i \rangle_+^{q-1} \cdot \boldsymbol{\xi}'_i]) \right) \end{aligned}$$

526 Taking inner product with feature vector \mathbf{v} and noise patch $\boldsymbol{\xi}_l$ and note that \mathbf{v} is orthogonal to $\boldsymbol{\xi}_l$
527 according to the data model, we have

$$\begin{aligned} \langle \mathbf{w}_j^{(t+1)}, \mathbf{v} \rangle &= (1 - \eta\lambda) \cdot \langle \mathbf{w}_j^{(t)}, \mathbf{v} \rangle + \frac{q\eta u_j}{n_1 + n_u} \left(\sum_{i=1}^{n_u} c_i^{(t)} \widehat{y}_i ([\langle \mathbf{w}_j, y_i \cdot \mathbf{v} \rangle_+^{q-1} y_i \|\mathbf{v}\|_2^2 + [\langle \mathbf{w}_j, \boldsymbol{\xi}_i \rangle_+^{q-1} \langle \boldsymbol{\xi}_i, \mathbf{v} \rangle]) \right. \\ &\quad \left. + \sum_{i=1}^{n_1} b_i^{(t)} y'_i ([\langle \mathbf{w}_j, y'_i \cdot \mathbf{v} \rangle_+^{q-1} y'_i \|\mathbf{v}\|_2^2 + [\langle \mathbf{w}_j, \boldsymbol{\xi}'_i \rangle_+^{q-1} \langle \boldsymbol{\xi}'_i, \mathbf{v} \rangle]) \right) \\ &= (1 - \eta\lambda) \cdot \langle \mathbf{w}_j^{(t)}, \mathbf{v} \rangle + \frac{q\eta u_j}{n_1 + n_u} \left(\sum_{i=1}^{n_u} y_i \widehat{y}_i c_i^{(t)} [\langle \mathbf{w}_j^{(t)}, y_i \cdot \mathbf{v} \rangle_+^{q-1} \|\mathbf{v}\|_2^2 \right. \\ &\quad \left. + \sum_{i=1}^{n_1} b_i^{(t)} [\langle \mathbf{w}_j^{(t)}, y'_i \cdot \mathbf{v} \rangle_+^{q-1} \|\mathbf{v}\|_2^2] \right), \end{aligned}$$

528

$$\begin{aligned} \langle \mathbf{w}_j^{(t+1)}, \boldsymbol{\xi}_l \rangle &= (1 - \eta\lambda) \cdot \langle \mathbf{w}_j^{(t)}, \boldsymbol{\xi}_l \rangle + \frac{q\eta u_j}{n_1 + n_u} \left(\sum_{i=1}^{n_u} c_i^{(t)} \widehat{y}_i ([\langle \mathbf{w}_j, y_i \cdot \mathbf{v} \rangle_+^{q-1} y_i \langle \mathbf{v}, \boldsymbol{\xi}_l \rangle + [\langle \mathbf{w}_j, \boldsymbol{\xi}_i \rangle_+^{q-1} \langle \boldsymbol{\xi}_i, \boldsymbol{\xi}_l \rangle]) \right. \\ &\quad \left. + \sum_{i=1}^{n_1} b_i^{(t)} y'_i ([\langle \mathbf{w}_j, y'_i \cdot \mathbf{v} \rangle_+^{q-1} y'_i \langle \mathbf{v}, \boldsymbol{\xi}_l \rangle + [\langle \mathbf{w}_j, \boldsymbol{\xi}'_i \rangle_+^{q-1} \langle \boldsymbol{\xi}'_i, \boldsymbol{\xi}_l \rangle]) \right) \\ &= (1 - \eta\lambda) \cdot \langle \mathbf{w}_j^{(t)}, \boldsymbol{\xi}_l \rangle + \frac{q\eta u_j}{n_1 + n_u} \left(\sum_{i=1}^{n_u} \widehat{y}_i c_i^{(t)} [\langle \mathbf{w}_j^{(t)}, \boldsymbol{\xi}_i \rangle_+^{q-1} \langle \boldsymbol{\xi}_i, \boldsymbol{\xi}_l \rangle \right. \\ &\quad \left. + \sum_{i=1}^{n_1} y'_i b_i^{(t)} [\langle \mathbf{w}_j^{(t)}, \boldsymbol{\xi}'_i \rangle_+^{q-1} \langle \boldsymbol{\xi}'_i, \boldsymbol{\xi}_l \rangle] \right), \end{aligned}$$

529 and

$$\begin{aligned} \langle \mathbf{w}_j^{(t+1)}, \boldsymbol{\xi}'_l \rangle &= (1 - \eta\lambda) \cdot \langle \mathbf{w}_j^{(t)}, \boldsymbol{\xi}'_l \rangle + \frac{q\eta u_j}{n_1 + n_u} \left(\sum_{i=1}^{n_u} \widehat{y}_i c_i^{(t)} [\langle \mathbf{w}_j^{(t)}, \boldsymbol{\xi}_i \rangle_+^{q-1} \langle \boldsymbol{\xi}_i, \boldsymbol{\xi}'_l \rangle] \right. \\ &\quad \left. + \sum_{i=1}^{n_1} y'_i b_i^{(t)} [\langle \mathbf{w}_j^{(t)}, \boldsymbol{\xi}'_i \rangle_+^{q-1} \langle \boldsymbol{\xi}'_i, \boldsymbol{\xi}'_l \rangle] \right), \end{aligned}$$

530 which completes the proof. \square

531 **C.4 Estimate** $\widehat{\Lambda}_r^{(0)}, \bar{\Lambda}_r^{(0)}, \Gamma_i^{(0)}, \Gamma_i^{\prime(0)}$

532 Let $\widehat{\Lambda}_1^{(t)} = \max_{1 \leq j \leq m} \langle \mathbf{w}_j^{(t)}, \mathbf{v} \rangle$, $\widehat{\Lambda}_{-1}^{(t)} = \max_{m+1 \leq j \leq 2m} -\langle \mathbf{w}_j^{(t)}, \mathbf{v} \rangle$, $\bar{\Lambda}_1^{(t)} =$
533 $\max_{m+1 \leq j \leq 2m} \langle \mathbf{w}_j^{(t)}, \mathbf{v} \rangle$, $\bar{\Lambda}_{-1}^{(t)} = \max_{1 \leq j \leq m} -\langle \mathbf{w}_j^{(t)}, \mathbf{v} \rangle$, which characterize the *feature*
534 *learning* aspect of training process. An easy way to distinguish between $\widehat{\Lambda}_r^{(t)}$ and $\bar{\Lambda}_r^{(t)}$ is that $\widehat{\Lambda}_r^{(t)}$
535 should be large while $\bar{\Lambda}_r^{(t)}$ should be small.

536 Let $\Gamma_i^{(t)} = \max_{1 \leq j \leq 2m} \langle \mathbf{w}_j, \boldsymbol{\xi}_i \rangle, i \in [n_u], \Gamma_i^{\prime(t)} = \max_{1 \leq j \leq 2m} \langle \mathbf{w}_j, \boldsymbol{\xi}'_i \rangle, i \in [n_l]$, which characterize
 537 the *noise memorization* aspect of training process with respect to a particular sample.

538 Let $\Gamma^{(t)} = \max \{ \max_{i \in [n_u]} \Gamma_i^{(t)}, \max_{i \in [n_l]} \Gamma_i^{\prime(t)} \}$, which characterize the *noise memorization* aspect
 539 of training process regardless of the sample index.

540 We first provide the concentration inequality for $\widehat{\Lambda}_r^{(0)}$ and $\bar{\Lambda}_r^{(0)}$ in the following lemma.

541 **Lemma C.9.** With probability at least $1 - 4\delta$ with respect to the randomness of initialization of \mathbf{w} ,
 542 we have

$$|\widehat{\Lambda}_r^{(0)} - \mathbb{E}[\widehat{\Lambda}_r^{(0)}]| < \sqrt{8 \log\left(\frac{1}{\delta}\right)} \sigma_0 \|\mathbf{v}\|_2,$$

543

$$|\bar{\Lambda}_r^{(0)} - \mathbb{E}[\bar{\Lambda}_r^{(0)}]| < \sqrt{8 \log\left(\frac{1}{\delta}\right)} \sigma_0 \|\mathbf{v}\|_2,$$

544 and

$$\mathbb{E}[\widehat{\Lambda}_r^{(0)}] \asymp \sqrt{\log(m)} \sigma_0 \|\mathbf{v}\|_2, \mathbb{E}[\bar{\Lambda}_r^{(0)}] \asymp \sqrt{\log(m)} \sigma_0 \|\mathbf{v}\|_2, r \in \{\pm 1\}.$$

545 *Proof of Lemma C.9.* Note that $\widehat{\Lambda}_1^{(0)} = \max_{1 \leq j \leq m} \langle \mathbf{w}_j^{(0)}, \mathbf{v} \rangle, \widehat{\Lambda}_{-1}^{(0)} = \max_{m+1 \leq j \leq 2m} -\langle \mathbf{w}_j^{(0)}, \mathbf{v} \rangle,$
 546 $\bar{\Lambda}_1^{(0)} = \max_{m+1 \leq j \leq 2m} \langle \mathbf{w}_j^{(0)}, \mathbf{v} \rangle$ and $\bar{\Lambda}_{-1}^{(0)} = \max_{m+1 \leq j \leq 2m} -\langle \mathbf{w}_j^{(0)}, \mathbf{v} \rangle, \mathbf{w}_j^{(0)} \sim \mathcal{N}(\mathbf{0}, \sigma_0^2 \mathbf{I})$ and

547 \mathbf{v} is a fixed vector. Therefore, $\langle \mathbf{w}_j^{(0)}, \mathbf{v} \rangle \sim \mathcal{N}(0, \sigma_0^2 \|\mathbf{v}\|_2^2), -\langle \mathbf{w}_j^{(0)}, \mathbf{v} \rangle \sim \mathcal{N}(0, \sigma_0^2 \|\mathbf{v}\|_2^2)$ for all
 548 $1 \leq j \leq 2m$ and $\widehat{\Lambda}_r^{(0)}, \bar{\Lambda}_r^{(0)}, r \in \{\pm 1\}$ are identically distributed. Therefore, without loss of
 549 generality, we only need to discuss the concentration of $\widehat{\Lambda}_1^{(0)}$. By applying Lemma E.1, we have

$$\mathbb{P}\left(|\widehat{\Lambda}_1^{(0)} - \mathbb{E}[\widehat{\Lambda}_1^{(0)}]| > t\right) \leq 2e^{-\frac{t^2}{2\sigma_0^2 \|\mathbf{v}\|_2^2}}.$$

550 By applying Lemma E.2, we have

$$\mathbb{E}[\widehat{\Lambda}_1^{(0)}] \asymp \sqrt{\log(m)} \sigma_0 \|\mathbf{v}\|_2,$$

551 which completes the proof. \square

552 Then we provide concentration inequality for $\Gamma_i^{(0)}$ in the following lemma.

553 **Lemma C.10.** Suppose that $d \geq \Omega(\log(m(n_u + n_l)/\delta)), m = \Omega(\log(1/\delta))$. Then with probability
 554 at least $1 - \delta$,

$$\begin{aligned} \frac{\sigma_0 \sigma_p \sqrt{d}}{4} &\leq \Gamma_i^{(0)} \leq 2\sqrt{\log(16m(n_u + n_l)/\delta)} \cdot \sigma_0 \sigma_p \sqrt{d}, \text{ for all } i \in [n_u], \\ \frac{\sigma_0 \sigma_p \sqrt{d}}{4} &\leq \Gamma_i^{\prime(0)} \leq 2\sqrt{\log(16m(n_u + n_l)/\delta)} \cdot \sigma_0 \sigma_p \sqrt{d}, \text{ for all } i \in [n_l]. \end{aligned}$$

555 *Proof of Lemma C.10.* By Lemma E.3, with probability at least $1 - \delta/4$,

$$\begin{aligned} \sigma_p \sqrt{d}/\sqrt{2} &\leq \|\boldsymbol{\xi}_i\|_2 \leq \sqrt{3/2} \cdot \sigma_p \sqrt{d}, \text{ for } i \in [n_u], \\ \sigma_p \sqrt{d}/\sqrt{2} &\leq \|\boldsymbol{\xi}'_i\|_2 \leq \sqrt{3/2} \cdot \sigma_p \sqrt{d}, \text{ for } i \in [n_l]. \end{aligned} \tag{C.3}$$

556 Therefore, by Gaussian tail bound and union bound, with probability at least $1 - \delta/4$,

$$\begin{aligned} \langle \mathbf{w}_j^{(0)}, \boldsymbol{\xi}_i \rangle &\leq |\langle \mathbf{w}_j^{(0)}, \boldsymbol{\xi}_i \rangle| \leq \sqrt{2 \log(8m/\delta)} \cdot \sigma_0 \|\boldsymbol{\xi}_i\|_2, \text{ for } i \in [n_u], \\ \langle \mathbf{w}_j^{(0)}, \boldsymbol{\xi}'_i \rangle &\leq |\langle \mathbf{w}_j^{(0)}, \boldsymbol{\xi}'_i \rangle| \leq \sqrt{2 \log(8m/\delta)} \cdot \sigma_0 \|\boldsymbol{\xi}'_i\|_2, \text{ for } i \in [n_l]. \end{aligned} \tag{C.4}$$

557 Note that $\mathbb{P}(\sigma_0 \sigma_p \sqrt{d}/4 > \langle \mathbf{w}_j^{(0)}, \boldsymbol{\xi}_i \rangle)$ is an absolute constant and therefore by the condition on m ,
 558 we have

$$\begin{aligned} \mathbb{P}\left(\frac{\sigma_0 \sigma_p \sqrt{d}}{4} \leq \Gamma_i^{(t)}\right) &= \mathbb{P}\left(\frac{\sigma_0 \sigma_p \sqrt{d}}{4} \leq \max_{j \in [2m]} \langle \mathbf{w}_j^{(0)}, \boldsymbol{\xi}_i \rangle\right) \\ &= 1 - \mathbb{P}\left(\frac{\sigma_0 \sigma_p \sqrt{d}}{4} > \max_{j \in [2m]} \langle \mathbf{w}_j^{(0)}, \boldsymbol{\xi}_i \rangle\right) \end{aligned}$$

$$\begin{aligned}
&= 1 - \left(\mathbb{P} \left(\frac{\sigma_0 \sigma_p \sqrt{d}}{4} > \langle \mathbf{w}_j^{(0)}, \boldsymbol{\xi}_i \rangle \right) \right)^{2m} \\
&\geq 1 - \frac{\delta}{4},
\end{aligned}$$

559 and

$$\mathbb{P} \left(\frac{\sigma_0 \sigma_p \sqrt{d}}{4} \leq \Gamma_i^{(t)} \right) \geq 1 - \frac{\delta}{4}.$$

560 On the other hand, according to (C.3) and (C.4), we have

$$\begin{aligned}
&\mathbb{P}(\Gamma_i^{(t)} \leq 2\sqrt{\log(16m(n_u + n_l)/\delta)} \cdot \sigma_0 \sigma_p \sqrt{d}) \\
&= \mathbb{P} \left(\max_{j \in [2m]} \langle \mathbf{w}_j^{(0)}, \boldsymbol{\xi}_i \rangle \leq 2\sqrt{\log(16m(n_u + n_l)/\delta)} \cdot \sigma_0 \sigma_p \sqrt{d} \right) \\
&\geq 1 - \frac{\delta}{4},
\end{aligned}$$

561 and

$$\mathbb{P}(\Gamma_i'^{(t)} \leq 2\sqrt{\log(16m(n_u + n_l)/\delta)} \cdot \sigma_0 \sigma_p \sqrt{d}) \geq 1 - \frac{\delta}{4},$$

562 which completes the proof. \square

563 C.5 Stage I of GD: On-diagonal feature learning

564 In this stage, $\widehat{\Lambda}_1^{(t)}$ and $\widehat{\Lambda}_{-1}^{(t)}$ respectively increase to magnitude $\Theta(1/m)$ and $\bar{\Lambda}_1^{(t)}$, $\bar{\Lambda}_{-1}^{(t)}$ and $\Gamma_j^{(t)}$
565 remain small, the same magnitude as initialization. In order to characterize the behaviour of feature
566 learning and noise memorization during Stage I, we decompose the analysis into following three
567 parts:

568 1. First, in Lemma C.15, we provide a lower bound of the update rules of on-diagonal feature learning
569 term of $\widehat{\Lambda}_1^{(t)}$, $\widehat{\Lambda}_{-1}^{(t)}$ to lower-bound their increasing speed, and an upper bound of off-diagonal
570 feature learning term $\bar{\Lambda}_1^{(t)}$, $\bar{\Lambda}_{-1}^{(t)}$ to indicate their decrease.

571 2. Second, in Lemma C.17, we provide an upper bound of the update rules of noise memorization
572 term $\Gamma^{(t)}$ to upper-bound its increasing speed.

573 3. Third, we provide a useful lemma, which is a derivation of Claim C.20 in Allen-Zhu and Li
574 (2020b), which is called tensor power method. By applying tensor power method, we will prove
575 that:

576 • When $\widehat{\Lambda}_1^{(t)}$ reaches $\Theta(1/m)$ at T_1 , $\bar{\Lambda}_1^{(t)}$ and $\Gamma^{(t)}$ remain a magnitude no more than initialization.

577 • When $\widehat{\Lambda}_{-1}^{(t)}$ reaches $\Theta(1/m)$ at T_{-1} , $\bar{\Lambda}_{-1}^{(t)}$ and $\Gamma^{(t)}$ remain a magnitude no more than initializa-
578 tion.

579 C.5.1 Upper bound and lower bound for $\widehat{\Lambda}_1^{(t)}$, $\widehat{\Lambda}_{-1}^{(t)}$ and $\bar{\Lambda}_1^{(t)}$, $\bar{\Lambda}_{-1}^{(t)}$

580 We first consider Stage I of GD when $\max_{r \in \{\pm 1\}} \{\widehat{\Lambda}_r^{(t)}, \bar{\Lambda}_r^{(t)}\} \leq \Theta(m^{-1})$.

581 In this stage, we first prove following lemma:

582 **Lemma C.11.** As long as $\max_{r \in \{\pm 1\}} \{\widehat{\Lambda}_r^{(t)}, \bar{\Lambda}_r^{(t)}\} \leq \Theta(m^{-1})$, we have $c_i^{(t)} := -\ell'(\widehat{y}_i \cdot f_{\mathbf{W}^{(t)}}(\mathbf{x}_i))$
583 and $b_i^{(t)} := -\ell'(y_i' \cdot f_{\mathbf{W}^{(t)}}(\mathbf{x}_i'))$ remains $1/2 \pm o(1)$.

584 *Proof of Lemma C.11.* Note that $\ell(z) = \log(1 + \exp(-z))$ and $-\ell'(z) = \exp(-z)/(1 + \exp(-z))$,
585 and without loss of generality assuming $\widehat{y}_i = y_i = 1$, we can express $c_i^{(t)}$ as follow:

$$c_i^{(t)} = -\ell'(f_{\mathbf{W}^{(t)}}(\mathbf{x}_i)) = \frac{e^{\sum_{j=m+1}^{2m} [\sigma(\langle \mathbf{w}_j^{(t)}, \mathbf{v} \rangle) + \sigma(\langle \mathbf{w}_j^{(t)}, \boldsymbol{\xi}_i \rangle)]}}{e^{\sum_{j=1}^m [\sigma(\langle \mathbf{w}_j^{(t)}, \mathbf{v} \rangle) + \sigma(\langle \mathbf{w}_j^{(t)}, \boldsymbol{\xi}_i \rangle)]} + e^{\sum_{j=m+1}^{2m} [\sigma(\langle \mathbf{w}_j^{(t)}, \mathbf{v} \rangle) + \sigma(\langle \mathbf{w}_j^{(t)}, \boldsymbol{\xi}_i \rangle)]}},$$

586 Since $\sigma(\langle \mathbf{w}_j^{(t)}, \mathbf{v} \rangle)$ dominates $\sigma(\langle \mathbf{w}_j^{(t)}, \boldsymbol{\xi} \rangle)$ for $j \in [m]$, which will be proved later by using *tensor*
 587 *power method*, we have

$$c_i^{(t)} = \frac{e^{\sum_{j=m+1}^{2m} [\sigma(\langle \mathbf{w}_j^{(t)}, \mathbf{v} \rangle) + \sigma(\langle \mathbf{w}_j^{(t)}, \boldsymbol{\xi}_i \rangle)]}}{e^{\sum_{j=1}^m \sigma(\langle \mathbf{w}_j^{(t)}, \mathbf{v} \rangle) + \{\text{lower order term}\}} + e^{\sum_{j=m+1}^{2m} [\sigma(\langle \mathbf{w}_j^{(t)}, \mathbf{v} \rangle) + \sigma(\langle \mathbf{w}_j^{(t)}, \boldsymbol{\xi}_i \rangle)]}}.$$

588 On the one side,

$$c_i^{(t)} \geq \frac{1}{e^{\sum_{j=1}^m \sigma(\langle \mathbf{w}_j^{(t)}, \mathbf{v} \rangle) + \{\text{lower order term}\}} + 1} \geq \frac{1}{e^{m(\bar{\Lambda}_1^{(t)})^{q-1} + 1}} \geq \frac{1}{e^{\Theta(m^{-(q-1)}) + 1}} = \frac{1}{2 + o(1)} = \frac{1}{2} - o(1).$$

589 On the other side, according to Lemma C.4, we have $\bar{\Lambda}_1^{(t)} = \tilde{O}(d^{-\frac{1}{4}})$ and $\Gamma^{(t)} = \tilde{O}(d^{-\frac{1}{4} + \epsilon})$, it
 590 follows that

$$\begin{aligned} c_i^{(t)} &\leq \frac{e^{m(\bar{\Lambda}_1^{(t)})^{q-1} + m(\Gamma^{(t)})^{q-1}}}{e^{\sum_{j=1}^m \sigma(\langle \mathbf{w}_j^{(t)}, \mathbf{v} \rangle) + \{\text{lower order term}\}} + e^{m(\bar{\Lambda}_1^{(t)})^{q-1} + m(\Gamma^{(t)})^{q-1}}} \\ &= \frac{1 + o(1)}{e^{\sum_{j=1}^m \sigma(\langle \mathbf{w}_j^{(t)}, \mathbf{v} \rangle) + \{\text{lower order term}\}} + 1 + o(1)} \\ &\leq \frac{1 + o(1)}{1 + 1 + o(1)} = \frac{1}{2} + o(1). \end{aligned}$$

591 Therefore, we have $c_i^{(t)} = 1/2 \pm o(1)$ if $\hat{y}_i = y_i = 1$ and other cases ($\hat{y}_i = y_i = 1, \hat{y}_i = -y_i, b_i^{(t)}$)
 592 can be proved in a similar way. \square

593 By applying above lemma, we can obtain following lemma:

594 **Lemma C.12.** For any $\delta < 1/2$, with probability at least $1 - 2\delta$ over pseudo-labels generated by the
 595 pseudo-labeler, we have

$$\left| \frac{1}{n_u} \sum_{i=1}^{n_u} \hat{y}_i y_i c_i^{(t)} - \left(p - \frac{1}{2}\right) \right| < \sqrt{\frac{1}{8n_u} \log \frac{1}{\delta}} + o(1),$$

596 where $o(1)$ is with respect to d .

597 If we denote $\{(\mathbf{x}_i, y_i) | y_i = 1, i \in [n_u]\}$ as S_1 , $\{(\mathbf{x}_i, y_i) | y_i = -1, i \in [n_u]\}$ as S_{-1} , $|S_1|$ as n_1 and
 598 $|S_{-1}|$ as n_{-1} , we have with probability at least $1 - 4\delta$ that

$$\left| \frac{1}{n_1} \sum_{i=1}^{n_1} \hat{y}_i y_i c_i^{(t)} - \left(p - \frac{1}{2}\right) \right| < \sqrt{\frac{1}{8n_1} \log \frac{1}{\delta}} + o(1),$$

599 and

$$\left| \frac{1}{n_{-1}} \sum_{i=1}^{n_{-1}} \hat{y}_i y_i c_i^{(t)} - \left(p - \frac{1}{2}\right) \right| < \sqrt{\frac{1}{8n_{-1}} \log \frac{1}{\delta}} + o(1).$$

600 *Proof of Lemma C.12.* First, according to Lemma C.11, we have

$$\frac{1}{n_u} \sum_{i=1}^{n_u} \hat{y}_i y_i c_i^{(t)} = \frac{1}{n_u} \sum_{i=1}^{n_u} \hat{y}_i y_i \left(c_i^{(t)} - \frac{1}{2}\right) + \frac{1}{2n_u} \sum_{i=1}^{n_u} \hat{y}_i y_i = \frac{1}{2n_u} \sum_{i=1}^{n_u} \hat{y}_i y_i \pm o(1) \quad (\text{C.5})$$

601 Then, according to Hoeffding's inequality when $a_i = -1, b_i = 1$, we have

$$\mathbb{P}\left(\left|\frac{1}{n_u} \sum_{i=1}^{n_u} \hat{y}_i y_i - \mathbb{E}\left[\frac{1}{n_u} \sum_{i=1}^{n_u} \hat{y}_i y_i\right]\right| \geq t\right) \leq 2 \exp\left(-\frac{2n_u^2 t^2}{\sum_{i=1}^{n_u} (a_i - b_i)^2}\right) = 2 \exp(-2n_u t^2).$$

602 Note that the pseudo-label \hat{y}_i generated by the pseudo-labeler takes y_i with probability p and $-y_i$
 603 with probability $1 - p$, we have $\mathbb{E}\left[\frac{1}{n_u} \sum_{i=1}^{n_u} \hat{y}_i y_i\right] = \frac{1}{n_u} \sum_{i=1}^{n_u} \mathbb{E}[\hat{y}_i y_i] = 2p - 1$. It follows that

$$\mathbb{P}\left(\left|\frac{1}{2n_u} \sum_{i=1}^{n_u} \hat{y}_i y_i - \left(p - \frac{1}{2}\right)\right| \geq t\right) \leq 2 \exp(-8n_u t^2),$$

604 and therefore

$$\left|\frac{1}{2n_u} \sum_{i=1}^{n_u} \hat{y}_i y_i - \left(p - \frac{1}{2}\right)\right| < \sqrt{\frac{1}{8n_u} \log \frac{1}{\delta}} \quad (\text{C.6})$$

605 holds with probability at least $1 - 2\delta$. According to (C.5) and (C.6), we have

$$\left|\frac{1}{2n_u} \sum_{i=1}^{n_u} \hat{y}_i y_i - \left(p - \frac{1}{2}\right)\right| < \sqrt{\frac{1}{8n_u} \log \frac{1}{\delta}} + o(1),$$

606 which verifies the first statement of the lemma. And the other part of the lemma can be proved in a
 607 similar way. \square

608 According to above lemma and note that $n_u, n_1, n_{-1} = \omega(1)$, we have further that

$$\left|\frac{1}{n_u} \sum_{i=1}^{n_u} \hat{y}_i y_i c_i^{(t)} - \left(p - \frac{1}{2}\right)\right| = o(1), \left|\frac{1}{n_r} \sum_{i=1}^{n_r} \hat{y}_i y_i c_i^{(t)} - \left(p - \frac{1}{2}\right)\right| = o(1), r \in \{\pm 1\}, \quad (\text{C.7})$$

609 with high probability.

610 Besides, we also need an approximation about n_1 and n_{-1} , which is given as the following lemma:

611 **Lemma C.13.** For $r \in \{\pm 1\}$, it holds with probability at least $1 - 2\delta$ that

$$\left|n_r - \frac{n_u}{2}\right| < \sqrt{\frac{n_u}{2} \log \frac{1}{\delta}},$$

612 where $n_r := |\{(x_i, y_i) | y_i = r, i \in [n_u]\}|$.

613 *Proof of Lemma C.13.* Note that $n_r = \sum_{i=1}^{n_u} \mathbb{1}[X_i = r]$, $r \in \{\pm 1\}$ where X_i takes label $+1$ or -1
 614 with equal probability $1/2$, according to Hoeffding's inequality, we have

$$\mathbb{P}\left(\left|\sum_{i=1}^{n_u} \mathbb{1}[X_i = r] - \mathbb{E}\left[\sum_{i=1}^{n_u} \mathbb{1}[X_i = r]\right]\right| \geq t\right) \leq 2 \exp\left(-\frac{2t^2}{n_u}\right), r \in \{\pm 1\},$$

615 and it follows that

$$\mathbb{P}\left(\left|n_r - \frac{n_u}{2}\right| \geq t\right) \leq 2 \exp\left(-\frac{2t^2}{n_u}\right), r \in \{\pm 1\},$$

616 leading to

$$\left|n_r - \frac{n_u}{2}\right| < \sqrt{\frac{n_u}{2} \log \frac{1}{\delta}},$$

617 with probability at least $1 - 2\delta$. \square

618 For labeled dataset $S' = \{(x'_i, y'_i)\}_{i=1}^{n_1}$, we also have

619 **Lemma C.14.** For $r \in \{\pm 1\}$, it holds with probability at least $1 - 2\delta$ that

$$\left|n'_r - \frac{n_1}{2}\right| < \sqrt{\frac{n_1}{2} \log \frac{1}{\delta}},$$

620 where $n'_r := |\{(x'_i, y'_i) | y'_i = r, i \in [n_1]\}|$.

621 Then we are prepared to estimate a lower bound of increasing speed of $\hat{\Lambda}^{(t)}$ and an upper bound of
 622 decreasing speed of $\bar{\Lambda}^{(t)}$ in the following lemma.

623 **Lemma C.15.** For $\widehat{\Lambda}_1^{(t)} := \max_{1 \leq j \leq m} \langle \mathbf{w}_j^{(t)}, \mathbf{v} \rangle$ and $\widehat{\Lambda}_{-1}^{(t)} := \max_{m+1 \leq j \leq 2m} \langle \mathbf{w}_j^{(t)}, -\mathbf{v} \rangle$, we have
 624 with high probability that

$$\widehat{\Lambda}_r^{(t+1)} \geq (1 - \eta\lambda) \cdot \widehat{\Lambda}_r^{(t)} + \eta \cdot \left(p - \frac{1}{2}\right) \cdot \Theta(d) \cdot (\widehat{\Lambda}_r^{(t)})^{q-1}, r \in \{\pm 1\}.$$

625 For $\bar{\Lambda}_1^{(t)} := \max_{m+1 \leq j \leq 2m} \langle \mathbf{w}_j^{(t)}, \mathbf{v} \rangle$ and $\bar{\Lambda}_{-1}^{(t)} := \max_{1 \leq j \leq m} \langle \mathbf{w}_j^{(t)}, -\mathbf{v} \rangle$, we have with high proba-
 626 bility that

$$\bar{\Lambda}_r^{(t+1)} \leq (1 - \eta\lambda) \cdot \bar{\Lambda}_r^{(t)}, r \in \{\pm 1\}.$$

627 *Proof of Lemma C.15.* We first prove the former inequality. Let $j^* = \arg \max_{1 \leq j \leq m} \langle \mathbf{w}_j^{(t)}, \mathbf{v} \rangle$ and
 628 note that $u_{j^*} = \mathbb{1}_{[1 \leq j \leq m]} - \mathbb{1}_{[m+1 \leq j \leq 2m]} = 1$, then we have

$$\begin{aligned} \widehat{\Lambda}_1^{(t+1)} &\geq \langle \mathbf{w}_{j^*}^{(t+1)}, \mathbf{v} \rangle \\ &= (1 - \eta\lambda) \cdot \langle \mathbf{w}_{j^*}^{(t)}, \mathbf{v} \rangle + \frac{q\eta}{n_1 + n_u} \left(\underbrace{\sum_{i=1}^{n_u} y_i \widehat{y}_i c_i^{(t)} [\langle \mathbf{w}_{j^*}^{(t)}, y_i \cdot \mathbf{v} \rangle]_+^{q-1} \|\mathbf{v}\|_2^2}_{\clubsuit} + \underbrace{\sum_{i=1}^{n_1} b_i^{(t)} [\langle \mathbf{w}_{j^*}^{(t)}, y'_i \cdot \mathbf{v} \rangle]_+^{q-1} \|\mathbf{v}\|_2^2}_{\star} \right) \end{aligned}$$

629 Then we respectively estimate terms \clubsuit and \star .

630 For \clubsuit , note the definition of j^* that $\widehat{\Lambda}_1^{(t)} = \langle \mathbf{w}_{j^*}^{(t)}, \mathbf{v} \rangle$ and note the increasing property of $\widehat{\Lambda}_1^{(t)}$ and
 631 $\widehat{\Lambda}_1^{(0)} > 0$ with high probability, we have $\langle \mathbf{w}_{j^*}^{(t)}, \mathbf{v} \rangle > 0$. It follows that

$$\begin{aligned} \underbrace{\sum_{i=1}^{n_u} y_i \widehat{y}_i c_i^{(t)} [\langle \mathbf{w}_{j^*}^{(t)}, y_i \cdot \mathbf{v} \rangle]_+^{q-1} \|\mathbf{v}\|_2^2}_{\clubsuit} &= \sum_{i \in S_1} y_i \widehat{y}_i c_i^{(t)} [\langle \mathbf{w}_{j^*}^{(t)}, \mathbf{v} \rangle]_+^{q-1} \|\mathbf{v}\|_2^2 + \sum_{i \in S_{-1}} y_i \widehat{y}_i c_i^{(t)} [-\langle \mathbf{w}_{j^*}^{(t)}, \mathbf{v} \rangle]_+^{q-1} \|\mathbf{v}\|_2^2 \\ &= \sum_{i \in S_1} y_i \widehat{y}_i c_i^{(t)} [\langle \mathbf{w}_{j^*}^{(t)}, \mathbf{v} \rangle]_+^{q-1} \|\mathbf{v}\|_2^2 \\ &= \left(\sum_{i \in S_1} y_i \widehat{y}_i c_i^{(t)} \right) \cdot \|\mathbf{v}\|_2^2 \cdot (\widehat{\Lambda}_1^{(t)})^{q-1} \\ &= n_1 \cdot \left(p - \frac{1}{2} \pm o(1)\right) \cdot \|\mathbf{v}\|_2^2 \cdot (\widehat{\Lambda}_1^{(t)})^{q-1}, \end{aligned} \quad (\text{C.8})$$

632 where $S_1 := \{(\mathbf{x}_i, y_i) | y_i = 1, i \in [n_u]\}$, $S_{-1} := \{(\mathbf{x}_i, y_i) | y_i = -1, i \in [n_u]\}$, $n_1 = |S_1|$ and the
 633 last equality is due to (C.7).

634 For \star , similarly we have

$$\begin{aligned} \underbrace{\sum_{i=1}^{n_1} b_i^{(t)} [\langle \mathbf{w}_{j^*}^{(t)}, y'_i \cdot \mathbf{v} \rangle]_+^{q-1} \|\mathbf{v}\|_2^2}_{\star} &= \sum_{i \in S'_1} b_i^{(t)} [\langle \mathbf{w}_{j^*}^{(t)}, \mathbf{v} \rangle]_+^{q-1} \|\mathbf{v}\|_2^2 + \sum_{i \in S'_{-1}} b_i^{(t)} [-\langle \mathbf{w}_{j^*}^{(t)}, \mathbf{v} \rangle]_+^{q-1} \|\mathbf{v}\|_2^2 \\ &= \sum_{i \in S'_1} b_i^{(t)} [\langle \mathbf{w}_{j^*}^{(t)}, \mathbf{v} \rangle]_+^{q-1} \|\mathbf{v}\|_2^2 \\ &= \left(\sum_{i \in S'_1} b_i^{(t)} \right) \cdot \|\mathbf{v}\|_2^2 \cdot (\widehat{\Lambda}_1^{(t)})^{q-1} \\ &= n'_1 \cdot \left(\frac{1}{2} \pm o(1)\right) \cdot \|\mathbf{v}\|_2^2 \cdot (\widehat{\Lambda}_1^{(t)})^{q-1}, \end{aligned} \quad (\text{C.9})$$

635 where $S'_1 = \{(\mathbf{x}'_i, y'_i) | y'_i = 1, i \in [n_1]\}$, $S'_{-1} = \{(\mathbf{x}'_i, y'_i) | y'_i = -1, i \in [n_1]\}$, $n'_1 = |S'_1|$ and the last
 636 equality is due to Lemma C.11.

637 According to (C.8) and (C.9), we have

$$\widehat{\Lambda}_1^{(t+1)}$$

$$\begin{aligned}
&\geq (1 - \eta\lambda) \cdot \widehat{\Lambda}_1^{(t)} + \frac{q\eta}{n_1 + n_u} \left(n_1 \cdot \left(p - \frac{1}{2} \pm o(1) \right) \cdot \|\mathbf{v}\|_2^2 \cdot (\widehat{\Lambda}_1^{(t)})^{q-1} + n'_1 \cdot \left(\frac{1}{2} \pm o(1) \right) \cdot \|\mathbf{v}\|_2^2 \cdot (\widehat{\Lambda}_1^{(t)})^{q-1} \right) \\
&= (1 - \eta\lambda) \cdot \widehat{\Lambda}_1^{(t)} + \frac{q\eta n_1}{n_1 + n_u} \cdot \left(p - \frac{1}{2} \pm o(1) \right) \cdot \|\mathbf{v}\|_2^2 \cdot (\widehat{\Lambda}_1^{(t)})^{q-1} + \frac{q\eta n'_1}{n_1 + n_u} \cdot \left(\frac{1}{2} \pm o(1) \right) \cdot \|\mathbf{v}\|_2^2 \cdot (\widehat{\Lambda}_1^{(t)})^{q-1} \\
&= (1 - \eta\lambda) \cdot \widehat{\Lambda}_1^{(t)} + q\eta \cdot \left(\frac{n_1}{n_1 + n_u} \cdot \left(p - \frac{1}{2} \pm o(1) \right) + \frac{n'_1}{n_1 + n_u} \cdot \left(\frac{1}{2} \pm o(1) \right) \right) \cdot \|\mathbf{v}\|_2^2 \cdot (\widehat{\Lambda}_1^{(t)})^{q-1} \\
&= (1 - \eta\lambda) \cdot \widehat{\Lambda}_1^{(t)} + q\eta \cdot \underbrace{\left(\frac{n_1}{n_1 + n_u} \cdot \left(p - \frac{1}{2} \right) + \frac{n'_1}{n_1 + n_u} \cdot \frac{1}{2} \pm o(1) \right)}_{\spadesuit} \cdot \|\mathbf{v}\|_2^2 \cdot (\widehat{\Lambda}_1^{(t)})^{q-1}.
\end{aligned} \tag{C.10}$$

638 According to Lemma C.13 and Lemma C.14, and note that $n_l = \widetilde{\Theta}(1)$, $n_u = \omega(d^{4\epsilon})$, we have for \spadesuit
639 that with probability at least $1 - 4\delta$

$$\begin{aligned}
&\left| \underbrace{\frac{n_1}{n_1 + n_u} \cdot \left(p - \frac{1}{2} \right) + \frac{n'_1}{n_1 + n_u} \cdot \frac{1}{2}}_{\spadesuit} - \frac{n_u}{2(n_1 + n_u)} \cdot \left(p - \frac{1}{2} \right) - \frac{n_1}{2(n_1 + n_u)} \cdot \frac{1}{2} \right| \\
&\leq \frac{|n_1 - \frac{n_u}{2}|}{n_1 + n_u} \cdot \left(p - \frac{1}{2} \right) + \frac{|n'_1 - \frac{n_1}{2}|}{n_1 + n_u} \cdot \frac{1}{2} \\
&\leq \frac{\sqrt{\frac{n_u}{2} \log \frac{1}{\delta}}}{n_1 + n_u} \cdot \left(p - \frac{1}{2} \right) + \frac{\sqrt{\frac{n_1}{2} \log \frac{1}{\delta}}}{n_1 + n_u} \cdot \frac{1}{2} \\
&= \Theta\left(\frac{1}{\sqrt{n_u}}\right) \\
&= o(1)
\end{aligned}$$

640 Therefore, note that $n_u = \omega(n_l)$ and $n_u = \omega(1)$, we have

$$\begin{aligned}
\underbrace{\frac{n_1}{n_1 + n_u} \cdot \left(p - \frac{1}{2} \right) + \frac{n'_1}{n_1 + n_u} \cdot \frac{1}{2}}_{\spadesuit} &= \frac{n_u}{2(n_1 + n_u)} \cdot \left(p - \frac{1}{2} \right) + \frac{n_1}{2(n_1 + n_u)} \cdot \frac{1}{2} \pm o(1) \\
&= \frac{1}{2} \cdot \left(p - \frac{1}{2} \right) \pm o(1)
\end{aligned} \tag{C.11}$$

641 Plugging (C.11) into (C.10), we have

$$\begin{aligned}
\widehat{\Lambda}_1^{(t+1)} &\geq (1 - \eta\lambda) \cdot \widehat{\Lambda}_1^{(t)} + q\eta \cdot \left(\frac{1}{2} \cdot \left(p - \frac{1}{2} \right) \pm o(1) \right) \cdot \|\mathbf{v}\|_2^2 \cdot (\widehat{\Lambda}_1^{(t)})^{q-1} \\
&= (1 - \eta\lambda) \cdot \widehat{\Lambda}_1^{(t)} + \eta \cdot \left(p - \frac{1}{2} \right) \cdot \Theta(d) \cdot (\widehat{\Lambda}_1^{(t)})^{q-1},
\end{aligned} \tag{C.12}$$

642 which verifies the first inequality of case $r = 1$ in the lemma.

643 Let $j^{**} = \operatorname{argmax}_{m+1 \leq j \leq 2m} \langle \mathbf{w}_j^{(t)}, -\mathbf{v} \rangle$ and note that $u_{j^{**}} = \mathbb{1}_{[1 \leq j \leq m]} - \mathbb{1}_{[m+1 \leq j \leq 2m]} = -1$, we
644 have

$$\begin{aligned}
\widehat{\Lambda}_{-1}^{(t+1)} &\geq \langle \mathbf{w}_{j^{**}}^{(t+1)}, -\mathbf{v} \rangle \\
&= (1 - \eta\lambda) \cdot \langle \mathbf{w}_{j^{**}}^{(t)}, -\mathbf{v} \rangle + \frac{q\eta}{n_1 + n_u} \underbrace{\left(\sum_{i=1}^{n_u} y_i \widehat{y}_i c_i^{(t)} [\langle \mathbf{w}_{j^{**}}^{(t)}, y_i \cdot \mathbf{v} \rangle]_+^{q-1} \|\mathbf{v}\|_2^2 \right)}_{\clubsuit} \\
&\quad + \underbrace{\left(\sum_{i=1}^{n_l} b_i^{(t)} [\langle \mathbf{w}_{j^{**}}^{(t)}, y'_i \cdot \mathbf{v} \rangle]_+^{q-1} \|\mathbf{v}\|_2^2 \right)}_{\star}
\end{aligned}$$

645 For \clubsuit , note the definition of j^{**} that $\widehat{\Lambda}_{-1}^{(t)} = \langle \mathbf{w}_{j^{**}}^{(t)}, -\mathbf{v} \rangle$ and note the increasing property of $\widehat{\Lambda}_{-1}^{(t)}$
 646 and $\widehat{\Lambda}_{-1}^{(0)} > 0$ with high probability, we have $\langle \mathbf{w}_{j^{**}}^{(t)}, -\mathbf{v} \rangle > 0$. According to (C.7), it follows that

$$\underbrace{\sum_{i=1}^{n_u} y_i \widehat{y}_i c_i^{(t)} [\langle \mathbf{w}_{j^{**}}^{(t)}, y_i \cdot \mathbf{v} \rangle]_+^{q-1} \|\mathbf{v}\|_2^2}_{\clubsuit} = \sum_{i \in S_{-1}} y_i \widehat{y}_i c_i^{(t)} [\langle \mathbf{w}_{j^{**}}^{(t)}, -\mathbf{v} \rangle]_+^{q-1} \|\mathbf{v}\|_2^2$$

$$= n_{-1} \cdot \left(p - \frac{1}{2} \pm o(1) \right) \cdot \|\mathbf{v}\|_2^2 \cdot (\widehat{\Lambda}_{-1}^{(t)})^{q-1}, \quad (\text{C.13})$$

647 where $S_{-1} := \{(\mathbf{x}_i, y_i) | y_i = -1, i \in [n_u]\}$, $n_{-1} = |S_{-1}|$.

648 For \star , according to Lemma C.11, similarly we have

$$\underbrace{\sum_{i=1}^{n_l} b_i^{(t)} [\langle \mathbf{w}_{j^{**}}^{(t)}, y'_i \cdot \mathbf{v} \rangle]_+^{q-1} \|\mathbf{v}\|_2^2}_{\star} = \sum_{i \in S'_{-1}} b_i^{(t)} [\langle \mathbf{w}_{j^{**}}^{(t)}, -\mathbf{v} \rangle]_+^{q-1} \|\mathbf{v}\|_2^2 = n'_{-1} \cdot \left(\frac{1}{2} \pm o(1) \right) \cdot \|\mathbf{v}\|_2^2 \cdot (\widehat{\Lambda}_{-1}^{(t)})^{q-1},$$

$$(\text{C.14})$$

649 where $S'_{-1} = \{(\mathbf{x}'_i, y'_i) | y'_i = -1, i \in [n_l]\}$ and $n'_{-1} = |S'_{-1}|$.

650 According to (C.13) and (C.14), we have

$$\widehat{\Lambda}_{-1}^{(t+1)} \geq (1 - \eta\lambda) \cdot \widehat{\Lambda}_{-1}^{(t)} + q\eta \cdot \underbrace{\left(\frac{n_{-1}}{n_l + n_u} \cdot \left(p - \frac{1}{2} \right) + \frac{n'_{-1}}{n_l + n_u} \cdot \frac{1}{2} \pm o(1) \right)}_{\spadesuit} \cdot \|\mathbf{v}\|_2^2 \cdot (\widehat{\Lambda}_{-1}^{(t)})^{q-1}.$$

$$(\text{C.15})$$

651 According to Lemma C.13 and Lemma C.14, and note that $n_l = \widetilde{\Theta}(1)$, $n_u = \omega(d^{4\epsilon})$, we have for \spadesuit
 652 that with probability at least $1 - 4\delta$

$$\left| \underbrace{\frac{n_{-1}}{n_l + n_u} \cdot \left(p - \frac{1}{2} \right) + \frac{n'_{-1}}{n_l + n_u} \cdot \frac{1}{2}}_{\spadesuit} - \frac{n_u}{2(n_l + n_u)} \cdot \left(p - \frac{1}{2} \right) - \frac{n_l}{2(n_l + n_u)} \cdot \frac{1}{2} \right|$$

$$\leq \frac{|n_{-1} - \frac{n_u}{2}|}{n_l + n_u} \cdot \left(p - \frac{1}{2} \right) + \frac{|n'_{-1} - \frac{n_l}{2}|}{n_l + n_u} \cdot \frac{1}{2}$$

$$\leq \frac{\sqrt{\frac{n_u}{2} \log \frac{1}{\delta}}}{n_l + n_u} \cdot \left(p - \frac{1}{2} \right) + \frac{\sqrt{\frac{n_l}{2} \log \frac{1}{\delta}}}{n_l + n_u} \cdot \frac{1}{2}$$

$$= \Theta\left(\frac{1}{\sqrt{n_u}} \right)$$

$$= o(1).$$

653 Therefore, note that $n_u = \omega(n_l)$ and $n_u = \omega(1)$, we have

$$\underbrace{\frac{n_{-1}}{n_l + n_u} \cdot \left(p - \frac{1}{2} \right) + \frac{n'_{-1}}{n_l + n_u} \cdot \frac{1}{2}}_{\spadesuit} = \frac{n_u}{2(n_l + n_u)} \cdot \left(p - \frac{1}{2} \right) + \frac{n_l}{2(n_l + n_u)} \cdot \frac{1}{2} \pm o(1)$$

$$= \frac{1}{2} \cdot \left(p - \frac{1}{2} \right) \pm o(1) \quad (\text{C.16})$$

654 Plugging (C.16) into (C.15), we have

$$\widehat{\Lambda}_{-1}^{(t+1)} \geq (1 - \eta\lambda) \cdot \widehat{\Lambda}_{-1}^{(t)} + q\eta \cdot \left(\frac{1}{2} \cdot \left(p - \frac{1}{2} \right) \pm o(1) \right) \cdot \|\mathbf{v}\|_2^2 \cdot (\widehat{\Lambda}_{-1}^{(t)})^{q-1}$$

$$= (1 - \eta\lambda) \cdot \widehat{\Lambda}_{-1}^{(t)} + \eta \cdot \left(p - \frac{1}{2} \right) \cdot \Theta(d) \cdot (\widehat{\Lambda}_{-1}^{(t)})^{q-1}, \quad (\text{C.17})$$

655 which verifies the first inequality of case $r = -1$ in the lemma.

656 Next, we prove the latter part of the lemma. Let $j^{\natural} = \arg \max_{m+1 \leq j \leq 2m} \langle \mathbf{w}_j^{(t+1)}, \mathbf{v} \rangle$, then we have:

$$\begin{aligned} \bar{\Lambda}_1^{(t+1)} &= \langle \mathbf{w}_{j^{\natural}}^{(t+1)}, \mathbf{v} \rangle \\ &= (1 - \eta\lambda) \cdot \langle \mathbf{w}_{j^{\natural}}^{(t)}, \mathbf{v} \rangle - \frac{q\eta}{n_1 + n_u} \underbrace{\left(\sum_{i=1}^{n_u} y_i \hat{y}_i c_i^{(t)} [\langle \mathbf{w}_{j^{\natural}}^{(t)}, y_i \cdot \mathbf{v} \rangle]_+^{q-1} \|\mathbf{v}\|_2^2 \right)}_{\clubsuit} \\ &\quad + \underbrace{\sum_{i=1}^{n_1} b_i^{(t)} [\langle \mathbf{w}_{j^{\natural}}^{(t)}, y'_i \cdot \mathbf{v} \rangle]_+^{q-1} \|\mathbf{v}\|_2^2}_{\star}. \end{aligned}$$

657 For \clubsuit , according to (C.7), we have

$$\begin{aligned} &\underbrace{\sum_{i=1}^{n_u} y_i \hat{y}_i c_i^{(t)} [\langle \mathbf{w}_{j^{\natural}}^{(t)}, y_i \cdot \mathbf{v} \rangle]_+^{q-1} \|\mathbf{v}\|_2^2}_{\clubsuit} \\ &= \sum_{i \in S_1} y_i \hat{y}_i c_i^{(t)} [\langle \mathbf{w}_{j^{\natural}}^{(t)}, \mathbf{v} \rangle]_+^{q-1} \|\mathbf{v}\|_2^2 + \sum_{i \in S_{-1}} y_i \hat{y}_i c_i^{(t)} [\langle \mathbf{w}_{j^{\natural}}^{(t)}, -\mathbf{v} \rangle]_+^{q-1} \|\mathbf{v}\|_2^2 \\ &= \left(\sum_{i \in S_1} y_i \hat{y}_i c_i^{(t)} \right) \cdot [\langle \mathbf{w}_{j^{\natural}}^{(t)}, \mathbf{v} \rangle]_+^{q-1} \|\mathbf{v}\|_2^2 + \left(\sum_{i \in S_{-1}} y_i \hat{y}_i c_i^{(t)} \right) \cdot [\langle \mathbf{w}_{j^{\natural}}^{(t)}, -\mathbf{v} \rangle]_+^{q-1} \|\mathbf{v}\|_2^2 \\ &= n_1 \cdot \left(p - \frac{1}{2} \pm o(1) \right) \cdot [\langle \mathbf{w}_{j^{\natural}}^{(t)}, \mathbf{v} \rangle]_+^{q-1} \|\mathbf{v}\|_2^2 + n_{-1} \cdot \left(p - \frac{1}{2} \pm o(1) \right) \cdot [\langle \mathbf{w}_{j^{\natural}}^{(t)}, -\mathbf{v} \rangle]_+^{q-1} \|\mathbf{v}\|_2^2 \geq 0, \end{aligned}$$

658 and for \star it's obvious that

$$\underbrace{\sum_{i=1}^{n_1} b_i^{(t)} [\langle \mathbf{w}_{j^{\natural}}^{(t)}, y'_i \cdot \mathbf{v} \rangle]_+^{q-1} \|\mathbf{v}\|_2^2}_{\star} \geq 0.$$

659 Therefore, it follows that

$$\bar{\Lambda}_1^{(t+1)} \leq (1 - \eta\lambda) \cdot \langle \mathbf{w}_{j^{\natural}}^{(t)}, \mathbf{v} \rangle \leq (1 - \eta\lambda) \bar{\Lambda}_1^{(t)}.$$

660 Let $j^{\natural\sharp} = \arg \max_{1 \leq j \leq m} \langle \mathbf{w}_j^{(t+1)}, -\mathbf{v} \rangle$, then we have:

$$\begin{aligned} \bar{\Lambda}_{-1}^{(t+1)} &= \langle \mathbf{w}_{j^{\natural\sharp}}^{(t+1)}, -\mathbf{v} \rangle \\ &= (1 - \eta\lambda) \cdot \langle \mathbf{w}_{j^{\natural\sharp}}^{(t)}, -\mathbf{v} \rangle - \frac{q\eta}{n_1 + n_u} \left(\sum_{i=1}^{n_u} y_i \hat{y}_i c_i^{(t)} [\langle \mathbf{w}_{j^{\natural\sharp}}^{(t)}, y_i \cdot \mathbf{v} \rangle]_+^{q-1} \|\mathbf{v}\|_2^2 \right. \\ &\quad \left. + \sum_{i=1}^{n_1} b_i^{(t)} [\langle \mathbf{w}_{j^{\natural\sharp}}^{(t)}, y'_i \cdot \mathbf{v} \rangle]_+^{q-1} \|\mathbf{v}\|_2^2 \right) \\ &\leq (1 - \eta\lambda) \cdot \langle \mathbf{w}_{j^{\natural\sharp}}^{(t)}, -\mathbf{v} \rangle \\ &\leq (1 - \eta\lambda) \cdot \bar{\Lambda}_{-1}^{(t)}, \end{aligned}$$

661 which verifies the second part of the lemma. \square

662 Although the accuracy of pseudo-labeler is larger than 1/2, which is used as an assumption in the
663 previous proof, we can also analyse the model with high label flipping probability and the accuracy of
664 pseudo-labeler p is smaller than 1/2. In this case, the neural network for pre-training will turn to fit
665 the opposite direction of feature vector, $\bar{\Lambda}_r^{(t)}$ will increase and $\hat{\Lambda}_r^{(t)}$ will decrease, which is formulated
666 as the following lemma.

667 **Lemma C.16.** For $\widehat{\Lambda}_1^{(t)} := \max_{1 \leq j \leq m} \langle \mathbf{w}_j^{(t)}, \mathbf{v} \rangle$ and $\widehat{\Lambda}_{-1}^{(t)} := \max_{m+1 \leq j \leq 2m} \langle \mathbf{w}_j^{(t)}, -\mathbf{v} \rangle$, we have
 668 with high probability that

$$\widehat{\Lambda}_r^{(t+1)} \leq (1 - \eta\lambda) \cdot \widehat{\Lambda}_r^{(t)}, r \in \{\pm 1\}.$$

669 For $\bar{\Lambda}_1^{(t)} := \max_{m+1 \leq j \leq 2m} \langle \mathbf{w}_j^{(t)}, \mathbf{v} \rangle$ and $\bar{\Lambda}_{-1}^{(t)} := \max_{1 \leq j \leq m} \langle \mathbf{w}_j^{(t)}, -\mathbf{v} \rangle$, we have with high proba-
 670 bility that

$$\bar{\Lambda}_r^{(t+1)} \geq (1 - \eta\lambda) \cdot \bar{\Lambda}_r^{(t)} + \eta \cdot \left(\frac{1}{2} - p\right) \cdot \Theta(d) \cdot (\bar{\Lambda}_r^{(t)})^{q-1}, r \in \{\pm 1\}.$$

671 *Proof of Lemma C.16.* First, we prove the former part of this lemma. Let $j^* =$
 672 $\arg \max_{1 \leq j \leq m} \langle \mathbf{w}_j^{(t+1)}, \mathbf{v} \rangle$ and note that $u_{j^*} = \mathbb{1}_{[1 \leq j \leq m]} - \mathbb{1}_{[m+1 \leq j \leq 2m]} = 1$, then we have

$$\begin{aligned} \widehat{\Lambda}_1^{(t+1)} &= \langle \mathbf{w}_{j^*}^{(t+1)}, \mathbf{v} \rangle \\ &= (1 - \eta\lambda) \cdot \langle \mathbf{w}_{j^*}^{(t)}, \mathbf{v} \rangle + \underbrace{\frac{q\eta}{n_1 + n_u} \left(\sum_{i=1}^{n_u} y_i \widehat{y}_i c_i^{(t)} [\langle \mathbf{w}_{j^*}^{(t)}, y_i \cdot \mathbf{v} \rangle]_+^{q-1} \|\mathbf{v}\|_2^2 \right)}_{\clubsuit} \\ &\quad + \underbrace{\sum_{i=1}^{n_1} b_i^{(t)} [\langle \mathbf{w}_{j^*}^{(t)}, y'_i \cdot \mathbf{v} \rangle]_+^{q-1} \|\mathbf{v}\|_2^2}_{\star}. \end{aligned}$$

673 For \clubsuit , according to (C.7), we have

$$\begin{aligned} &\underbrace{\sum_{i=1}^{n_u} y_i \widehat{y}_i c_i^{(t)} [\langle \mathbf{w}_{j^*}^{(t)}, y_i \cdot \mathbf{v} \rangle]_+^{q-1} \|\mathbf{v}\|_2^2}_{\clubsuit} \\ &= \sum_{i \in S_1} y_i \widehat{y}_i c_i^{(t)} [\langle \mathbf{w}_{j^*}^{(t)}, \mathbf{v} \rangle]_+^{q-1} \|\mathbf{v}\|_2^2 + \sum_{i \in S_{-1}} y_i \widehat{y}_i c_i^{(t)} [\langle \mathbf{w}_{j^*}^{(t)}, -\mathbf{v} \rangle]_+^{q-1} \|\mathbf{v}\|_2^2 \\ &= \left(\sum_{i \in S_1} y_i \widehat{y}_i c_i^{(t)} \right) \cdot [\langle \mathbf{w}_{j^*}^{(t)}, \mathbf{v} \rangle]_+^{q-1} \|\mathbf{v}\|_2^2 + \left(\sum_{i \in S_{-1}} y_i \widehat{y}_i c_i^{(t)} \right) \cdot [\langle \mathbf{w}_{j^*}^{(t)}, -\mathbf{v} \rangle]_+^{q-1} \|\mathbf{v}\|_2^2 \\ &= n_1 \cdot \left(p - \frac{1}{2} \pm o(1) \right) \cdot [\langle \mathbf{w}_{j^*}^{(t)}, \mathbf{v} \rangle]_+^{q-1} \|\mathbf{v}\|_2^2 + n_{-1} \cdot \left(p - \frac{1}{2} \pm o(1) \right) \cdot [\langle \mathbf{w}_{j^*}^{(t)}, -\mathbf{v} \rangle]_+^{q-1} \|\mathbf{v}\|_2^2, \end{aligned}$$

674 For \star , according to (C.7), we have

$$\begin{aligned} &\underbrace{\sum_{i=1}^{n_1} b_i^{(t)} [\langle \mathbf{w}_{j^*}^{(t)}, y'_i \cdot \mathbf{v} \rangle]_+^{q-1} \|\mathbf{v}\|_2^2}_{\star} \\ &= \sum_{i \in S'_1} b_i^{(t)} [\langle \mathbf{w}_{j^*}^{(t)}, \mathbf{v} \rangle]_+^{q-1} \|\mathbf{v}\|_2^2 + \sum_{i \in S'_{-1}} b_i^{(t)} [\langle \mathbf{w}_{j^*}^{(t)}, -\mathbf{v} \rangle]_+^{q-1} \|\mathbf{v}\|_2^2 \\ &= n'_1 \cdot \left(\frac{1}{2} \pm o(1) \right) \cdot [\langle \mathbf{w}_{j^*}^{(t)}, \mathbf{v} \rangle]_+^{q-1} \|\mathbf{v}\|_2^2 + n'_{-1} \cdot \left(\frac{1}{2} \pm o(1) \right) \cdot [\langle \mathbf{w}_{j^*}^{(t)}, -\mathbf{v} \rangle]_+^{q-1} \|\mathbf{v}\|_2^2, \end{aligned}$$

675 It follows that

$$\begin{aligned} &\underbrace{\sum_{i=1}^{n_u} y_i \widehat{y}_i c_i^{(t)} [\langle \mathbf{w}_{j^*}^{(t)}, y_i \cdot \mathbf{v} \rangle]_+^{q-1} \|\mathbf{v}\|_2^2}_{\clubsuit} + \underbrace{\sum_{i=1}^{n_1} b_i^{(t)} [\langle \mathbf{w}_{j^*}^{(t)}, y'_i \cdot \mathbf{v} \rangle]_+^{q-1} \|\mathbf{v}\|_2^2}_{\star} \\ &= \left(n_1 \cdot \left(p - \frac{1}{2} \pm o(1) \right) + n'_1 \cdot \left(\frac{1}{2} \pm o(1) \right) \right) \cdot [\langle \mathbf{w}_{j^*}^{(t)}, \mathbf{v} \rangle]_+^{q-1} \|\mathbf{v}\|_2^2 \end{aligned}$$

$$+ \left(n_{-1} \cdot \left(p - \frac{1}{2} \pm o(1) \right) + n'_{-1} \cdot \left(\frac{1}{2} \pm o(1) \right) \right) \cdot [\langle \mathbf{w}_{j^*}^{(t)}, \mathbf{v} \rangle]_+^{q-1} \|\mathbf{v}\|_2^2.$$

676 According to Lemma C.13 and note that $n_u = \omega(n_1)$, it holds with probability at least $1 - 8\delta$ that

$$\begin{aligned} n'_1 \cdot \left(\frac{1}{2} \pm o(1) \right) &\leq \left(\frac{n_1}{2} + \sqrt{\frac{n_1}{2} \log \frac{1}{\delta}} \right) \cdot \left(\frac{1}{2} \pm o(1) \right) = \Theta(n_1) = o(n_u) \\ &\leq \left(\frac{n_u}{2} + \sqrt{\frac{n_u}{2} \log \frac{1}{\delta}} \right) \cdot \left(\frac{1}{2} - p \pm o(1) \right) \leq n_1 \cdot \left(\frac{1}{2} - p \pm o(1) \right), \end{aligned}$$

677

$$\begin{aligned} n'_{-1} \cdot \left(\frac{1}{2} \pm o(1) \right) &\leq \left(\frac{n_1}{2} + \sqrt{\frac{n_1}{2} \log \frac{1}{\delta}} \right) \cdot \left(\frac{1}{2} \pm o(1) \right) = \Theta(n_1) = o(n_u) \\ &\leq \left(\frac{n_u}{2} + \sqrt{\frac{n_u}{2} \log \frac{1}{\delta}} \right) \cdot \left(\frac{1}{2} - p \pm o(1) \right) \leq n_{-1} \cdot \left(\frac{1}{2} - p \pm o(1) \right), \end{aligned}$$

678 leading to $\clubsuit + \star \leq 0$. Therefore,

$$\widehat{\Lambda}_1^{(t+1)} \leq (1 - \eta\lambda) \langle \mathbf{w}_{j^*}^{(t)}, \mathbf{v} \rangle \leq (1 - \eta\lambda) \cdot \widehat{\Lambda}_1^{(t)}.$$

679 And we can prove in a similar way that $\widehat{\Lambda}_{-1}^{(t+1)} \leq (1 - \eta\lambda) \cdot \widehat{\Lambda}_{-1}^{(t)}$.

680 Next, we prove the second part of the lemma. Let $j^\ddagger = \arg \max_{m+1 \leq j \leq 2m} \langle \mathbf{w}_j^{(t)}, \mathbf{v} \rangle$ and note that

681 $u_{j^\ddagger} = \mathbb{1}_{[1 \leq j \leq m]} - \mathbb{1}_{[m+1 \leq j \leq 2m]} = -1$, then we have

$$\begin{aligned} \bar{\Lambda}_1^{(t+1)} &\geq \langle \mathbf{w}_{j^\ddagger}^{(t+1)}, \mathbf{v} \rangle \\ &= (1 - \eta\lambda) \cdot \langle \mathbf{w}_{j^\ddagger}^{(t)}, \mathbf{v} \rangle - \frac{q\eta}{n_1 + n_u} \underbrace{\left(\sum_{i=1}^{n_u} y_i \widehat{y}_i c_i^{(t)} [\langle \mathbf{w}_{j^\ddagger}^{(t)}, y_i \cdot \mathbf{v} \rangle]_+^{q-1} \|\mathbf{v}\|_2^2 \right)}_{\clubsuit} \\ &\quad + \underbrace{\left(\sum_{i=1}^{n_1} b_i^{(t)} [\langle \mathbf{w}_{j^\ddagger}^{(t)}, y'_i \cdot \mathbf{v} \rangle]_+^{q-1} \|\mathbf{v}\|_2^2 \right)}_{\star}. \end{aligned}$$

682 For \clubsuit , note the definition of j^\ddagger that $\bar{\Lambda}_1^{(t)} = \langle \mathbf{w}_{j^\ddagger}^{(t)}, \mathbf{v} \rangle$ and note the increasing property of $\bar{\Lambda}_1^{(t)}$ in this

683 case and $\bar{\Lambda}_1^{(0)} > 0$ with high probability, we have $\langle \mathbf{w}_{j^\ddagger}^{(t)}, \mathbf{v} \rangle > 0$. It follows that

$$\begin{aligned} \underbrace{\sum_{i=1}^{n_u} y_i \widehat{y}_i c_i^{(t)} [\langle \mathbf{w}_{j^\ddagger}^{(t)}, y_i \cdot \mathbf{v} \rangle]_+^{q-1} \|\mathbf{v}\|_2^2}_{\clubsuit} &= \sum_{i \in S_1} y_i \widehat{y}_i c_i^{(t)} [\langle \mathbf{w}_{j^\ddagger}^{(t)}, \mathbf{v} \rangle]_+^{q-1} \|\mathbf{v}\|_2^2 + \sum_{i \in S_{-1}} y_i \widehat{y}_i c_i^{(t)} [-\langle \mathbf{w}_{j^\ddagger}^{(t)}, \mathbf{v} \rangle]_+^{q-1} \|\mathbf{v}\|_2^2 \\ &= \sum_{i \in S_1} y_i \widehat{y}_i c_i^{(t)} [\langle \mathbf{w}_{j^\ddagger}^{(t)}, \mathbf{v} \rangle]_+^{q-1} \|\mathbf{v}\|_2^2 \\ &= \left(\sum_{i \in S_1} y_i \widehat{y}_i c_i^{(t)} \right) \cdot \|\mathbf{v}\|_2^2 \cdot (\bar{\Lambda}_1^{(t)})^{q-1} \\ &= n_1 \cdot \left(p - \frac{1}{2} \pm o(1) \right) \cdot \|\mathbf{v}\|_2^2 \cdot (\bar{\Lambda}_1^{(t)})^{q-1}, \end{aligned} \tag{C.18}$$

684 For \star , similarly we have

$$\underbrace{\sum_{i=1}^{n_1} b_i^{(t)} [\langle \mathbf{w}_{j^\ddagger}^{(t)}, y'_i \cdot \mathbf{v} \rangle]_+^{q-1} \|\mathbf{v}\|_2^2}_{\star} = \sum_{i \in S'_1} b_i^{(t)} [\langle \mathbf{w}_{j^\ddagger}^{(t)}, \mathbf{v} \rangle]_+^{q-1} \|\mathbf{v}\|_2^2 + \sum_{i \in S'_{-1}} b_i^{(t)} [-\langle \mathbf{w}_{j^\ddagger}^{(t)}, \mathbf{v} \rangle]_+^{q-1} \|\mathbf{v}\|_2^2$$

$$\begin{aligned}
&= \sum_{i \in S'_1} b_i^{(t)} [\langle \mathbf{w}_{j^*}^{(t)}, \mathbf{v} \rangle]_+^{q-1} \|\mathbf{v}\|_2^2 \\
&= \left(\sum_{i \in S'_1} b_i^{(t)} \right) \cdot \|\mathbf{v}\|_2^2 \cdot (\bar{\Lambda}_1^{(t)})^{q-1} \\
&= n'_1 \cdot \left(\frac{1}{2} \pm o(1) \right) \cdot \|\mathbf{v}\|_2^2 \cdot (\bar{\Lambda}_1^{(t)})^{q-1}. \tag{C.19}
\end{aligned}$$

685 According to Lemma C.13, (C.18) and (C.19), we have $n'_1 = o(n_1)$ with high probability, therefore

$$\clubsuit + \star = n_1 \cdot \left(p - \frac{1}{2} \pm o(1) \right) \cdot \|\mathbf{v}\|_2^2 \cdot (\bar{\Lambda}_1^{(t)})^{q-1},$$

686 leading to

$$\begin{aligned}
\bar{\Lambda}_1^{(t+1)} &\geq (1 - \eta\lambda) \cdot \langle \mathbf{w}_{j^*}^{(t)}, \mathbf{v} \rangle - \frac{q\eta n_1}{n_1 + n_u} \cdot \left(p - \frac{1}{2} \pm o(1) \right) \cdot \|\mathbf{v}\|_2^2 \cdot (\bar{\Lambda}_1^{(t)})^{q-1} \\
&= (1 - \eta\lambda) \cdot \bar{\Lambda}_1^{(t)} + \eta \cdot \left(\frac{1}{2} - p \right) \cdot \Theta(d) \cdot (\bar{\Lambda}_1^{(t)})^{q-1}.
\end{aligned}$$

687 And we can prove in a similar way that

$$\bar{\Lambda}_1^{(t+1)} \geq (1 - \eta\lambda) \cdot \bar{\Lambda}_1^{(t)} + \eta \cdot \left(\frac{1}{2} - p \right) \cdot \Theta(d) \cdot (\bar{\Lambda}_1^{(t)})^{q-1}.$$

688

□

689 In this case ($p < 1/2$), given a small amount of labeled data, downstream task parameter \mathbf{a} will learn
690 the negative direction and the main theorems still hold.

691 C.5.2 Uniform upper bound for $\Gamma^{(t)}$

692 The following lemma provides an upper bound for the increasing rate of $\Gamma^{(t)}$.

693 **Lemma C.17.** For $\Gamma_i^{(t)} := \max_{j \in [2m]} \langle \mathbf{w}_j, \boldsymbol{\xi}_i \rangle, i \in [n_u], \Gamma'_i{}^{(t)} := \max_{j \in [2m]} \langle \mathbf{w}_j, \boldsymbol{\xi}'_i \rangle, i \in [n_l],$
694 $\Gamma^{(t)} := \max\{\max_{i \in [n_u]} \Gamma_i^{(t)}, \max_{i \in [n_l]} \Gamma'_i{}^{(t)}\}$, we have with high probability that

$$\begin{aligned}
\Gamma_i^{(t+1)} &\leq (1 - \eta\lambda) \cdot \Gamma_i^{(t)} + \eta \cdot \max \left\{ \tilde{\Theta}(d^{\frac{1}{2}+2\epsilon}), \tilde{\Theta}\left(\frac{d^{1+2\epsilon}}{n_u}\right) \right\} \cdot (\Gamma^{(t)})^{q-1}, i \in [n_l], \\
\Gamma'_i{}^{(t+1)} &\leq (1 - \eta\lambda) \cdot \Gamma'_i{}^{(t)} + \eta \cdot \max \left\{ \tilde{\Theta}(d^{\frac{1}{2}+2\epsilon}), \tilde{\Theta}\left(\frac{d^{1+2\epsilon}}{n_u}\right) \right\} \cdot (\Gamma^{(t)})^{q-1}, i \in [n_l],
\end{aligned}$$

696 and

$$\Gamma^{(t+1)} \leq (1 - \eta\lambda) \cdot \Gamma^{(t)} + \eta \cdot \max \left\{ \tilde{\Theta}(d^{\frac{1}{2}+2\epsilon}), \tilde{\Theta}\left(\frac{d^{1+2\epsilon}}{n_u}\right) \right\} \cdot (\Gamma^{(t)})^{q-1},$$

697 where $\epsilon < 1/8$.

698 *Proof of Lemma C.17.* We first prove the former inequality. Let $j^* = \arg \max_{1 \leq j \leq 2m} \langle \mathbf{w}_j^{(t+1)}, \boldsymbol{\xi}_l \rangle,$
699 where $l \in [n_u]$ is fixed. According to Lemma C.8, we have

$$\begin{aligned}
\Gamma_l^{(t+1)} &= \langle \mathbf{w}_{j^*}^{(t+1)}, \boldsymbol{\xi}_l \rangle \\
&= (1 - \eta\lambda) \cdot \langle \mathbf{w}_{j^*}^{(t)}, \boldsymbol{\xi}_l \rangle + \frac{q\eta u_{j^*}}{n_1 + n_u} \left(\sum_{i=1}^{n_u} \hat{y}_i c_i^{(t)} [\langle \mathbf{w}_{j^*}^{(t)}, \boldsymbol{\xi}_i \rangle]_+^{q-1} \langle \boldsymbol{\xi}_i, \boldsymbol{\xi}_l \rangle + \sum_{i=1}^{n_l} y'_i b_i^{(t)} [\langle \mathbf{w}_{j^*}^{(t)}, \boldsymbol{\xi}'_i \rangle]_+^{q-1} \langle \boldsymbol{\xi}'_i, \boldsymbol{\xi}_l \rangle \right) \\
&\leq (1 - \eta\lambda) \cdot \langle \mathbf{w}_{j^*}^{(t)}, \boldsymbol{\xi}_l \rangle + \frac{q\eta}{n_1 + n_u} \left(\underbrace{\sum_{i=1}^{n_u} c_i^{(t)} [\langle \mathbf{w}_{j^*}^{(t)}, \boldsymbol{\xi}_i \rangle]_+^{q-1} |\langle \boldsymbol{\xi}_i, \boldsymbol{\xi}_l \rangle|}_{\clubsuit} + \underbrace{\sum_{i=1}^{n_l} b_i^{(t)} [\langle \mathbf{w}_{j^*}^{(t)}, \boldsymbol{\xi}'_i \rangle]_+^{q-1} |\langle \boldsymbol{\xi}'_i, \boldsymbol{\xi}_l \rangle|}_{\star} \right), \tag{C.20}
\end{aligned}$$

700 where the last inequality is due to triangle inequality.

701 For \clubsuit , note that $l \in [n_u]$ and there exists an $i \in [n_u]$ equivalent to l , it follows that

$$\begin{aligned}
& \underbrace{\sum_{i=1}^{n_u} c_i^{(t)} [\langle \mathbf{w}_{j^*}^{(t)}, \boldsymbol{\xi}_i \rangle_+^{q-1} |\langle \boldsymbol{\xi}_i, \boldsymbol{\xi}_l \rangle|]}_{\clubsuit} \\
&= \sum_{i \in [n_u], i \neq l} c_i^{(t)} [\langle \mathbf{w}_{j^*}^{(t)}, \boldsymbol{\xi}_i \rangle_+^{q-1} |\langle \boldsymbol{\xi}_i, \boldsymbol{\xi}_l \rangle| + c_l^{(t)} [\langle \mathbf{w}_{j^*}^{(t)}, \boldsymbol{\xi}_l \rangle_+^{q-1} \|\boldsymbol{\xi}_l\|_2^2] \\
&\leq (n_u - 1) \cdot \left(\frac{1}{2} + o(1)\right) \cdot \tilde{\Theta}(d^{\frac{1}{2}+2\epsilon}) \cdot (\Gamma^{(t)})^{q-1} + \left(\frac{1}{2} + o(1)\right) \cdot \tilde{\Theta}(d^{1+2\epsilon}) \cdot (\Gamma^{(t)})^{q-1} \\
&= (n_u - 1) \cdot \tilde{\Theta}(d^{\frac{1}{2}+2\epsilon}) \cdot (\Gamma^{(t)})^{q-1} + \tilde{\Theta}(d^{1+2\epsilon}) \cdot (\Gamma^{(t)})^{q-1},
\end{aligned} \tag{C.21}$$

702 where the inequality is due to Lemma C.11, $\|\boldsymbol{\xi}_l\|_2^2 = \tilde{\Theta}(d\sigma_p^2) = \tilde{\Theta}(d^{1+2\epsilon})$, $|\langle \boldsymbol{\xi}_i, \boldsymbol{\xi}_l \rangle| = \tilde{\Theta}(d^{\frac{1}{2}}\sigma_p^2) =$
703 $\tilde{\Theta}(d^{\frac{1}{2}+2\epsilon})$ according to Lemma E.3 and the definition of $\Gamma^{(t)}$.

704 For \star , we have

$$\underbrace{\sum_{i=1}^{n_1} b_i^{(t)} [\langle \mathbf{w}_{j^*}^{(t)}, \boldsymbol{\xi}'_i \rangle_+^{q-1} |\langle \boldsymbol{\xi}'_i, \boldsymbol{\xi}_l \rangle|]}_{\star} \leq n_1 \cdot \left(\frac{1}{2} + o(1)\right) \cdot \tilde{\Theta}(d^{\frac{1}{2}+2\epsilon}) \cdot (\Gamma^{(t)})^{q-1} = n_1 \cdot \tilde{\Theta}(d^{\frac{1}{2}+2\epsilon}) \cdot (\Gamma^{(t)})^{q-1}, \tag{C.22}$$

705 Plugging (C.21) and (C.22) into (C.20), we have

$$\begin{aligned}
\Gamma_l^{(t+1)} &\leq (1 - \eta\lambda) \cdot \Gamma_l^{(t)} + \eta \cdot \left(\frac{q}{n_1 + n_u} \cdot \left((n_u + n_1 - 1) \cdot \tilde{\Theta}(d^{\frac{1}{2}+2\epsilon}) + \tilde{\Theta}(d^{1+2\epsilon}) \right) \right) \cdot (\Gamma^{(t)})^{q-1} \\
&\leq (1 - \eta\lambda) \cdot \Gamma_l^{(t)} + \eta \cdot \max \left\{ \tilde{\Theta}(d^{\frac{1}{2}+2\epsilon}), \tilde{\Theta}\left(\frac{d^{1+2\epsilon}}{n_u}\right) \right\} \cdot (\Gamma^{(t)})^{q-1},
\end{aligned}$$

706 which is the first part of this lemma.

707 Let $j^* = \operatorname{argmax}_{1 \leq j \leq 2m} \langle \mathbf{w}_j^{(t+1)}, \boldsymbol{\xi}'_l \rangle$, where $l \in [n_1]$ is fixed. According to Lemma C.8, we have

$$\begin{aligned}
\Gamma_l^{(t+1)} &= \langle \mathbf{w}_{j^*}^{(t+1)}, \boldsymbol{\xi}'_l \rangle \\
&= (1 - \eta\lambda) \cdot \langle \mathbf{w}_{j^*}^{(t)}, \boldsymbol{\xi}'_l \rangle + \frac{q\eta u_{j^*}}{n_1 + n_u} \left(\sum_{i=1}^{n_u} \hat{y}_i c_i^{(t)} [\langle \mathbf{w}_{j^*}^{(t)}, \boldsymbol{\xi}_i \rangle_+^{q-1} \langle \boldsymbol{\xi}_i, \boldsymbol{\xi}'_l \rangle] + \sum_{i=1}^{n_1} y'_i b_i^{(t)} [\langle \mathbf{w}_{j^*}^{(t)}, \boldsymbol{\xi}'_i \rangle_+^{q-1} \langle \boldsymbol{\xi}'_i, \boldsymbol{\xi}'_l \rangle] \right) \\
&\leq (1 - \eta\lambda) \cdot \langle \mathbf{w}_{j^*}^{(t)}, \boldsymbol{\xi}'_l \rangle + \frac{q\eta}{n_1 + n_u} \left(\underbrace{\sum_{i=1}^{n_u} c_i^{(t)} [\langle \mathbf{w}_{j^*}^{(t)}, \boldsymbol{\xi}_i \rangle_+^{q-1} |\langle \boldsymbol{\xi}_i, \boldsymbol{\xi}'_l \rangle|]}_{\clubsuit} + \underbrace{\sum_{i=1}^{n_1} b_i^{(t)} [\langle \mathbf{w}_{j^*}^{(t)}, \boldsymbol{\xi}'_i \rangle_+^{q-1} |\langle \boldsymbol{\xi}'_i, \boldsymbol{\xi}'_l \rangle|]}_{\star} \right),
\end{aligned} \tag{C.23}$$

708 For \clubsuit , we have

$$\underbrace{\sum_{i=1}^{n_u} c_i^{(t)} [\langle \mathbf{w}_{j^*}^{(t)}, \boldsymbol{\xi}_i \rangle_+^{q-1} |\langle \boldsymbol{\xi}_i, \boldsymbol{\xi}'_l \rangle|]}_{\clubsuit} \leq \sum_{i=1}^{n_u} \left(\frac{1}{2} \pm o(1)\right) \cdot \tilde{\Theta}(d^{\frac{1}{2}+2\epsilon}) \cdot (\Gamma^{(t)})^{q-1} = n_u \cdot \tilde{\Theta}(d^{\frac{1}{2}+2\epsilon}) \cdot (\Gamma^{(t)})^{q-1}, \tag{C.24}$$

709 where the inequality is due to Lemma C.11, $|\langle \boldsymbol{\xi}_i, \boldsymbol{\xi}'_l \rangle| = \tilde{\Theta}(d^{\frac{1}{2}}\sigma_p^2) = \tilde{\Theta}(d^{\frac{1}{2}+2\epsilon})$ and the definition of
710 $\Gamma^{(t)}$.

711 For \star , note that $l \in [n_l]$ and there exists an $i \in [n_l]$ equivalent to l , it follows that

$$\begin{aligned}
& \underbrace{\sum_{i=1}^{n_l} b_i^{(t)} [\langle \mathbf{w}_{j^*}^{(t)}, \boldsymbol{\xi}'_i \rangle_+^{q-1} |\langle \boldsymbol{\xi}'_i, \boldsymbol{\xi}'_l \rangle|]}_{\star} \\
&= \sum_{i \in [n_l], i \neq l} b_i^{(t)} [\langle \mathbf{w}_{j^*}^{(t)}, \boldsymbol{\xi}'_i \rangle_+^{q-1} |\langle \boldsymbol{\xi}'_i, \boldsymbol{\xi}'_l \rangle| + b_l^{(t)} [\langle \mathbf{w}_{j^*}^{(t)}, \boldsymbol{\xi}'_l \rangle_+^{q-1} \|\boldsymbol{\xi}'_l\|_2^2] \\
&\leq (n_l - 1) \cdot \left(\frac{1}{2} + o(1)\right) \cdot \tilde{\Theta}(d^{\frac{1}{2}+2\epsilon}) \cdot (\Gamma^{(t)})^{q-1} + \left(\frac{1}{2} + o(1)\right) \cdot \tilde{\Theta}(d^{1+2\epsilon}) \cdot (\Gamma^{(t)})^{q-1} \\
&= (n_l - 1) \cdot \tilde{\Theta}(d^{\frac{1}{2}+2\epsilon}) + \tilde{\Theta}(d^{1+2\epsilon}) \cdot (\Gamma^{(t)})^{q-1}
\end{aligned} \tag{C.25}$$

712 Plugging (C.24) and (C.25) into (C.23), we have

$$\begin{aligned}
\Gamma_l^{(t+1)} &\leq (1 - \eta\lambda) \cdot \Gamma_l^{(t+1)} + \eta \cdot \left(\frac{q}{n_l + n_u} \cdot ((n_u + n_l - 1) \cdot \tilde{\Theta}(d^{\frac{1}{2}+2\epsilon}) + \tilde{\Theta}(d^{1+2\epsilon})) \right) \cdot (\Gamma^{(t)})^{q-1} \\
&\leq (1 - \eta\lambda) \cdot \Gamma_l^{(t+1)} + \eta \cdot \max \left\{ \tilde{\Theta}(d^{\frac{1}{2}+2\epsilon}), \tilde{\Theta}\left(\frac{d^{1+2\epsilon}}{n_u}\right) \right\} \cdot (\Gamma^{(t)})^{q-1},
\end{aligned}$$

713 which verifies the second inequality in this lemma.

714 Note that $\Gamma^{(t)} = \max\{\max_{l \in [n_u]} \Gamma_l^{(t)}, \max_{l \in [n_l]} \Gamma_l^{(t)}\}$, without loss of generality, we assume $\Gamma^{(t)} =$

715 $\max_{l \in [n_u]} \Gamma_l^{(t)}$ and assume $l^* = \operatorname{argmax}_{l \in [n_u]} \Gamma_l^{(t+1)}$, we have

$$\begin{aligned}
\Gamma^{(t+1)} = \Gamma_{l^*}^{(t+1)} &\leq (1 - \eta\lambda) \cdot \Gamma_{l^*}^{(t)} + \eta \cdot \max \left\{ \tilde{\Theta}(d^{\frac{1}{2}+2\epsilon}), \tilde{\Theta}\left(\frac{d^{1+2\epsilon}}{n_u}\right) \right\} \cdot (\Gamma^{(t)})^{q-1} \\
&\leq (1 - \eta\lambda) \cdot \Gamma^{(t)} + \eta \cdot \max \left\{ \tilde{\Theta}(d^{\frac{1}{2}+2\epsilon}), \tilde{\Theta}\left(\frac{d^{1+2\epsilon}}{n_u}\right) \right\} \cdot (\Gamma^{(t)})^{q-1},
\end{aligned}$$

716 which verifies the third inequality in this lemma.

717 □

718 C.5.3 Tensor Power Method: Proving $\Gamma^{(t)} = O(\Gamma^{(0)})$ during $[0, T_r]$ and computing the 719 magnitude of T_r

720 In this section, we first show that off-diagonal correlation ($\bar{\Lambda}_r^{(t)}$ for $p > 1/2$ and $\hat{\Lambda}_r^{(t)}$ for $p < 1/2$)
721 remains initialization magnitude during $[0, T_r]$. If the accuracy of pseudo-labeler $p > 1/2$, we have

722 off-diagonal correlation $\bar{\Lambda}_r^{(t+1)} \leq (1 - \eta\lambda) \cdot \bar{\Lambda}_r^{(t)}$ for $r \in \{\pm 1\}$, therefore, $\bar{\Lambda}_r^{(t)} = O(\bar{\Lambda}_r^{(0)}) = \tilde{O}(d^{-\frac{1}{4}})$.

723 If $p < 1/2$, we have off-diagonal correlation $\hat{\Lambda}_r^{(t+1)} \leq (1 - \eta\lambda) \cdot \hat{\Lambda}_r^{(t)}$ for $r \in \{\pm 1\}$, therefore,
724 $\hat{\Lambda}_r^{(t)} = O(\hat{\Lambda}_r^{(0)}) = \tilde{O}(d^{-\frac{1}{4}})$. In this paper, we mainly focus on $p > 1/2$.

725 According to Sections C.5.1 and C.5.2, we have obtained following upper bounds and lower bounds
726 for *feature learning* term $\hat{\Lambda}_r^{(t)}, \bar{\Lambda}_r^{(t)}, r \in \{\pm 1\}$ and *noise memorization* term $\Gamma^{(t)}$: When $t \in [0, T_r]$,
727 we have

$$\begin{aligned}
\hat{\Lambda}_r^{(t+1)} &\geq \hat{\Lambda}_r^{(t)} + \eta \cdot (2p - 1) \cdot \Theta(d) \cdot (\hat{\Lambda}_r^{(t)})^{q-1} \text{ and } \bar{\Lambda}_r^{(t+1)} \leq (1 - \eta\lambda) \cdot \bar{\Lambda}_r^{(t)}, \text{ for } r \in \{\pm 1\}; \\
\Gamma^{(t+1)} &\leq (1 - \eta\lambda) \cdot \Gamma^{(t)} + \eta \cdot \max \left\{ \tilde{\Theta}(d^{\frac{1}{2}+2\epsilon}), \tilde{\Theta}\left(\frac{d^{1+2\epsilon}}{n_u}\right) \right\} \cdot (\Gamma^{(t)})^{q-1}.
\end{aligned} \tag{C.26}$$

728 According to Condition 3.1, assume $n_u = \Omega(d^{4\epsilon})$ and note that $\epsilon < 1/8$, we have

$$\max \left\{ \tilde{\Theta}(d^{\frac{1}{2}+2\epsilon}), \tilde{\Theta}\left(\frac{d^{1+2\epsilon}}{n_u}\right) \right\} = \max \left\{ \tilde{\Theta}(d^{\frac{1}{2}+2\epsilon}), \tilde{O}(d^{1-2\epsilon}) \right\} = \tilde{O}(d^{1-2\epsilon}),$$

729 leading to

$$\Gamma^{(t+1)} \leq (1 - \eta\lambda) \cdot \Gamma^{(t)} + \eta \cdot \tilde{\Theta}(d^{1-2\epsilon}) \cdot (\Gamma^{(t)})^{q-1}.$$

730 By leveraging tensor power method introduced in Lemma E.4, we can prove following lemma about
 731 the magnitude of $\Gamma^{(t)}$:

732 **Lemma C.18.** $\Gamma^{(t)}$ remains initialization magnitude during $[0, \max_{r \in \{\pm 1\}} \{T_r\}]$.

733 *Proof of Lemma C.18.* Let T_r^* be the first iteration t in which $\widehat{\Lambda}_r^{(t)} \geq A$ for $r \in \{\pm 1\}$, let T^* be the
 734 first iteration t in which $\Gamma^{(t)} \geq A'$, then according to Lemma E.4, we know

$$\sum_{t \geq 0, x_t \leq A} \eta \leq \frac{\delta}{(1 - (1 + \delta)^{-(q-2)})x_0 C_1} + \eta \cdot \frac{C_2}{C_1} (1 + \delta)^{q-1} \left(1 + \frac{\log(A/x_0)}{\log(1 + \delta)}\right),$$

735

$$\sum_{t \geq 0, x_t \leq A} \eta \geq \frac{\delta(1 - (x_0/A)^{q-2})}{(1 + \delta)^{q-1}(1 - (1 + \delta)^{-(q-2)})x_0 C_2} - \eta \cdot (1 + \delta)^{-(q-1)} \left(1 + \frac{\log(A/x_0)}{\log(1 + \delta)}\right).$$

736 And it follows that

$$\eta \cdot T_r^* \leq \frac{\delta}{(1 - (1 + \delta)^{-(q-2)})\widehat{\Lambda}_r^{(0)} C_1} + \eta \cdot \frac{C_2}{C_1} (1 + \delta)^{q-1} \left(1 + \frac{\log(A/\widehat{\Lambda}_r^{(0)})}{\log(1 + \delta)}\right),$$

737

$$\eta \cdot T^* \geq \frac{\delta'(1 - (x_0/A')^{q-2})}{(1 + \delta)^{q-1}(1 - (1 + \delta)^{-(q-2)})\Gamma^{(0)} C_2} - \eta \cdot (1 + \delta')^{-(q-1)} \left(1 + \frac{\log(A'/\Gamma^{(0)})}{\log(1 + \delta')}\right),$$

738 where $C_1, C_2 = (2p - 1) \cdot \widetilde{\Theta}(d)$ and $C_1', C_2' = \widetilde{\Theta}(d^{1-2\epsilon})$ according to (C.26).

739 Taking $A = \Theta(1/m)$, $A' = C \cdot \Gamma^{(t)}$ where C is a large constant and $C = \Theta(1)$, $\delta = \delta' = \frac{1}{2}$ and

740 note that $\widehat{\Lambda}_r^{(0)} = \widetilde{\Theta}(\sigma_0 d^{\frac{1}{2}}) = \widetilde{\Theta}(d^{-\frac{1}{4}})$, $\Gamma^{(0)} = \widetilde{\Theta}(\sigma_0 \sigma_p d^{\frac{1}{2}}) = \widetilde{\Theta}(d^{-\frac{1}{4} + \epsilon})$, we have

$$\eta \cdot T_r^* \leq \widetilde{\Theta}(d^{-\frac{3}{4}}) + \eta \cdot \widetilde{\Theta}(1) = \widetilde{\Theta}(d^{-\frac{3}{4}}), \quad (\text{C.27})$$

741 and

$$\eta \cdot T^* \geq \widetilde{\Theta}(d^{-\frac{3}{4} + \epsilon}) - \eta \cdot \widetilde{\Theta}(1) = \widetilde{\Theta}(d^{-\frac{3}{4} + \epsilon}). \quad (\text{C.28})$$

742 Therefore, combining (C.27) and (C.28), we have $\eta \cdot T^* \geq \widetilde{\Theta}(d^{-\frac{3}{4} + \epsilon}) > \widetilde{\Theta}(d^{-\frac{3}{4}}) \geq \eta \cdot T_r^*$, leading
 743 to $T^* > T_r^*$ for both $r \in \{-1, +1\}$. This indicates that when $\widehat{\Lambda}_1^{(t)}, \widehat{\Lambda}_{-1}^{(t)}$ reach $\Theta(1/m)$, $\Gamma^{(t)}$ remain
 744 the same magnitude as initialization. \square

745 By leveraging tensor power method, we can also estimate the length of Stage I, i.e. T_1, T_{-1} , by
 746 applying tensor power method. To use tensor power method, we need to upper-bound the increasing
 747 speed of $\widehat{\Lambda}_r^{(t)}$. We have the following lemma:

748 **Lemma C.19.** For $r \in \{\pm 1\}$, we have with high probability that

$$\widehat{\Lambda}_r^{(t+1)} \geq (1 - \eta\lambda) \cdot \widehat{\Lambda}_r^{(t)} + \eta \cdot q \left(p - \frac{1}{2} - o(1)\right) \cdot \|\mathbf{v}\|_2^2 \cdot (\widehat{\Lambda}_r^{(t)})^{q-1},$$

749

$$\widehat{\Lambda}_r^{(t+1)} \leq (1 - \eta\lambda) \cdot \widehat{\Lambda}_r^{(t)} + \eta \cdot q \left(p - \frac{1}{2} + o(1)\right) \cdot \|\mathbf{v}\|_2^2 \cdot (\widehat{\Lambda}_r^{(t)})^{q-1}.$$

750 *Proof of Lemma C.19.* Let $j^* = \arg \max_{1 \leq j \leq m} \langle \mathbf{w}_j^{(t+1)}, \mathbf{v} \rangle$ and note that $u_{j^*} = \mathbf{1}_{[1 \leq j \leq m]} =$
 751 $\mathbf{1}_{[m+1 \leq j \leq 2m]} = 1$, then we have

$$\begin{aligned} \widehat{\Lambda}_1^{(t+1)} &= \langle \mathbf{w}_{j^*}^{(t+1)}, \mathbf{v} \rangle \\ &= (1 - \eta\lambda) \cdot \langle \mathbf{w}_{j^*}^{(t)}, \mathbf{v} \rangle + \frac{q\eta}{n_l + n_u} \left(\underbrace{\sum_{i \in S_1} y_i \widehat{y}_i c_i^{(t)} [\langle \mathbf{w}_{j^*}^{(t)}, \mathbf{v} \rangle]_+^{q-1} \|\mathbf{v}\|_2^2 + \sum_{i \in S_{-1}} y_i \widehat{y}_i c_i^{(t)} [-\langle \mathbf{w}_{j^*}^{(t)}, \mathbf{v} \rangle]_+^2 \|\mathbf{v}\|_2^{q-1}}_{\clubsuit} \right) \\ &\quad + \frac{q\eta}{n_l + n_u} \left(\underbrace{\sum_{i \in S'_1} b_i^{(t)} [\langle \mathbf{w}_{j^*}^{(t)}, \mathbf{v} \rangle]_+^{q-1} \|\mathbf{v}\|_2^2 + \sum_{i \in S'_1} b_i^{(t)} [-\langle \mathbf{w}_{j^*}^{(t)}, \mathbf{v} \rangle]_+^{q-1} \|\mathbf{v}\|_2^2}_{\star} \right). \end{aligned} \quad (\text{C.29})$$

752 For ♣, according to Lemma C.12, we have

$$\begin{aligned}
& \underbrace{\sum_{i \in S_1} y_i \widehat{y}_i c_i^{(t)} [\langle \mathbf{w}_{j^*}^{(t)}, \mathbf{v} \rangle]_+^{q-1} \|\mathbf{v}\|_2^2 + \sum_{i \in S_{-1}} y_i \widehat{y}_i c_i^{(t)} [-\langle \mathbf{w}_{j^*}^{(t)}, \mathbf{v} \rangle]_+^{q-1} \|\mathbf{v}\|_2^2}_{\clubsuit} \\
&= n_1 \cdot \left(p - \frac{1}{2} \pm o(1)\right) \cdot [\langle \mathbf{w}_{j^*}^{(t)}, \mathbf{v} \rangle]_+^{q-1} \|\mathbf{v}\|_2^2 + n_{-1} \cdot \left(p - \frac{1}{2} \pm o(1)\right) \cdot [-\langle \mathbf{w}_{j^*}^{(t)}, \mathbf{v} \rangle]_+^{q-1} \|\mathbf{v}\|_2^2 \\
&\leq n_1 \cdot \left(p - \frac{1}{2} \pm o(1)\right) \cdot \|\mathbf{v}\|_2^2 \cdot (\widehat{\Lambda}_1^{(t)})^{q-1} + n_{-1} \cdot \left(p - \frac{1}{2} \pm o(1)\right) \cdot \|\mathbf{v}\|_2^2 \cdot (\bar{\Lambda}_{-1}^{(t)})^{q-1} \\
&= n_1 \cdot \left(p - \frac{1}{2} \pm o(1)\right) \cdot \|\mathbf{v}\|_2^2 \cdot (\widehat{\Lambda}_1^{(t)})^{q-1},
\end{aligned} \tag{C.30}$$

753 where the last equality is due to $\widehat{\Lambda}_1^{(t)} = \omega(\bar{\Lambda}_{-1}^{(t)})$.

754 For ★, according to Lemma C.11, we have

$$\begin{aligned}
& \underbrace{\sum_{i \in S'_1} b_i^{(t)} [\langle \mathbf{w}_{j^*}^{(t)}, \mathbf{v} \rangle]_+^{q-1} \|\mathbf{v}\|_2^2 + \sum_{i \in S'_{-1}} b_i^{(t)} [-\langle \mathbf{w}_{j^*}^{(t)}, \mathbf{v} \rangle]_+^{q-1} \|\mathbf{v}\|_2^2}_{\star} \\
&= n'_1 \cdot \left(\frac{1}{2} \pm o(1)\right) \cdot [\langle \mathbf{w}_{j^*}^{(t)}, \mathbf{v} \rangle]_+^{q-1} \|\mathbf{v}\|_2^2 + n'_{-1} \cdot \left(\frac{1}{2} \pm o(1)\right) \cdot [-\langle \mathbf{w}_{j^*}^{(t)}, \mathbf{v} \rangle]_+^{q-1} \|\mathbf{v}\|_2^2 \tag{C.31} \\
&\leq n'_1 \cdot \left(\frac{1}{2} \pm o(1)\right) \cdot \|\mathbf{v}\|_2^2 \cdot (\widehat{\Lambda}_1^{(t)})^{q-1} + n'_{-1} \cdot \left(\frac{1}{2} \pm o(1)\right) \cdot \|\mathbf{v}\|_2^2 \cdot (\bar{\Lambda}_{-1}^{(t)})^{q-1} \\
&= n'_1 \cdot \left(\frac{1}{2} \pm o(1)\right) \cdot \|\mathbf{v}\|_2^2 \cdot (\widehat{\Lambda}_1^{(t)})^{q-1},
\end{aligned}$$

755 where the last equality is due to $\widehat{\Lambda}_1^{(t)} = \omega(\bar{\Lambda}_{-1}^{(t)})$.

756 Plugging (C.30) and (C.31) into (C.29), we have

$$\begin{aligned}
& \widehat{\Lambda}_1^{(t+1)} \\
&\leq (1 - \eta\lambda) \cdot \widehat{\Lambda}_1^{(t)} + \frac{q\eta}{n_1 + n_u} \left(n_1 \cdot \left(p - \frac{1}{2} \pm o(1)\right) \cdot \|\mathbf{v}\|_2^2 \cdot (\widehat{\Lambda}_1^{(t)})^{q-1} + n'_1 \cdot \left(\frac{1}{2} \pm o(1)\right) \cdot \|\mathbf{v}\|_2^2 \cdot (\widehat{\Lambda}_1^{(t)})^{q-1} \right) \\
&= (1 - \eta\lambda) \cdot \widehat{\Lambda}_1^{(t)} + \frac{q\eta n_1}{n_1 + n_u} \cdot \left(p - \frac{1}{2} \pm o(1)\right) \cdot \|\mathbf{v}\|_2^2 \cdot (\widehat{\Lambda}_1^{(t)})^{q-1} + \frac{q\eta n'_1}{n_1 + n_u} \cdot \left(\frac{1}{2} \pm o(1)\right) \cdot \|\mathbf{v}\|_2^2 \cdot (\widehat{\Lambda}_1^{(t)})^{q-1} \\
&= (1 - \eta\lambda) \cdot \widehat{\Lambda}_1^{(t)} + q\eta \cdot \left(\frac{n_1}{n_1 + n_u} \cdot \left(p - \frac{1}{2} \pm o(1)\right) + \frac{n'_1}{n_1 + n_u} \cdot \left(\frac{1}{2} \pm o(1)\right) \right) \cdot \|\mathbf{v}\|_2^2 \cdot (\widehat{\Lambda}_1^{(t)})^{q-1} \\
&= (1 - \eta\lambda) \cdot \widehat{\Lambda}_1^{(t)} + q\eta \cdot \underbrace{\left(\frac{n_1}{n_1 + n_u} \cdot \left(p - \frac{1}{2}\right) + \frac{n'_1}{n_1 + n_u} \cdot \frac{1}{2} \pm o(1) \right)}_{\spadesuit} \cdot \|\mathbf{v}\|_2^2 \cdot (\widehat{\Lambda}_1^{(t)})^{q-1}.
\end{aligned} \tag{C.32}$$

757 Note that we have already proved in (C.10) that

$$\widehat{\Lambda}_1^{(t+1)} \leq (1 - \eta\lambda) \cdot \widehat{\Lambda}_1^{(t)} + q\eta \cdot \underbrace{\left(\frac{n_1}{n_1 + n_u} \cdot \left(p - \frac{1}{2}\right) + \frac{n'_1}{n_1 + n_u} \cdot \frac{1}{2} \pm o(1) \right)}_{\spadesuit} \cdot \|\mathbf{v}\|_2^2 \cdot (\widehat{\Lambda}_1^{(t)})^{q-1}. \tag{C.33}$$

758 Note we have already prove in (C.11) that

$$\underbrace{\frac{n_1}{n_1 + n_u} \cdot \left(p - \frac{1}{2}\right) + \frac{n'_1}{n_1 + n_u} \cdot \frac{1}{2}}_{\spadesuit} = \frac{1}{2} \cdot \left(p - \frac{1}{2}\right) \pm o(1)$$

759 Therefore, we have

$$\widehat{\Lambda}_1^{(t+1)} \geq (1 - \eta\lambda) \cdot \widehat{\Lambda}_1^{(t)} + q\eta \cdot \left(p - \frac{1}{2} - o(1)\right) \cdot \|\mathbf{v}\|_2^2 \cdot (\widehat{\Lambda}_1^{(t)})^{q-1},$$

760

$$\widehat{\Lambda}_1^{(t+1)} \leq (1 - \eta\lambda) \cdot \widehat{\Lambda}_1^{(t)} + q\eta \cdot \left(p - \frac{1}{2} + o(1)\right) \cdot \|\mathbf{v}\|_2^2 \cdot (\widehat{\Lambda}_1^{(t)})^{q-1}.$$

761 In a similar way, we can prove that

$$\widehat{\Lambda}_{-1}^{(t+1)} \geq (1 - \eta\lambda) \cdot \widehat{\Lambda}_{-1}^{(t)} + q\eta \cdot \left(p - \frac{1}{2} - o(1)\right) \cdot \|\mathbf{v}\|_2^2 \cdot (\widehat{\Lambda}_{-1}^{(t)})^{q-1},$$

762

$$\widehat{\Lambda}_{-1}^{(t+1)} \leq (1 - \eta\lambda) \cdot \widehat{\Lambda}_{-1}^{(t)} + q\eta \cdot \left(p - \frac{1}{2} + o(1)\right) \cdot \|\mathbf{v}\|_2^2 \cdot (\widehat{\Lambda}_{-1}^{(t)})^{q-1},$$

763 which completes the proof of this lemma. \square

764 **Lemma C.20** (Length of pre-training). For $r \in \{\pm 1\}$, let T_r be the first iteration that $\widehat{\Lambda}_r^{(t)}$ reaches
765 $\Theta(1/m)$ respectively. Then $T_r = \widetilde{\Theta}(d^{\frac{q}{4} - \frac{3}{2}})/\eta$ for all $r \in \{\pm 1\}$.

766 *Proof of Lemma C.20.* By leveraging tensor power method given in Lemma E.4,

$$\sum_{t \geq 0, x_t \leq A} \eta \leq \frac{\delta}{(1 - (1 + \delta)^{-(q-2)})x_0^{q-2}C_1} + \eta \cdot \frac{C_2}{C_1}(1 + \delta)^{q-1} \left(1 + \frac{\log(A/x_0)}{\log(1 + \delta)}\right),$$

767

$$\sum_{t \geq 0, x_t \leq A} \eta \geq \frac{\delta(1 - (x_0/A)^{q-2})}{(1 + \delta)^{q-1}(1 - (1 + \delta)^{-(q-2)})x_0^{q-2}C_2} - \eta \cdot (1 + \delta)^{-(q-1)} \left(1 + \frac{\log(A/x_0)}{\log(1 + \delta)}\right),$$

768 we have for $r \in \{\pm 1\}$ that

$$\eta \cdot T_r^* = \sum_{t \geq 0, \widehat{\Lambda}_r^{(t)} \leq A} \eta \leq \underbrace{\frac{\delta}{(1 - (1 + \delta)^{-(q-2)}) (\widehat{\Lambda}_r^{(0)})^{q-2} C_1}}_{(i)} + \underbrace{\eta \cdot \frac{C_2}{C_1} (1 + \delta)^{q-1} \left(1 + \frac{\log(A/\widehat{\Lambda}_r^{(0)})}{\log(1 + \delta)}\right)}_{(ii)},$$

$$\eta \cdot T_r^* = \sum_{t \geq 0, \widehat{\Lambda}_r^{(t)} \leq A} \eta \geq \underbrace{\frac{\delta(1 - (x_0/A)^{q-2})}{(1 + \delta)^{q-1}(1 - (1 + \delta)^{-(q-2)}) (\widehat{\Lambda}_r^{(0)})^{q-2} C_2}}_{(iii)} - \underbrace{\eta \cdot (1 + \delta)^{-(q-1)} \left(1 + \frac{\log(A/\widehat{\Lambda}_r^{(0)})}{\log(1 + \delta)}\right)}_{(iv)},$$

769 where C_1 is taken as $q(p - \frac{1}{2} - o(1)) \cdot \|\mathbf{v}\|_2^2$ and C_2 is taken as $q(p - \frac{1}{2} + o(1)) \cdot \|\mathbf{v}\|_2^2$ according to
770 Lemma C.19. Taking $\delta = \frac{1}{k}$, $A = \Theta(1/m)$ and note that terms (ii), (iv) are respectively dominated
771 by terms (i), (iii) when η is sufficiently small and letting $k \rightarrow \infty$, we have

$$\frac{1}{(\widehat{\Lambda}_r^{(0)})^{q-2}C_2} - \{\text{lower order terms}\} \leq \eta \cdot T_r^* \leq \frac{1}{(\widehat{\Lambda}_r^{(0)})^{q-2}C_1} + \{\text{lower order terms}\},$$

772 for $r \in \{\pm 1\}$. It follows that

$$\eta \cdot T_r^* = \frac{1}{q(p - \frac{1}{2})\|\mathbf{v}\|_2^2 \cdot (\widehat{\Lambda}_r^{(0)})^{q-2}} \pm \{\text{lower order terms}\}. \quad (\text{C.34})$$

773 And by Lemma C.9, we have $\eta \cdot T_r^* = \Theta(1/q(p - \frac{1}{2}))\|\mathbf{v}\|_2^2 \cdot (\sqrt{\log(m)}\sigma_0\|\mathbf{v}\|_2)^{q-2} = \widetilde{\Theta}(d^{q/4 - 3/2})$,
774 which completes the proof. \square

775 The discussion in this section verifies Lemma C.4 and provides a clear understanding about how
776 $\widehat{\Lambda}_r^{(t)}$, $\bar{\Lambda}_r^{(t)}$ varies within the iteration range $[0, T_r]$ for $r \in \{\pm 1\}$. Note that the iteration numbers
777 when $\widehat{\Lambda}_1^{(t)}$ and $\widehat{\Lambda}_{-1}^{(t)}$ reaches $\Theta(1/m)$ (T_1 and T_{-1}) are different, however, since T_{-1} and T_1 have the
778 same magnitude, it remains clear that although $T_1 \neq T_{-1}$ (wlog, assume $T_1 < T_{-1}$), we still have
779 $\widehat{\Lambda}_1^{(t)} = \widetilde{\Theta}(1)$ and $\bar{\Lambda}_1^{(t)} = \widetilde{O}(d^{-\frac{1}{4}})$ within the iteration range $[T_1, T_{-1}]$, since off-diagonal feature
780 learning also costs time no less than order $\Theta(1/\eta\sigma_0\|\mathbf{v}\|_2^q(\log m)^{(q-2)/2})$, which is higher order than

781 $|T_1 - T_{-1}| = \Theta(1/\eta\sigma_0\|\mathbf{v}\|_2^q(\log m)^{(q-1)/2})$, according to (C.34) and Lemma C.9. Therefore, at
782 time $T_0 := \max\{T_1, T_{-1}\}$, off-diagonal $\bar{\Lambda}_1^{(t)}, \bar{\Lambda}_{-1}^{(t)}$ still remain initialization magnitude $\tilde{O}(d^{-\frac{1}{4}})$,
783 $\Gamma_1^{(t)}, \Gamma_{-1}^{(t)}$ remain initialization magnitude $\tilde{O}(d^{-\frac{1}{4}+\epsilon})$, while on-diagonal $\hat{\Lambda}_1^{(t)}, \hat{\Lambda}_{-1}^{(t)}$ reach and then
784 remain $\tilde{\Theta}(1)$.

785 C.6 Proof of Lemma C.3

786 If we only use labeled data S' for the optimization of CNN, according to Lemma D.1, we have

$$\begin{aligned} \mathbf{w}_j^{(t+1)} &= \mathbf{w}_j^{(t)} - \nabla_{\mathbf{w}_j} L_{S'}(\mathbf{W}) \\ &= (1 - \eta\lambda) \cdot \mathbf{w}_j^{(t)} + \frac{q\eta u_j}{n_1} \sum_{i=1}^{n_1} b_i^{(t)} y'_i ([\langle \mathbf{w}_j, y'_i \cdot \mathbf{v} \rangle]_+^{q-1} \cdot y'_i \cdot \mathbf{v} + [\langle \mathbf{w}_j, \boldsymbol{\xi}'_i \rangle]_+^{q-1} \cdot \boldsymbol{\xi}'_i), \end{aligned}$$

787 where $u_j := \mathbb{1}_{[1 \leq j \leq m]} - \mathbb{1}_{[m+1 \leq 2m]}$, $b_i^{(t)} = -\ell'(y'_i \cdot f_{\mathbf{W}}(\mathbf{x}'_i)) = \exp[-y'_i \cdot f_{\mathbf{W}}(\mathbf{x}'_i)] / (1 + \exp[-y'_i \cdot$
788 $f_{\mathbf{W}}(\mathbf{x}'_i)])$.

789 Notice that \mathbf{v} and $\boldsymbol{\xi}'_i$ are orthogonal to each other, we have

$$\begin{aligned} \langle \mathbf{w}_j^{(t+1)}, \mathbf{v} \rangle &= (1 - \eta\lambda) \cdot \langle \mathbf{w}_j^{(t)}, \mathbf{v} \rangle + \frac{q\eta u_j}{n_1} \sum_{i=1}^{n_1} b_i^{(t)} \cdot [\langle \mathbf{w}_j^{(t)}, y'_i \cdot \mathbf{v} \rangle]_+^{q-1} \cdot \|\mathbf{v}\|_2^2, \\ \langle \mathbf{w}_j^{(t+1)}, \boldsymbol{\xi}'_i \rangle &= (1 - \eta\lambda) \cdot \langle \mathbf{w}_j^{(t)}, \boldsymbol{\xi}'_i \rangle + \frac{q\eta u_j}{n_1} \sum_{i=1}^{n_1} b_i^{(t)} y'_i \cdot [\langle \mathbf{w}_j^{(t)}, \boldsymbol{\xi}'_i \rangle]_+^{q-1} \cdot \langle \boldsymbol{\xi}'_i, \boldsymbol{\xi}'_i \rangle, i \in [n_1]. \end{aligned}$$

790 Let T'_i be the first iteration that $\Gamma_i^{(t)}$ reaches $\Theta(1/m)$, then we have following lemma:

791 **Lemma C.21.** As long as $\Gamma_i^{(t)} \leq \Theta(1/m)$, $b_i^{(t)} := -\ell'(y'_i \cdot f_{\mathbf{W}^{(t)}}(\mathbf{x}'_i))$ will remain $1/2 \pm o(1)$.

792 *Proof of Lemma C.21.* Note that $\ell(z) = \log(1 + \exp(-z))$ and $-\ell'(z) = \exp(-z) / (1 + \exp(-z))$,
793 and without loss of generality assuming $y'_i = 1$, we can express $b_i^{(t)}$ as follow:

$$b_i^{(t)} = -\ell'(f_{\mathbf{W}^{(t)}}(\mathbf{x}'_i)) = \frac{e^{\sum_{j=m+1}^{2m} [\sigma(\langle \mathbf{w}_j^{(t)}, \mathbf{v} \rangle) + \sigma(\langle \mathbf{w}_j^{(t)}, \boldsymbol{\xi}'_i \rangle)]}}{e^{\sum_{j=1}^m [\sigma(\langle \mathbf{w}_j^{(t)}, \mathbf{v} \rangle) + \sigma(\langle \mathbf{w}_j^{(t)}, \boldsymbol{\xi}'_i \rangle)]} + e^{\sum_{j=m+1}^{2m} [\sigma(\langle \mathbf{w}_j^{(t)}, \mathbf{v} \rangle) + \sigma(\langle \mathbf{w}_j^{(t)}, \boldsymbol{\xi}'_i \rangle)]}},$$

794 Since $\sigma(\langle \mathbf{w}_j^{(t)}, \boldsymbol{\xi}'_i \rangle)$ will dominate $\sigma(\langle \mathbf{w}_j^{(t)}, \mathbf{v} \rangle)$, which will be proved later by using *tensor power*
795 *method*, we have

$$b_i^{(t)} = -\ell'(f_{\mathbf{W}^{(t)}}(\mathbf{x}'_i)) = \frac{e^{\sum_{j=m+1}^{2m} [\sigma(\langle \mathbf{w}_j^{(t)}, \mathbf{v} \rangle) + \sigma(\langle \mathbf{w}_j^{(t)}, \boldsymbol{\xi}'_i \rangle)]}}{e^{\sum_{j=1}^m \sigma(\langle \mathbf{w}_j^{(t)}, \boldsymbol{\xi}'_i \rangle) + \{\text{lower order term}\}} + e^{\sum_{j=m+1}^{2m} [\sigma(\langle \mathbf{w}_j^{(t)}, \mathbf{v} \rangle) + \sigma(\langle \mathbf{w}_j^{(t)}, \boldsymbol{\xi}'_i \rangle)]}},$$

796 On the one side,

$$\begin{aligned} b_i^{(t)} &\geq \frac{1}{e^{\sum_{j=1}^m \sigma(\langle \mathbf{w}_j^{(t)}, \boldsymbol{\xi}'_i \rangle) + \{\text{lower order term}\}} + 1} \\ &\geq \frac{1}{e^{m(\Gamma_i^{(t)})^q + \{\text{lower order term}\}} + 1} \\ &\geq \frac{1}{e^{\Theta(m^{-(q-1)})} + 1} = \frac{1}{2 + o(1)} = \frac{1}{2} - o(1). \end{aligned}$$

797 On the other side, according to Lemma C.5, we have $\bar{\Lambda}_1^{(t)} = \tilde{O}(d^{-\frac{1}{4}})$, it follows that

$$\begin{aligned} b_i^{(t)} &\leq \frac{e^{m(\bar{\Lambda}_1^{(t)})^q + o(1)}}{e^{\sum_{j=1}^m \sigma(\langle \mathbf{w}_j^{(t)}, \boldsymbol{\xi}'_i \rangle) + \{\text{lower order term}\}} + e^{m(\bar{\Lambda}_1^{(t)})^q + o(1)}} \\ &= \frac{1 + o(1)}{e^{\sum_{j=1}^m \sigma(\langle \mathbf{w}_j^{(t)}, \boldsymbol{\xi}'_i \rangle) + \{\text{lower order term}\}} + 1 + o(1)} \end{aligned}$$

$$\leq \frac{1 + o(1)}{1 + 1 + o(1)} = \frac{1}{2} + o(1).$$

798 Therefore, we have $b_i^{(t)} = 1/2 \pm o(1)$ and the other case of $y_i = -1$ can be proved in a similar
799 way. \square

800 With the help of above lemma, we are now ready to prove Lemma C.3.

801 *Proof of Lemma C.3.* Let $j^* = \arg \max_{1 \leq j \leq m} \langle \mathbf{w}_j^{(t+1)}, \mathbf{v} \rangle$ and note that $u_j = 1$, according to
802 Lemma C.21, we have

$$\begin{aligned} \widehat{\Lambda}_1^{(t+1)} &= \langle \mathbf{w}_{j^*}^{(t+1)}, \mathbf{v} \rangle \\ &= (1 - \eta\lambda) \cdot \langle \mathbf{w}_{j^*}^{(t)}, \mathbf{v} \rangle + \frac{q\eta}{n_1} \sum_{i=1}^{n_1} b_i^{(t)} [\langle \mathbf{w}_{j^*}^{(t)}, y_i \cdot \mathbf{v} \rangle_+^{q-1} \|\mathbf{v}\|_2^2 \\ &= (1 - \eta\lambda) \cdot \langle \mathbf{w}_{j^*}^{(t)}, \mathbf{v} \rangle + \frac{q\eta}{n_1} \sum_{i \in S'_1} b_i^{(t)} [\langle \mathbf{w}_{j^*}^{(t)}, \mathbf{v} \rangle_+^{q-1} \|\mathbf{v}\|_2^2 + \frac{q\eta}{n_1} \sum_{i \in S'_{-1}} b_i^{(t)} [\langle \mathbf{w}_{j^*}^{(t)}, -\mathbf{v} \rangle_+^{q-1} \|\mathbf{v}\|_2^2 \\ &= (1 - \eta\lambda) \cdot \langle \mathbf{w}_{j^*}^{(t)}, \mathbf{v} \rangle + \underbrace{\frac{q\eta}{n_1} \sum_{i \in S'_1} \left(\frac{1}{2} \pm o(1) \right) [\langle \mathbf{w}_{j^*}^{(t)}, \mathbf{v} \rangle_+^{q-1} \|\mathbf{v}\|_2^2}_{\clubsuit} \\ &\quad + \underbrace{\frac{q\eta}{n_1} \sum_{i \in S'_{-1}} \left(\frac{1}{2} \pm o(1) \right) [\langle \mathbf{w}_{j^*}^{(t)}, -\mathbf{v} \rangle_+^{q-1} \|\mathbf{v}\|_2^2}_{\star} \end{aligned} \tag{C.35}$$

803 For \clubsuit , we have

$$\begin{aligned} \underbrace{\sum_{i \in S'_1} \left(\frac{1}{2} \pm o(1) \right) [\langle \mathbf{w}_{j^*}^{(t)}, \mathbf{v} \rangle_+^{q-1} \|\mathbf{v}\|_2^2}_{\clubsuit} &= n'_1 \cdot \left(\frac{1}{2} \pm o(1) \right) \cdot [\langle \mathbf{w}_{j^*}^{(t)}, \mathbf{v} \rangle_+^{q-1} \|\mathbf{v}\|_2^2 \\ &\leq n'_1 \cdot \left(\frac{1}{2} \pm o(1) \right) \cdot \|\mathbf{v}\|_2^2 \cdot (\widehat{\Lambda}_1^{(t)})^{q-1}. \end{aligned} \tag{C.36}$$

804 For \star , we have

$$\begin{aligned} \underbrace{\sum_{i \in S'_{-1}} \left(\frac{1}{2} \pm o(1) \right) [\langle \mathbf{w}_{j^*}^{(t)}, -\mathbf{v} \rangle_+^{q-1} \|\mathbf{v}\|_2^2}_{\star} &= n'_{-1} \cdot \left(\frac{1}{2} \pm o(1) \right) \cdot [\langle \mathbf{w}_{j^*}^{(t)}, -\mathbf{v} \rangle_+^{q-1} \|\mathbf{v}\|_2^2 \\ &\leq n'_{-1} \cdot \left(\frac{1}{2} \pm o(1) \right) \cdot \|\mathbf{v}\|_2^2 \cdot (\bar{\Lambda}_{-1}^{(t)})^{q-1}. \end{aligned} \tag{C.37}$$

805 By plugging (C.36) and (C.37) in (C.35), and according to Lemma C.14, we have with probability at
806 least $1 - 4\delta$ that

$$\begin{aligned} \widehat{\Lambda}_1^{(t+1)} &\leq (1 - \eta\lambda) \cdot \widehat{\Lambda}_1^{(t)} + \frac{q\eta}{n_1} \left(n'_1 \cdot \left(\frac{1}{2} \pm o(1) \right) \cdot \|\mathbf{v}\|_2^2 \cdot (\widehat{\Lambda}_1^{(t)})^{q-1} + n'_{-1} \cdot \left(\frac{1}{2} \pm o(1) \right) \cdot \|\mathbf{v}\|_2^2 \cdot (\bar{\Lambda}_{-1}^{(t)})^{q-1} \right) \\ &\leq (1 - \eta\lambda) \cdot \widehat{\Lambda}_1^{(t)} + \frac{q\eta}{n_1} \left(\left(\frac{n_1}{2} + \sqrt{\frac{n_1}{2} \log \frac{1}{\delta}} \right) \cdot \left(\frac{1}{2} \pm o(1) \right) \cdot \|\mathbf{v}\|_2^2 \cdot (\widehat{\Lambda}_1^{(t)})^{q-1} \right. \\ &\quad \left. + \left(\frac{n_1}{2} + \sqrt{\frac{n_1}{2} \log \frac{1}{\delta}} \right) \cdot \left(\frac{1}{2} \pm o(1) \right) \cdot \|\mathbf{v}\|_2^2 \cdot (\bar{\Lambda}_{-1}^{(t)})^{q-1} \right) \\ &= (1 - \eta\lambda) \cdot \widehat{\Lambda}_1^{(t)} + q\eta \left(\left(\frac{1}{4} \pm o(1) \right) \cdot \|\mathbf{v}\|_2^2 \cdot (\widehat{\Lambda}_1^{(t)})^{q-1} + \left(\frac{1}{4} \pm o(1) \right) \cdot \|\mathbf{v}\|_2^2 \cdot (\bar{\Lambda}_{-1}^{(t)})^{q-1} \right) \end{aligned}$$

$$= (1 - \eta\lambda) \cdot \widehat{\Lambda}_1^{(t)} + \eta \cdot \Theta(d) \cdot \left((\widehat{\Lambda}_1^{(t)})^2 + (\bar{\Lambda}_{-1}^{(t)})^{q-1} \right).$$

807 And we can prove in the same way that with probability at least $1 - 4\delta$ we have

$$\widehat{\Lambda}_{-1}^{(t+1)} \leq (1 - \eta\lambda) \cdot \widehat{\Lambda}_{-1}^{(t)} + \eta \cdot \Theta(d) \cdot \left((\widehat{\Lambda}_{-1}^{(t)})^{q-1} + (\bar{\Lambda}_1^{(t)})^{q-1} \right).$$

808 Let $j^* = \arg \max_{m+1 \leq j \leq 2m} \langle \mathbf{w}_j^{(t+1)}, \mathbf{v} \rangle$ and note that $u_j = -1$, we have

$$\begin{aligned} \bar{\Lambda}_1^{(t+1)} &= \langle \mathbf{w}_{j^*}^{(t+1)}, \mathbf{v} \rangle \\ &= (1 - \eta\lambda) \cdot \langle \mathbf{w}_{j^*}^{(t)}, \mathbf{v} \rangle - \frac{q\eta}{n_1} \sum_{i=1}^{n_1} b_i^{(t)} [\langle \mathbf{w}_{j^*}^{(t)}, \mathbf{y}'_i \cdot \mathbf{v} \rangle]_+^{q-1} \|\mathbf{v}\|_2^2 \\ &\leq (1 - \eta\lambda) \cdot \langle \mathbf{w}_{j^*}^{(t)}, \mathbf{v} \rangle \\ &\leq (1 - \eta\lambda) \cdot \bar{\Lambda}_1^{(t)}. \end{aligned} \tag{C.38}$$

809 And we can prove in the same way that $\bar{\Lambda}_{-1}^{(t+1)} \leq (1 - \eta\lambda) \cdot \bar{\Lambda}_{-1}^{(t)}$.

810 Next, we consider the increasing rate of $\Gamma_l^{(t)}$ where $l \in [n_1]$ is fixed. If $y_l = 1$, let $j^\natural =$

811 $\arg \max_{1 \leq j \leq m} \langle \mathbf{w}_j^{(t)}, \boldsymbol{\xi}'_l \rangle$ and note that $u_j = 1$, we have

$$\begin{aligned} \Gamma_l^{(t+1)} &\geq \langle \mathbf{w}_{j^\natural}^{(t+1)}, \boldsymbol{\xi}'_l \rangle \\ &= (1 - \eta\lambda) \cdot \langle \mathbf{w}_{j^\natural}^{(t)}, \boldsymbol{\xi}'_l \rangle + \frac{q\eta}{n_1} \sum_{i=1}^{n_1} b_i^{(t)} y'_i \cdot [\langle \mathbf{w}_{j^\natural}^{(t)}, \boldsymbol{\xi}'_i \rangle]_+^{q-1} \cdot \langle \boldsymbol{\xi}'_i, \boldsymbol{\xi}'_l \rangle \\ &= (1 - \eta\lambda) \cdot \langle \mathbf{w}_{j^\natural}^{(t)}, \boldsymbol{\xi}'_l \rangle + \frac{q\eta}{n_1} b_l^{(t)} [\langle \mathbf{w}_{j^\natural}^{(t)}, \boldsymbol{\xi}'_l \rangle]_+^{q-1} \|\boldsymbol{\xi}'_l\|_2^2 + \frac{q\eta}{n_1} \sum_{i \in [n_1], i \neq l} b_i^{(t)} y'_i [\langle \mathbf{w}_{j^\natural}^{(t)}, \boldsymbol{\xi}'_i \rangle]_+^{q-1} \langle \boldsymbol{\xi}'_i, \boldsymbol{\xi}'_l \rangle \\ &= (1 - \eta\lambda) \cdot \langle \mathbf{w}_{j^\natural}^{(t)}, \boldsymbol{\xi}'_l \rangle + \frac{q\eta}{n_1} b_l^{(t)} [\langle \mathbf{w}_{j^\natural}^{(t)}, \boldsymbol{\xi}'_l \rangle]_+^{q-1} \|\boldsymbol{\xi}'_l\|_2^2 \{\pm \text{lower order terms}\} \\ &\geq (1 - \eta\lambda) \cdot \Gamma_l^{(t)} + \frac{q\eta}{n_1} \cdot \left(\frac{1}{2} - o(1) \right) \cdot \|\boldsymbol{\xi}'_l\|_2^2 \cdot (\Gamma_l^{(t)})^{q-1} \\ &= (1 - \eta\lambda) \cdot \Gamma_l^{(t)} + \eta \cdot \widetilde{\Theta}(d^{1+2\epsilon}) \cdot (\Gamma_l^{(t)})^{q-1}, \end{aligned} \tag{C.39}$$

812 where the third equality holds if we properly choose the order of λ .

813 If $y_l = -1$, let $j^\natural = \arg \max_{m+1 \leq j \leq 2m} \langle \mathbf{w}_j^{(t)}, \boldsymbol{\xi}'_l \rangle$ and note that $u_j = -1$, we have

$$\begin{aligned} \Gamma_l^{(t+1)} &\geq \langle \mathbf{w}_{j^\natural}^{(t+1)}, \boldsymbol{\xi}'_l \rangle \\ &= (1 - \eta\lambda) \cdot \langle \mathbf{w}_{j^\natural}^{(t)}, \boldsymbol{\xi}'_l \rangle - \frac{q\eta}{n_1} \sum_{i=1}^{n_1} b_i^{(t)} y'_i \cdot [\langle \mathbf{w}_{j^\natural}^{(t)}, \boldsymbol{\xi}'_i \rangle]_+^{q-1} \cdot \langle \boldsymbol{\xi}'_i, \boldsymbol{\xi}'_l \rangle \\ &= (1 - \eta\lambda) \cdot \langle \mathbf{w}_{j^\natural}^{(t)}, \boldsymbol{\xi}'_l \rangle + \frac{q\eta}{n_1} b_l^{(t)} [\langle \mathbf{w}_{j^\natural}^{(t)}, \boldsymbol{\xi}'_l \rangle]_+^{q-1} \|\boldsymbol{\xi}'_l\|_2^2 - \frac{q\eta}{n_1} \sum_{i \in [n_1], i \neq l} b_i^{(t)} y'_i [\langle \mathbf{w}_{j^\natural}^{(t)}, \boldsymbol{\xi}'_i \rangle]_+^{q-1} \langle \boldsymbol{\xi}'_i, \boldsymbol{\xi}'_l \rangle \\ &= (1 - \eta\lambda) \cdot \langle \mathbf{w}_{j^\natural}^{(t)}, \boldsymbol{\xi}'_l \rangle + \frac{q\eta}{n_1} b_l^{(t)} [\langle \mathbf{w}_{j^\natural}^{(t)}, \boldsymbol{\xi}'_l \rangle]_+^{q-1} \|\boldsymbol{\xi}'_l\|_2^2 \{\pm \text{lower order terms}\} \\ &\geq (1 - \eta\lambda) \cdot \Gamma_l^{(t)} + \frac{q\eta}{n_1} \cdot \left(\frac{1}{2} - o(1) \right) \cdot \|\boldsymbol{\xi}'_l\|_2^2 \cdot (\Gamma_l^{(t)})^{q-1} \\ &= (1 - \eta\lambda) \cdot \Gamma_l^{(t)} + \eta \cdot \widetilde{\Theta}(d^{1+2\epsilon}) \cdot (\Gamma_l^{(t)})^{q-1}, \end{aligned} \tag{C.40}$$

814 where the third equality holds if we properly choose the order of λ .

815 According to (C.39) and (C.40), we always have

$$\Gamma_l^{(t+1)} \geq (1 - \eta\lambda) \cdot \Gamma_l^{(t)} + \eta \cdot \widetilde{\Theta}(d^{1+2\epsilon}) \cdot (\Gamma_l^{(t)})^{q-1}.$$

816

□

817 **C.7 Proof of Lemma C.5**

818 By applying Lemma E.4 to $\Gamma_i^{(t)}$ and taking $C_1 = \tilde{\Theta}(d^{1+2\epsilon})$, $\delta = 1/2$, $A = \Theta(1/m)$, we have

$$\sum_{t \geq 0, \Gamma_i^{(t)} \leq A} \eta \leq \Theta(1/C_1(\Gamma_i^{(0)})^{q-2}) = \tilde{\Theta}(d^{(\frac{1}{4}-\epsilon)q-\frac{3}{2}}).$$

819 And note the definition of T'_i , we have

$$\eta \cdot T'_i = \tilde{\Theta}(d^{(\frac{1}{4}-\epsilon)q-\frac{3}{2}}). \quad (\text{C.41})$$

820 In Lemma C.3, we have already prove that

$$\begin{aligned} \widehat{\Lambda}_r^{(t+1)} &\leq (1 - \eta\lambda) \cdot \widehat{\Lambda}_r^{(t)} + \eta \cdot \Theta(d) \cdot \left((\widehat{\Lambda}_r^{(t)})^{q-1} + (\bar{\Lambda}_{-r}^{(t)})^{q-1} \right), \\ \widehat{\Lambda}_r^{(t+1)} &\leq (1 - \eta\lambda) \cdot \widehat{\Lambda}_r^{(t+1)}, r \in \{\pm 1\}. \end{aligned} \quad (\text{C.42})$$

821 Define $\Lambda^{(t)} := \max_{r \in \{\pm 1\}} \{\widehat{\Lambda}_r^{(t)}, \bar{\Lambda}_r^{(t)}\}$, according to (C.42), we have

$$\Lambda^{(t+1)} \leq (1 - \eta\lambda) \cdot \Lambda^{(t)} + \eta \cdot \Theta(d) \cdot (\Lambda^{(t)})^{q-1}.$$

822 By applying Lemma E.4 to $\Lambda^{(t)}$, and taking $C_1 = \Theta(d)$, $\delta = 1/2$, $A = C \cdot \Lambda^{(0)}$, where A is a large
823 constant, we have

$$\sum_{t \geq 0, \Lambda^{(t)} \leq A} \eta \geq \Theta(1/C_1(\Lambda^{(0)})^{q-2}) = \tilde{\Theta}(d^{\frac{q}{4}-\frac{3}{2}}).$$

824 Let T' be the first iteration that $\Lambda^{(t)}$ reaches $C \cdot \Lambda^{(0)}$, then we have

$$\eta \cdot T' = \tilde{\Theta}(d^{\frac{q}{4}-\frac{3}{2}}). \quad (\text{C.43})$$

825 According to (C.41) and (C.43), we have $T' = \omega(T'_i)$, which indicates that when $\Gamma_i^{(t)}$ reaches
826 $\Theta(1/m)$, $\Lambda^{(t)}$ remains initialization magnitude $\tilde{\Theta}(d^{-\frac{1}{4}})$.

827 **C.8 Empirical, test error and loss for early stopped classifier**

828 Assume the accuracy of pseudo-labeler p is larger than $1/2$. We first estimate the empirical loss
829 for early stopped classifier $f_{\mathbf{W}(T_0)}$, where $T_0 = \max_{r \in \{\pm 1\}} \{T_r\}$ and T_r is defined as the first
830 iteration that $\widehat{\Lambda}_r^{(t)}$ reaches $\Theta(1/m)$. According to Section C.5.3 and Lemma C.18, we have $\widehat{\Lambda}_r^{(T_0)} =$
831 $\tilde{\Theta}(1)$, $\bar{\Lambda}_r^{(T_0)} = \tilde{O}(d^{-\frac{1}{4}})$, $\Gamma^{(t)} = \tilde{O}(d^{-\frac{1}{4}+\epsilon})$, for $r \in \{\pm 1\}$. We have the following lemma:

832 **Lemma C.22.** Early stopped classifier $f_{\mathbf{W}(T_0)}(\mathbf{x})$ possesses following properties:

- 833 1. Training error of early stopped classifier $f_{\mathbf{W}(T_0)}(\mathbf{x})$ is asymptotically $1 - p$: $\frac{1}{n_u+n_l} (\sum_{i=1}^{n_u} \mathbb{1}[\widehat{y}_i \cdot$
834 $f_{\mathbf{W}(T_0)}(\mathbf{x}_i) \leq 0] + \sum_{i=1}^{n_l} \mathbb{1}[y'_i \cdot f_{\mathbf{W}(T_0)}(\mathbf{x}'_i) \leq 0]) = 1 - p \pm o(1)$.
- 835 2. Test error is nearly $1 - p$, if we use pseudo-label \widehat{y} generated by pseudo-labeler as target:
836 $\mathbb{P}_{(\mathbf{x}, y) \sim \mathcal{D}, \widehat{y} \sim y \cdot \mathcal{B}(p)}[\widehat{y} \cdot f_{\mathbf{W}(T_0)}(\mathbf{x}) \leq 0] = 1 - p \pm o(1)$.
- 837 3. Test error is nearly 0, if we use true label y as target: $\mathbb{P}_{(\mathbf{x}, y) \sim \mathcal{D}}[y \cdot f_{\mathbf{W}(T_0)}(\mathbf{x}) \leq 0] = o(1)$ and
838 hence $\text{sign } f_{\mathbf{W}(T_0)}(\mathbf{x}) = \text{sign}(y)$ with high probability,

839 where p is the accuracy of the pseudo-labeler. We can regard p as the probability that \mathbf{x}_i is paired
840 with true label y_i , $1 - p$ is the probability that \mathbf{x}_i is paired with wrong label $-y_i$.

841 *Proof of Lemma C.22.* Recall the definition of $f_{\mathbf{W}}$ in (2.1) that

$$\begin{aligned} f_{\mathbf{W}(T_0)}(\mathbf{x}_i) &= \sum_{j=1}^m \left[\sigma(\langle \mathbf{w}_j^{(T_0)}, y_i \cdot \mathbf{v} \rangle) + \sigma(\langle \mathbf{w}_j^{(T_0)}, \boldsymbol{\xi}_i \rangle) \right] \\ &\quad - \sum_{j=m+1}^{2m} \left[\sigma(\langle \mathbf{w}_j^{(T_0)}, y_i \cdot \mathbf{v} \rangle) + \sigma(\langle \mathbf{w}_j^{(T_0)}, \boldsymbol{\xi}_i \rangle) \right]. \end{aligned}$$

842 According to Section C.5.3 and Lemma C.18, we have $\widehat{\Lambda}_r^{(T_0)} = \widetilde{\Theta}(1)$, $\bar{\Lambda}_r^{(T_0)} = \widetilde{O}(d^{-\frac{1}{4}})$, $\Gamma^{(t)} =$
843 $\max\{\max_{i \in [n_u]} \Gamma_i^{(t)}, \max_{i \in [n_l]} \Gamma_i^{(t)}\} = \widetilde{O}(d^{-\frac{1}{4} + \epsilon})$, for $r \in \{\pm 1\}$. If $y_i = 1$, we have following
844 lower bound for $f_{\mathbf{W}^{(T_0)}}(\mathbf{x}_i)$

$$\begin{aligned} f_{\mathbf{W}^{(T_0)}}(\mathbf{x}_i) &= \sum_{j=1}^m \left[\sigma(\langle \mathbf{w}_j^{(T_0)}, \mathbf{v} \rangle) + \sigma(\langle \mathbf{w}_j^{(T_0)}, \boldsymbol{\xi}_i \rangle) \right] - \sum_{j=m+1}^{2m} \left[\sigma(\langle \mathbf{w}_j^{(T_0)}, \mathbf{v} \rangle) + \sigma(\langle \mathbf{w}_j^{(T_0)}, \boldsymbol{\xi}_i \rangle) \right] \\ &\geq (\widehat{\Lambda}_1^{(T_0)})^q + (\Gamma_i^{(T_0)})^q - m(\bar{\Lambda}_1^{(T_0)})^q - m(\Gamma_i^{(T_0)})^q \\ &\geq (\widehat{\Lambda}_1^{(T_0)})^q \{- \text{lower order terms}\}, \end{aligned}$$

845 and following upper bound for $f_{\mathbf{W}^{(T_0)}}(\mathbf{x}_i)$:

$$\begin{aligned} f_{\mathbf{W}^{(T_0)}}(\mathbf{x}_i) &= \sum_{j=1}^m \left[\sigma(\langle \mathbf{w}_j^{(T_0)}, \mathbf{v} \rangle) + \sigma(\langle \mathbf{w}_j^{(T_0)}, \boldsymbol{\xi}_i \rangle) \right] - \sum_{j=m+1}^{2m} \left[\sigma(\langle \mathbf{w}_j^{(T_0)}, \mathbf{v} \rangle) + \sigma(\langle \mathbf{w}_j^{(T_0)}, \boldsymbol{\xi}_i \rangle) \right] \\ &\leq m(\widehat{\Lambda}_1^{(T_0)})^q + m(\Gamma_i^{(T_0)})^q - (\bar{\Lambda}_1^{(T_0)})^q - (\Gamma_i^{(T_0)})^q \\ &\leq (\widehat{\Lambda}_1^{(T_0)})^q \{+ \text{lower order terms}\}. \end{aligned}$$

846 If $y_i = -1$, we have following upper bound for $f_{\mathbf{W}^{(T_0)}}(\mathbf{x}_i)$:

$$\begin{aligned} f_{\mathbf{W}^{(T_0)}}(\mathbf{x}_i) &= \sum_{j=1}^m \left[\sigma(-\langle \mathbf{w}_j^{(T_0)}, \mathbf{v} \rangle) + \sigma(\langle \mathbf{w}_j^{(T_0)}, \boldsymbol{\xi}_i \rangle) \right] - \sum_{j=m+1}^{2m} \left[\sigma(-\langle \mathbf{w}_j^{(T_0)}, \mathbf{v} \rangle) + \sigma(\langle \mathbf{w}_j^{(T_0)}, \boldsymbol{\xi}_i \rangle) \right] \\ &\leq m(\bar{\Lambda}_{-1}^{(T_0)})^q + m(\Gamma_i^{(T_0)})^q - (\widehat{\Lambda}_{-1}^{(T_0)})^q - (\Gamma_i^{(T_0)})^q \\ &\leq -(\widehat{\Lambda}_{-1}^{(T_0)})^q \{+ \text{lower order terms}\}, \end{aligned}$$

847 and following lower bound for $f_{\mathbf{W}^{(T_0)}}(\mathbf{x}_i)$:

$$\begin{aligned} f_{\mathbf{W}^{(T_0)}}(\mathbf{x}_i) &= \sum_{j=1}^m \left[\sigma(-\langle \mathbf{w}_j^{(T_0)}, \mathbf{v} \rangle) + \sigma(\langle \mathbf{w}_j^{(T_0)}, \boldsymbol{\xi}_i \rangle) \right] - \sum_{j=m+1}^{2m} \left[\sigma(-\langle \mathbf{w}_j^{(T_0)}, \mathbf{v} \rangle) + \sigma(\langle \mathbf{w}_j^{(T_0)}, \boldsymbol{\xi}_i \rangle) \right] \\ &\geq (\bar{\Lambda}_{-1}^{(T_0)})^q + (\Gamma_i^{(T_0)})^q - m(\widehat{\Lambda}_{-1}^{(T_0)})^q - m(\Gamma_i^{(T_0)})^q \\ &\geq -m(\bar{\Lambda}_{-1}^{(T_0)})^q \{- \text{lower order terms}\}. \end{aligned}$$

848 Therefore, for unlabeled data, we have $y_i \cdot f_{\mathbf{W}^{(T_0)}}(\mathbf{x}_i) \in [(1-o(1)) \cdot (\widehat{\Lambda}_{y_i}^{(T_0)})^q, (m+o(1)) \cdot (\widehat{\Lambda}_{y_i}^{(T_0)})^q]$
849 and hence $\text{sign}(f_{\mathbf{W}^{(T_0)}}(\mathbf{x}_i)) = \text{sign}(y_i)$ holds with high probability. We can also prove for labeled
850 data (\mathbf{x}'_i, y'_i) that $y'_i \cdot f_{\mathbf{W}^{(T_0)}}(\mathbf{x}'_i) \in [(1-o(1)) \cdot (\widehat{\Lambda}_{y'_i}^{(T_0)})^q, (m+o(1)) \cdot (\widehat{\Lambda}_{y'_i}^{(T_0)})^q]$, $\text{sign}(f_{\mathbf{W}^{(T_0)}}(\mathbf{x}'_i)) =$
851 $\text{sign}(y'_i)$ in the same way.

852 Note that \widehat{y}_i takes y_i with probability p , $-y_i$ with probability p and $n_l = o(n_u)$, the first statement in
853 this lemma follows obviously.

854 To prove the other two statement, we need to give an upper bound for the norm of \mathbf{w}_j . According to
855 the update rule of $\mathbf{w}_j^{(t)}$, we have

$$\begin{aligned} \mathbf{w}_j^{(t+1)} &= (1 - \eta\lambda) \cdot \mathbf{w}_j^{(t)} + \frac{q\eta u_j}{n_l + n_u} \cdot \left(\sum_{i=1}^{n_u} c_i \widehat{y}_i ([\langle \mathbf{w}_j^{(t)}, y_i \cdot \mathbf{v} \rangle]_+^{q-1} \cdot y_i \cdot \mathbf{v} + [\langle \mathbf{w}_j^{(t)}, \boldsymbol{\xi}_i \rangle]_+^{q-1} \cdot \boldsymbol{\xi}_i) \right. \\ &\quad \left. + \sum_{i=1}^{n_l} b_i y'_i ([\langle \mathbf{w}_j^{(t)}, y'_i \cdot \mathbf{v} \rangle]_+^{q-1} \cdot y'_i \cdot \mathbf{v} + [\langle \mathbf{w}_j^{(t)}, \boldsymbol{\xi}'_i \rangle]_+^{q-1} \cdot \boldsymbol{\xi}'_i) \right), \end{aligned}$$

856 leading to

$$\begin{aligned}
\|\mathbf{w}_j^{(t+1)}\|_2 &\leq (1 - \eta\lambda) \cdot \|\mathbf{w}_j^{(t)}\|_2 + \frac{q\eta}{n_l + n_u} \cdot \left(\sum_{i=1}^{n_u} ([\langle \mathbf{w}_j^{(t)}, y_i \cdot \mathbf{v} \rangle]_+^{q-1} \cdot \|\mathbf{v}\|_2 + [\langle \mathbf{w}_j^{(t)}, \boldsymbol{\xi}_i \rangle]_+^{q-1} \cdot \|\boldsymbol{\xi}_i\|_2) \right. \\
&\quad \left. + \sum_{i=1}^{n_l} ([\langle \mathbf{w}_j^{(t)}, y'_i \cdot \mathbf{v} \rangle]_+^{q-1} \cdot \|\mathbf{v}\|_2 + [\langle \mathbf{w}_j^{(t)}, \boldsymbol{\xi}'_i \rangle]_+^{q-1} \cdot \|\boldsymbol{\xi}'_i\|_2) \right) \\
&\leq (1 - \eta\lambda) \cdot \|\mathbf{w}_j^{(t)}\|_2 + \frac{q\eta}{n_l + n_u} \cdot \left((n_l + n_u) \cdot \|\mathbf{v}\|_2 \cdot \left(\max_{r \in \{\pm 1\}} \{\widehat{\Lambda}_r^{(t)}, \bar{\Lambda}_r^{(t)}\} \right)^{q-1} \right. \\
&\quad \left. + \left(\sum_{i \in [n_u]} \|\boldsymbol{\xi}_i\|_2 + \sum_{i \in [n_l]} \|\boldsymbol{\xi}'_i\|_2 \right) \cdot (\Gamma^{(t)})^{q-1} \right) \\
&\leq \|\mathbf{w}_j^{(t)}\|_2 + \eta \cdot \left(\Theta(d^{\frac{1}{2}}) \cdot \tilde{\Theta}(1) + \Theta(d^{\frac{1}{2}+\epsilon}) \cdot \tilde{O}(d^{(q-1)(-\frac{1}{4}+\epsilon)}) \right) \\
&= \|\mathbf{w}_j^{(t)}\|_2 + \eta \cdot \tilde{\Theta}(d^{\frac{1}{2}}), \tag{C.44}
\end{aligned}$$

857 where the first inequality is by triangle inequality; the second inequality is due to the definition of
858 $\widehat{\Lambda}_r^{(t)}, \bar{\Lambda}_r^{(t)}, \Gamma^{(t)}$, the last inequality is due to Lemma C.4.

859 According to Lemma C.20, we know that $T_r \cdot \eta = \tilde{\Theta}(d^{-\frac{3}{4}})$, $r \in \{\pm 1\}$ and $T_0 \cdot \eta = \max_{r \in \{\pm 1\}} \{T_r \cdot$
860 $\eta\} = \tilde{\Theta}(d^{-\frac{3}{4}})$. Note that $\mathbf{w}_j^{(0)} \sim \mathcal{N}(\mathbf{0}, \sigma_0^2 \mathbf{I}_d)$, $\sigma_0 = \Theta(d^{-\frac{3}{4}})$ and hence $\|\mathbf{w}_j^{(0)}\|_2 = \tilde{\Theta}(d^{-\frac{1}{4}})$, we
861 know that

$$\|\mathbf{w}_j^{(T_0)}\|_2 \leq \|\mathbf{w}_j^{(0)}\|_2 + \eta \cdot T_0 \cdot \tilde{\Theta}(d^{-\frac{1}{4}}) = \tilde{\Theta}(d^{-\frac{1}{4}}) + \tilde{\Theta}(d^{-\frac{1}{4}}) = \tilde{\Theta}(d^{-\frac{1}{4}}).$$

862 Therefore, for any (\mathbf{x}, y) sampled from distribution \mathcal{D} where $\mathbf{x} = [y \cdot \mathbf{v}^\top, \boldsymbol{\xi}^\top]^\top$ and $\boldsymbol{\xi} \sim \mathcal{N}(0, \sigma_p^2)$,
863 we have

$$\langle \mathbf{w}_j^{(T_0)}, \boldsymbol{\xi} \rangle \sim \mathcal{N}(0, \sigma_p^2 \|\mathbf{w}_j^{(T_0)}\|_2^2), |\langle \mathbf{w}_j^{(T_0)}, \boldsymbol{\xi} \rangle| = \Theta(\sigma_p \|\mathbf{w}_j^{(T_0)}\|_2) = \tilde{O}(d^{-\frac{1}{4}+\epsilon}). \tag{C.45}$$

864 And this indicates that $\langle \mathbf{w}_j^{(T_0)}, \boldsymbol{\xi} \rangle$ will still be dominated by $\langle \mathbf{w}_j^{(T_0)}, \mathbf{v} \rangle$, therefore it holds for newly
865 sampled (\mathbf{x}, y) that

$$y \cdot f_{\mathbf{W}^{(T_0)}}(\mathbf{x}) \in [(1 - o(1)) \cdot (\widehat{\Lambda}_{y_i}^{(T_0)})^q, (m + o(1)) \cdot (\widehat{\Lambda}_{y_i}^{(T_0)})^q],$$

866 which means that

$$\mathbb{P}_{(\mathbf{x}, y) \sim \mathcal{D}} [y \cdot f_{\mathbf{W}^{(T_0)}}(\mathbf{x}) \leq 0] = o(1).$$

867 This verifies the third statement that test error is nearly zero.

868 For the second statement, note that

$$\begin{aligned}
&\mathbb{P}_{(\mathbf{x}, y) \sim \mathcal{D}, \widehat{y} \sim y \cdot \mathcal{B}(p)} [\widehat{y} \cdot f_{\mathbf{W}^{(T_0)}}(\mathbf{x}) \leq 0] \\
&= \mathbb{P}_{(\mathbf{x}, y) \sim \mathcal{D}} [\widehat{y} \cdot f_{\mathbf{W}^{(T_0)}}(\mathbf{x}) \leq 0 | \widehat{y} = y] \cdot \mathbb{P}_{\widehat{y} \sim y \cdot \mathcal{B}(p)} (\widehat{y} = y) \\
&\quad + \mathbb{P}_{(\mathbf{x}, y) \sim \mathcal{D}} [\widehat{y} \cdot f_{\mathbf{W}^{(T_0)}}(\mathbf{x}) \leq 0 | \widehat{y} = -y] \cdot \mathbb{P}_{\widehat{y} \sim y \cdot \mathcal{B}(p)} (\widehat{y} = -y) \\
&= p \cdot \mathbb{P}_{(\mathbf{x}, y) \sim \mathcal{D}} [y \cdot f_{\mathbf{W}^{(T_0)}}(\mathbf{x}) \leq 0] + (1 - p) \cdot \mathbb{P}_{(\mathbf{x}, y) \sim \mathcal{D}} [y \cdot f_{\mathbf{W}^{(T_0)}}(\mathbf{x}) \geq 0] \\
&= p \cdot o(1) + (1 - p) \cdot (1 - o(1)) \\
&= 1 - p \pm o(1),
\end{aligned}$$

869 which verifies the second statement. \square

870 C.9 Downstream task

871 For downstream tasks, we use early stopped classifiers, which are stopped when on-diagonal feature
872 $\widehat{\Lambda}_r^{(t)}$ are learned while off-diagonal feature $\bar{\Lambda}_r^{(t)}$ and noise $\Gamma^{(t)}$ are not memorized. Assume we have
873 learned K early stopped classifiers $f_{\mathbf{W}_1^{(T_0^1)}}(\mathbf{x}), \dots, f_{\mathbf{W}_K^{(T_0^K)}}(\mathbf{x})$ by using n_u pseudo-labeled data
874 generated by pseudo-labeler f_1^w, \dots, f_K^w and n_l labeled data.

875 Then, we want to design a classifier on the learned representation $f_{\mathbf{W}_1^{(T_0^1)}}(\mathbf{x}), \dots, f_{\mathbf{W}_K^{(T_0^K)}}(\mathbf{x})$ to fit
 876 y . Here we consider training a downstream linear model

$$g_{\mathbf{a}}(\mathbf{x}) = \sum_{k=1}^K a_k f_{\mathbf{W}_k^{(T_0^k)}}(\mathbf{x}),$$

877 where $a_k \in \mathbb{R}$ denotes the weight as the k -th pre-trained model. Given labeled training data
 878 $S' = \{(\mathbf{x}'_i, y'_i)\}_{i=1}^{n_1}$, we want to optimize the empirical loss function

$$L_{S'}(\mathbf{a}) = \frac{1}{n_1} \sum_{i=1}^{n_1} \ell(y'_i \cdot g_{\mathbf{a}}(\mathbf{x}'_i)),$$

879 where $\ell(z) = \log(1 + \exp(-z))$ denotes the cross entropy loss. We initialize \mathbf{a} as zero and optimize
 880 empirical loss function by gradient descent, i.e.

$$\mathbf{a}^{(t+1)} = \mathbf{a}^{(t)} - \eta \cdot \nabla_{\mathbf{a}} L_{S'}(\mathbf{a}^{(t)}), \mathbf{a}^{(0)} = \mathbf{0}.$$

881 In order to estimate the training error and test error for downstream task, we first introduce following
 882 lemma about the increasing rate of $\|\mathbf{a}^{(t)}\|_1$.

883 **Lemma C.23** (Logarithmic increasing rate). For any learning rate $\eta > 0$, $a_k^{(t)}$ will always increase
 884 for any $k \in [K]$ and hence $\|\mathbf{a}^{(t)}\|_1 = \sum_{k=1}^K a_k^{(t)}$. And it holds that $\|\mathbf{a}^{(t)}\|_1 = \Theta(\log(t))$.

885 In order to give the increasing rate of $\|\mathbf{a}^{(t)}\|_1$, we introduce and prove the following lemma:

886 **Lemma C.24.** Consider following sequence $\{x_t\}_{t=1}^{\infty}$ with

$$x_{t+1} = x_t + C \cdot a^{-x_t}, x_0 = 0,$$

887 where $a > 1$ and $C > 0$ are constants, and it follows that

$$\log_a(\ln a \cdot C \cdot t + 1) \leq x_t \leq \log_a(\ln a \cdot C \cdot t + 1) + C,$$

888 and

$$x_{t+1} - x_t \leq \frac{C}{C \cdot \ln a \cdot t + 1}.$$

889 *Proof of Lemma C.24.* Note that

$$x_{i+1} - x_i = C \cdot a^{-x_i} \iff a^{x_i}(x_{i+1} - x_i) = C,$$

890 by adding up above equation from $i = 0$ to $i = t - 1$, we have

$$\begin{aligned} & \sum_{i=0}^{t-1} a^{x_i}(x_{i+1} - x_i) = C \cdot t & (C.46) \\ \implies & \int_{x_0}^{x_t} a^x dx \geq C \cdot t \\ \implies & \frac{a^{x_t} - a^{x_0}}{\ln a} \geq C \cdot t \\ \implies & a^{x_t} \geq C \cdot \ln a \cdot t + 1 \\ \implies & \begin{cases} x_t \geq \log_a(C \cdot \ln a \cdot t + 1), \\ x_{t+1} - x_t = C \cdot a^{-x_t} \leq \frac{C}{C \cdot \ln a \cdot t + 1}, \end{cases} \end{aligned}$$

891 where the first arrow is due to a^x is monotone increasing.

892 On the other hand,

$$a^{x_{i+1}} = a^{x_i + C \cdot a^{-x_i}} = a^{x_i} \cdot a^{C \cdot a^{-x_i}} \leq a^{x_i} \cdot a^{C / (C \cdot \ln a \cdot i + 1)} \leq a^{x_i} \cdot a^C,$$

893 which implies

$$\begin{aligned}
& \sum_{i=0}^{t-1} a^{x_{i+1}} \cdot (x_{i+1} - x_i) \leq a^C \sum_{i=0}^{t-1} a^{x_i} \cdot (x_{i+1} - x_i) \\
\implies & \sum_{i=0}^{t-1} a^{x_{i+1}} \cdot (x_{i+1} - x_i) \leq a^C \cdot Ct \\
\implies & \int_{x_0}^{x_t} a^x dx \leq a^C \cdot Ct,
\end{aligned}$$

894 where the first arrow is due to (C.46) and the last arrow is due to a^x is monotone increasing.

895 This leads to

$$\begin{aligned}
x_t & \leq \log_a (\ln a \cdot C \cdot a^C \cdot n + 1) \\
& \leq \log_a (\ln a \cdot C \cdot a^C \cdot n + a^C) \\
& = \log_a (\ln a \cdot C \cdot t + 1) + C
\end{aligned}$$

896 Therefore, we have

$$\log_a (\ln a \cdot C \cdot t + 1) \leq x_t \leq \log_a (\ln a \cdot C \cdot t + 1) + C,$$

897 and

$$x_{t+1} - x_t \leq \frac{C}{\ln a \cdot C \cdot t + 1}.$$

898

□

899 Now we are ready to prove Lemma C.23.

900 *Proof of Lemma C.23.* Note that we take downstream task linear model $g_{\mathbf{a}}(\mathbf{x})$ as

$$\begin{aligned}
g_{\mathbf{a}}(\mathbf{x}) & = \sum_{k=1}^d a_k \left\{ \sum_{j=1}^m \left[\sigma(\langle \mathbf{w}_{k,j}^{(T_0^k)}, y \cdot \mathbf{v} \rangle) + \sigma(\langle \mathbf{w}_{k,j}^{(T_0^k)}, \boldsymbol{\xi} \rangle) \right] \right. \\
& \quad \left. - \sum_{j=m+1}^{2m} \left[\sigma(\langle \mathbf{w}_{k,j}^{(T_0^k)}, y \cdot \mathbf{v} \rangle) + \sigma(\langle \mathbf{w}_{k,j}^{(T_0^k)}, \boldsymbol{\xi} \rangle) \right] \right\} \\
& = \sum_{k=1}^d a_k f_{\mathbf{w}_k^{(T_0^k)}}(\mathbf{x}).
\end{aligned}$$

901 Then, we have following update rule for model parameter \mathbf{a} :

$$a_k^{(t+1)} = a_k^{(t)} - \eta \cdot \frac{1}{n_1} \sum_{i=1}^{n_1} \ell'(y'_i \cdot g_{\mathbf{a}^{(t)}}(\mathbf{x}'_i)) \cdot y'_i f_{\mathbf{w}_k^{(T_0^k)}}(\mathbf{x}'_i),$$

902 where we initialize $a_k^{(0)}$ as zero for all $k \in [K]$.

903 Next, we prove following statement by using induction method: when $t \geq 1$,

904 • $a_k^{(t)}, \forall k \in [K]$ is non-negative and increasing.

905 • $\|\mathbf{a}^{(t)}\|_1 = \sum_{i=1}^K a_k^{(t)}$.

906 • $a_k^{(t+1)} = a_k^{(t)} + \eta \cdot \tilde{\Theta}(1) \cdot \left(\exp(-\|\mathbf{a}^{(1)}\|_1 \cdot \tilde{\Theta}(1)) \right), \forall k \in [K]$.

907 Note that $a_k^{(0)} = 0$ for all $k \in [d]$ and therefore $g_{\mathbf{a}^{(0)}}(\mathbf{x}'_i) = 0, \ell'(y'_i \cdot g_{\mathbf{a}^{(0)}}(\mathbf{x}'_i)) = \ell'(0) = -1/2$,

$$a_k^{(1)} = a_k^{(0)} - \eta \cdot \frac{1}{n_1} \sum_{i=1}^{n_1} \ell'(y'_i \cdot g_{\mathbf{a}^{(0)}}(\mathbf{x}'_i)) \cdot y'_i f_{\mathbf{w}_k^{(T_0^k)}}(\mathbf{x}'_i)$$

$$= a_k^{(0)} + \eta \cdot \frac{1}{2n_1} \sum_{i=1}^{n_1} y'_i f_{\mathbf{w}_k^{(T_0^k)}}(\mathbf{x}'_i) = \eta \cdot \frac{1}{2n_1} \sum_{i=1}^{n_1} y'_i f_{\mathbf{w}_k^{(T_0^k)}}(\mathbf{x}'_i) \text{ for all } k \in [K].$$

908 Note that the accuracy of the k -th pseudo-labeler $p_k > 1/2$, according to the proof of Lemma C.22,
909 we have

$$\begin{aligned} f_{\mathbf{w}_k^{(T_0^k)}}(\mathbf{x}'_i) &= \sum_{j=1}^m \left[\sigma(\langle \mathbf{w}_{k,j}^{(T_0^k)}, y'_i \cdot \mathbf{v} \rangle) + \sigma(\langle \mathbf{w}_{k,j}^{(T_0^k)}, \boldsymbol{\xi}'_i \rangle) \right] \\ &\quad - \sum_{j=m+1}^{2m} \left[\sigma(\langle \mathbf{w}_{k,j}^{(T_0^k)}, y'_i \cdot \mathbf{v} \rangle) + \sigma(\langle \mathbf{w}_{k,j}^{(T_0^k)}, \boldsymbol{\xi}'_i \rangle) \right] \\ &= y'_i \cdot \tilde{\Theta}((\hat{\Lambda}_{y'_i}^{(T_0^k)})^q), \end{aligned}$$

910 for all $k \in [K]$. Therefore

$$a_k^{(1)} = \eta \cdot \frac{1}{2n_1} \sum_{i=1}^{n_1} y'_i f_{\mathbf{w}_k^{(T_0^k)}}(\mathbf{x}'_i) \geq \frac{\eta}{2} \cdot \tilde{\Theta}((\hat{\Lambda}_{y'_i}^{(T_0^k)})^q) > 0, \forall k \in [K].$$

911 It follows that

$$\|\mathbf{a}^{(t)}\|_1 = \sum_{i=1}^K |a_i^{(t)}| = \sum_{i=1}^K a_i^{(t)}.$$

912 Note that

$$\begin{aligned} y'_i \cdot g_{\mathbf{a}^{(1)}}(\mathbf{x}'_i) &= y'_i \cdot \sum_{k=1}^K a_k^{(1)} f_{\mathbf{w}_k^{(T_0^k)}}(\mathbf{x}'_i) \\ &= \sum_{k=1}^K a_k^{(1)} \cdot \left(y'_i \cdot f_{\mathbf{w}_k^{(T_0^k)}}(\mathbf{x}'_i) \right) \\ &= \sum_{k=1}^K a_k^{(1)} \cdot \tilde{\Theta}((\hat{\Lambda}_{y'_i}^{(T_0^k)})^q) \\ &= \sum_{k=1}^K a_k^{(1)} \cdot \tilde{\Theta}(1) \\ &= \|\mathbf{a}^{(1)}\|_1 \cdot \tilde{\Theta}(1). \end{aligned} \tag{C.47}$$

913 This leads to

$$\begin{aligned} \ell'(y'_i \cdot g_{\mathbf{a}^{(1)}}(\mathbf{x}'_i)) &= -\frac{\exp(-y'_i \cdot g_{\mathbf{a}^{(1)}}(\mathbf{x}'_i))}{1 + \exp(-y'_i \cdot g_{\mathbf{a}^{(1)}}(\mathbf{x}'_i))} \\ &= -c \cdot \left(\exp(-y'_i \cdot g_{\mathbf{a}^{(1)}}(\mathbf{x}'_i)) \right) \\ &= -c \cdot \left(\exp(-\|\mathbf{a}^{(1)}\|_1 \cdot \tilde{\Theta}(1)) \right), \end{aligned}$$

914 where the second equality is due to $y'_i \cdot g_{\mathbf{a}^{(1)}}(\mathbf{x}'_i) > 0$, $\exp(-y'_i \cdot g_{\mathbf{a}^{(1)}}(\mathbf{x}'_i)) < 1$ and $c \in (1/2, 1)$;
915 the last equality is due to (C.47). It follows that

$$\begin{aligned} a_k^{(2)} &= a_k^{(1)} - \eta \cdot \frac{1}{n_1} \sum_{i=1}^{n_1} \ell'(y'_i \cdot g_{\mathbf{a}^{(1)}}(\mathbf{x}'_i)) \cdot y'_i f_{\mathbf{w}_k^{(T_0^k)}}(\mathbf{x}'_i) \\ &= a_k^{(1)} + \eta \cdot c \cdot \tilde{\Theta}(1) \cdot \left(\exp(-\|\mathbf{a}^{(1)}\|_1 \cdot \tilde{\Theta}(1)) \right), \forall k \in [K] \end{aligned}$$

916 where $c \in (1/2, 1)$. By then, we have already proved the induction hypothesis of $t = 1$.

917 Next, assume the induction hypotheses hold for t . For $t + 1$, we have

$$a_k^{(t+1)} = a_k^{(t)} - \eta \cdot \frac{1}{n_1} \sum_{i=1}^{n_1} \underbrace{\ell'(y_i' \cdot g_{\mathbf{a}^{(t)}}(\mathbf{x}_i'))}_{<0} \cdot \underbrace{y_i' f_{\mathbf{W}_k^{(T_0^k)}}(\mathbf{x}_i')}_{>0} > a_k^{(t)} > 0.$$

918 And it follows that

$$\|\mathbf{a}^{(t+1)}\|_1 = \sum_{i=1}^K a_k^{(t+1)} \text{ and } y_i' \cdot g_{\mathbf{a}^{(t+1)}}(\mathbf{x}_i') = \|\mathbf{a}^{(t+1)}\|_1 \cdot \tilde{\Theta}(1), \quad (\text{C.48})$$

919 leading to

$$\ell'(y_i' \cdot g_{\mathbf{a}^{(t+1)}}(\mathbf{x}_i')) = -c \cdot \left(\exp(-\|\mathbf{a}^{(t+1)}\|_1 \cdot \tilde{\Theta}(1)) \right), c \in (1/2, 1),$$

920 and

$$a_k^{(t+2)} = a_k^{(t+1)} + \eta \cdot \tilde{\Theta}(1) \cdot \left(\exp(-\|\mathbf{a}^{(t+1)}\|_1 \cdot \tilde{\Theta}(1)) \right), \forall k \in [K].$$

921 This indicates that if induction hypotheses hold for t , then they holds for $t + 1$.

922 Adding up $k \in [K]$, we can obtain

$$\|\mathbf{a}^{(t+1)}\|_1 = \|\mathbf{a}^{(t)}\|_1 + \eta \cdot \tilde{\Theta}(1) \cdot \exp(-\tilde{\Theta}(1) \cdot \|\mathbf{a}^{(t)}\|_1) \quad (\text{C.49})$$

923 According to Lemma C.24, we know that $\|\mathbf{a}^{(t)}\|_1 = \log t / \tilde{\Theta}(1) \{\pm \text{lower order terms w.r.t. } t\}$. \square

924 The following lemma gives the convergence guarantee of downstream task:

925 **Lemma C.25.** (Convergence Guarantee) For any learning rate $\eta > 0$,

$$\|\nabla_{\mathbf{a}} L_{S'}(\mathbf{a}^{(t)})\|_1 \leq \frac{\tilde{\Theta}(1)}{\eta \cdot \tilde{\Theta}(1) \cdot t + 1} \text{ and } \nabla_{\mathbf{a}}^2 L_S(\mathbf{a}) \succeq 0 \text{ for any } \mathbf{a} \in \mathbb{R}^d,$$

926 which means within polynomial steps, gradient descent is guaranteed to find a point with small
927 gradient.

928 *Proof of Lemma C.25.* Note that

$$\begin{aligned} \|\nabla_{\mathbf{a}} L_{S'}(\mathbf{a}^{(t)})\|_1 &= \sum_{k=1}^K |\partial_{a_k} L_{S'}(\mathbf{a}^{(t)})| \\ &= - \sum_{k=1}^K \partial_{a_k} L_{S'}(\mathbf{a}^{(t)}) \\ &= \sum_{k=1}^K \frac{a_k^{(t+1)} - a_k^{(t)}}{\eta} \\ &= \frac{\|\mathbf{a}^{(t+1)}\|_1 - \|\mathbf{a}^{(t)}\|_1}{\eta}, \end{aligned}$$

929 then according to Lemma C.24 and (C.49), we know

$$\|\mathbf{a}^{(t+1)}\|_1 - \|\mathbf{a}^{(t)}\|_1 \leq \frac{\eta \cdot \tilde{\Theta}(1)}{\eta \cdot \tilde{\Theta}(1) \cdot t + 1}. \quad (\text{C.50})$$

930 And it follows that

$$\|\nabla_{\mathbf{a}} L_{S'}(\mathbf{a}^{(t)})\|_1 \leq \frac{\tilde{\Theta}(1)}{\eta \cdot \tilde{\Theta}(1) \cdot t + 1},$$

931 which shows that within polynomial steps, gradient descent is guaranteed to find a point with small
932 gradient.

933 Note that

$$\partial_{a_k} L_{S'}(\mathbf{a}) = \frac{1}{n_1} \sum_{i=1}^{n_1} \ell'(y'_i \cdot g_{\mathbf{a}^{(t)}}(\mathbf{x}'_i)) \cdot y'_i f_{\mathbf{w}_k^{(T_0^k)}}(\mathbf{x}'_i),$$

934

$$\partial_{a_k} \partial_{a_j} L_{S'}(\mathbf{a}) = \frac{1}{n_1} \sum_{i=1}^{n_1} \ell''(y'_i \cdot g_{\mathbf{a}^{(t)}}(\mathbf{x}'_i)) \cdot \left(f_{\mathbf{w}_k^{(T_0^k)}}(\mathbf{x}_i) \cdot f_{\mathbf{w}_j^{(T_0^j)}}(\mathbf{x}_i) \right) \text{ for all } k, j \in [K],$$

935 Denote $\left[f_{\mathbf{w}_1^{(T_0^1)}}(\mathbf{x}'_i), \dots, f_{\mathbf{w}_K^{(T_0^K)}}(\mathbf{x}'_i) \right]^\top$ as $\mathbf{f}_{\mathbf{w}^*}(\mathbf{x}'_i)$, then

$$\nabla_{\mathbf{a}}^2 L_S(\mathbf{a}) = \frac{1}{n_1} \sum_{i=1}^{n_1} \ell''(y'_i \cdot g_{\mathbf{a}^{(t)}}(\mathbf{x}'_i)) \cdot \left(\mathbf{f}_{\mathbf{w}^*}(\mathbf{x}'_i) \cdot \mathbf{f}_{\mathbf{w}^*}(\mathbf{x}'_i)^\top \right).$$

936 Note that $\mathbf{f}_{\mathbf{w}^*}(\mathbf{x}'_i) \cdot \mathbf{f}_{\mathbf{w}^*}(\mathbf{x}'_i)^\top$ is a non-negative definite matrix, $\ell''(z) = \exp(-z)/(1 +$
 937 $\exp(-z))^2 > 0$ and the fact that sum of non-negative definite matrices is still a non-negative
 938 definite matrix, it follows that $\nabla_{\mathbf{a}}^2 L_S(\mathbf{a}) \succeq 0$. \square

939 **Theorem C.26** (Restatement of Theorem 3.3). Under semi-supervised learning setting, for down-
 940 stream task, suppose K early stopped classifiers $\{f_{\mathbf{w}_k^*}\}_{k=1}^K$ are obtained after the pre-training of
 941 KK CNN models finished, and after $T_{\text{dt}} = \Theta(d^{0.1}/\eta)$ iterations with learning rate $\eta = \Theta(1)$,
 942 then we can find a linear model $\mathbf{a}^{(T_{\text{dt}})}$, which satisfies: Both test error and loss are nearly 0, i.e.
 943 $\mathbb{P}_{(\mathbf{x}, y) \sim \mathcal{D}}[y \cdot g_{\mathbf{a}^{(T_{\text{dt}})}}(\mathbf{x}) \leq 0] = o(1)$, $L_{\mathcal{D}}(\ell(y \cdot g_{\mathbf{a}^{(T_{\text{dt}})}}(\mathbf{x}))) = o(1)$.

944 *Proof of Theorem C.26.* For test error, we have

$$\begin{aligned} \mathbb{P}_{(\mathbf{x}, y) \sim \mathcal{D}}[y \cdot g_{\mathbf{a}^{(T_{\text{dt}})}}(\mathbf{x}) \leq 0] &= \mathbb{P}_{(\mathbf{x}, y) \sim \mathcal{D}} \left[\sum_{k=1}^K a_k^{(T_{\text{dt}})} \cdot (y \cdot f_{\mathbf{w}_k^*}(\mathbf{x})) \leq 0 \right] \\ &= \mathbb{P}_{(\mathbf{x}, y) \sim \mathcal{D}} \left[\sum_{k=1}^K a_k^{(T_{\text{dt}})} \cdot \tilde{\Theta}(1) \leq 0 \right] = o(1) \end{aligned}$$

945 where the last equality is due to $a_k^{(T_{\text{dt}})} > 0$ according to Lemma C.23.

946 For test loss, we have

$$L_{\mathcal{D}}(\ell(y \cdot g_{\mathbf{a}^{(T_{\text{dt}})}}(\mathbf{x}))) = \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}}[\ell(y \cdot g_{\mathbf{a}^{(T_{\text{dt}})}}(\mathbf{x}))],$$

947 i.e., we estimate for newly generated data (\mathbf{x}, y) the magnitude of $\ell(y \cdot g_{\mathbf{a}^{(t)}}(\mathbf{x}))$. In order to do so,
 948 we will first estimate $\ell(y'_i \cdot g_{\mathbf{a}^{(t)}}(\mathbf{x}'_i))$. Then, we will show that $\ell(y \cdot g_{\mathbf{a}^{(t)}}(\mathbf{x}))$ and $\ell(y'_i \cdot g_{\mathbf{a}^{(t)}}(\mathbf{x}'_i))$
 949 nearly equal to each other.

950 According to the update rule of $a_k^{(t)}$, we have

$$a_k^{(t+1)} = a_k^{(t)} - \eta \cdot \frac{1}{n_1} \sum_{i=1}^{n_1} \ell'(y'_i \cdot g_{\mathbf{a}^{(t)}}(\mathbf{x}'_i)) \cdot y'_i f_{\mathbf{w}_k^{(T_0^k)}}(\mathbf{x}'_i).$$

951 Adding up the above equation for $k \in [K]$, we obtain

$$\|\mathbf{a}^{(t+1)}\|_1 = \|\mathbf{a}^{(t)}\|_1 - \eta \cdot \frac{1}{n_1} \sum_{i=1}^{n_1} \ell'(y'_i \cdot g_{\mathbf{a}^{(t)}}(\mathbf{x}'_i)) \cdot y'_i \sum_{k=1}^K f_{\mathbf{w}_k^{(T_0^k)}}(\mathbf{x}'_i).$$

952 And according to (C.50), we have

$$\|\mathbf{a}^{(t+1)}\|_1 - \|\mathbf{a}^{(t)}\|_1 \leq \frac{\eta \cdot \tilde{\Theta}(1)}{\eta \cdot \tilde{\Theta}(1) \cdot t + 1},$$

953 therefore it follows that

$$-\frac{1}{n_1} \sum_{i=1}^{n_1} \ell'(y'_i \cdot g_{\mathbf{a}^{(t)}}(\mathbf{x}'_i)) \cdot y'_i \sum_{k=1}^K f_{\mathbf{w}_k^{(T_0^k)}}(\mathbf{x}'_i) \leq \frac{\tilde{\Theta}(1)}{\eta \cdot \tilde{\Theta}(1) \cdot t + 1}.$$

954 Note that $K = \Theta(1)$ and for all $k \in [K]$ we have $y'_i \cdot f_{\mathbf{w}_k^{(T_0^k)}}(\mathbf{x}'_i) = \tilde{\Theta}(1)$, it follows that

$$-\frac{1}{n_1} \sum_{i=1}^{n_1} \ell'(y'_i \cdot g_{\mathbf{a}^{(t)}}(\mathbf{x}'_i)) \leq \frac{\tilde{\Theta}(1)}{\eta \cdot \tilde{\Theta}(1) \cdot t + 1}.$$

955 Note that $n_1 = \tilde{\Theta}(1)$ and according to Lemma C.14, there exists a positive sample $(\mathbf{x}_{i_1}, y_{i_1})$ and a
956 negative sample $(\mathbf{x}_{i_2}, y_{i_2})$ with the property that

$$-\ell'(y'_{i_1} \cdot g_{\mathbf{a}^{(t)}}(\mathbf{x}'_{i_1})) \leq \frac{\tilde{\Theta}(1)}{\eta \cdot \tilde{\Theta}(1) \cdot t + 1}, \quad -\ell'(y'_{i_2} \cdot g_{\mathbf{a}^{(t)}}(\mathbf{x}'_{i_2})) \leq \frac{\tilde{\Theta}(1)}{\eta \cdot \tilde{\Theta}(1) \cdot t + 1}.$$

957 Note that $\ell(z) = \log(1 + \exp(-z))$ and $\ell'(z) = -\exp(-z)/(1 + \exp(-z))$, we know that for
958 $z > 0$,

$$\begin{aligned} -\ell'(z) &= c \cdot \exp(-z), \\ \ell(z) &< \exp(-z) = -\ell'(z)/c, c \in (1/2, 1). \end{aligned}$$

959 It follows that

$$\ell(y'_{i_1} \cdot g_{\mathbf{a}^{(t)}}(\mathbf{x}'_{i_1})) \leq \frac{\tilde{\Theta}(1)}{\eta \cdot \tilde{\Theta}(1) \cdot t + 1}, \quad \ell(y'_{i_2} \cdot g_{\mathbf{a}^{(t)}}(\mathbf{x}'_{i_2})) \leq \frac{\tilde{\Theta}(1)}{\eta \cdot \tilde{\Theta}(1) \cdot t + 1}.$$

960 Note that $\ell(z)$ is 1-Lipschitz, we have

$$\begin{aligned} |\ell(y \cdot g_{\mathbf{a}^{(t)}}(\mathbf{x})) - \ell(y'_{i_1} \cdot g_{\mathbf{a}^{(t)}}(\mathbf{x}'_{i_1}))| &\leq |y \cdot g_{\mathbf{a}^{(t)}}(\mathbf{x}) - y'_{i_1} \cdot g_{\mathbf{a}^{(t)}}(\mathbf{x}'_{i_1})|, \\ |\ell(y \cdot g_{\mathbf{a}^{(t)}}(\mathbf{x})) - \ell(y'_{i_2} \cdot g_{\mathbf{a}^{(t)}}(\mathbf{x}'_{i_2}))| &\leq |y \cdot g_{\mathbf{a}^{(t)}}(\mathbf{x}) - y'_{i_2} \cdot g_{\mathbf{a}^{(t)}}(\mathbf{x}'_{i_2})|. \end{aligned} \quad (\text{C.51})$$

961 If $y = 1$, we have

$$\begin{aligned} |y \cdot g_{\mathbf{a}^{(t)}}(\mathbf{x}) - y'_{i_1} \cdot g_{\mathbf{a}^{(t)}}(\mathbf{x}'_{i_1})| &= |g_{\mathbf{a}^{(t)}}(\mathbf{x}) - g_{\mathbf{a}^{(t)}}(\mathbf{x}'_{i_1})| \\ &= \left| \sum_{k=1}^K a_k^{(t)} f_{\mathbf{w}_k^{(T_0^k)}}(\mathbf{x}) - \sum_{k=1}^K a_k^{(t)} f_{\mathbf{w}_k^{(T_0^k)}}(\mathbf{x}'_{i_1}) \right| \\ &= \left| \sum_{k=1}^K a_k^{(t)} \left(f_{\mathbf{w}_k^{(T_0^k)}}(\mathbf{x}) - f_{\mathbf{w}_k^{(T_0^k)}}(\mathbf{x}'_{i_1}) \right) \right|, \end{aligned} \quad (\text{C.52})$$

962 and

$$\begin{aligned}
f_{\mathbf{w}_k^{(T_0^k)}}(\mathbf{x}) - f_{\mathbf{w}_k^{(T_0^k)}}(\mathbf{x}'_{i_1}) &= \sum_{j=1}^m \left[\sigma(\langle \mathbf{w}_j^{(T_0^k)}, \mathbf{v} \rangle) + \sigma(\langle \mathbf{w}_j^{(T_0^k)}, \boldsymbol{\xi} \rangle) \right] \\
&\quad - \sum_{j=m+1}^{2m} \left[\sigma(\langle \mathbf{w}_j^{(T_0^k)}, \mathbf{v} \rangle) + \sigma(\langle \mathbf{w}_j^{(T_0^k)}, \boldsymbol{\xi} \rangle) \right] \\
&\quad - \sum_{j=1}^m \left[\sigma(\langle \mathbf{w}_j^{(T_0^k)}, \mathbf{v} \rangle) + \sigma(\langle \mathbf{w}_j^{(T_0^k)}, \boldsymbol{\xi}'_{i_1} \rangle) \right] \\
&\quad + \sum_{j=m+1}^{2m} \left[\sigma(\langle \mathbf{w}_j^{(T_0^k)}, \mathbf{v} \rangle) + \sigma(\langle \mathbf{w}_j^{(T_0^k)}, \boldsymbol{\xi}'_{i_1} \rangle) \right] \tag{C.53} \\
&= \sum_{j=1}^m \left[\sigma(\langle \mathbf{w}_j^{(T_0^k)}, \boldsymbol{\xi} \rangle) - \sigma(\langle \mathbf{w}_j^{(T_0^k)}, \boldsymbol{\xi}'_{i_1} \rangle) \right] \\
&\quad + \sum_{j=m+1}^{2m} \left[\sigma(\langle \mathbf{w}_j^{(T_0^k)}, \boldsymbol{\xi}'_{i_1} \rangle) - \sigma(\langle \mathbf{w}_j^{(T_0^k)}, \boldsymbol{\xi} \rangle) \right] \\
&= \tilde{O}(d^{-\frac{1}{4}+\epsilon}),
\end{aligned}$$

963 where the last equality is due to (C.45) and Lemma C.4.

964 Plugging (C.53) into (C.52), we have

$$|y \cdot g_{\mathbf{a}^{(t)}}(\mathbf{x}) - y'_{i_1} \cdot g_{\mathbf{a}^{(t)}}(\mathbf{x}'_{i_1})| = \tilde{O}(d^{-\frac{1}{4}+\epsilon}) \cdot \|\mathbf{a}^{(t)}\|_1. \tag{C.54}$$

965 If $y = -1$, we can prove in a similar way that

$$|y \cdot g_{\mathbf{a}^{(t)}}(\mathbf{x}) - y'_{i_2} \cdot g_{\mathbf{a}^{(t)}}(\mathbf{x}'_{i_2})| = \tilde{O}(d^{-\frac{1}{4}+\epsilon}) \cdot \|\mathbf{a}^{(t)}\|_1. \tag{C.55}$$

966 Plugging (C.54) and (C.55) into (C.51), we have

$$\ell(y \cdot g_{\mathbf{a}^{(t)}}(\mathbf{x})) \leq \max \{y'_{i_1} \cdot g_{\mathbf{a}^{(t)}}(\mathbf{x}'_{i_1}), y'_{i_2} \cdot g_{\mathbf{a}^{(t)}}(\mathbf{x}'_{i_2})\} + \tilde{O}(d^{-\frac{1}{4}+\epsilon}) \cdot \|\mathbf{a}^{(t)}\|_1$$

967 According to Lemma C.24 and (C.49), we have $\|\mathbf{a}^{(t)}\|_1 = \log t / \tilde{\Theta}(1) \{\pm \text{lower order terms w.r.t. } t\}$,

968 therefore

$$\ell(y \cdot g_{\mathbf{a}^{(t)}}(\mathbf{x})) \leq \frac{\tilde{\Theta}(1)}{\eta \cdot \tilde{\Theta}(1) \cdot t + 1} + \tilde{O}(d^{-\frac{1}{4}+\epsilon}) \cdot \log t \{\pm \text{lower order terms w.r.t. } t\}$$

969 Taking $\eta = \Theta(1)$ and $T_{dt} = \Theta(d^\alpha / \eta)$ where $\alpha > 0$ is a sufficiently small constant, we know that

$$\begin{aligned}
&L_{\mathcal{D}}(\ell(y \cdot g_{\mathbf{a}^{(T_{dt})}}(\mathbf{x}))) \\
&= \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}}[\ell(y \cdot g_{\mathbf{a}^{(T_{dt})}}(\mathbf{x}))] \\
&\leq \frac{\tilde{\Theta}(1)}{\eta \cdot \tilde{\Theta}(1) \cdot T_{dt} + 1} + \tilde{O}(d^{-\frac{1}{4}+\epsilon}) \cdot \log T_{dt} \{\pm \text{lower order terms w.r.t. } T_{dt}\} + o(1) \\
&= o(1),
\end{aligned}$$

970 which completes the proof. □

971 **D Proof of supervised learning setting**

972 Here we prove Theorem 3.4. First, we give following lemma to facilitate the proof.

973 **Lemma D.1** (Gradient Calculation). The gradient of loss function $L_S(\mathbf{W})$ with respect to weight
 974 parameter \mathbf{w}_j is

$$\nabla_{\mathbf{w}_j} L_{S'}(\mathbf{W}) = -\frac{qu_j}{n_1} \cdot \sum_{i=1}^{n_1} b_i y'_i ([\langle \mathbf{w}_j, y'_i \cdot \mathbf{v} \rangle]_+^{q-1} \cdot y'_i \cdot \mathbf{v} + [\langle \mathbf{w}_j, \boldsymbol{\xi}'_i \rangle]_+^{q-1} \cdot \boldsymbol{\xi}'_i),$$

975 where $u_j := (\mathbb{1}_{[1 \leq j \leq m]} - \mathbb{1}_{[m+1 \leq j \leq 2m]})$ and $-\ell'(y'_i \cdot f_{\mathbf{W}}(\mathbf{x}'_i)) = \exp[-y'_i \cdot f_{\mathbf{W}}(\mathbf{x}'_i)] / (1 +$
 976 $\exp[-y'_i \cdot f_{\mathbf{W}}(\mathbf{x}'_i)])$ is denoted as b_i .

977 *Proof of Lemma D.1.* When $1 \leq j \leq m$,

$$\begin{aligned} \nabla_{\mathbf{w}_j} \ell(y'_i \cdot f_{\mathbf{W}}(\mathbf{x}'_i)) &= \ell'(y'_i \cdot f_{\mathbf{W}}(\mathbf{x}'_i)) \cdot y'_i \cdot \nabla_{\mathbf{w}_j} f_{\mathbf{W}}(\mathbf{x}'_i) \\ &= -b_i \cdot y'_i \cdot \nabla_{\mathbf{w}_j} f_{\mathbf{W}}(\mathbf{x}'_i) \\ &= -b_i y'_i \cdot (\sigma'(\langle \mathbf{w}_j, y'_i \cdot \mathbf{v} \rangle) \cdot y'_i \cdot \mathbf{v} + \sigma'(\langle \mathbf{w}_j, \boldsymbol{\xi}'_i \rangle) \cdot \boldsymbol{\xi}'_i) \\ &= -qb_i y'_i ([\langle \mathbf{w}_j, y'_i \cdot \mathbf{v} \rangle]_+^{q-1} \cdot y'_i \cdot \mathbf{v} + [\langle \mathbf{w}_j, \boldsymbol{\xi}'_i \rangle]_+^{q-1} \cdot \boldsymbol{\xi}'_i) \end{aligned}$$

978 and when $m+1 \leq j \leq 2m$,

$$\nabla_{\mathbf{w}_j} \ell(y'_i \cdot f_{\mathbf{W}}(\mathbf{x}'_i)) = qb_i y'_i ([\langle \mathbf{w}_j, y'_i \cdot \mathbf{v} \rangle]_+^{q-1} \cdot y'_i \cdot \mathbf{v} + [\langle \mathbf{w}_j, \boldsymbol{\xi}'_i \rangle]_+^{q-1} \cdot \boldsymbol{\xi}'_i)$$

979 Combining above two cases, we have

$$\begin{aligned} \nabla_{\mathbf{w}_j} \ell(y'_i \cdot f_{\mathbf{W}}(\mathbf{x}'_i)) &= -q(\mathbb{1}_{[1 \leq j \leq m]} - \mathbb{1}_{[m+1 \leq j \leq 2m]}) b_i y'_i ([\langle \mathbf{w}_j, y'_i \cdot \mathbf{v} \rangle]_+^{q-1} \cdot y'_i \cdot \mathbf{v} + [\langle \mathbf{w}_j, \boldsymbol{\xi}'_i \rangle]_+^{q-1} \cdot \boldsymbol{\xi}'_i) \\ &= -qu_j b_i y'_i ([\langle \mathbf{w}_j, y'_i \cdot \mathbf{v} \rangle]_+^{q-1} \cdot y'_i \cdot \mathbf{v} + [\langle \mathbf{w}_j, \boldsymbol{\xi}'_i \rangle]_+^{q-1} \cdot \boldsymbol{\xi}'_i) \end{aligned}$$

980 and therefore

$$\begin{aligned} \nabla_{\mathbf{w}_j} L_{S'}(\mathbf{W}) &= \frac{1}{n_1} \sum_{i=1}^{n_1} \nabla_{\mathbf{w}_j} L_i(\mathbf{W}) = \frac{1}{n_1} \sum_{i=1}^{n_1} \nabla_{\mathbf{w}_j} \ell(y'_i \cdot f_{\mathbf{W}}(\mathbf{x}'_i)) \\ &= -\frac{qu_j}{n_1} \sum_{i=1}^{n_1} b_i y'_i ([\langle \mathbf{w}_j, y'_i \cdot \mathbf{v} \rangle]_+^{q-1} \cdot y'_i \cdot \mathbf{v} + [\langle \mathbf{w}_j, \boldsymbol{\xi}'_i \rangle]_+^{q-1} \cdot \boldsymbol{\xi}'_i). \end{aligned}$$

981 □

982 *Proof of Theorem 3.4.* Recall the definition of $f_{\mathbf{W}}$ in (2.1) that

$$f_{\mathbf{W}}(\mathbf{x}) = \sum_{j=1}^m [\sigma(\langle \mathbf{w}_j, y \cdot \mathbf{v} \rangle) + \sigma(\langle \mathbf{w}_j, \boldsymbol{\xi} \rangle)] - \sum_{j=m+1}^{2m} [\sigma(\langle \mathbf{w}_j, y \cdot \mathbf{v} \rangle) + \sigma(\langle \mathbf{w}_j, \boldsymbol{\xi} \rangle)].$$

983 Define $\tilde{\mathbf{w}}_j := m^{1/q} \cdot \mathbf{w}_j$, we have

$$\begin{aligned} f_{\mathbf{W}}(\mathbf{x}) &= \sum_{j=1}^m [\sigma(\langle m^{-1/q} \cdot \tilde{\mathbf{w}}_j, y \cdot \mathbf{v} \rangle) + \sigma(\langle m^{-1/q} \cdot \tilde{\mathbf{w}}_j, \boldsymbol{\xi} \rangle)] \\ &\quad - \sum_{j=m+1}^{2m} [\sigma(\langle m^{-1/q} \cdot \tilde{\mathbf{w}}_j, y \cdot \mathbf{v} \rangle) + \sigma(\langle m^{-1/q} \cdot \tilde{\mathbf{w}}_j, \boldsymbol{\xi} \rangle)] \\ &= \frac{1}{m} \sum_{j=1}^m [\sigma(\langle \tilde{\mathbf{w}}_j, y \cdot \mathbf{v} \rangle) + \sigma(\langle \tilde{\mathbf{w}}_j, \boldsymbol{\xi} \rangle)] - \frac{1}{m} \sum_{j=m+1}^{2m} [\sigma(\langle \tilde{\mathbf{w}}_j, y \cdot \mathbf{v} \rangle) + \sigma(\langle \tilde{\mathbf{w}}_j, \boldsymbol{\xi} \rangle)] \\ &:= f_{\tilde{\mathbf{W}}}(\mathbf{x}). \end{aligned}$$

984 Since the standard deviation of Gaussian initialization of \mathbf{w}_j is σ_0 and note that $\tilde{\mathbf{w}}_j := m^{1/q} \cdot \mathbf{w}_j$,
 985 the standard deviation of Gaussian initialization of $\tilde{\mathbf{w}}_j$ is $m^{1/q} \sigma_0 := \tilde{\sigma}_0$.

986 On the other hand, note that the update rule of $\mathbf{w}_j^{(t)}$ is $\mathbf{w}_j^{(t+1)} = \mathbf{w}_j^{(t)} - \eta \cdot \nabla_{\mathbf{w}_j} L_{S'}(\mathbf{W}^{(t)})$, and in
 987 Lemma D.1, we have

$$\nabla_{\mathbf{w}_j} L_{S'}(\mathbf{W}) = -\frac{qu_j}{n_1} \cdot \sum_{i=1}^{n_1} b_i y'_i ([\langle \mathbf{w}_j, y'_i \cdot \mathbf{v} \rangle]_+^{q-1} \cdot y'_i \cdot \mathbf{v} + [\langle \mathbf{w}_j, \boldsymbol{\xi}'_i \rangle]_+^{q-1} \cdot \boldsymbol{\xi}'_i).$$

988 It follows that

$$\mathbf{w}_j^{(t+1)} = \mathbf{w}_j^{(t)} + \frac{q\eta u_j}{n_1} \cdot \sum_{i=1}^{n_1} b_i^{(t)} y'_i ([\langle \mathbf{w}_j^{(t)}, y'_i \cdot \mathbf{v} \rangle]_+^{q-1} \cdot y'_i \cdot \mathbf{v} + [\langle \mathbf{w}_j^{(t)}, \boldsymbol{\xi}'_i \rangle]_+^{q-1} \cdot \boldsymbol{\xi}'_i). \quad (\text{D.1})$$

989 By plugging $\mathbf{w}_j = m^{-1/q} \cdot \tilde{\mathbf{w}}_j$ into (D.1), we have

$$\tilde{\mathbf{w}}_j^{(t+1)} = \tilde{\mathbf{w}}_j^{(t)} + \frac{q\eta m^{-\frac{1}{q}} u_j}{n_1} \cdot \sum_{i=1}^{n_1} b_i^{(t)} y'_i ([\langle \tilde{\mathbf{w}}_j^{(t)}, y'_i \cdot \mathbf{v} \rangle]_+^{q-1} \cdot y'_i \cdot \mathbf{v} + [\langle \tilde{\mathbf{w}}_j^{(t)}, \boldsymbol{\xi}'_i \rangle]_+^{q-1} \cdot \boldsymbol{\xi}'_i)$$

990 Assume $\tilde{\eta} = m^{-\frac{1}{q}} \eta$, we have $\tilde{\mathbf{w}}_j^{(t+1)} = \tilde{\mathbf{w}}_j^{(t)} - \tilde{\eta} \cdot \nabla_{\tilde{\mathbf{w}}_j} L_{S'}(\tilde{\mathbf{W}}^{(t)})$. Therefore, our data model and
 991 training algorithm is equivalent to the model and algorithm below:

$$\begin{aligned} f_{\tilde{\mathbf{W}}+1}(\mathbf{x}) &= \frac{1}{m} \sum_{j=1}^m \left[\sigma(\langle \tilde{\mathbf{w}}_j, y \cdot \mathbf{v} \rangle) + \sigma(\langle \tilde{\mathbf{w}}_j, \boldsymbol{\xi} \rangle) \right], \\ f_{\tilde{\mathbf{W}}-1}(\mathbf{x}) &= \frac{1}{m} \sum_{j=m+1}^{2m} \left[\sigma(\langle \tilde{\mathbf{w}}_j, y \cdot \mathbf{v} \rangle) + \sigma(\langle \tilde{\mathbf{w}}_j, \boldsymbol{\xi} \rangle) \right], \\ f_{\tilde{\mathbf{W}}}(\mathbf{x}) &= f_{\tilde{\mathbf{W}}+1}(\mathbf{x}) - f_{\tilde{\mathbf{W}}-1}(\mathbf{x}), \end{aligned}$$

992 and we use gradient decent with learning rate $\tilde{\eta}$ and cross-entropy loss to optimize such a data model,
 993 i.e.

$$\tilde{\mathbf{w}}_0^{(t)} \sim \mathcal{N}(\mathbf{0}, \tilde{\sigma}_0^2 \mathbf{I}_d), \tilde{\mathbf{w}}_j^{(t+1)} = \tilde{\mathbf{w}}_j^{(t)} - \tilde{\eta} \cdot \nabla_{\tilde{\mathbf{w}}_j} L_{S'}(\tilde{\mathbf{W}}^{(t)}), L_{S'}(\tilde{\mathbf{W}}^{(t)}) = \sum_{i=1}^{n_1} \ell(y'_i \cdot f_{\tilde{\mathbf{W}}}(\mathbf{x}'_i)),$$

994 where $\ell(z) = \log(1 + \exp(-z))$, $\tilde{\sigma}_0 = m^{1/q} \sigma_0$. Note that the new model meets the one used in Cao
 995 et al. (2022). To leverage their result, we introduce condition 4.3 from Cao et al. (2022) and verify
 996 that the new model meets the new condition.

997 **Condition D.2** (Condition 4.2 in Cao et al. (2022)). Dimension d is sufficiently large that
 998 $d = \tilde{\Omega}(m^{2\vee[4/(q-2)]} n^{4\vee[(2q-2)/(q-2)]})$. Training sample size n and neural network width m
 999 satisfy $n, m = \Omega(\text{polylog}(d))$. Learning rate η satisfies $\eta \leq \tilde{O}(\min\{\|\mathbf{v}\|_2^{-2}, \sigma_p^{-2} d^{-1}\})$. The
 1000 standard deviation of Gaussian initialization σ_0 is approximately chosen such that $\tilde{O}(nd^{-\frac{1}{2}}) \cdot$
 1001 $\min\{(\sigma_p \sqrt{d})^{-1}, \|\mathbf{v}\|_2^{-1}\} \leq \sigma_0 \leq \tilde{O}(m^{-2/(q-2)} n^{-[1/(q-2)]\vee 1}) \cdot \min\{(\sigma_p \sqrt{d})^{-1}, \|\mathbf{v}\|_2^{-1}\}$.

1002 **Theorem D.3** (Theorem 4.4 in Cao et al. (2022)). For any $\epsilon > 0$, let $T = \tilde{\Theta}(\eta^{-1} m \cdot n (\sigma_p \sqrt{d})^{-q} \cdot$
 1003 $\sigma_0^{-(q-2)} + \eta^{-1} \epsilon^{-1} n m^3 d^{-1} \sigma_p^{-2})$. Under Condition D.2, if $n^{-1} \cdot \text{SNR}^{-q} = \tilde{\Omega}(1)$, $\text{SNR} =$
 1004 $\|\mathbf{v}\|_2 / \sigma_p \sqrt{d}$, then with probability at least $1 - d^{-1}$, there exists $0 \leq t \leq T$ such that:

- 1005 1. The training loss converges to δ , i.e., $L_S(\mathbf{W}^{(t)}) \leq \delta$.
- 1006 2. The trained CNN has a constant order test loss: $L_{\mathcal{D}}(\mathbf{W}^{(t)}) = \Theta(1)$.

1007 Note that in our setting, $m = \Theta(\text{polylog}(d))$, $n_1 = \tilde{\Theta}(1)$, $\|\mathbf{v}\|_2 = \Theta(d^{\frac{1}{2}})$, $\tilde{\sigma}_0 = m^{1/q} \sigma_0$, $\sigma_0 =$
 1008 $\Theta(d^{-\frac{3}{4}}) \sigma_p = \Theta(d^{0.01})$, $\tilde{\eta} = m^{-\frac{1}{q}} \eta$ and $\eta = O(d^{-1-2\epsilon})$, it's not difficult to verify that Condition
 1009 D.2 holds. Besides, $\text{SNR} = d^{-0.01}$, $n^{-1} \cdot \text{SNR}^{-q} = \tilde{\Theta}(d^{q\epsilon}) = \tilde{\Omega}(1)$. Therefore, the conclusion of
 1010 Theorem D.3 holds for

$$\begin{aligned} T &= \tilde{\Theta}(\tilde{\eta}^{-1} m \cdot n (\sigma_p \sqrt{d})^{-q} \cdot \sigma_0^{-(q-2)} + \tilde{\eta}^{-1} \epsilon^{-1} n m^3 d^{-1} \sigma_p^{-2}) \\ &= \tilde{\Theta}(\tilde{\eta}^{-1} \cdot (d^{1/2+\epsilon})^{-q} \cdot (d^{-3/4})^{-(q-2)} + \tilde{\eta}^{-1} \epsilon^{-1} d^{-1} d^{-2\epsilon}) \end{aligned}$$

$$\begin{aligned}
&= \tilde{\Theta}(\tilde{\eta}^{-1} \cdot d^{(1/4-\epsilon)q-3/2} + \tilde{\eta}^{-1}\epsilon^{-1}d^{-1-2\epsilon}) \\
&= \tilde{\Theta}(\eta^{-1} \cdot d^{(1/4-\epsilon)q-3/2}).
\end{aligned}$$

1011

□

1012 E Auxiliary Lemmas

1013 For the estimation of $\bar{\Lambda}^{(0)}$ and $\hat{\Lambda}^{(0)}$, we introduce the following lemma.

1014 **Lemma E.1** (Borell-TIS inequality). Let X be a centered Gaussian on \mathbb{R}^m and set $\sigma_X^2 :=$
1015 $\max_{i \in [m]} \mathbb{E}(X_i^2)$. Then for each $t > 0$,

$$\mathbb{P}\left(\left|\max_{i \in [m]} X_i - \mathbb{E}\left(\max_{i \in [m]} X_i\right)\right| > t\right) \leq 2e^{-\frac{t^2}{2\sigma_X^2}}.$$

1016 For the expectation of $\hat{\Lambda}_r^{(0)}$ and $\bar{\Lambda}_r^{(0)}$, we give the following lemma.

1017 **Lemma E.2.** Let $Y = \max_{1 \leq i \leq m} X_i$, where $X_i \sim \mathcal{N}(0, \sigma^2)$ are i.i.d. random variables. Then

$$\frac{1}{\sqrt{\pi \log 2}} \sigma \sqrt{\log m} \leq \mathbb{E}[Y] \leq \sqrt{2} \sigma \sqrt{\log m}.$$

1018 For the estimation of $\|\xi_i\|_2^2$ and $\langle \xi_i, \xi_l \rangle$, we introduce following lemma.

1019 **Lemma E.3** (Lemma B.2 in Cao et al. (2022)). Suppose that $\delta > 0$ and $d = \Omega(\log(4n/\delta))$. Then
1020 with probability at least $1 - \delta$,

$$\begin{aligned}
\sigma_p^2 d/2 &\leq \|\xi_i\|_2^2 \leq 3\sigma_p^2 d/2, \\
|\langle \xi_i, \xi_l \rangle| &\leq 2\sigma_p^2 \cdot \sqrt{d \log(4n^2/\delta)},
\end{aligned}$$

1021 for all $i, l \in [n], i \neq l$.

1022 Besides, we introduce the following lemma about the tensor power method.

1023 **Lemma E.4.** Consider an increasing sequence $x_t \geq 0$ defined as $x_{t+1} = x_t + \eta \cdot C_t x_t^{q-1}$, and
1024 $C_1 \leq C_t \leq C_2$ for all $t > 0$, then we have for $A > x_0$, every $\delta > 0$, and every $\eta > 0$:

$$\sum_{t \geq 0, x_t \leq A} \eta \leq \frac{\delta}{(1 - (1 + \delta)^{-(q-2)})x_0^{q-2}C_1} + \eta \cdot \frac{C_2}{C_1}(1 + \delta)^{q-1} \left(1 + \frac{\log(A/x_0)}{\log(1 + \delta)}\right),$$

1025

$$\sum_{t \geq 0, x_t \leq A} \eta \geq \frac{\delta(1 - (x_0/A)^{q-2})}{(1 + \delta)^{q-1}(1 - (1 + \delta)^{-(q-2)})x_0^{q-2}C_2} - \eta \cdot (1 + \delta)^{-(q-1)} \left(1 + \frac{\log(A/x_0)}{\log(1 + \delta)}\right).$$

1026 *Proof of Lemma E.4.* For every $g = 0, 1, 2, \dots$, let τ_g be the first iteration such that $x_t \geq (1 + \delta)^g x_0$.
1027 Let b be the smallest integer such that $(1 + \delta)^b x_0 \geq A$. By the definition of τ_g , we have $x_t \in$
1028 $[(1 + \delta)^g x_0, (1 + \delta)^{g+1} x_0)$ for all $t \in [\tau_g, \tau_{g+1})$ and $x_{\tau_{g+1}} \geq (1 + \delta)^{g+1} x_0$, $x_{\tau_g - 1} < (1 + \delta)^g x_0$,
1029 leading to

$$\begin{aligned}
\sum_{t \in [\tau_g, \tau_{g+1})} \eta \cdot C_t [(1 + \delta)^g x_0]^{q-1} &\leq x_{\tau_{g+1}} - x_{\tau_g} = \sum_{t \in [\tau_g, \tau_{g+1})} (x_{t+1} - x_t) \\
&= \sum_{t \in [\tau_g, \tau_{g+1})} \eta \cdot C_t x_t^{q-1} \leq \sum_{t \in [\tau_g, \tau_{g+1})} \eta \cdot C_t [(1 + \delta)^{g+1} x_0]^{q-1},
\end{aligned}$$

1030 following lower bound for $x_{\tau_{g+1}} - x_{\tau_g}$:

$$\begin{aligned}
x_{\tau_{g+1}} - x_{\tau_g} &= x_{\tau_{g+1}} - x_{\tau_g - 1} - \eta \cdot C_{\tau_g - 1} x_{\tau_g - 1}^{q-1} \\
&\geq (1 + \delta)^{g+1} x_0 - (1 + \delta)^g x_0 - \eta \cdot C_{\tau_g - 1} [(1 + \delta)^g x_0]^{q-1} \\
&= \delta(1 + \delta)^g x_0 - \eta \cdot C_{\tau_g - 1} (1 + \delta)^{(q-1)g} x_0^{q-1},
\end{aligned}$$

1031 and following upper bound for $x_{\tau_{g+1}} - x_{\tau_g}$:

$$\begin{aligned} x_{\tau_{g+1}} - x_{\tau_g} &= x_{\tau_{g+1}-1} + \eta \cdot C_{\tau_{g+1}-1} x_{\tau_{g+1}-1}^{q-1} - x_{\tau_g} \\ &\leq (1 + \delta)^{g+1} x_0 + \eta \cdot C_{\tau_{g+1}-1} [(1 + \delta)^{(g+1)} x_0]^{q-1} - (1 + \delta)^g x_0 \\ &= \delta(1 + \delta)^g x_0 + \eta \cdot C_{\tau_{g+1}-1} (1 + \delta)^{(q-1)(g+1)} x_0^{q-1}. \end{aligned}$$

1032 Therefore,

$$\sum_{t \in [\tau_g, \tau_{g+1})} \eta \cdot C_t [(1 + \delta)^g x_0]^{q-1} \leq \delta(1 + \delta)^g x_0 + \eta \cdot C_{\tau_{g+1}-1} (1 + \delta)^{(q-1)(g+1)} x_0^{q-1},$$

1033

$$\sum_{t \in [\tau_g, \tau_{g+1})} \eta \cdot C_t [(1 + \delta)^{g+1} x_0]^{q-1} \geq \delta(1 + \delta)^g x_0 - \eta \cdot C_{\tau_g-1} (1 + \delta)^{(q-1)g} x_0^{q-1}.$$

1034 These imply that

$$\sum_{t \in [\tau_g, \tau_{g+1})} \eta \cdot C_t \leq \frac{\delta}{(1 + \delta)^{(q-2)g} x_0^{q-2}} + \eta \cdot C_{\tau_{g+1}-1} (1 + \delta)^{q-1} \leq \frac{\delta}{(1 + \delta)^{(q-2)g} x_0^{q-2}} + \eta \cdot C_2 (1 + \delta)^{q-1},$$

1035

$$\begin{aligned} \sum_{t \in [\tau_g, \tau_{g+1})} \eta \cdot C_t &\geq \frac{\delta}{(1 + \delta)^{(q-2)g+(q-1)} x_0^{q-2}} - \eta \cdot C_{\tau_g-1} (1 + \delta)^{-(q-1)} \\ &\geq \frac{\delta}{(1 + \delta)^{(q-2)g+(q-1)} x_0^{q-2}} - \eta \cdot C_2 (1 + \delta)^{-(q-1)}. \end{aligned}$$

1036 Recall b is the smallest integer such that $(1 + \delta)^b x_0 \geq A$, so we can calculate that

$$\begin{aligned} \sum_{t \geq 0, x_t \leq A} \eta \cdot C_t &\leq \sum_{g=0}^{b-1} \frac{\delta}{(1 + \delta)^{(q-2)g} x_0^{q-2}} + \eta \cdot C_2 (1 + \delta)^{q-1} b \\ &= \frac{\delta(1 - (1 + \delta)^{-(q-2)b})}{(1 - (1 + \delta)^{-(q-2)}) x_0^{q-2}} + \eta \cdot C_2 (1 + \delta)^{q-1} b \\ &\leq \frac{\delta}{(1 - (1 + \delta)^{-(q-2)}) x_0^{q-2}} + \eta \cdot C_2 (1 + \delta)^{q-1} b, \end{aligned}$$

1037 and

$$\begin{aligned} \sum_{t \geq 0, x_t \leq A} \eta \cdot C_t &\geq \sum_{g=0}^{b-1} \frac{\delta}{(1 + \delta)^{(q-2)g+(q-1)} x_0^{q-2}} - \eta \cdot C_2 (1 + \delta)^{-(q-1)} b \\ &= \frac{\delta(1 - (1 + \delta)^{-(q-2)b})}{(1 + \delta)^{q-1} (1 - (1 + \delta)^{-(q-2)}) x_0^{q-2}} - \eta \cdot C_2 (1 + \delta)^{-(q-1)} b \\ &\geq \frac{\delta(1 - (x_0/A)^{q-2})}{(1 + \delta)^{q-1} (1 - (1 + \delta)^{-(q-2)}) x_0^{q-2}} - \eta \cdot C_2 (1 + \delta)^{-(q-1)} b, \end{aligned}$$

1038 where the last inequality is due to $(1 + \delta)^b x_0 \geq A$.

1039 Note that $(1 + \delta)^{b-1} x_0 < A$, i.e. $b \leq 1 + \frac{\log(A/x_0)}{\log(1+\delta)}$, therefore

$$\sum_{t \geq 0, x_t \leq A} \eta \cdot C_t \leq \frac{\delta}{(1 - (1 + \delta)^{-(q-2)}) x_0^{q-2}} + \eta \cdot C_2 (1 + \delta)^{q-1} \left(1 + \frac{\log(A/x_0)}{\log(1 + \delta)}\right),$$

1040

$$\sum_{t \geq 0, x_t \leq A} \eta \cdot C_t \geq \frac{\delta(1 - x_0/A)}{(1 + \delta)^{q-1} (1 - (1 + \delta)^{-(q-2)}) x_0^{q-2}} - \eta \cdot C_2 (1 + \delta)^{-(q-1)} \left(1 + \frac{\log(A/x_0)}{\log(1 + \delta)}\right),$$

1041 Note that $C_1 \leq C_t \leq C_2$, we have

$$\sum_{t \geq 0, x_t \leq A} \eta \leq \frac{\delta}{(1 - (1 + \delta)^{-(q-2)})x_0^{q-2}C_1} + \eta \cdot \frac{C_2}{C_1}(1 + \delta)^{q-1} \left(1 + \frac{\log(A/x_0)}{\log(1 + \delta)}\right),$$

1042

$$\sum_{t \geq 0, x_t \leq A} \eta \geq \frac{\delta(1 - (x_0/A)^{q-2})}{(1 + \delta)^{q-1}(1 - (1 + \delta)^{-(q-2)})x_0^{q-2}C_2} - \eta \cdot (1 + \delta)^{-(q-1)} \left(1 + \frac{\log(A/x_0)}{\log(1 + \delta)}\right).$$

1043

□