
Information-Theoretic Generalization Bounds for Transductive Learning and its Applications

Huayi Tang Yong Liu*

Gaoling School of Artificial Intelligence, Renmin University of China, Beijing, China
Beijing Key Laboratory of Big Data Management and Analysis Methods, Beijing, China

Abstract

We develop data-dependent and algorithm-dependent generalization bounds for transductive learning algorithms based on information theory. We first show that the generalization error for a transductive learning algorithm can be controlled by the mutual information between training labels and the hypothesis. By introducing the concept of transductive supersamples, we go beyond inductive learning setting further derive upper bounds in terms of various information measure. From this results, we derive PAC-Bayesian bounds with faster rate and build connection between generalization and the sharpness of loss landscape under transductive learning setting. Finally, we present the upper bounds for adaptive optimization algorithms and demonstrate the applications of our results on semi-supervised and transductive graph learning scenarios.

1 Introduction

In the standard supervised learning paradigm [107], we receive a set of instances sampling independently from unknown distribution that are composed of features and targets. Our task is to build a learner (or model) by specific optimization algorithm that maps features to corresponding targets based on limited instances. A modern and popular practice of this paradigm is training a deep neural network for image classification via stochastic gradient descent [59]. Generalization, referring to the prediction performance of learner on unseen examples, is one of the fundamental questions in machine learning theory. The past decades have witnessed efforts devoted to characterizing and understanding the generalization of machine learning models. In the category of classical learning theory, generalization is connected to the complexity of hypothesis space [58, 8, 7], the stability of algorithm [13, 92] or the divergence between two probability measure on hypothesis space [95, 70]. Recently, mutual information and its variants are shown to serve as an ideal generalization metric, since they reflect both the impact of dataset and optimization on generalization. Research of this viewpoint originates from work [85, 86, 117], and is further developed by subsequent studies [73, 15, 45, 99, 46, 43, 89, 114]. Despite the diverse forms of these results, they possess a common key insight: the less information on training data (or its selection) a hypothesis (or the variables it induced) reveals, the better generalization it will have.

The aforementioned supervised learning setting are far not enough to cover all machine learning scenarios. Data collected from real-world scenarios could come from various domains, and some examples could lack of targets due to the expensive cost of annotations. This raises the need of exploring new theory to measure the generalization, where the key challenge is relaxing the identical and independent assumptions of instances. We move towards this direction by analyzing a classical but important regime termed as transductive learning [106]. In this learning paradigm, we are provided with a fixed set of instances containing both labeled example and unlabeled examples, and our task is to build a learner that make prediction for those unlabeled ones. Note that under this setting, the examples to be predicted are available to learner. As a comparison, the learning paradigm of supervised learning is inductive, since the examples to be predicted are not available

to learner during training. Popular practices of transductive learning are semi-supervised learning [91, 11, 52, 116] and graph learning [39, 87, 36, 57].

Existing results for transductive generalization bounds fall into three categories: complexity-based bounds derived by VC dimension [22] or transductive Rademacher complexity [30, 105], stability-based bounds [29, 23] and PAC-Bayesian bounds [26, 10]. These results are sufficient to provide learning guarantee for classical learner or algorithms such as transductive support vector machine and unlabeled-labeled representation. However, they are far from the optimal approaches to understanding the generalization behavior of deep transductive model such as Graph Neural Networks (GNNs). The reasons are three folds. First, [32] shows that VC dimension results in trivial generalization error bounds, and transductive Rademacher complexity is algorithm-independent and fail to characterize the impact of optimization algorithms on generalization [101]. Second, stability-based bounds depend on Lipschitz and smoothness constant [21], which is difficulty to compute for deep models [74]. Third, existing transductive PAC-Bayesian bounds are of slow order and not sufficient to provide non-vacuous bounds. Furthermore, it is unclear to what extent they reflect the impact of optimization algorithms. In a nutshell, effort to establish data-dependent and algorithm-dependent generalization bounds for transductive learning is limited.

In this paper, we comprehensively study the generalization theory of transductive learning in the context of information theory. First, we derive upper bounds of transductive generalization gap in expectation and high probability. These results demonstrate that the dependence of output hypothesis on the randomness of training labels serve as a measure to quantify the generalization performance of transductive learner. Second, we introduce the concept of transductive supersamples and establish upper bounds include the conditional mutual information between the indicator random variables and output hypothesis. These bounds are non-vacuous and convenient in computation. Third, by observing the connection between information theory and PAC-Bayesian theory, we give novel transductive PAC-Bayesian bounds that has weaker assumption and sharper bounds. With this result, we further show that the sharpness of loss landscape affects generalization still holds in transductive learning setting, as supported by the empirical evidence in recent study [17]. Forth, we provide the upper bounds for adaptive optimization algorithms composed of trajectory adaptive gradient norm and the sensitivity of the final parameter to noise under transductive setting. Fifth, we illustrates the applications of the established results on semi-supervised learning and graph learning scenarios. Our main contributions are summarized as follows:

- For the first time, we systematically establish information-theoretic generalization bounds for transductive learning and reveal their connection with PAC-Bayesian bounds, which provide new perspective to understanding the generalization transductive learning.
- We propose the concept of transductive supersamples and use it to bridge the gap between supersample setting in inductive learning and transductive learning setting. This also help us obtain the first non-vacuous bound for deep transductive learner.

2 Related Work

Information-theoretic Generalization Theory. [85, 86, 117] are the pioneer work that associates expected generalization error with the mutual information between training examples and algorithm output. The subsequent studies mainly fall into four categories: (i) deriving sharper upper bounds by introducing new information measure [46, 49, 114], problem setting [99, 81, 44] or proof techniques [6, 15, 41, 34, 125, 20], (ii) establishing bounds described by various divergences [64, 109, 31, 4, 3] (iii) applying existing results to establish upper bounds for optimization algorithms such as SGD [74, 113] or SDLD [78, 73, 110], and (iv) extending the theoretical results to diverse scenarios such as meta-learning [54, 83, 18, 55], transfer learning [115, 53, 67, 14], semi-supervised learning [2, 47], self-supervised learning [123] and domain adaption [114]. However, the training and test examples are independent in existing studies, which makes them not applicable on transductive learning. Another related topic is information bottleneck theory [102] and its applications to explaining the representation [103, 82] and generalization [98, 42, 112, 56] of deep neural networks, which are parallel to our work. Interested readers are referred to a recent monograph [50] for more illustrations.

PAC-Bayesian Generalization Theory. The classical results in PAC-Bayesian generalization theory include McAllester’s bound [70], Seeger’s bound [88], Catoni’s bound [16] and Maurer’s

bound [68]. Based on this, there have been numerous studies that applying or extending these results to the analysis of deep neural networks generalization, including computing non-vacuous bounds for deep neural networks [28, 126, 79, 65] and establishing upper bounds for optimization algorithms [63, 75, 84, 5, 72, 121, 61, 66]. We refer interested readers to a comprehensive survey [1] for more detail on this topic. All the above results require that training examples are independent to test examples and could not be applied on transductive learning setting.

Transductive Learning Generalization Theory. The concept of transductive learning and the earliest bounds are presented to [106]. The authors in [29] study the stability of transductive learning algorithms. They further propose another tool named transductive Rademacher Complexity [30] as a complexity of hypothesis space under transductive setting. Permutational Rademacher Complexity is latter presented in [105], which is shown to be more suitable for transductive learning setting than transductive Rademacher Complexity. By considering the variance of functions, the authors in [104] establish new concentration inequalities and sharper bounds. [26] is the first work to analyze the generalization of transductive learning in the context of PAC-Bayesian, and their results are improved latter in [10]. Furthermore, the above results have been applied on the theoretical analysis in graph learning [96, 97, 24, 76, 32, 21, 101], semi-supervised learning [69, 38, 118], matrix completion [37, 94], distributed optimization [93] and collaborative filtering [119, 25], among other areas. Our results also apply to these fields, particularly when the learner is deep neural network.

3 Preliminaries

3.1 Notations

We stipulate that random variables and their realization are denoted by uppercase and lowercase letters, respectively. For given random variable X , we denote its distribution measure by P_X . Similarly, the conditional distribution measure of X given Y is given by $P_{X|Y}$. We use $D_{\text{KL}}(P||Q)$ to denote the Kullback–Leibler (KL) divergence between two probability measure P and Q , where we have assumed that they are on the same measure space and the Radon-Nikodym derivative of P with respect to Q is well defined. With this notation, the mutual information between X and Y is defined as $I(X; Y) \triangleq D_{\text{KL}}(P_{X,Y}||P_X P_Y)$. Furthermore, the disintegrated mutual information is denoted by $I^Z(X; Y) \triangleq D_{\text{KL}}(P_{X,Y|Z=z}||P_{X|Z=z} P_{Y|Z=z})$, whose expectation taking over $Z \sim P_Z$ is the conditional mutual information $I(X; Y|Z) = \mathbb{E}_{P_Z}[I^Z(X; Y)]$.

3.2 Transductive Learning

Let $D = \{\mathbf{z}_1, \dots, \mathbf{z}_n\}$ be a given set with finite cardinality, where $\mathbf{z} = (\mathbf{x}, y)$ is an instance composed of attribute $\mathbf{x} \in \mathcal{X}$ and target $y \in \mathcal{Y}$ from $\mathcal{Z} \triangleq \mathcal{X} \times \mathcal{Y}$. We use $\text{Perm}(D)$ to denote the set containing all bijections $\pi : D \rightarrow D$. Here each bijection $\pi \in \text{Perm}(D)$ could be regarded as a permutation on D . Note that sampling without replacement from D is equivalent to firstly sampling a permutation from $\text{Perm}(D)$ with equal probability and then applying it on D . Denote by Π a random variable follows uniform distribution over $\text{Perm}(D)$, namely $\mathbb{P}\{\Pi = \pi\} = \frac{1}{n!}$ holds for any $\pi \in \text{Perm}(D)$. With this notation, we use $Z \triangleq (Z_1, \dots, Z_n)$ to denote the random permutation vector induced by Π , where $Z_j = \Pi(\mathbf{z}_j)$ represents the j -th element of the sequence after permutation. Once the permutation Z is determined, the training set is defined as $D_{\text{train}} \triangleq \{Z_1, \dots, Z_m, X_{m+1}, \dots, X_{m+u}\}$, where m and u are the number of training and test instances, respectively. Let \mathcal{W} be the space of parameter, the transductive learning algorithm takes D_{train} as input and outputs a random element $W \in \mathcal{W}$ as the hypothesis, which is characterized by a Markov kernel $P_{W|Z}$. Let $\ell : \mathcal{W} \times \mathcal{Z} \rightarrow \mathbb{R}_{\geq 0}$ be the objective function, the transductive training and test error of a hypothesis W are defined as $R_m(W, Z) \triangleq \frac{1}{m} \sum_{i=1}^m \ell(W, Z_i)$ and $R_u(W, Z) = \frac{1}{u} \sum_{i=m+1}^{m+u} \ell(W, Z_i)$, respectively. The *transductive generalization error* is then defined as $\mathcal{E}(W, Z) \triangleq R_u(W, Z) - R_m(W, Z)$. Furthermore, we use $\mathbb{E}_{W,Z}[\mathcal{E}(W, Z)]$ to denote the expectation of $\mathcal{E}(W, Z)$ taking over $P_{W,Z} = P_Z \otimes P_{W|Z}$, which represents the the average performance difference of the hypothesis W between testing and training instances over all permutation Z . Under certain circumstances, the objective $\ell(W, Z)$ is represented as $l(f_W(X), Y)$, where $f(\cdot) : \mathcal{X} \times \mathcal{W} \rightarrow \hat{\mathcal{Y}}$ is the learner parameterized by W and $l : \hat{\mathcal{Y}} \times \mathcal{Y} \rightarrow \mathbb{R}_{\geq 0}$ is the criterion.

4 Information-Theoretic Generalization Bounds for Transductive Learning

4.1 Establishing Upper Bounds by Mutual Information

Different from supervised learning, the randomness of training and testing examples in transductive learning come from the partition determined by permutation rather than from sampling. This also brings another challenge, namely the dependence of training and testing examples, since the testing examples are uniquely determined once training examples are chosen. The most widely adopt technique to tackle dependence is the martingales method, which enables us to derive similar “sub-Gaussian” property for the transductive generalization error. Together with the Donsker-Varadhan’s variational formula, we establish the following transductive generalization bounds.

Theorem 1. Suppose $\ell(\mathbf{w}, \mathbf{z}) \leq B$ holds for any $\mathbf{w} \in \mathcal{W}$ and $\mathbf{z} \in D$, then

$$|\mathbb{E}_{W,Z} [R_u(W, Z) - R_m(W, Z)]| \leq \sqrt{\frac{C_{m,u}}{2} \left(\frac{1}{m} + \frac{1}{u} \right) I(Z; W)}, \quad (1)$$

$$\mathbb{E}_{W,Z} [(R_u(W, Z) - R_m(W, Z))^2] \leq C_{m,u} \left(\frac{1}{m} + \frac{1}{u} \right) (I(Z; W) + \log 3). \quad (2)$$

where $C_{m,u} \triangleq \frac{2B^2(m+u) \max(m,u)}{(m+u-1/2)(2 \max(m,u)-1)}$.

Theorem 1 shows that the expectation of transductive generalization error is upper bounded by the mutual information between permutation Z and hypothesis W . This result implies that the less dependence the output hypothesis has on the selection of training data, the better generalization a transductive learning algorithms will have. One can image that if the algorithm only “memorize” the partition of training and testing examples (or heavily depends on the training samples it sees), we could not expect that it will have strong generalization ability. As a comparison, similar results (Theorem 1 in [117]) under supervised learning setting say that the generalization error is upper bounded by the mutual information between training set S and hypothesis W , which could be regarded as a special case of our result that u is infinite. In this event, the transductive training and test error corresponding to the supervised training and test error, respectively. Also, m is the number of examples in S , and the term $\frac{1}{u}$ vanish. Note that since the test error is computed over infinite examples, there is no dependence between it and the training error. Furthermore, the assumption of Theorem 1 is slightly stronger than that in supervised learning setting, where they only requires the loss to be sub-Gaussian while we require the loss to be bounded. However, we believe that our result could be extend to the unbounded loss setting under proper assumptions.

The result presented in Theorem 1 is a expectation bound over all possible selection of training data. In real-world application particularly deep learning scenario, only a few partitions (determined by random seed) are adopted to verify the quality of a transductive learning algorithm, and the empirical results show that whose performance could generally be guaranteed. This urges us to establish the high probability bound in order to better describe the generalization behavior of deep transductive learners. Achieve this relies on the monitor technique proposed in [9]. Another derivant is the expectation bound on the absolute value of transductive generalization error, which serves as a supplement of Theorem 1. The aforementioned results are summarized in Theorem 2.

Theorem 2. Suppose $\ell(\mathbf{w}, \mathbf{z}) \in [0, B]$ holds for any $\mathbf{w} \in \mathcal{W}$ and $\mathbf{z} \in D$, with probability at least $1 - \delta$ over the randomness of Z, W :

$$|R_u(W, Z) - R_m(W, Z)| \leq 2\sqrt{\frac{C_{m,u}}{2} \left(\frac{1}{m} + \frac{1}{u} \right) \left(\log \left(\frac{1}{\delta} \right) + \frac{I(Z; W)}{\delta} \right)}, \quad (3)$$

where $C_{m,u}$ follows the definition in Theorem 1. Furthermore, we have

$$\mathbb{E} |R_u(W, Z) - R_m(W, Z)| \leq \sqrt{\frac{C_{m,u}}{2} \left(\frac{1}{m} + \frac{1}{u} \right) (I(Z; W) + \log 2)}. \quad (4)$$

Since $C_{m,u} \approx B^2$ when $m + u$ is large, the high probability bound presented in Eq. 3 is of order $(1/m + 1/u)^{\frac{1}{2}}$. Despite a degenerated constant factor from $\log(1/\delta)$ to $1/\delta$, our bound is sharper

than that in existing studies [76, 32]. Although the mutual information term $I(Z; W)$ could not be easily computed, we will show in Sec. that it has a unique advantage when the learner is optimized by stochastic algorithms such as stochastic gradient descent and its variants. Besides, result similar to Eq. (4) can be derived from Eq. (2), despite the constant factor is slightly larger.

4.2 Establishing Upper Bounds by Conditional Information

So far, all the bounds we have established contains the mutual information term $I(Z; W)$, either in expectation or high probability. One unsatisfied property of mutual information is that it does not have a finite upper bound, which may lead to vacuous bounds under some circumstances. Fortunately, this issue could be addressed by adopting “supersamples” setting proposed in [99] under supervised learning setting. The key insight is that introducing another random variable to control the randomness of sampling training and test examples, which is independent of the instances. However, directly applying this technique on transductive learning setting is not feasible. The reason is that the training and test examples given returned by this setting is independent, which is yet dependent in transductive learning setting. To bridge this gap, we propose the following transductive supersamples under a specific condition that the number of training examples is equal to that of test examples.

Definition 1 (Transductive Supersamples). Let $D = \{\mathbf{z}_i\}_{i=1}^n$ be a given set where n is a finite even number. Denote by $m = \frac{1}{2}n$, the transductive supersamples is a sequence $\tilde{Z}^m \triangleq (\tilde{Z}_1, \dots, \tilde{Z}_m)$ generated by sampling without replacement from D , where $\tilde{Z}_i \triangleq \{\tilde{Z}_{i,0}, \tilde{Z}_{i,1}\}$ for $i \in [m]$ is an unordered set with cardinality 2.

Definition 1 shows that the transductive supersamples is obtained by continuously sampling an unordered instance pairs from a fixed set until there is no remained instances. Please refer to Appendix E for a illustrated example of this definition. As a comparison, transductive samples Z^n is obtained by each time sampling an instance from D . A deeper relationship between Z^m and \tilde{Z}^m is given by the following lemma.

Proposition 1. Denote by \tilde{Z} the set containing all \tilde{Z}^m . Let $S = (S_1, \dots, S_m) \sim \text{Unif}(\{0, 1\})^m$ be the sequence of random variables that is independent of \tilde{Z}^m . Sampling without replacement from D is equivalent to firstly sampling \tilde{Z} from \tilde{Z} and applying S to permute \tilde{Z} .

Proposition 1 implies that there is another way to obtain the random permutation vector Z^n with the help of transductive supersamples. Let \tilde{Z}^m and S be the random variables described in Proposition 1, Z can be expressed by $Z^n = (\tilde{Z}_{1,S_1}, \dots, \tilde{Z}_{m,S_m}, \tilde{Z}_{1,1-S_1}, \dots, \tilde{Z}_{m,1-S_m})$. Let $\mathcal{E}(W, \tilde{Z}, S)$ be the transductive generalization error under supersampling setting defined by

$$\mathcal{E}(W, \tilde{Z}, S) \triangleq \frac{1}{m} \sum_{i=1}^m (\ell(W, \tilde{Z}_{i,S_i}) - \ell(W, \tilde{Z}_{i,1-S_i})), \quad (5)$$

we have $\mathbb{E}[\mathcal{E}(W, \tilde{Z}, S)] = \mathbb{E}[\mathcal{E}(W, Z)]$. This enables us to characterize the generalization bounds using conditional mutual information, as presented in Theorem 3.

Theorem 3. Suppose that $\ell(\mathbf{w}, \mathbf{z}) \in [0, B]$ holds for any $\mathbf{w} \in \mathcal{W}$ and $\mathbf{z} \in D$, then

$$|\mathbb{E}_{Z,W} [R_u(W, Z) - R_m(W, Z)]| \leq \mathbb{E}_{\tilde{Z}} \sqrt{\frac{2B^2}{m} I^{\tilde{Z}}(S; W)} \quad (6)$$

$$\mathbb{E}_{S,W} [(R_u(W, Z) - R_m(W, Z))^2] \leq \frac{4B^2}{m} (I(S; W|\tilde{Z}) + \log 3). \quad (7)$$

By [99], we have $I(S; W|\tilde{Z}) \leq H(S; Z) \leq m \log 2$ holds, suggesting that the conditional mutual information has a finite upper bound and thus provides a non-vacuous generalization bound. Eq. (6) is consistent with the results in supervised learning setting [99] in form, and the only difference is that \tilde{Z} should be interpreted as the transductive supersamples. Also, we can recover the result in supervised learning (Theorem 1.2 in [99]) whereas the full sample set D has a infinite cardinality. In this event, D is exactly the space containing all instances, implying that entries in the sequence \tilde{Z}^m are independent to each other.

Although the mutual information term $I(S; W|\tilde{Z})$ in Theorem 3 is bounded, computing its numerical value is still difficult, due to W is always a high-dimensional random variable. Thanks to the transductive supersampling setting, various information-theoretical measures [46, 49, 114] adopted in supervised learning setting can be extended to transductive learning setting, as shown in the following Corollary.

Corollary 1. *Suppose the criterion l satisfies $l(\hat{y}, y) \in [0, B]$ holds for any $\hat{y} \in \hat{\mathcal{Y}}, y \in \mathcal{Y}$. Let $f_w(\mathbf{x}) \in \mathbb{R}^d$ be the prediction given by the learner. Denote by $F \in \mathbb{R}^{m \times 2d}$ the prediction matrix where the i -th row is given by $F_{i,:} \triangleq (f_w(\tilde{X}_{i,0}), f_w(\tilde{X}_{i,1}))$, where \tilde{X} is the feature of $\tilde{Z} = (\tilde{X}, \tilde{Y})$. We have*

$$|\mathbb{E}_{Z,W} [R_u(W, Z) - R_m(W, Z)]| \leq \frac{B}{m} \sum_{i=1}^m \mathbb{E}_{\tilde{Z}} \sqrt{2I(F_i; S_i|\tilde{Z})}. \quad (8)$$

Denote by $L \in \{0, 1\}^{m \times 2}$ the loss value matrix, where the i -th row is $L_{i,:} \triangleq (\ell(W, \tilde{z}_{i,0}), \ell(W, \tilde{z}_{i,1}))$. Let $\Delta_i \triangleq \ell(W, \tilde{z}_{i,1}) - \ell(W, \tilde{z}_{i,0})$ be the difference of loss value. We have

$$|\mathbb{E}_{S,W} [R_u(W, Z) - R_m(W, Z)]| \leq \frac{B}{m} \sum_{i=1}^m \mathbb{E}_{\tilde{Z}} \sqrt{2I^{\tilde{Z}}(L_i; S_i)}, \quad (9)$$

$$|\mathbb{E}_{S,W} [R_u(W, Z) - R_m(W, Z)]| \leq \frac{B}{m} \sum_{i=1}^m \mathbb{E}_{\tilde{Z}} \sqrt{2I^{\tilde{Z}}(\Delta_i; S_i)}. \quad (10)$$

Compared with previous studies, we termed the bounds in Eqs. (8,9,10) as transductive f -CMI, transductive e-CMI and transductive Id-CMI bounds, respectively. In applications, the prediction of learner is a low-dimension vector and thus reduce the difficulty of computing the conditional mutual information $I(W; S|\tilde{Z})$. Note that L_i in Eq. (9) and Δ_i in Eq. (10) are two-dimensional and one-dimensional random variable, yielding more convenient computation and sharper bounds. Here we point out that each result has its advantage. The vanilla bound in Eq. (6) is more informative to understanding generalization, and the expense is its difficulty of calculating numerical value. In contrast, the other bounds in Corollary 1 has computation convenience, yet they are inferior in reflecting factors that affect generalization. Despite the exist of this trade-off, these results are sufficient for us to understand the generalization behavior of transductive learner or establish non-vacuous bounds for them.

We close this part by briefly discuss how to extend the aforementioned results to more ordinary cases. According to Definition 1, Theorem 3 and Corollary 1 only apply to the case that $m = u$. Here we point out that they can be extended to cases that $m = ku$ or $u = km$ where $k \in \mathbb{N}_+$. Considering the symmetric, it is sufficient to discuss the case that $u = km$. To this end, the transductive supersamples are extended to the following k -transductive supersamples.

Definition 2 (k -Transductive Supersamples). *Let $D = \{\mathbf{z}_i\}_{i=1}^n$ be a given set where n is a finite even number. Let $k \in \{1, \dots, n-1\}$ be a given integer. Denote by $m = \frac{n}{k+1} \in \mathbb{N}_+$, the k -transductive supersamples is a sequence $\tilde{Z}^m \triangleq (\tilde{Z}_1, \dots, \tilde{Z}_m)$ generated by sampling without replacement from D , where $\tilde{Z}_i \triangleq \{\tilde{z}_{i,0}, \dots, \tilde{z}_{i,k}\}$ is an unordered set with cardinality $k+1$.*

Note that Definition 1 is a special case of Definition 2 where $k = 1$. Similarly we need to extend the definition of the indicator variable S . Let $S = (S_1, \dots, S_m) \sim \text{Unif}(\{0, \dots, k\})^m$ be the sequence of random variables that is independent of \tilde{Z}^m . Let $U = (U_1, \dots, U_m) \sim \text{Unif}(\{0, \dots, (k-1)!\})^m$ be the sequence of random variables that is independent of \tilde{Z}^m and S . For a fixed set \mathcal{S} , we use $\Pi(\mathcal{S}, U)$ to represent the permutation of \mathcal{S} induced by random variable U . For example, $\Pi(\tilde{Z}_i \setminus \tilde{z}_{i,0}, U_1)$ denote the result of drawing a permutation of $\tilde{Z}_i \setminus \tilde{z}_{i,0}$ according to U_1 , which is essentially equivalent to sampling without replacement from $\tilde{Z}_i \setminus \tilde{z}_{i,0}$. With this definition, the permutation vector Z can be expressed by $Z^n = (\tilde{Z}_{1,S_1}, \dots, \tilde{Z}_{m,S_m}, \Pi_{U_1}(\tilde{Z}_1 \setminus \tilde{z}_{1,S_1}), \dots, \Pi_{U_m}(\tilde{Z}_m \setminus \tilde{z}_{m,S_m}))$. By this way, all results in Theorem 3 and Corollary 1 can be extended to the case that $u = km$, and we place the details in Appendix E. One main differences between the case that $k = 1$ and $k > 1$ is the increased cost of estimating the mutual information, since each entry of S has more possible values to take. In other words, with the increase of k , we need to accordingly increase the samples of S to reduce the estimated error of the mutual information.

4.3 Connection with Transductive PAC-Bayesian Bounds

PAC-Bayesian methods and Information-theoretic methods are closely related, since both of them are based on Donsker-Varadhan’s variational formulation. Borrowing the proof of Theorem 1, we obtain the following new transductive PAC-Bayesian bounds.

Theorem 4. *Suppose that $\ell(w, \mathbf{z}) \leq B$ holds for any $w \in \mathcal{W}$ and $\mathbf{z} \in D$. Let P be a prior distribution P on \mathcal{W} . With probability at least $1 - \delta$ over the randomness of Z , for any distribution Q on \mathcal{W} we have*

$$|\mathbb{E}_{W \sim Q} [R_u(W, Z) - R_m(W, Z)]| \leq \sqrt{\frac{C_{m,u}}{2} \left(\frac{1}{m} + \frac{1}{u} \right) \left(D_{\text{KL}}(Q||P) + \log \left(\frac{1}{\delta} \right) \right)}, \quad (11)$$

where $C_{m,u}$ follows the definition in Theorem 1.

Compared with previous transductive PAC-Bayesian bound (Corollary 7(b) in [10]), the result in Theorem 4 are as follows. For one hand, the assumptions of our bound is weaker. Previous result only apply to zero-one loss, and the value of m and n are required to satisfy $m \geq 40$ and $20 \leq m \leq n - 20$. In contrast, our result apply to any bounded loss and there are no constraints on the value of m and n . For the other hand, our result is strictly sharper than previous result by removing the term $\log \left(\log(m) \sqrt{\frac{mu}{m+u}} \right)$. Furthermore, incorporating the technique in [75], the results in [62] could be directly extend to transductive learning setting and thus providing generalization guarantee for many GNNs models on node classification and link prediction task.

One of the most important insights provided by PAC-Bayesian bounds in that the generalization performance is closely related with the sharpness of the loss landscape, and a flat minimum is beneficial for generalization [75]. With the help of Theorem 4, this result can be extended to transductive learning setting when the criterion is 0-1 loss.

Corollary 2. *Suppose that $R_u(\mathbf{w}, Z) \leq \mathbb{E}_{\epsilon \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})} [R_u(\mathbf{w} + \epsilon, Z)]$ holds for any Z , where $\mathbf{w} \in \mathbb{R}^d$ is the parameter returned by given transductive learning algorithm and $\epsilon \in \mathbb{R}^d$ is a random noise. With probability at least $1 - \delta$ over the randomness of Z ,*

$$\begin{aligned} & R_u(\mathbf{w}, Z) \\ & \leq \max_{\|\epsilon\|_2 \leq \rho} R_m(\mathbf{w} + \epsilon, Z) \\ & \quad + \sqrt{\frac{C_{m,u}(m+u) \left(1 + \frac{d}{2} \log \left(1 + \frac{(1+\tilde{C}_{m,u})^2 \|\mathbf{w}\|_2^2}{\rho^2} \right) + \log \left(\frac{1}{6\delta} \right) + 2 \log \left(\frac{4\pi mu}{m+u} \right) \right)}{mu}}, \end{aligned} \quad (12)$$

where $\tilde{C}_{m,u} \triangleq \sqrt{\log(4mu/(m+u))/d}$ and $C_{m,u}$ follows the definition in Theorem 1.

The term $\max_{\|\epsilon\|_2 \leq \rho} R_m(\mathbf{w} + \epsilon, Z)$ characterizes the change of loss landscape within a ball with \mathbf{w} as the center and ρ as the radius. Formally, we call \mathbf{w} as sharp minima if the loss values around it differ significantly from itself, namely $R_m(\mathbf{w} + \epsilon, Z) > R_m(\mathbf{w}, Z)$. Therefore, Corollary 2 suggests that a flat optima could have better transductive generalization performance. This theoretical result has been verify by recent work [17], which applies sharpness-aware minimization [33] on recommendation task to achieve better generalization performance for GNNs. We believe that this result could shed light in understanding the correlation between sharpness and generalization in transductive learning.

4.4 Upper Bounds for Adaptive Optimization Algorithms

As we have mentioned, one advantages of our theoretical results against previous one is that the effect of optimization algorithm on generalization can be fully considered. In this work we focus on AdaGrad, which is one of the most widely adopted optimization algorithms in practice. Different from SGD, the learning rate is adaptively adjusted during training. Denote by $\{W_t\}_{t \in [T]}$ the weights along the training trajectory of AdaGrad. Following [113], we consider the setting that mini-batches examples are fixed. Denote by (B_1, \dots, B_T) the sequence of mini-batches where B_i is the of examples used in the i -th epoch. For concise we assume that the learner only minimizes the loss

on labeled examples, and the number of each mini-batch examples are equal to b . Then the average gradient on the B_t is defined as

$$g(w, B_t(Z)) \triangleq \frac{1}{b} \sum_{\mathbf{z} \in B_t} \nabla_w \ell(w, \mathbf{z}), \quad (13)$$

where $B_t(Z) \subseteq \{\mathbf{z} | \mathbf{z} \in Z^m\}$. For $t \in [T]$, the update rule of AdaGrad can be formulated as

$$v_t = \sum_{k=1}^{t-1} g(W_k, B_k(Z)) \odot g(W_k, B_k(Z)), W_t = W_{t-1} - \frac{\eta}{\sqrt{v_t} + \epsilon} \odot g(W_{t-1}, B_t(Z)), \quad (14)$$

where W_0 is the initial parameter and \odot is the Hadamard product. η and ϵ are two hyper-parameters. Note that v_t is a random variable determined by $W^{[t-1]} \triangleq (W_0, \dots, W_{t-1})$. For the concise of notations we use $\Psi(W^{[t-1]}, Z) \triangleq (\eta/(\sqrt{v_t} + \epsilon)) \odot g(W_{t-1}, B_t(Z))$ to denote the ‘‘adaptive gradient’’, which is computed by normalizing the current gradient with accumulate squared gradient. Inspired by [74, 113], we introduce the following auxiliary weight process $\{\widetilde{W}_t\}_{t \in [T]}$ for analysis:

$$\widetilde{W}_0 = W_0, \widetilde{W}_t = \widetilde{W}_{t-1} - \Psi(W^{[t-1]}, Z) + N_t, \quad (15)$$

where $N_t \triangleq \sigma_t N$. Here $\{\sigma_t\}_{t \in [T]}$ are predefined hyperparameters and N is a Gaussian random variable that is independent to $W^{[T]}$ and Z . For concise, we define $U_t \triangleq \sum_{k=1}^t N_k$. The upper bound for a transductive learner trained by AdaGrad is presented in the following theorem.

Theorem 5. Suppose that (i) $\ell(w, \mathbf{z}) \leq B$ holds for any $w \in \mathcal{W}$ and $\mathbf{z} \in D$ and (ii) $\mathbb{E}[R_u(w_T + U_T, Z) - R_u(w_T, Z)] \geq 0$ holds for any realization of $W_T = w_T$. Then we have

$$\begin{aligned} & \mathbb{E}_{Z, W_T} [R_u(W_T, Z) - R_m(W_T, Z)] \\ & \leq \sqrt{\frac{C_{m,u}}{2} \left(\frac{1}{m} + \frac{1}{u} \right) \sum_{t=1}^T \frac{d}{2} \log \left(\frac{1}{d\sigma_t^2} \mathbb{E} \left[\left\| \Psi(W^{[t-1]}, Z) - \mathbb{E}_{W^{[t-1]}, Z} [\Psi(W^{[t-1]}, Z)] \right\|_2^2 \right] + 1 \right)} \\ & \quad + \mathbb{E}_{Z, W_T, U_T} [R_m(W_T + U_T, Z) - R_m(W_T, Z)]. \end{aligned} \quad (16)$$

Assumption (ii) in Theorem 5 is also used in Corollary 2 to establish the PAC-Bayesian bound. This assumption requires that adding random noise to the final parameter does not decrease the risk on unlabeled examples in expectation. The first term in Eq. (16) records the variance of ‘‘adaptive gradient’’ along the training trajectory, and the second term measures the expected change of training risk after adding random noise. Compared with stability based [21] or complexity based [101] methods, our results do not contain any Lipschitz or smoothness constants and thus more easier to computed. Besides, the smoothness or Hölder smoothness assumption limit the application scope of previous results, e.g., they could not be applied to neural networks with ReLU as activation function. As a comparison, our result does not rely on these assumptions and has a wider applicability. Also, since the Lipschitz constant is the upper bound of the norm of gradient, our results can more finely depict the impact of optimization trajectories on generalization. Furthermore, the second term characterizes the flatness of the final parameter, which conveys the same insight as Theorem 4, namely a flat optima is beneficial to achieve smaller generalization gap. If u is infinity and v_t is set to all-ones vector, the result for SGD under supervised learning setting [74, 113] are recovered from our result. Considering the popularity of Adam in real-world applications, we also derive the corresponding results. The reflected insights are similar yet the formulations are more tedious, and we place the details in Appendix I.

5 Applications

In this section, we demonstrate the application of theoretical results through two specific examples: semi-supervised learning and graph learning.

5.1 Semi-supervised Learning

Due to the expensive cost of collecting high-quality labeled data, semi-supervised learning aims to train a learner with a few labeled examples and a large amount of unlabeled data. The analysis for the

generalization of semi-supervised learner has been widely explored [71], and the theoretical results differ by the problem setting and assumptions. Here we focus on the transductive setting, which is also termed as setting 2 of transductive learning [106]. Formally, the labeled and unlabeled data is represented as $D_m \triangleq \{(X_i, Y_i)\}_{i \in [m]}$ and $D_u \triangleq \{(X_i, Y_i)\}_{i \in [u]}$, which are samples independently from certain distribution. The semi-supervised learner takes $D_m \cup D_u^X$ as input and outputs the hypothesis characterized by $W \in P_{W|D_m \cup D_u^X}$, where $D_u^X \triangleq \{X_i\}_{i \in [u]}$. Different from the setting we present in Subsection 3.2, here each example (X, Y) should be regarded as random variable rather than constant pair. Furthermore, the training and test risk are defined as $R(W, D_m) \triangleq \frac{1}{m} \sum_{i=1}^m \ell(W, (X_i, Y_i))$ and $R(W, D_u) \triangleq \frac{1}{u} \sum_{i=m+1}^{m+u} \ell(W, (X_i, Y_i))$, respectively. Leveraging the theoretical results established in Section 4, the generalization of semi-supervised learner is described as follows.

Proposition 2. Suppose $\ell(\mathbf{w}, \mathbf{z}) \in [0, B]$ holds for any $\mathbf{w} \in \mathcal{W}$ and $\mathbf{z} \in D$, then we have

$$|\mathbb{E}_{D_m, D_u^X, W} [R(W, D_u) - R(W, D_m)]| \leq \mathbb{E}_{D_{m+u}} \sqrt{\frac{C_{m,u}}{2} \left(\frac{1}{m} + \frac{1}{u} \right) I(Z; W | D_{m+u})}, \quad (17)$$

where $C_{m,u}$ follows the definition in Theorem 1. Denote by $R_\mu(W) \triangleq \mathbb{E}_{(X,Y) \sim \mu} \ell(W, (X, Y))$ the population risk of the semi-supervised learner, under the same assumption we have

$$\mathbb{E}_{D_m, D_u^X, W} R_\mu(W) \leq R(W, D_m) + \sqrt{\frac{2B^2 I(D_{m+u}, Z; W)}{m+u}} + \mathbb{E}_{D_{m+u}} \sqrt{\frac{C_{m,u} I(Z; W | D_{m+u})}{2m(m/u + 1)}}. \quad (18)$$

Compared with Theorem 1, the upper bound in Eq. (17) includes the conditional mutual information $I(Z, W | D_{m+u})$ rather than mutual information $I(Z, W)$. The reason is that each element in D_{m+u} is random variables rather than constants, and the randomness of D_{m+u} should also be taken into consideration. Eq. (18) reveals that if the unlabeled data follows the same distribution as that of labeled data, the population risk decreases with the number of u , namely increasing the number of unlabeled data is beneficial for generalization. However, label/data noise and distribution shift could occur in unlabeled data collected from real-world scenario. Thus, developing robust semi-supervised learners is more urgent.

5.2 Transductive Graph Learning

Composed of several objects and their relationship, graph-structured data plays an important role in real-world applications, *e.g.*, recommendation system [111, 48, 51] drug discovery [100, 12] and traffic prediction. Recent years have witnessed the success of GNNs in various learning and inference tasks on graph-structured data. The graph learning tasks can be divided into transductive task and inductive task, and we only discuss the first one. Transductive graph learning tasks include graph-level task and node/edge-level task. The goal for node-level task is to predict the label of nodes. As for edge-level task, the learner is required to predict whether there is a link between two nodes. Both of them fall into the category of transductive learning. Taking node classification as an example, partial nodes features together with labels randomly selected from the graph are used for training, and the GNN model need to predict the labels of remaining nodes. During this process, all node features are visible. Let D be the set containing all nodes and all edges respectively, our results can be applied on analyzing the generalization gap of GNNs on node classification task and link prediction task. Notably, the upper bound in Theorem 5 includes the norms along trajectory, which has close relation to both the architecture of GNN model and property of graph-structured information. Following the technique in [21, 101], we believe that it is possible to derive fine-grained upper bounds for specific GNN models and obtain more insights on its generalization behavior.

6 Experiments

In this section, we conduct experiments on semi-supervised learning and transductive graph learning scenarios to validate our theoretical findings.

6.1 Semi-supervised Learning

Datasets and Learning Algorithms. We choose semi-supervised image classification on MNIST and CIFAR-10 as the learning tasks. The semi-supervised learning loss for unlabeled images is defined as the mean square error between the prediction of the augmented images and the vanilla images by the model, which is also termed as consistency regularization in semi-supervised learning. Following [46, 40], we adopt a four-layer CNN and Wide ResNet-28-10 [124] as the model for MNIST and CIFAR-10, respectively. Please refer to Appendix for the detail of network architecture. For both these two experiments, we train the model on xx mini batches using Adam optimizer with learning rate 0.001, and the number per mini batch is fixed to xx . The loss is set to zero-one loss. Following [46], we make the training process be deterministic by fixing the sequence of mini batch and the initialization of parameters via random seed.

Estimating the Expected Transductive Generalization Gaps and the Derived Bounds. Notice that compute the accurate value of the expected transductive generalization gap (and also the derived upper bounds) is not applicable since we need to run the algorithm on $(m+u)!$ partitions in total. Therefore we use Monte Carlo simulation to estimate these expectations based on finite samples. The sampling process are as follows: (i) randomly draw t_1 full samples set s_{m+u} by each time sampling $m+u$ images from the raw images set, (ii) randomly draw t_2 transductive supersamples \tilde{z}^m based on Definition 2, (iii) randomly draw t_3 train/test split variables s and obtain the training and test samples set according to Subsection 4.2. Notice that here we do not consider the randomness of U for the case $k \geq 2$. The reason is that U control the permutation of test samples, and the learning algorithm we consider is independent of this permutation. Now we discuss the estimation of transductive generalization gap and the upper bounds established in Corollary 1. Taking Eq. () as an example, for each (s_{m+u}, \tilde{z}) we use the mean value over t_3 samples of S to estimate the conditional term expectation term $\frac{1}{m} \sum_{i=1}^m \mathbb{E}_{F_i, U_i | \tilde{z}, s_{m+u}} g(F_i, S_i, \tilde{y}_i)$. After that we use $t_1 t_2$ samples of S_{m+u} and \tilde{Z} to estimate the expected generalization gap, whose mean and standard deviation are shown in Figure. Similarly, we use a plug-in estimator [77] to estimate the disentangled mutual information $I^{s_{m+u}, \tilde{z}}(F_i; S_i)$ over the t_3 samples of S . Then the upper bound in Eq. () is estimated by the $t_1 t_2$ samples of S_{m+u} and \tilde{Z} , whose mean and standard deviation are shown in Figure.

6.2 Transductive Graph Learning

Datasets and Learning Algorithms. We choose semi-supervised node classification on synthetic and real world datasets as the learning tasks. Specifically, we select cSBMs [27] to be the synthetic data, and Cora, CiteSeer, PubMed [90, 122] to be the real world dataset. For each of these dataset, we adopt GCN [57], GAT [108], JKNet [120], APPNP [35] and GPR-GNN [19] as the learner, which are popular baselines in graph learning literature. Please refer to Appendix for the detail of network architecture. Following , we train the model on all labeled nodes for 200 epochs with Adam optimizer with learning rate

Estimating the Expected Transductive Generalization Gaps and the Derived Bounds. The estimation process is generally follows that of semi-supervised learning, expect that we do not need to consider the sample of S_{m+u} . Concretely, the sampling process is only composed of (ii) and (iii). Accordingly, we use t_2 samples of \tilde{Z} to estimate the expected bounds and the conditional mutual information.

7 Conclusion

In this work, we study the generalization of transductive learning under the viewpoint of information theory, and establish upper bounds for general transductive algorithms and iterative algorithms in terms of different information measure. Furthermore, we demonstrate their applications in semi-supervised learning and transductive graph learning. Promising future directions include apply our results to other scenarios and develop the lower bounds that matched the upper bounds.

References

- [1] Pierre Alquier. User-friendly introduction to pac-bayes bounds. *arXiv preprint arXiv:2110.11216*, 2021.

- [2] Gholamali Aminian, Mahed Abroshan, Mohammad Mahdi Khalili, Laura Toni, and Miguel Rodrigues. An information-theoretical approach to semi-supervised learning under covariate-shift. In *Proceedings of The 25th International Conference on Artificial Intelligence and Statistics*, pages 7433–7449, 2022.
- [3] Gholamali Aminian, Laura Toni, and Miguel R. D. Rodrigues. Information-theoretic bounds on the moments of the generalization error of learning algorithms. In *2021 IEEE International Symposium on Information Theory*, pages 682–687, 2021.
- [4] Gholamali Aminian, Laura Toni, and Miguel R. D. Rodrigues. Jensen-shannon information based characterization of the generalization error of learning algorithms. In *2020 IEEE Information Theory Workshop*, pages 1–5, 2021.
- [5] Sanjeev Arora, Rong Ge, Behnam Neyshabur, and Yi Zhang. Stronger generalization bounds for deep nets via a compression approach. In *Proceedings of the 35th International Conference on Machine Learning*, pages 254–263, 2018.
- [6] Amir R. Asadi, Emmanuel Abbe, and Sergio Verdú. Chaining mutual information and tightening generalization bounds. In *Advances in Neural Information Processing Systems*, page 7245–7254, 2018.
- [7] Peter L. Bartlett, Olivier Bousquet, and Shahar Mendelson. Local Rademacher complexities. *The Annals of Statistics*, 33(4):1497–1537, 2005.
- [8] Peter L. Bartlett and Shahar Mendelson. Rademacher and gaussian complexities: Risk bounds and structural results. *Journal of Machine Learning Research*, 3:463–482, 2002.
- [9] Raef Bassily, Kobbi Nissim, Adam Smith, Thomas Steinke, Uri Stemmer, and Jonathan Ullman. Algorithmic stability for adaptive data analysis. In *Proceedings of the 48th Annual ACM SIGACT Symposium on Theory of Computing*, pages 1046–1059, 2016.
- [10] Luc Bégin, Pascal Germain, François Laviolette, and Jean-François Roy. PAC-bayesian theory for transductive learning. In *Proceedings of the 17th International Conference on Artificial Intelligence and Statistics*, volume 33, pages 105–113, 2014.
- [11] Avrim Blum and Tom Mitchell. Combining labeled and unlabeled data with co-training. In *Proceedings of the 8th Annual Conference on Computational Learning Theory*, pages 92–100, 1998.
- [12] Pietro Bongini, Monica Bianchini, and Franco Scarselli. Molecular generative graph neural networks for drug discovery. *Neurocomputing*, 450:242–252, 2021.
- [13] Olivier Bousquet and André Elisseeff. Stability and generalization. *Journal of Machine Learning Research*, 2:499–526, 2002.
- [14] Yuheng Bu, Gholamali Aminian, Laura Toni, Gregory W. Wornell, and Miguel Rodrigues. Characterizing and understanding the generalization error of transfer learning with gibbs algorithm. In *Proceedings of the 25th International Conference on Artificial Intelligence and Statistics*, pages 8673–8699, 2022.
- [15] Yuheng Bu, Shaofeng Zou, and Venugopal V. Veeravalli. Tightening mutual information-based bounds on generalization error. *IEEE Journal on Selected Areas in Information Theory*, 1(1):121–130, 2020.
- [16] Olivier Catoni. Pac-bayesian supervised classification: The thermodynamics of statistical learning. *Institute of Mathematical Statistics Lecture Notes Monograph Series*, 56:1–163, 2007.
- [17] Huiyuan Chen, Chin-Chia Michael Yeh, Yujie Fan, Yan Zheng, Junpeng Wang, Vivian Lai, Mahashweta Das, and Hao Yang. Sharpness-aware graph collaborative filtering. In *Proceedings of the 46th International ACM SIGIR Conference on Research and Development in Information Retrieval*, page 2369–2373, 2023.
- [18] Qi Chen, Changjian Shui, and Mario Marchand. Generalization bounds for meta-learning: An information-theoretic analysis. In *Advances in Neural Information Processing Systems*, 2021.
- [19] Eli Chien, Jianhao Peng, Pan Li, and Olgica Milenkovic. Adaptive universal generalized pagerank graph neural network. In *International Conference on Learning Representations*, 2021.
- [20] Eugenio Clerico, Amitis Shidani, George Deligiannidis, and Arnaud Doucet. Chained generalisation bounds. In *Proceedings of 35th Conference on Learning Theory*, pages 4212–4257, 2022.
- [21] Weilin Cong, Morteza Ramezani, and Mehrdad Mahdavi. On provable benefits of depth in training graph convolutional networks. In *Advances in Neural Information Processing Systems*, 2021.

- [22] Corinna Cortes and Mehryar Mohri. On transductive regression. In *Advances in Neural Information Processing Systems*, pages 305–312, 2006.
- [23] Corinna Cortes, Mehryar Mohri, Dmitry Pechyony, and Ashish Rastogi. Stability of transductive regression algorithms. In *Proceedings of the 25th International Conference on Machine Learning*, page 176–183, 2008.
- [24] Jaydeep De, Xiaowei Zhang, Feng Lin, and Li Cheng. Transduction on directed graphs via absorbing random walks. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 40(7):1770–1784, 2018.
- [25] Leyan Deng, Defu Lian, Chenwang Wu, and Enhong Chen. Graph convolution network based recommender systems: Learning guarantee and item mixture powered strategy. In *Advances in Neural Information Processing Systems*, 2022.
- [26] Philip Derbeko, Ran El-Yaniv, and Ron Meir. Explicit learning curves for transduction and application to clustering and compression algorithms. *Journal of Artificial Intelligence Research*, 22:117–142, 2004.
- [27] Yash Deshpande, Subhabrata Sen, Andrea Montanari, and Elchanan Mossel. Contextual stochastic block models. In *Advances in Neural Information Processing Systems*, pages 8590–8602, 2018.
- [28] Gintare Karolina Dziugaite and Daniel M. Roy. Computing nonvacuous generalization bounds for deep (stochastic) neural networks with many more parameters than training data. In *Uncertainty in Artificial Intelligence*, 2017.
- [29] Ran El-Yaniv and Dmitry Pechyony. Stable transductive learning. In *Conference on Learning Theory*, pages 35–49, 2006.
- [30] Ran El-Yaniv and Dmitry Pechyony. Transductive rademacher complexity and its applications. In *Conference on Learning Theory*, pages 157–171, 2007.
- [31] Amedeo Roberto Esposito, Michael Gastpar, and Ibrahim Issa. Generalization error bounds via rényi-, f-divergences and maximal leakage. *IEEE Transactions on Information Theory*, 67(8):4986–5004, 2021.
- [32] Pascal Mattia Esser, Leena C. Vankadara, and Debarghya Ghoshdastidar. Learning theory can (sometimes) explain generalisation in graph neural networks. In *Advances in Neural Information Processing Systems*, pages 27043–27056, 2021.
- [33] Pierre Foret, Ariel Kleiner, Hossein Mobahi, and Behnam Neyshabur. Sharpness-aware minimization for efficiently improving generalization. In *International Conference on Learning Representations*, 2021.
- [34] Borja Rodríguez Gálvez, Germán Bassi, Ragnar Thobaben, and Mikael Skoglund. On random subset generalization error bounds and the stochastic gradient langevin dynamics algorithm. In *2020 IEEE Information Theory Workshop*, page 1–5, 2021.
- [35] Johannes Gasteiger, Aleksandar Bojchevski, and Stephan Günnemann. Predict then propagate: Graph neural networks meet personalized pagerank. In *International Conference on Learning Representations*, 2019.
- [36] Justin Gilmer, Samuel S. Schoenholz, Patrick F. Riley, Oriol Vinyals, and George E. Dahl. Neural message passing for quantum chemistry. In *Proceedings of the 34th International Conference on Machine Learning*, pages 1263–1272, 2017.
- [37] Pere Giménez-Febrer, Alba Pagès-Zamora, and Georgios B. Giannakis. Generalization error bounds for kernel matrix completion and extrapolation. *IEEE Signal Processing Letters*, 27:326–330, 2020.
- [38] Chen Gong, Xiaojun Chang, Meng Fang, and Jian Yang. Teaching semi-supervised classifier via generalized distillation. In Jérôme Lang, editor, *Proceedings of the 27th International Joint Conference on Artificial Intelligence*, pages 2156–2162, 2018.
- [39] M. Gori, G. Monfardini, and F. Scarselli. A new model for learning in graph domains. In *Proceedings. 2005 IEEE International Joint Conference on Neural Networks*, 2005., pages 729–734, 2005.
- [40] Lan-Zhe Guo, Zhenyu Zhang, Yuan Jiang, Yu-Feng Li, and Zhi-Hua Zhou. Safe deep semi-supervised learning for unseen-class unlabeled data. In *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pages 3897–3906, 2020.
- [41] Hassan Hafez-Kolahi, Zeinab Golgooni, Shohreh Kasaei, and Mahdieh Soleymani. Conditioning and processing: Techniques to improve information-theoretic generalization bounds. In *Advances in Neural Information Processing Systems*, pages 16457–16467, 2020.

- [42] Hassan Hafez-Kolahi, Shohreh Kasaei, and Mahdiyeh Soleymani-Baghshah. Sample complexity of classification with compressed input. *Neurocomputing*, 415:286–294, 2020.
- [43] Mahdi Haghifam, Gintare Karolina Dziugaite, Shay Moran, and Daniel M. Roy. Towards a unified information-theoretic framework for generalization. In *Advances in Neural Information Processing Systems*, pages 26370–26381, 2021.
- [44] Mahdi Haghifam, Shay Moran, Daniel M. Roy, and Gintare Karolina Dziugaite. Understanding generalization via leave-one-out conditional mutual information. *arXiv preprint arXiv:2206.14800*, 2022.
- [45] Mahdi Haghifam, Jeffrey Negrea, Ashish Khisti, Daniel M. Roy, and Gintare Karolina Dziugaite. Sharpened generalization bounds based on conditional mutual information and an application to noisy, iterative algorithms. In *Advances in Neural Information Processing Systems*, 2020.
- [46] Hrayr Harutyunyan, Maxim Raginsky, Greg Ver Steeg, and Aram Galstyan. Information-theoretic generalization bounds for black-box learning algorithms. In *Advances in Neural Information Processing Systems*, pages 24670–24682, 2021.
- [47] Haiyun He, Hanshu Yan, and Vincent Y. F. Tan. Information-theoretic characterization of the generalization error for iterative semi-supervised learning. *Journal of Machine Learning Research*, 23(287):1–52, 2022.
- [48] Xiangnan He, Kuan Deng, Xiang Wang, Yan Li, Yongdong Zhang, and Meng Wang. LightGCN: Simplifying and powering graph convolution network for recommendation. In *Proceedings of the 43th International ACM SIGIR Conference on Research and Development in Information Retrieval*, page 639–648, 2020.
- [49] Fredrik Hellström and Giuseppe Durisi. A new family of generalization bounds using samplewise evaluated CMI. In *Advances in Neural Information Processing Systems*, 2022.
- [50] Fredrik Hellström, Giuseppe Durisi, Benjamin Guedj, and Maxim Raginsky. Generalization bounds: Perspectives from information theory and pac-bayes. *arXiv preprint arXiv:2309.04381*, 2023.
- [51] Tinglin Huang, Yuxiao Dong, Ming Ding, Zhen Yang, Wenzheng Feng, Xinyu Wang, and Jie Tang. Mixgcf: An improved training method for graph neural network-based recommender systems. In *Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery & Data Mining*, page 665–674, 2021.
- [52] Thorsten Joachims. Transductive inference for text classification using support vector machines. In *Proceedings of the 16th International Conference on Machine Learning*, pages 200–209, 1999.
- [53] Sharu Theresa Jose and Osvaldo Simeone. Information-theoretic bounds on transfer generalization gap based on jensen-shannon divergence. In *2021 29th European Signal Processing Conference*, pages 1461–1465, 2021.
- [54] Sharu Theresa Jose and Osvaldo Simeone. Information-theoretic generalization bounds for meta-learning and applications. *Entropy*, 23(1), 2021.
- [55] Sharu Theresa Jose, Osvaldo Simeone, and Giuseppe Durisi. Transfer meta-learning: Information-theoretic bounds and information meta-risk minimization. *IEEE Transactions on Information Theory*, 68(1):474–501, 2022.
- [56] Kenji Kawaguchi, Zhun Deng, Xu Ji, and Jiaoyang Huang. How does information bottleneck help deep learning? In *Proceedings of the 40th International Conference on Machine Learning*, pages 16049–16096, 2023.
- [57] Thomas N. Kipf and Max Welling. Semi-supervised classification with graph convolutional networks. In *International Conference on Learning Representations*, 2017.
- [58] Vladimir Koltchinskii and Dmitriy Panchenko. Rademacher processes and bounding the risk of function learning. In *High Dimensional Probability II*, page 443–457, 2000.
- [59] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E. Hinton. Imagenet classification with deep convolutional neural networks. In *Proceedings of the 25th International Conference on Neural Information Processing Systems*, page 1097–1105, 2012.
- [60] Beatrice Laurent and Pascal Massart. Adaptive estimation of a quadratic functional by model selection. *The Annals of Statistics*, 28(5):1302 – 1338, 2000.

- [61] Jian Li, Xuanyuan Luo, and Mingda Qiao. On generalization error bounds of noisy gradient methods for non-convex learning. In *International Conference on Learning Representations*, 2020.
- [62] Renjie Liao, Raquel Urtasun, and Richard Zemel. A PAC-bayesian approach to generalization bounds for graph neural networks. In *International Conference on Learning Representations*, 2021.
- [63] Ben London. A pac-bayesian analysis of randomized learning with application to stochastic gradient descent. In *Advances in Neural Information Processing Systems*, page 2935–2944, 2017.
- [64] Adrian Tovar Lopez and Varun Jog. Generalization error bounds using wasserstein distances. In *2018 IEEE Information Theory Workshop*, pages 1–5, 2018.
- [65] Sanae Lotfi, Marc Anton Finzi, Sanyam Kapoor, Andres Potapczynski, Micah Goldblum, and Andrew Gordon Wilson. PAC-bayes compression bounds so tight that they can explain generalization. In *Advances in Neural Information Processing Systems*, 2022.
- [66] Xuanyuan Luo, Bei Luo, and Jian Li. Generalization bounds for gradient methods via discrete and continuous prior. In *Advances in Neural Information Processing Systems*, 2022.
- [67] Mohammad Saeed Masiha, Amin Gohari, Mohammad Hossein Yassaee, and Mohammad Reza Aref. Learning under distribution mismatch and model misspecification. In *2021 IEEE International Symposium on Information Theory*, page 2912–2917, 2021.
- [68] Andreas Maurer. A note on the pac bayesian theorem. *arXiv preprint arXiv:cs/0411099*, 2004.
- [69] Yury Maximov, Massih-Reza Amini, and Zaid Harchaoui. Rademacher complexity bounds for a penalized multi-class semi-supervised algorithm (extended abstract). In *Proceedings of the 27th International Joint Conference on Artificial Intelligence*, pages 5637–5641, 2018.
- [70] David A. McAllester. Some pac-bayesian theorems. *Maching Learning*, 37(3):355–363, 1999.
- [71] Alexander Mey and Marco Loog. Improved generalization in semi-supervised learning: A survey of theoretical results. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(4):4747–4767, 2023.
- [72] Wenlong Mou, Liwei Wang, Xiyu Zhai, and Kai Zheng. Generalization bounds of sgld for non-convex learning: Two theoretical viewpoints. In *Proceedings of the 31st Conference On Learning Theory*, pages 605–638, 2018.
- [73] Jeffrey Negrea, Mahdi Haghifam, Gintare Karolina Dziugaite, Ashish Khisti, and Daniel M. Roy. Information-theoretic generalization bounds for SGLD via data-dependent estimates. In *Advances in Neural Information Processing Systems*, pages 11013–11023, 2019.
- [74] Gergely Neu, Gintare Karolina Dziugaite, Mahdi Haghifam, and Daniel M. Roy. Information-theoretic generalization bounds for stochastic gradient descent. In *Conference on Learning Theory*, volume 134, pages 3526–3545, 2021.
- [75] Behnam Neyshabur, Srinadh Bhojanapalli, and Nathan Srebro. A PAC-bayesian approach to spectrally-normalized margin bounds for neural networks. In *International Conference on Learning Representations*, 2018.
- [76] Kenta Oono and Taiji Suzuki. Optimization and generalization analysis of transduction through gradient boosting and application to multi-scale graph neural networks. In *Advances in Neural Information Processing Systems*, 2020.
- [77] Liam Paninski. Estimation of entropy and mutual information. *Neural Comput.*, 15(6):1191–1253, 2003.
- [78] Ankit Pensia, Varun Jog, and Po-Ling Loh. Generalization error bounds for noisy, iterative algorithms. In *2018 IEEE International Symposium on Information Theory*, pages 546–550, 2018.
- [79] María Pérez-Ortiz, Omar Rivasplata, John Shawe-Taylor, and Csaba Szepesvári. Tighter risk certificates for neural networks. *Journal of Machine Learning Research*, 22(227):1–40, 2021.
- [80] Yury Polyanskiy and Yihong Wu. *Information Theory: From Coding to Learning*. Cambridge University Press, 2022.
- [81] Mohamad Rida Rammal, Alessandro Achille, Aditya Golatkar, Suhas Diggavi, and Stefano Soatto. On leave-one-out conditional mutual information for generalization. In *Advances in Neural Information Processing Systems*, pages 10179–10190, 2022.

- [82] Naftali Tishby Ravid Shwartz-Ziv. Opening the black box of deep neural networks via information. *arXiv preprint arXiv:1703.00810*, 2017.
- [83] Arezou Rezazadeh, Sharu Theresa Jose, Giuseppe Durisi, and Osvaldo Simeone. Conditional mutual information-based generalization bound for meta learning. In *2021 IEEE International Symposium on Information Theory*, pages 1176–1181, 2021.
- [84] Omar Rivasplata, Emilio Parrado-Hernández, John Shawe-Taylor, Shiliang Sun, and Csaba Szepesvári. Pac-bayes bounds for stable algorithms with instance-dependent priors. In *Advances in Neural Information Processing Systems*, page 9234–9244, 2018.
- [85] Daniel Russo and James Zou. Controlling bias in adaptive data analysis using information theory. In *Proceedings of the 19th International Conference on Artificial Intelligence and Statistics*, volume 51, pages 1232–1240, 2016.
- [86] Daniel Russo and James Zou. How much does your data exploration overfit? controlling bias via information usage. *IEEE Transactions on Information Theory*, 66(1):302–323, 2020.
- [87] Franco Scarselli, Marco Gori, Ah Chung Tsoi, Markus Hagenbuchner, and Gabriele Monfardini. The graph neural network model. *IEEE Transactions on Neural Networks*, 20(1):61–80, 2009.
- [88] Matthias Seeger. Pac-bayesian generalisation error bounds for gaussian process classification. *Journal of Machine Learning Research*, 3:233–269, 2002.
- [89] Milad Sefidgaran, Amin Gohari, Gaël Richard, and Umut Simsekli. Rate-distortion theoretic generalization bounds for stochastic learning algorithms. In *Proceedings of 35th Conference on Learning Theory*, pages 4416–4463, 2022.
- [90] Prithviraj Sen, Galileo Namata, Mustafa Bilgic, Lise Getoor, Brian Gallagher, and Tina Eliassi-Rad. Collective classification in network data. *AI magazine*, 29(3):93–106, 2008.
- [91] B.M. Shahshahani and D.A. Landgrebe. The effect of unlabeled samples in reducing the small sample size problem and mitigating the hughes phenomenon. *IEEE Transactions on Geoscience and Remote Sensing*, 32(5):1087–1095, 1994.
- [92] Shai Shalev-Shwartz, Ohad Shamir, Nathan Srebro, and Karthik Sridharan. Learnability, stability and uniform convergence. *Journal of Machine Learning Research*, 11:2635–2670, 2010.
- [93] Ohad Shamir. Without-replacement sampling for stochastic gradient methods. In *Advances in Neural Information Processing Systems*, pages 46–54, 2016.
- [94] Ohad Shamir and Shai Shalev-Shwartz. Matrix completion with the trace norm: Learning, bounding, and transducing. *Journal of Machine Learning Research*, 15(98):3401–3423, 2014.
- [95] John Shawe-Taylor and Robert C. Williamson. A pac analysis of a bayesian estimator. In *Proceedings of the Tenth Annual Conference on Computational Learning Theory*, page 2–9, 1997.
- [96] Rakesh Shivanna and Chiranjib Bhattacharyya. Learning on graphs using orthonormal representation is statistically consistent. In *Advances in Neural Information Processing Systems*, pages 3635–3643, 2014.
- [97] Rakesh Shivanna, Bibaswan K. Chatterjee, Raman Sankaran, Chiranjib Bhattacharyya, and Francis R. Bach. Spectral norm regularization of orthonormal representations for graph transduction. In *Advances in Neural Information Processing Systems*, pages 2215–2223, 2015.
- [98] Ravid Shwartz-Ziv, Amichai Painsky, and Naftali Tishby. Representation compression and generalization in deep neural networks, 2019.
- [99] Thomas Steinke and Lydia Zakyntinou. Reasoning about generalization via conditional mutual information. In *Conference on Learning Theory*, volume 125, pages 3437–3452, 2020.
- [100] Mengying Sun, Sendong Zhao, Coryandar Gilvary, Olivier Elemento, Jiayu Zhou, and Fei Wang. Graph convolutional networks for computational drug development and discovery. *Briefings in bioinformatics*, 21(3):919–935, 2020.
- [101] Huayi Tang and Yong Liu. Towards understanding generalization of graph neural networks. In *Proceedings of the 40th International Conference on Machine Learning*, pages 33674–33719, 2023.
- [102] Naftali Tishby, Fernando C. Pereira, and William Bialek. The information bottleneck method. *arXiv preprint arXiv:physics/0004057*, 2000.

- [103] Naftali Tishby and Noga Zaslavsky. Deep learning and the information bottleneck principle. In *2015 IEEE Information Theory Workshop*, pages 1–5, 2015.
- [104] Ilya Tolstikhin, Gilles Blanchard, and Marius Kloft. Localized complexities for transductive learning. In *Proceedings of The 27th Conference on Learning Theory*, volume 35, pages 857–884, 2014.
- [105] Ilya O. Tolstikhin, Nikita Zhivotovskiy, and Gilles Blanchard. Permutational rademacher complexity - A new complexity measure for transductive learning. In Kamalika Chaudhuri, Claudio Gentile, and Sandra Zilles, editors, *Proceedings of the 26th International Conference on Algorithmic Learning Theory - 26th International Conference*, pages 209–223, 2015.
- [106] V. N. Vapnik. *Estimation of Dependences Based on Empirical Data: Empirical Inference Science*. Springer, New York, 1982.
- [107] V. N. Vapnik. *Statistical Learning Theory*. Wiley, New York, 1998.
- [108] Petar Velickovic, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Liò, and Yoshua Bengio. Graph attention networks. In *International Conference on Learning Representations*, 2018.
- [109] Hao Wang, Mario Diaz, José Cândido S. Santos Filho, and Flavio P. Calmon. An information-theoretic view of generalization via wasserstein distance. In *2019 IEEE International Symposium on Information Theory*, pages 577–581, 2019.
- [110] Hao Wang, Yizhe Huang, Rui Gao, and Flavio Calmon. Analyzing the generalization capability of SGLD using properties of gaussian channels. In *Advances in Neural Information Processing Systems*, 2021.
- [111] Xiang Wang, Xiangnan He, Meng Wang, Fuli Feng, and Tat-Seng Chua. Neural graph collaborative filtering. In *Proceedings of the 42nd International ACM SIGIR Conference on Research and Development in Information Retrieval*, page 165–174, 2019.
- [112] Zifeng Wang, Shao-Lun Huang, Ercan Engin Kuruoglu, Jimeng Sun, Xi Chen, and Yefeng Zheng. PAC-bayes information bottleneck. In *International Conference on Learning Representations*, 2022.
- [113] Ziqiao Wang and Yongyi Mao. On the generalization of models trained with SGD: Information-theoretic bounds and implications. In *International Conference on Learning Representations*, 2022.
- [114] Ziqiao Wang and Yongyi Mao. Tighter information-theoretic generalization bounds from supersamples. In *Proceedings of the 40th International Conference on Machine Learning*, pages 36111–36137, 2023.
- [115] Xuetong Wu, Jonathan H. Manton, Uwe Aickelin, and Jingge Zhu. Information-theoretic analysis for transfer learning. In *2020 IEEE International Symposium on Information Theory*, pages 2819–2824, 2020.
- [116] John Lafferty Xiaojin Zhu, Zoubin Ghahramani. Semi-supervised learning using gaussian fields and harmonic functions. In *Proceedings of the 20th International Conference on Machine Learning*, pages 912–919, 2003.
- [117] Aolin Xu and Maxim Raginsky. Information-theoretic analysis of generalization capability of learning algorithms. In *Advances in Neural Information Processing Systems*, pages 2524–2533, 2017.
- [118] Chao Xu, Hong Tao, Jing Zhang, Dewen Hu, and Chenping Hou. Label distribution changing learning with sample space expanding. *Journal of Machine Learning Research*, 24(36):1–48, 2023.
- [119] Da Xu, Chuanwei Ruan, Evren Körpeoglu, Sushant Kumar, and Kannan Achan. Rethinking neural vs. matrix-factorization collaborative filtering: the theoretical perspectives. In *Proceedings of the 38th International Conference on Machine Learning*, pages 11514–11524, 2021.
- [120] Keyulu Xu, Chengtao Li, Yonglong Tian, Tomohiro Sonobe, Ken-ichi Kawarabayashi, and Stefanie Jegelka. Representation learning on graphs with jumping knowledge networks. In *Proceedings of the 35th International Conference on Machine Learning*, volume 80, pages 5449–5458, 2018.
- [121] Jun Yang, Shengyang Sun, and Daniel M. Roy. Fast-rate pac-bayes generalization bounds via shifted rademacher processes. In *Advances in Neural Information Processing Systems*, 2019.
- [122] Zhilin Yang, William W. Cohen, and Ruslan Salakhutdinov. Revisiting semi-supervised learning with graph embeddings. In *Proceedings of the 33rd International Conference on Machine Learning*, pages 40–48, 2016.

- [123] Yige Yuan, Bingbing Xu, Huawei Shen, Qi Cao, Keting Cen, Wen Zheng, and Xueqi Cheng. Towards generalizable graph contrastive learning: An information theory perspective. *arXiv preprint arXiv:2211.10929*, 2022.
- [124] Sergey Zagoruyko and Nikos Komodakis. Wide residual networks. In *Proceedings of the British Machine Vision Conference 2016*, 2016.
- [125] Ruida Zhou, Chao Tian, and Tie Liu. Stochastic chaining and strengthened information-theoretic generalization bounds. In *2022 IEEE International Symposium on Information Theory*, pages 690–695, 2022.
- [126] Wenda Zhou, Victor Veitch, Morgane Austern, Ryan P. Adams, and Peter Orbanz. Non-vacuous generalization bounds at the imagenet scale: a PAC-bayesian compression approach. In *International Conference on Learning Representations*, 2019.

A Notations and Lemma

We introduce additional notations used throughout the paper. We use $[m]$ to denote the set $\{1, 2, \dots, m\}$. The combination number is denoted as $C_m^n = \frac{m!}{n!(m-n)!}$, where $!$ represents factorial. Furthermore, we use \mathbb{N} to denote the set of all non-negative integers, and \mathbb{N}_+ to denote the set of all positive integers. Also, \mathbb{R} is the set of real number and $\mathbb{R}_{\geq 0}$ is the set of non-negative real numbers. The gamma function is denoted as $\Gamma(\cdot)$. We finish this section by introducing the following lemma, which is also termed as the Donsker-Varadhan dual characterization of KL divergence or Gibbs variational principle in the literature. This lemma is the foundation of most information-theoretic and PAC-Bayesian theoretical results.

Lemma 1 (Theorem 4.6 in [80]). *Let P , and Q be two probability measure on \mathcal{X} and $\mathcal{F} \triangleq \{f : \mathcal{X} \rightarrow \mathbb{R}\}$ the family of bounded measurable function. Then we have*

$$D_{\text{KL}}(P||Q) = \sup_{f \in \mathcal{F}} \mathbb{E}_P[f(X)] - \log \mathbb{E}_Q[\exp\{f(X)\}]. \quad (19)$$

B Proof of Theorem 1

We firstly show that $\mathcal{E}(w, Z)$ satisfies sub-gaussian property for a fixed realization w of W . Inspired by [23, 30], we construct the following martingale difference sequences:

$$V_i \triangleq \mathbb{E}[\mathcal{E}(w, Z)|Z_1, \dots, Z_i] - \mathbb{E}[\mathcal{E}(w, Z)|Z_1, \dots, Z_{i-1}], i \in [n]. \quad (20)$$

With this definition, one can verify that $\mathcal{E}(w, Z) - \mathbb{E}[\mathcal{E}(w, Z)] = \sum_{i=1}^n V_i$. Note that V_i is a function of Z_1, \dots, Z_i . Define

$$\begin{aligned} L_i &\triangleq \inf_z \mathbb{E}[\mathcal{E}(w, Z)|Z_1, \dots, Z_{i-1}, Z_i = z] - \mathbb{E}[\mathcal{E}(w, Z)|Z_1, \dots, Z_{i-1}], \\ U_i &\triangleq \sup_z \mathbb{E}[\mathcal{E}(w, Z)|Z_1, \dots, Z_{i-1}, Z_i = z] - \mathbb{E}[\mathcal{E}(w, Z)|Z_1, \dots, Z_{i-1}], \end{aligned}$$

we have $L_i \leq V_i \leq U_i$. Then we show that $U_i - L_i$ is a bounded random variable when Z_1, \dots, Z_{i-1} are given:

$$\begin{aligned} U_i - L_i &= \sup_z \mathbb{E}[\mathcal{E}(w, Z)|Z_1, \dots, Z_{i-1}, Z_i = z] - \inf_z \mathbb{E}[\mathcal{E}(w, Z)|Z_1, \dots, Z_{i-1}, Z_i = z] \\ &= \sup_{z, \tilde{z}} \left\{ \mathbb{E}[\mathcal{E}(w, Z)|Z_1, \dots, Z_{i-1}, Z_i = z] - \mathbb{E}[\mathcal{E}(w, Z)|Z_1, \dots, Z_{i-1}, Z_i = \tilde{z}] \right\} \\ &= \frac{u!(m-i)!C_{n-i-1}^{m-i}}{(n-i)!} \cdot \frac{(m+u)B}{mu} \\ &= \frac{(m+u)B}{m(m+u-i)} \triangleq c_i. \end{aligned} \quad (21)$$

Since $\mathbb{E}[V_i|Z_1, \dots, Z_{i-1}] = 0$, by Hoeffding's lemma, $\mathbb{E}[e^{\lambda V_i}|Z_1, \dots, Z_{i-1}] \leq e^{\frac{\lambda^2 c_i^2}{8}}$ holds for any $\lambda \in \mathbb{R}$. Note that

$$\begin{aligned} \mathbb{E} \left[\exp \left\{ \lambda \sum_{i=1}^n V_i \right\} \right] &= \mathbb{E} \left[\mathbb{E} \left[\exp \left\{ \lambda \sum_{i=1}^{n-1} V_i \right\} \exp \{V_n\} \right] \middle| Z_1, \dots, Z_{n-1} \right] \\ &= \mathbb{E} \left[\exp \left\{ \lambda \sum_{i=1}^{n-1} V_i \right\} \mathbb{E}[\exp \{V_n\}] \middle| Z_1, \dots, Z_{n-1} \right] \\ &\leq \exp \left\{ \frac{\lambda^2 c_n^2}{8} \right\} \mathbb{E} \left[\exp \left\{ \lambda \sum_{i=1}^{n-1} V_i \right\} \right]. \end{aligned} \quad (22)$$

By recursively repeating the above process, we obtain

$$\begin{aligned} &\mathbb{E} \left[\exp \left\{ \lambda \sum_{i=1}^n V_i \right\} \right] \\ &\leq \exp \left\{ \frac{\lambda^2}{8} \sum_{i=1}^n c_i^2 \right\} = \exp \left\{ \frac{\lambda^2 (m+u)^2 B^2}{8m^2} \sum_{i=1}^n \frac{1}{(m+u-i)^2} \right\} \\ &\leq \exp \left\{ \frac{\lambda^2 B^2 (m+u)^2}{8m(u-1/2)(m+u-1/2)} \right\}. \end{aligned} \quad (23)$$

Due to the symmetric of training and test partition, the final bound is obtained by taking the largest one among them:

$$\mathbb{E} \left[\exp \left\{ \lambda \sum_{i=1}^n V_i \right\} \right] \leq \exp \left\{ \frac{\lambda^2 B^2 (m+u)^2}{8mu(m+u-1/2)} \cdot \frac{2 \max(m, u)}{2 \max(m, u) - 1} \right\} \quad (24)$$

Combining Eq. (24) and the facts that $\mathcal{E}(w, Z) - \mathbb{E}[\mathcal{E}(w, Z)] = \sum_{i=1}^n V_i$ and $\mathbb{E}[\mathcal{E}(w, Z)] = 0$, we obtain

$$\begin{aligned} & \mathbb{E}_Z [\exp \{ \lambda (R_u(w, Z) - R_m(w, Z)) \}] \\ & \leq \exp \left\{ \frac{\lambda^2 B^2 (m+u)^2}{8mu(m+u-1/2)} \cdot \frac{2 \max(m, u)}{2 \max(m, u) - 1} \right\} \\ & = \exp \left\{ \frac{\lambda^2 (m+u) C_{m,u}}{8mu} \right\} \end{aligned} \quad (25)$$

where $C_{m,u} \triangleq \frac{2B^2(m+u) \max(m, u)}{(m+u-1/2)(2 \max(m, u)-1)}$. Denote by Z' the independent copy of Z , which is independent from W and has the same distribution as Z . Then we have

$$\begin{aligned} & \log \mathbb{E}_{W, Z'} [\exp \{ \lambda (R_u(W, Z') - R_m(W, Z')) \}] \\ & = \log \left(\int_w \mathbb{E}_{Z'} [\exp \{ \lambda (R_u(w, Z') - R_m(w, Z')) \}] dP_W(w) \right) \\ & \leq \log \left(\int_w \exp \left\{ \frac{\lambda^2 (m+u) C_{m,u}}{8mu} \right\} dP_W(w) \right) \\ & = \frac{\lambda^2 C_{m,u}}{8} \left(\frac{1}{m} + \frac{1}{u} \right). \end{aligned} \quad (26)$$

By Donsker-Varadhan's variational formula, for any $\lambda \in \mathbb{R}$:

$$\begin{aligned} & D_{\text{KL}}(P_{Z,W} || P_{Z',W}) \\ & \geq \mathbb{E}_{Z,W} [\lambda (R_u(W, Z) - R_m(W, Z))] - \log \mathbb{E}_{Z',W} [\exp \{ \lambda (R_u(W, Z') - R_m(W, Z')) \}] \\ & \geq \mathbb{E}_{Z,W} [\lambda (R_u(W, Z) - R_m(W, Z))] - \frac{\lambda^2 C_{m,u}}{8} \left(\frac{1}{m} + \frac{1}{u} \right), \end{aligned} \quad (27)$$

which implies that

$$|\mathbb{E}_{Z,W} [R_u(W, Z) - R_m(W, Z)]| \leq \sqrt{\frac{C_{m,u}}{2} \left(\frac{1}{m} + \frac{1}{u} \right) I(Z; W)}. \quad (28)$$

This finishes the proof for the first part. For the second part, note that Eq. (25) can be rewritten as $\mathbb{E}_Z [\exp \{ \lambda \mathcal{E}(w, Z) \}] \leq \exp \{ \lambda^2 \sigma_{m,u} \}$, where $\sigma_{m,u} \triangleq C_{m,u} (1/m + 1/u)/8$. Similarly we have $\mathbb{E}_Z [\exp \{ -\lambda \mathcal{E}(w, Z) \}] \leq \exp \{ \lambda^2 \sigma_{m,u} \}$. Then we have

$$\mathbb{P} \{ |\mathcal{E}(w, Z)| \geq t \} \leq \mathbb{P} \{ \mathcal{E}(w, Z) \geq t \} + \mathbb{P} \{ \mathcal{E}(w, Z) \leq -t \} \leq 2 \exp \left\{ -\frac{t^2}{4\sigma_{m,u}} \right\}. \quad (29)$$

where the first and the second inequality are due to the Boole's inequality and the Chernoff technique, respectively. For any $k \in \mathbb{N}_+$, we have

$$\begin{aligned} \mathbb{E} [|\mathcal{E}(w, Z)|^k] &= \int_0^\infty \mathbb{P} \{ |\mathcal{E}(w, Z)|^k \geq u \} du \\ &= k \int_0^\infty \mathbb{P} \{ |\mathcal{E}(w, Z)| \geq t \} t^{k-1} dt \\ &\leq 2k \int_0^\infty \exp \left\{ -\frac{t^2}{4\sigma_{m,u}} \right\} t^{k-1} dt = (4\sigma_{m,u})^{\frac{k}{2}} k \Gamma(k/2), \end{aligned} \quad (30)$$

which implies that

$$\mathbb{E} [\exp \{ \lambda \mathcal{E}^2(w, Z) \}] = 1 + \sum_{k=1}^\infty \frac{\lambda^k}{k!} \mathbb{E} [|\mathcal{E}(w, Z)|^{2k}] \leq 1 + 2 \sum_{k=1}^\infty (4\lambda\sigma_{m,u})^k. \quad (31)$$

By Donsker-Varadhan's variational formula, for any $\lambda \in \mathbb{R}$:

$$\begin{aligned} & D_{\text{KL}}(P_{Z,W} || P_{Z,W'}) \\ & \geq \mathbb{E}_{Z,W} [\lambda (R_u(W, Z) - R_m(W, Z))^2] - \log \mathbb{E}_{Z,W'} [\exp \{ \lambda (R_u(W', Z) - R_m(W', Z))^2 \}] \\ & \geq \mathbb{E}_{Z,W} [\lambda (R_u(W, Z) - R_m(W, Z))^2] - \log \left(1 + 2 \sum_{k=1}^\infty (4\lambda\sigma_{m,u})^k \right), \end{aligned} \quad (32)$$

Let $\lambda = 1/8\sigma_{m,u}$ and plugging into $\sigma_{m,u} \triangleq C_{m,u} (1/m + 1/u)/8$, we obtain

$$\mathbb{E}_{Z,W} [(R_u(W, Z) - R_m(W, Z))^2] \leq C_{m,u} \left(\frac{1}{m} + \frac{1}{u} \right) (I(Z; W) + \log 3). \quad (33)$$

C Proof of Theorem 2

Denote by $Z^{(1)}, \dots, Z^{(k)}$ the k independent copy of Z . By running a transductive algorithm \mathcal{A} on each $Z^{(j)}$ respectively, we obtain the corresponding output $W^{(j)} = \mathcal{A}(Z^{(j)})$ for $j \in [k]$. By this way, $(Z^{(j)}, W^{(j)})$ can be regarded as independent copy of (Z, W) for $j \in [k]$. Now assume that there is a monitor that returns

$$(J^*, R^*) \triangleq \operatorname{argmax}_{j \in [k], r \in \{\pm 1\}} r \mathcal{E}(W^{(j)}, Z^{(j)}), \quad W^* = W_{J^*}.$$

One can verify that

$$R^* \mathcal{E}(W^{(J^*)}, Z^{(J^*)}) = \max_{j \in [k]} |\mathcal{E}(Z^{(j)}, W^{(j)})|.$$

Now taking expectation on both side, we have

$$\mathbb{E}_{Z^{(1)}, \dots, Z^{(k)}, J^*, R^*, W^*} [R^* \mathcal{E}(W^{(J^*)}, Z^{(J^*)})] = \mathbb{E}_{Z^{(1)}, \dots, Z^{(k)}, W_1, \dots, W_k} \left[\max_{j \in [k]} |\mathcal{E}(Z^{(j)}, W^{(j)})| \right].$$

Using the same procedure used in the proof of Theorem 1, we have

$$\log \mathbb{E}_{J^*, R^*, W^*} \mathbb{E}_{Z^{(1)}, \dots, Z^{(k)}} \left[\exp \left\{ \lambda R^* \mathcal{E}(W^{(J^*)}, Z^{(J^*)}) \right\} \right] \leq \frac{\lambda^2 C_{m,u}}{8} \left(\frac{1}{m} + \frac{1}{u} \right).$$

By Donsker-Varadhan's variational formula, the following inequality holds for any $\lambda \in \mathbb{R}$:

$$\begin{aligned} & D(P_{Z^{(1)}, \dots, Z^{(k)}, J^*, R^*, W^*} \| P_{Z^{(1)}, \dots, Z^{(k)}} \otimes P_{J^*, R^*, W^*}) \\ & \geq \mathbb{E}_{Z^{(1)}, \dots, Z^{(k)}, J^*, R^*, W^*} \left[\lambda R^* \mathcal{E}(W^{(J^*)}, Z^{(J^*)}) \right] \\ & \quad - \log \mathbb{E}_{J^*, R^*, W^*} \mathbb{E}_{Z^{(1)}, \dots, Z^{(k)}} \left[\exp \left\{ \lambda R^* \mathcal{E}(W^{(J^*)}, Z^{(J^*)}) \right\} \right] \\ & \geq \lambda \mathbb{E}_{Z^{(1)}, \dots, Z^{(k)}, J^*, R^*, W^*} \left[R^* \mathcal{E}(W^{(J^*)}, Z^{(J^*)}) \right] - \frac{\lambda^2 C_{m,u}}{8} \left(\frac{1}{m} + \frac{1}{u} \right). \end{aligned}$$

which implies

$$\begin{aligned} & \mathbb{E}_{Z^{(1)}, \dots, Z^{(k)}, J^*, R^*, W^*} \left[R^* \mathcal{E}(W^{(J^*)}, Z^{(J^*)}) \right] \\ & \leq \sqrt{\frac{C_{m,u}}{2} \left(\frac{1}{m} + \frac{1}{u} \right) I(Z^{(1)}, \dots, Z^{(k)}; J^*, R^*, W^*)}. \end{aligned} \tag{34}$$

Next we provide a upper bound for the mutual information. Note that

$$\begin{aligned} & I(Z^{(1)}, \dots, Z^{(k)}; J^*, R^*, W^*) \\ & \leq I(S(\pi_1), \dots, S(\pi_k); J^*, R^*, W^*, W_1, \dots, W_k) \\ & = I(Z^{(1)}, \dots, Z^{(k)}; W_1, \dots, W_k) + I(S(\pi_1), \dots, S(\pi_k); J^*, R^*, W^* | W_1, \dots, W_k) \\ & = \sum_{j=1}^k I(S(\pi_j); W_j) + I(Z^{(1)}, \dots, Z^{(k)}; J^*, R^*, W^* | W_1, \dots, W_k) \\ & \leq k I(Z; W) + \log(2k). \end{aligned} \tag{35}$$

where we have use that $(S(\pi_j), W_j), j \in [k]$ are independent. Plugging Eq. (35) into Eq. (34) yields

$$\mathbb{E}_{Z^{(1)}, \dots, Z^{(k)}, W^{(1)}, \dots, W^{(k)}} \left[\max_{j \in [k]} |\mathcal{E}(Z^{(j)}, W^{(j)})| \right] \leq \sqrt{\frac{C_{m,u}}{2} \left(\frac{1}{m} + \frac{1}{u} \right) (\log(2k) + k I(Z, W))}.$$

Since $(Z^{(j)}, W^{(j)})$ are independent copy of (Z, W) , for any $\alpha > 0$ we have

$$\mathbb{P}_{Z^{(1)}, W^{(1)}, \dots, Z^{(k)}, W^{(k)}} \left\{ \max_{j \in [k]} |\mathcal{E}(Z^{(j)}, W^{(j)})| < \alpha \right\} = (\mathbb{P}_{Z, W} \{ |\mathcal{E}(Z, W)| < \alpha \})^k.$$

By Markov's inequality:

$$\begin{aligned} & \mathbb{P}_{Z^{(1)}, W^{(1)}, \dots, Z^{(k)}, W^{(k)}} \left\{ \max_{j \in [k]} |\mathcal{E}(Z^{(j)}, W^{(j)})| \geq \alpha \right\} \\ & \leq \frac{1}{\alpha} \mathbb{E}_{Z^{(1)}, \dots, Z^{(k)}, W^{(1)}, \dots, W^{(k)}} \left[\max_{j \in [k]} |\mathcal{E}(Z^{(j)}, W^{(j)})| \right] \\ & \leq \frac{1}{\alpha} \sqrt{\frac{C_{m,u}}{2} \left(\frac{1}{m} + \frac{1}{u} \right) (\log(2k) + k I(Z, W))}. \end{aligned}$$

Therefore,

$$\begin{aligned}
& \mathbb{P}_{Z,W} \{|\mathcal{E}(Z, W)| \geq \alpha\} \\
&= 1 - \mathbb{P}_{Z,W} \{|\mathcal{E}(Z, W)| < \alpha\} \\
&= 1 - \left(\mathbb{P}_{Z^{(1)}, W^{(1)}, \dots, Z^{(k)}, W^{(k)}} \left\{ \max_{j \in [k]} |\mathcal{E}(Z^{(j)}, W^{(j)})| < \alpha \right\} \right)^{\frac{1}{k}} \\
&= 1 - \left(1 - \mathbb{P}_{Z^{(1)}, W^{(1)}, \dots, Z^{(k)}, W^{(k)}} \left\{ \max_{j \in [k]} |\mathcal{E}(Z^{(j)}, W^{(j)})| < \alpha \right\} \right)^{\frac{1}{k}} \\
&\leq 1 - \left(1 - \frac{1}{\alpha} \sqrt{\frac{C_{m,u}}{2} \left(\frac{1}{m} + \frac{1}{u} \right) (\log(2k) + kI(Z, W))} \right)^{\frac{1}{k}}.
\end{aligned}$$

Let $\alpha = 2\sqrt{\frac{C_{m,u}}{2} \left(\frac{1}{m} + \frac{1}{u} \right) (\log(2k) + kI(Z, W))}$ and $k = \lfloor \frac{1}{\delta} \rfloor$, we have obtained the result. Let $k = 1$, we obtain

$$\mathbb{E}_{Z,W} [|\mathcal{E}(Z, W)|] \leq \sqrt{\frac{C_{m,u}}{2} \left(\frac{1}{m} + \frac{1}{u} \right) (\log 2 + I(Z, W))}.$$

D Proof of Proposition 1

Denote by \mathcal{S} the set containing all value of S . For any $\tilde{z} \in \tilde{\mathcal{Z}}$ and $s \in \mathcal{S}$, we use

$$\bar{z}(\tilde{z}, s) \triangleq (\tilde{z}_{1,s_1}, \dots, \tilde{z}_{m,s_m}, \tilde{z}_{1,1-s_1}, \dots, \tilde{z}_{m,1-s_m}) \quad (36)$$

to denote the random permutation vector induced by \tilde{z} and s . Let $\bar{\mathcal{Z}} : \{\bar{z}(\tilde{z}, s) | \tilde{z} \in \tilde{\mathcal{Z}}, s \in \mathcal{S}\}$, it is sufficient to proof that $\bar{\mathcal{Z}} = \mathfrak{Z}$, where \mathfrak{Z} is the set includes all possible values of Z^n defined in Subsection 3.2. Note that each element \tilde{z} in $\tilde{\mathcal{Z}}$ differs by each other, since \tilde{z} is a partitions of $m + u$ elements into m subsets, where each subset contains 2 elements. Thus, the cardinality of $\tilde{\mathcal{Z}}$ is $\frac{(2m)!}{2^m}$. Furthermore, for fixed \tilde{z} and $s_1, s_2 \in \mathcal{S}$, it is clear that $s_1 \neq s_2$ implies $\bar{z}(\tilde{z}, s_1) \neq \bar{z}(\tilde{z}, s_2)$, due to the fact that $\bar{z}(\tilde{z}, s_1) = \bar{z}(\tilde{z}, s_2)$ if and only if $\bar{z}(\tilde{z}, s_1)_j = \bar{z}(\tilde{z}, s_2)_j$ for $j \in [2m]$. Here $\bar{z}(\tilde{z}, s)_j$ represents the j -th entry in the sequence. Now we claim that for any two element $\tilde{z}_1(\tilde{z}_1, s_1), \tilde{z}_2(\tilde{z}_2, s_2) \in \tilde{\mathcal{Z}}$, if $\bar{z}_1(\tilde{z}_1, s_1) = \bar{z}_2(\tilde{z}_2, s_2)$, then $\tilde{z}_1 = \tilde{z}_2$ holds. To see this, we show that \tilde{z} can be uniquely determined when seeing $\bar{z}(\tilde{z}, s)$. Recall that \tilde{z} is a sequence containing m element, and each element \tilde{z}_j is a set. By Eq. (36), we conclude that \tilde{z}_{1,s_1} comes from \tilde{z}_1 , \tilde{z}_{2,s_2} comes from \tilde{z}_2 , and so on. Similarly, $\tilde{z}_{1,1-s_1}$ comes from \tilde{z}_1 , $\tilde{z}_{2,1-s_2}$ comes from \tilde{z}_2 , and so on. By this way, we have recovered \tilde{z} from $\bar{z}(\tilde{z}, s)$ and it is unique. Together with the fact $\bar{z}(\tilde{z}, s_1) = \bar{z}(\tilde{z}, s_2) \implies s_1 = s_2$ that we have just shown, we conclude that $\bar{z}_1(\tilde{z}_1, s_1) = \bar{z}_2(\tilde{z}_2, s_2) \implies \tilde{z}_1 = \tilde{z}_2, s_1 = s_2$, which suggest that $(\tilde{z}, s) \rightarrow \bar{z}(\tilde{z}, s)$ is a one-to-one mapping. Since $|\mathcal{S}| = 2^m$, we have $|\bar{\mathcal{Z}}| = |\tilde{\mathcal{Z}}||\mathcal{S}| = (2m)!$. Note that $\bar{\mathcal{Z}} \subseteq \mathfrak{Z}$ and $|\bar{\mathcal{Z}}| = |\mathfrak{Z}|$, we conclude that $\bar{\mathcal{Z}} = \mathfrak{Z}$.

E Proof of Theorem 3

We firstly present a warm-up example for illustrating the proposed transductive supersamples. Suppose $m = 2$, all entries in $\tilde{\mathcal{Z}}$ are as follows:

- $(\{\mathbf{z}_1, \mathbf{z}_2\}, \{\mathbf{z}_3, \mathbf{z}_4\})$
- $(\{\mathbf{z}_1, \mathbf{z}_3\}, \{\mathbf{z}_2, \mathbf{z}_4\})$
- $(\{\mathbf{z}_1, \mathbf{z}_4\}, \{\mathbf{z}_2, \mathbf{z}_3\})$
- $(\{\mathbf{z}_2, \mathbf{z}_3\}, \{\mathbf{z}_1, \mathbf{z}_4\})$
- $(\{\mathbf{z}_2, \mathbf{z}_4\}, \{\mathbf{z}_1, \mathbf{z}_3\})$
- $(\{\mathbf{z}_3, \mathbf{z}_4\}, \{\mathbf{z}_1, \mathbf{z}_2\})$

Here we use set $\{\cdot\}$ and tuple (\cdot) to indicate whether we consider the order or not. Now we give the formal proof. Denote by w and \tilde{z} the fixed realizations of W and $\tilde{\mathcal{Z}}$. For any $\lambda \in \mathbb{R}$, by Hoeffding's Lemma:

$$\begin{aligned}
& \mathbb{E}_S [\exp \{ \lambda \mathcal{E}(w, \tilde{z}, S) \}] \\
&= \mathbb{E}_S \left[\exp \left\{ \frac{\lambda}{m} \sum_{i=1}^m \ell(w, \tilde{z}_{i,s_i}) - \ell(w, \tilde{z}_{i,1-s_i}) \right\} \right] \leq \exp \left\{ \frac{\lambda^2 B^2}{2m} \right\}. \quad (37)
\end{aligned}$$

Let S' be the independent copy of S , we have

$$\begin{aligned} & \log \mathbb{E}_{S', W | \tilde{Z} = \tilde{z}} [\exp \{ \lambda \mathcal{E}(W, \tilde{z}, S') \}] \\ &= \log \left(\int_w \mathbb{E}_{S'} [\exp \{ \lambda \mathcal{E}(w, \tilde{z}, S') \}] dP_{W | \tilde{Z} = \tilde{z}}(w) \right) \leq \frac{\lambda^2 B^2}{2m}. \end{aligned} \quad (38)$$

where we have used the fact that $P_{S', W | \tilde{Z} = \tilde{z}} = P_{W | \tilde{Z} = \tilde{z}} P_{S'}$, due to S' is independent from both \tilde{Z} and W . By Donsker-Varadhan's variational formula, for any $\lambda \in \mathbb{R}$:

$$\begin{aligned} I(S; W | \tilde{Z} = \tilde{z}) &= D_{\text{KL}}(P_{S, W | \tilde{Z} = \tilde{z}} \| P_{S', W | \tilde{Z} = \tilde{z}}) \\ &\geq \mathbb{E}_{S, W | \tilde{Z} = \tilde{z}} [\lambda \mathcal{E}(W, \tilde{z}, S)] - \log \mathbb{E}_{S', W | \tilde{Z} = \tilde{z}} [\exp \{ \lambda \mathcal{E}(W, \tilde{z}, S') \}] \\ &\geq \lambda \mathbb{E}_{S, W | \tilde{Z} = \tilde{z}} [\mathcal{E}(W, \tilde{z}, S)] - \frac{\lambda^2 B^2}{2m}, \end{aligned} \quad (39)$$

which implies that

$$\left| \mathbb{E}_{S, W | \tilde{Z} = \tilde{z}} [\mathcal{E}(W, \tilde{z}, S)] \right| \leq \sqrt{\frac{2B^2}{m} I(S; W | \tilde{Z} = \tilde{z})}.$$

Taking expectation over \tilde{Z} on both side, we have obtain

$$\begin{aligned} & \left| \mathbb{E}_{Z, W} [R_u(W, Z) - R_m(W, Z)] \right| = \left| \mathbb{E}_{\tilde{Z}, S, W} [\mathcal{E}(W, \tilde{Z}, S)] \right| \\ & \leq \mathbb{E}_{\tilde{Z}} \left| \mathbb{E}_{S, W | \tilde{Z}} [\mathcal{E}(W, \tilde{Z}, S)] \right| \leq \mathbb{E}_{\tilde{Z}} \sqrt{\frac{2B^2}{m} I(S; W | \tilde{Z})}. \end{aligned} \quad (40)$$

For the second part, note that Eq. (37) can be rewritten as $\mathbb{E}_S [\exp \{ \lambda \mathcal{E}(w, \tilde{z}, S) \}] \leq \exp \{ \lambda^2 B^2 / 2m \}$. Similarly we have $\mathbb{E}_S [\exp \{ -\lambda \mathcal{E}(w, \tilde{z}, S) \}] \leq \exp \{ \lambda^2 B^2 / 2m \}$. By the same technique in Section B, we have

$$\mathbb{E}_{S, W | \tilde{Z} = \tilde{z}} [\mathcal{E}^2(W, \tilde{z}, S)] \leq \frac{4B^2}{m} (I(S; W | \tilde{Z} = \tilde{z}) + \log 3). \quad (41)$$

Taking expectation on both side, we have obtain

$$\mathbb{E}_{S, W} [(R_u(W, Z) - R_m(W, Z))^2] = \mathbb{E}_{\tilde{Z}, S, W} [\mathcal{E}^2(W, \tilde{Z}, S)] \leq \frac{4B^2}{m} (I(S; W | \tilde{Z}) + \log 3). \quad (42)$$

We close this proof by presenting the results for more ordinary cases where $u = km$ with $k \in \mathbb{N}_+$. By Definition 2, the transductive generalization under this setting is defined by

$$\mathcal{E}(W, \tilde{Z}, S) \triangleq \frac{1}{m} \sum_{i=1}^m \ell(W, \tilde{Z}_{i, S_i}) - \frac{1}{km} \sum_{i=1}^m \sum_{j=0}^k \ell(W, \tilde{Z}_{i, \tilde{S}_{i,j}}). \quad (43)$$

One can verify that $\mathbb{E}_{S, W} [R_u(W, Z) - R_m(W, Z)] = \mathbb{E}_{\tilde{Z}, S, W} [\mathcal{E}(W, \tilde{Z}, S)]$ holds. Following the same procedure we have

$$\begin{aligned} & \left| \mathbb{E}_{S, W} [R_u(W, Z) - R_m(W, Z)] \right| \leq \mathbb{E}_{\tilde{Z}} \sqrt{\frac{2(k+1)B^2}{n} I(S; W | \tilde{Z})} \\ & \mathbb{E}_{S, W} [(R_u(W, Z) - R_m(W, Z))^2] \leq \frac{4(k+1)B^2}{n} (I(S; W | \tilde{Z}) + \log 3). \end{aligned} \quad (44)$$

F Proof of Corollary 1

Denote by $g(F_i, S_i, \tilde{Y}_i) \triangleq \ell(F_i, S_i, \tilde{Y}_i) - \ell(F_i, 1-S_i, \tilde{Y}_i)$ the function of (F_i, S_i, \tilde{Y}_i) . Let f_i and $\tilde{z}_i = (\tilde{x}_i, \tilde{y}_i)$ be the fixed realizations of F_i and \tilde{Z}_i . For any $\lambda \in \mathbb{R}$ and $i \in [m]$, by Hoeffding's Lemma:

$$\mathbb{E}_{S_i} [\exp \{ \lambda g(f_i, S_i, \tilde{y}_i) \}] \leq \exp \left\{ \frac{\lambda^2 B^2}{2} \right\}. \quad (45)$$

Let S'_i be the independent copy of S_i , by Donsker-Varadhan's variational formula:

$$\begin{aligned} I(F_i; S_i | \tilde{Z} = \tilde{z}) &\geq \lambda \mathbb{E}_{F_i, S_i | \tilde{Z} = \tilde{z}} [g(F_i, S_i, \tilde{y}_i)] - \log \mathbb{E}_{F_i, S'_i | \tilde{Z} = \tilde{z}} [\exp \{ \lambda g(F_i, S'_i, \tilde{y}_i) \}] \\ &\geq \lambda \mathbb{E}_{F_i, S_i | \tilde{Z} = \tilde{z}} [g(F_i, S_i, \tilde{y}_i)] - \frac{\lambda^2 B^2}{2}. \end{aligned} \quad (46)$$

Then we have

$$\left| \mathbb{E}_{S, W | \tilde{Z} = \tilde{z}} [\ell(W, \tilde{z}_{i, S_i}) - \ell(W, \tilde{z}_{i, 1-S_i})] \right| = \left| \mathbb{E}_{F_i, S_i | \tilde{Z} = \tilde{z}} [g(F_i, S_i, \tilde{y}_i)] \right| \leq B \sqrt{2I(F_i; S_i | \tilde{Z} = \tilde{z})}, \quad (47)$$

which implies that

$$\begin{aligned}
& |\mathbb{E}_{Z,W} [R_u(W, Z) - R_m(W, Z)]| = |\mathbb{E}_{\tilde{Z}, S, W} [\mathcal{E}(W, \tilde{Z}, S)]| \\
& \leq \mathbb{E}_{\tilde{Z}} \left| \frac{1}{m} \sum_{i=1}^m \mathbb{E}_{S, W | \tilde{Z}} [\ell(W, \tilde{Z}_{i, S_i}) - \ell(W, \tilde{Z}_{i, 1-S_i})] \right| \\
& \leq \frac{1}{m} \sum_{i=1}^m \mathbb{E}_{\tilde{Z}} |\mathbb{E}_{S, W | \tilde{Z}} [\ell(W, \tilde{Z}_{i, S_i}) - \ell(W, \tilde{Z}_{i, 1-S_i})]| \\
& = \frac{1}{m} \sum_{i=1}^m \mathbb{E}_{\tilde{Z}} |\mathbb{E}_{F_i, S_i | \tilde{Z}} [g(F_i, S_i, \tilde{Y}_i)]| \leq \frac{B}{m} \sum_{i=1}^m \mathbb{E}_{\tilde{Z}} \sqrt{2I(\ell(F_i; S_i | \tilde{Z}))}.
\end{aligned} \tag{48}$$

Denote by $g(L_i, S_i) = L_{i, S_i} - L_{i, 1-S_i}$, by the same technique we have

$$|\mathbb{E}_{S, W} [R_u(W, Z) - R_m(W, Z)]| \leq \frac{1}{m} \sum_{i=1}^m \mathbb{E}_{\tilde{Z}} \sqrt{2I(L_i; S_i | \tilde{Z})}. \tag{49}$$

Denote by $g(\Delta_i, S_i) \triangleq (-1)^{S_i} \Delta_i$, by the same technique we have

$$|\mathbb{E}_{S, W} [R_u(W, Z) - R_m(W, Z)]| \leq \frac{1}{m} \sum_{i=1}^m \mathbb{E}_{\tilde{Z}} \sqrt{2I(\Delta_i; S_i | \tilde{Z})}. \tag{50}$$

G Proof of Theorem 4

By Markov inequality, for any distribution P that independent to Z :

$$\mathbb{P} \left\{ \mathbb{E}_{W \sim P} [e^{\lambda(R_u(W, Z) - R_m(W, Z))}] \geq \frac{1}{\delta} \mathbb{E}_Z \mathbb{E}_{W \sim P} [e^{\lambda(R_u(W, Z) - R_m(W, Z))}] \right\} \leq \delta. \tag{51}$$

By Donsker-Varadhan's variational formula and Eq. (25), for any distribution Q , with probability at least $1 - \delta$ over the randomness of Z :

$$\begin{aligned}
& \lambda \mathbb{E}_{W \sim Q} [R_u(W, Z) - R_m(W, Z)] \\
& \leq \text{D}_{\text{KL}}(Q \| P) + \log \left(\mathbb{E}_{W \sim P} [e^{\lambda(R_u(W, Z) - R_m(W, Z))}] \right) \\
& \leq \text{D}_{\text{KL}}(Q \| P) + \log \left(\frac{1}{\delta} \right) + \log \mathbb{E}_Z \mathbb{E}_{W \sim P} [e^{\lambda(R_u(W, Z) - R_m(W, Z))}] \\
& = \text{D}_{\text{KL}}(Q \| P) + \log \left(\frac{1}{\delta} \right) + \log \left(\int_w \mathbb{E}_{Z'} [\exp \{ \lambda(R_u(w, Z') - R_m(w, Z')) \}] dP(w) \right) \\
& \leq \text{D}_{\text{KL}}(Q \| P) + \log \left(\frac{1}{\delta} \right) + \frac{\lambda^2 C_{m, u}}{8} \left(\frac{1}{m} + \frac{1}{u} \right),
\end{aligned} \tag{52}$$

which implies that: for any distribution Q , with probability at least $1 - \delta$ over the randomness of Z ,

$$|\mathbb{E}_{W \sim Q} [R_u(W, Z) - R_m(W, Z)]| \leq \sqrt{\frac{C_{m, u}}{2} \left(\frac{1}{m} + \frac{1}{u} \right) \left(\text{D}_{\text{KL}}(Q \| P) + \log \left(\frac{1}{\delta} \right) \right)}. \tag{53}$$

Note that by setting $Q = P_{W|S}$ and $P = P_W$, we recover a degenerated version of Theorem 1 holds with probability $1 - \delta$:

$$\begin{aligned}
& |\mathbb{E}_{W, Z} [R_u(W, Z) - R_m(W, Z)]| \leq \mathbb{E}_Z [|\mathbb{E}_{W \sim Q} [R_u(W, Z) - R_m(W, Z)]|] \\
& \leq \mathbb{E}_Z \left[\sqrt{\frac{C_{m, u}}{2} \left(\frac{1}{m} + \frac{1}{u} \right) \left(\text{D}_{\text{KL}}(P_{W|S} \| P_W) + \log \left(\frac{1}{\delta} \right) \right)} \right] \\
& \leq \sqrt{\frac{C_{m, u}}{2} \left(\frac{1}{m} + \frac{1}{u} \right) \left(\mathbb{E}_Z [\text{D}_{\text{KL}}(P_{W|S} \| P_W)] + \log \left(\frac{1}{\delta} \right) \right)} \\
& = \sqrt{\frac{C_{m, u}}{2} \left(\frac{1}{m} + \frac{1}{u} \right) \left(I(S; W) + \log \left(\frac{1}{\delta} \right) \right)}.
\end{aligned} \tag{54}$$

It is worth mentioning that we can also recover Theorem 1 from PAC-Bayesian perspective. Denote by $R_{m+u}(W, Z) \triangleq \frac{1}{m+u} \sum_{i=1}^{m+u} \ell(W, Z_i) = \frac{m}{m+u} R_m(W, Z) + \frac{u}{m+u} R_u(W, Z)$ the error on Z . Denote by

$\mathcal{D}(p, q) \triangleq p \log \frac{p}{q} + (1-p) \log \frac{1-p}{1-q}$ the KL divergence between two Bernoulli distributions with success probability p and q . Then the \mathcal{D} -function introduced in [10] is expressed by $\mathcal{D}_\beta^*(p, q) \triangleq \mathcal{D}(p, q) + \frac{1-\beta}{\beta} \mathcal{D}(\frac{q-\beta p}{1-\beta}, q)$. By Theorem 5 and Theorem 6 in [10], for fixed realization w of W :

$$\mathbb{E}_Z [\exp \{m \mathcal{D}_\beta^*(R_m(w, Z), R_{m+u}(w, Z))\}] \leq 3 \log(m) \sqrt{\frac{mu}{m+u}}, \quad (55)$$

which implies that

$$\begin{aligned} & \log \mathbb{E}_{W \otimes Z} [\exp \{m \mathcal{D}_\beta^*(R_m(W, Z), R_{m+u}(W, Z))\}] \\ &= \log \left(\int_w \mathbb{E}_Z [\exp \{m \mathcal{D}_\beta^*(R_m(w, Z), R_{m+u}(w, Z))\}] dP_W(w) \right) \\ &\leq \log \left(3 \log(m) \sqrt{\frac{mu}{m+u}} \right). \end{aligned} \quad (56)$$

By Donsker-Varadhan's variational formula we have

$$\begin{aligned} & D_{\text{KL}}(P_{Z,W} || P_{Z,W'}) \\ &\geq \mathbb{E}_{Z,W} [m \mathcal{D}_\beta^*(R_m(W, Z), R_{m+u}(W, Z))] - \log \mathbb{E}_{W \otimes Z} [e^{m \mathcal{D}_\beta^*(R_m(W, Z), R_{m+u}(W, Z))}] \\ &\geq m \mathbb{E}_{Z,W} [\mathcal{D}_\beta^*(R_m(W, Z), R_{m+u}(W, Z))] - \log \left(3 \log(m) \sqrt{\frac{mu}{m+u}} \right). \end{aligned} \quad (57)$$

By the Pinsker's inequality and plugging in $\beta = \frac{m}{m+u}$, the expectation term can be lower bounded by

$$\begin{aligned} & \mathbb{E}_{Z,W} [\mathcal{D}_\beta^*(R_m(W, Z), R_{m+u}(W, Z))] \\ &= \mathbb{E}_{Z,W} [\mathcal{D}(R_m(W, Z), R_{m+u}(W, Z))] \\ &\quad + \frac{u}{m} \mathbb{E}_{Z,W} \left[\mathcal{D} \left(\frac{m+u}{u} R_{m+u}(W, Z) - \frac{m}{u} R_m(W, Z), R_{m+u}(W, Z) \right) \right] \\ &\geq 2 \mathbb{E}_{Z,W} [(R_m(W, Z) - R_{m+u}(W, Z))^2] + 2 \frac{m}{u} \mathbb{E}_{Z,W} [(R_m(W, Z) - R_{m+u}(W, Z))^2] \\ &= 2 \frac{m+u}{u} \mathbb{E}_{Z,W} \left[\left(R_m(W, Z) - \frac{m}{m+u} R_m(W, Z) + \frac{u}{m+u} R_u(W, Z) \right)^2 \right] \\ &= \frac{2u}{m+u} \mathbb{E}_{Z,W} [(R_m(W, Z) - R_u(W, Z))^2] \\ &\geq \frac{2u}{m+u} (\mathbb{E}_{Z,W} [R_m(W, Z) - R_u(W, Z)])^2. \end{aligned} \quad (58)$$

Combining Eq. (57) and Eq. (58) we obtain a degenerated bound compared with that provided in Theorem 1 with extra factors $\log \left(3 \log(m) \sqrt{\frac{mu}{m+u}} \right)$:

$$|\mathbb{E}_{Z,W} [R_m(W, Z) - R_u(W, Z)]| \leq \sqrt{\frac{1}{2} \left(\frac{1}{m} + \frac{1}{u} \right) \left[I(Z; W) + \log \left(3 \log(m) \sqrt{\frac{mu}{m+u}} \right) \right]}.$$

H Proof of Corollary 2

The proof generally follows the proof of Theorem 2 in [33]. For given posterior distribution Q , we need to properly select the optimal prior distribution P^* such that the KL divergence term $D_{\text{KL}}(Q||P)$ can be minimized. However, this solution is not applicable. The reason is that P will depend on Q and Z , yet we require that P should be chosen before observing Z . Therefore, the most widely adopted method is to construct a predefined set of prior distribution $\mathcal{P} = \{P_j\}_{j \in \mathbb{N}}$, and then establishing a high probability guarantee for each $P \in \mathcal{P}$. After that, we can establish a high probability guarantee for the optimal prior P^* by using union bound inequality.

Formally, denote by $\mathbf{w} \in \mathbb{R}^d$ the parameter return by the learning algorithm and σ a predefined hyper-parameter, we define the posterior distribution as $Q \triangleq \mathcal{N}(\mathbf{w} + \epsilon, \sigma^2 \mathbf{I}_d)$. Let $c = \sigma_Q^2(1 + e^{4n/d})$, the predefined set \mathcal{P} is constructed as

$$\mathcal{P} \triangleq \left\{ \mathcal{N}(\epsilon, \sigma_j^2 \mathbf{I}) | \sigma_j = ce^{(1-j)/d} \right\}.$$

Here c is a constant depends on m, n, d, σ , whose value will be discussed later. For any $P \triangleq \mathcal{N}(\epsilon, \sigma_P^2 \mathbf{I}) \in \mathcal{P}$, by calculating the KL divergence term, we have

$$D_{\text{KL}}(Q||P) = \frac{1}{2} \left[\frac{d\sigma^2 + \|\mathbf{w}\|_2^2}{\sigma_P^2} - d + d \log \left(\frac{\sigma_P^2}{\sigma^2} \right) \right],$$

which implies that

$$\operatorname{argmin}_{\sigma_P} D_{\text{KL}}(Q||P) = \sqrt{\sigma^2 + \|\mathbf{w}\|_2^2/d}.$$

Therefore, we can define the optimal prior distribution as $P^* = \mathcal{N}(\boldsymbol{\epsilon}, \sigma_{j^*}^2 \mathbf{I})$ where

$$j^* = \left\lfloor 1 - d \log \left(\frac{\sigma^2 + \|\mathbf{w}\|_2^2/d}{c} \right) \right\rfloor,$$

which implies that

$$-d \log \left(\frac{\sigma^2 + \|\mathbf{w}\|_2^2/d}{c} \right) \leq j^* \leq 1 - d \log \left(\frac{\sigma^2 + \|\mathbf{w}\|_2^2/d}{c} \right), \quad (59)$$

and

$$\sigma^2 + \frac{\|\mathbf{w}\|_2^2}{d} \leq \sigma_{j^*}^2 \leq e^{1/d} \left(\sigma^2 + \frac{\|\mathbf{w}\|_2^2}{d} \right). \quad (60)$$

Here we have used a fact that $\sigma^2 + \|\mathbf{w}\|_2^2/d < c$, which will be shown later. Therefore, we have

$$\begin{aligned} D_{\text{KL}}(Q||P^*) &= \frac{1}{2} \left[\frac{d\sigma_Q^2 + \|\mathbf{w}\|_2^2}{\sigma_{j^*}^2} - d + d \log \left(\frac{\sigma_{j^*}^2}{\sigma_Q^2} \right) \right] \\ &\leq \frac{1}{2} \left[\frac{d(\sigma_Q^2 + \|\mathbf{w}\|_2^2/d)}{\sigma_Q^2 + \|\mathbf{w}\|_2^2/d} - d + d \log \left(\frac{e^{1/d} (\sigma_Q^2 + \|\mathbf{w}\|_2^2/d)}{\sigma^2} \right) \right] \\ &= \frac{1}{2} \left[d \log \left(\frac{e^{\frac{1}{d}} (\sigma^2 + \|\mathbf{w}\|_2^2/d)}{\sigma^2} \right) \right] = \frac{1}{2} \left[1 + d \log \left(1 + \frac{\|\mathbf{w}\|_2^2}{d\sigma^2} \right) \right]. \end{aligned} \quad (61)$$

Denote by A_j the event that

$$A_j \triangleq \left\{ |\mathbb{E}_{W \sim Q} [R_u(W, Z) - R_m(W, Z)]| \geq \sqrt{\frac{C_{m,u}}{2} \left(\frac{1}{m} + \frac{1}{u} \right) \left(D_{\text{KL}}(Q||P_j) + \log \left(\frac{1}{\delta_j} \right) \right)} \right\}.$$

Let $\delta_j \triangleq \frac{6\delta}{\pi^2 j^2}$, by Theorem 4, for any distribution Q we have

$$\mathbb{P}\{A_{j^*}\} \leq \mathbb{P}\{\cup_{j=1}^{\infty} A_j\} \leq \sum_{j=1}^{\infty} \mathbb{P}\{A_j\} = \sum_{j=1}^{\infty} \delta_j = \sum_{j=1}^{\infty} \frac{6\delta}{\pi^2 j^2} = \delta. \quad (62)$$

Therefore, with probability at least $1 - \delta$,

$$\begin{aligned} &\mathbb{E}_{W \sim Q} [R_u(W, Z) - R_m(W, Z)] \\ &= \mathbb{E}_{\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})} [R_u(\mathbf{w} + \boldsymbol{\epsilon}, Z) - R_m(\mathbf{w} + \boldsymbol{\epsilon}, Z)] \\ &\leq \sqrt{\frac{C_{m,u}}{2} \left(\frac{1}{m} + \frac{1}{u} \right) \left(D_{\text{KL}}(Q||P_{j^*}) + \log \left(\frac{1}{\delta_j} \right) \right)} \\ &\leq \sqrt{\frac{C_{m,u}}{2} \left(\frac{1}{m} + \frac{1}{u} \right) \left(\frac{1}{2} \left[1 + d \log \left(1 + \frac{\|\mathbf{w}\|_2^2}{d\sigma^2} \right) \right] + \log \left(\frac{1}{6\delta} \right) + 2 \log(\pi j^*) \right)}. \end{aligned} \quad (63)$$

Since $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$, Lemma 1 in [60] suggests that

$$\mathbb{P}\{\|\boldsymbol{\epsilon}\|_2^2 \geq d\sigma^2 + 2\sigma^2 \sqrt{dt} + 2t\sigma^2\} \leq e^{-t}. \quad (64)$$

Let $\tilde{C}_{m,u} \triangleq \sqrt{\log(4mu/(m+u))/d}$, with probability at least $1 - \sqrt{(m+u)/4mu}$ we have

$$\begin{aligned} \|\boldsymbol{\epsilon}\|_2^2 &\leq d\sigma^2 + 2\sigma^2 \sqrt{d \log(\sqrt{4mu/(m+u)})} + \sigma^2 \log(4mu/(m+u)) \\ &\leq d\sigma^2 + 2\sigma^2 \sqrt{d \log(4mu/(m+u))} + \sigma^2 \log(4mu/(m+u)) \\ &= \sigma^2 d \left(1 + \sqrt{\frac{\log(4mu/(m+u))}{d}} \right)^2 = \sigma^2 d \left(1 + \tilde{C}_{m,u} \right)^2 \triangleq \rho^2, \end{aligned} \quad (65)$$

Denote by $A = \{\|\epsilon\|_2^2 \leq \rho^2\}$ the event that the Euclidean norm of ϵ is not larger than ρ , with probability at least $1 - \delta$:

$$\begin{aligned}
& R_u(\mathbf{w}, Z) \\
& \leq \mathbb{E}_{\epsilon \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})} [R_u(\mathbf{w} + \epsilon, Z)] \\
& = \mathbb{P}\{A\} \mathbb{E}_{\epsilon \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})} [R_u(\mathbf{w} + \epsilon, Z)|A] + \mathbb{P}\{\bar{A}\} \mathbb{E}_{\epsilon \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})} [R_u(\mathbf{w} + \epsilon, Z)|\bar{A}] \\
& = \left(1 - \sqrt{\frac{m+u}{4mu}}\right) \mathbb{E}_{\epsilon \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})} [R_u(\mathbf{w} + \epsilon, Z)|A] + \sqrt{\frac{m+u}{4mu}} \mathbb{E}_{\epsilon \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})} [R_u(\mathbf{w} + \epsilon, Z)|\bar{A}] \\
& \leq \mathbb{E}_{\epsilon \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})} [R_u(\mathbf{w} + \epsilon, Z)|A] + \sqrt{\frac{m+u}{4mu}} \\
& \leq \mathbb{E}_{W \sim Q} R_m(W, Z) + \sqrt{\frac{m+u}{4mu}} \\
& \quad + \sqrt{\frac{C_{m,u}}{2} \left(\frac{1}{m} + \frac{1}{u}\right) \left(\frac{1}{2} \left[1 + d \log \left(1 + \frac{\|\mathbf{w}\|_2^2}{d\sigma^2}\right)\right] + \log \left(\frac{1}{6\delta}\right) + 2 \log(\pi j^*)\right)} \\
& \leq \max_{\|\epsilon\|_2 \leq \rho} R_m(\mathbf{w} + \epsilon, Z) \\
& \quad + \sqrt{C_{m,u} \left(\frac{1}{m} + \frac{1}{u}\right) \left(1 + \frac{d}{2} \log \left(1 + \frac{\|\mathbf{w}\|_2^2}{\rho^2} (1 + \tilde{C}_{m,u})^2\right) + \log \left(\frac{1}{6\delta}\right) + 2 \log(\pi j^*)\right)}.
\end{aligned} \tag{66}$$

Here we use the assumption to obtain the first inequality, and the second line is due to law of total expectation. The second inequality is due to the fact that $R_u(\mathbf{w}, Z) \leq 1$ for any \mathbf{w} and Z since the loss is 0-1 loss. We use the fact $\sqrt{a} + \sqrt{a+b} \leq \sqrt{2(2a+b)}$ in the last inequality. The remaining step is to specify the value of c . Note that if

$$\|\mathbf{w}\|_2^2 \geq \frac{\rho^2}{(1 + \tilde{C}_{m,u})^2} \left(\exp \left\{ \frac{2mu}{(m+u)d} \right\} - 1 \right), \tag{67}$$

the slack term in Eq. (66) will exceed 1 and the inequality holds trivially. Therefore, we only need to consider the case that

$$\|\mathbf{w}\|_2^2 < \frac{\rho^2}{(1 + \tilde{C}_{m,u})^2} \left(\exp \left\{ \frac{2mu}{(m+u)d} \right\} - 1 \right). \tag{68}$$

which implies that

$$\sigma^2 + \frac{\|\mathbf{w}\|_2^2}{d} < \frac{\rho^2}{(1 + \tilde{C}_{m,u})^2 d} \exp \left\{ \frac{2mu}{(m+u)d} \right\} = \sigma^2 \exp \left\{ \frac{2mu}{(m+u)d} \right\} \triangleq c. \tag{69}$$

Here we have used Eq. (65) since we only need to consider the case that $\|\epsilon\|_2^2 \leq \rho^2$. One can verify that j^* is an valid integer under this definition. Note that

$$\begin{aligned}
\log(j^*) & \leq \log \left(1 + d \log \left(\frac{c}{\sigma^2 + \|\mathbf{w}\|_2^2/d} \right) \right) \\
& \leq \log \left(1 + d \log \left(\frac{c}{\sigma^2} \right) \right) = \log \left(1 + \frac{2mu}{(m+u)} \right) \\
& \leq \log \left(\frac{4mu}{(m+u)} \right).
\end{aligned} \tag{70}$$

Plugging Eq. (70) into Eq. (66), with probability at least $1 - \delta$ over the randomness of Z :

$$\begin{aligned}
& R_u(\mathbf{w}, Z) \\
& \leq \max_{\|\epsilon\|_2 \leq \rho} R_m(\mathbf{w} + \epsilon, Z) \\
& \quad + \sqrt{\frac{C_{m,u} (m+u) \left(2 + d \log \left(1 + \frac{\|\mathbf{w}\|_2^2}{\rho^2} (1 + \tilde{C}_{m,u})^2 \right) + 2 \log \left(\frac{1}{6\delta} \right) + 4 \log \left(\frac{4\pi mu}{m+u} \right) \right)}{2mu}}.
\end{aligned} \tag{71}$$

I Proof of Theorem 5

This proof is inspired by [74, 113]. Since $\mathbb{E}[R_u(w_T + U_T, Z) - R_u(w_T, Z)] \geq 0$ holds for any realization of $W_T = w_T$ we have

$$\begin{aligned} & \mathbb{E}_{Z, W_T, U_T} [R_u(\widetilde{W}_T, Z) - R_u(W_T, Z)] \\ &= \int_{w_T} [\mathbb{E}_{Z, U_T} [R_u(w_T + U_T, Z) - R_u(w_T, Z) | W_T = w_T]] dP_{W_T}(w_T) \geq 0. \end{aligned} \quad (72)$$

Therefore, the transductive generalization error can be bounded by

$$\begin{aligned} & \mathbb{E}_{Z, W_T} [R_u(W_T, Z) - R_m(W_T, Z)] \\ &= \mathbb{E}_{Z, W_T, U_T} [R_m(\widetilde{W}_T, Z) - R_m(W_T, Z)] - \mathbb{E}_{Z, W_T, U_T} [R_u(\widetilde{W}_T, Z) - R_u(W_T, Z)] \\ & \quad + \mathbb{E}_{Z, W_T, U_T} [R_u(\widetilde{W}_T, Z) - R_m(\widetilde{W}_T, Z)] \\ & \leq \mathbb{E}_{Z, W_T, U_T} [R_m(\widetilde{W}_T, Z) - R_m(W_T, Z)] + \sqrt{\frac{C_{m,u}}{2} \left(\frac{1}{m} + \frac{1}{u} \right) I(Z; \widetilde{W}_T)}. \end{aligned} \quad (73)$$

Now the last step is to provide an upper bound for $I(Z; \widetilde{W}_T)$. Following [113], the mutual information term is decomposed by

$$\begin{aligned} & I(Z; \widetilde{W}_T) \\ &= I\left(Z; \widetilde{W}_{T-1} - \frac{\eta}{\sqrt{v_T} + \epsilon} \odot g(W_{T-1}, B_T(Z)) + N_T\right) \\ & \leq I\left(Z; \widetilde{W}_{T-1}, -\frac{\eta}{\sqrt{v_T} + \epsilon} \odot g(W_{T-1}, B_T(Z)) + N_T\right) \\ &= I(Z; \widetilde{W}_{T-1}) + I\left(-\frac{\eta}{\sqrt{v_T} + \epsilon} \odot g(W_{T-1}, B_T(Z)) + N_T; Z \middle| \widetilde{W}_{T-1}\right) \end{aligned}$$

Recursively repeating the above process, we obtain

$$\begin{aligned} I(Z; \widetilde{W}_T) & \leq \sum_{t=1}^T I\left(-\frac{\eta}{\sqrt{v_t} + \epsilon} \odot g(W_{t-1}, B_t(Z)) + N_t; Z \middle| \widetilde{W}_{t-1}\right) \\ &= \sum_{t=1}^T I\left(-\frac{\eta}{\sqrt{v_t} + \epsilon} \odot g(\widetilde{W}_{t-1} - U_{t-1}, B_t(Z)) + N_t; Z \middle| \widetilde{W}_{t-1}\right) \end{aligned} \quad (74)$$

Then we need to provide an upper bound for the conditional mutual information. Let V, X, U be random variables that are independent of $N \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$. Define Ψ as a function of random variables U, V, X, Y . Denote by $h(\cdot)$ the differential entropy, then

$$\begin{aligned} & I(\Psi(V, y - U, X) + \sigma N; X | Y = y) \\ &= h(\Psi(V, y - U, X) + \sigma N | Y = y) - h(\Psi(V, y - U, X) + \sigma N | X, Y = y) \end{aligned} \quad (75)$$

For the first term, using the fact that Gaussian minimizes entropy, we have

$$\begin{aligned} & h(\Psi(V, y - U, X) + \sigma N | Y = y) \\ & \leq \frac{d}{2} \log \left(2\pi e \frac{\mathbb{E}[\|\Psi(V, X, y - U) + \sigma N\|_2^2 | Y = y]}{d} \right) \\ &= \frac{d}{2} \log \left(2\pi e \frac{\mathbb{E}[\|\Psi(U, V, X, y - U)\|_2^2 | Y = y] + \sigma^2 \mathbb{E}[\|N\|_2^2]}{d} \right) \\ &= \frac{d}{2} \log \left(2\pi e \frac{\mathbb{E}[\|\Psi(V, X, y - U)\|^2 | Y = y] + d\sigma^2}{d} \right). \end{aligned} \quad (76)$$

For the second term, we have:

$$\begin{aligned} & h(\Psi(V, y - U, X) + \sigma N | X, Y = y) \geq h(\Psi(V, y - U, X) + \sigma N | U, V, X, Y = y) \\ &= h(\sigma N) = \frac{d}{2} \log 2\pi e \sigma^2. \end{aligned} \quad (77)$$

Let $V = W^{[t-2]} \triangleq (W_0, \dots, W_{t-2})$, $X = Z$, $Y = \widetilde{W}_{t-1}$, $U = U_{t-1}$ and

$$\begin{aligned}\Psi(V, y - U, X) &= \Psi(W^{[t-2]}, \widetilde{w}_{t-1} - U_{t-1}, Z) \\ &= -\frac{\eta}{\sqrt{v_t(W^{[t-2]}, \widetilde{w}_{t-1} - U_{t-1})} + \epsilon} \odot g(\widetilde{w}_{t-1} - U_{t-1}, B_t(Z)),\end{aligned}$$

plugging Eqs. (76,77) into Eq. (75), we have

$$\begin{aligned}& I\left(\Psi(W^{[t-2]}, \widetilde{w}_{t-1} - U_{t-1}, Z) + N_t; Z \middle| \widetilde{W}_{T-1} = \widetilde{w}_{t-1}\right) \\ &= I\left(\Psi(W^{[t-2]}, \widetilde{w}_{t-1} - U_{t-1}, Z) + N_t; Z \middle| \widetilde{W}_{t-1} = \widetilde{w}_{t-1}\right) \\ &\leq \frac{d}{2} \log \left(\frac{1}{d\sigma_t^2} \mathbb{E} \left[\left\| \Psi(W^{[t-2]}, \widetilde{w}_{t-1} - U_{t-1}, Z) \right\|_2^2 \middle| \widetilde{W}_{t-1} = \widetilde{w}_{t-1} \right] + 1 \right),\end{aligned}$$

which implies that

$$\begin{aligned}& I\left(\Psi(W^{[t-2]}, \widetilde{w}_{t-1} - U_{t-1}, Z) + N_t; Z \middle| \widetilde{W}_{t-1}\right) \\ &\leq \int_{\widetilde{w}_{t-1}} \frac{d}{2} \log \left(\frac{1}{d\sigma_t^2} \mathbb{E} \left[\left\| \Psi(W^{[t-2]}, \widetilde{w}_{t-1} - U_{t-1}, Z) \right\|_2^2 \middle| \widetilde{W}_{t-1} = \widetilde{w}_{t-1} \right] + 1 \right) dP_{\widetilde{W}_{T-1}}(\widetilde{w}_{T-1}) \quad (78) \\ &\leq \frac{d}{2} \log \left(\frac{1}{d\sigma_t^2} \mathbb{E} \left[\left\| \Psi(W^{[t-2]}, \widetilde{W}_{t-1} - U_{t-1}, Z) \right\|_2^2 \right] + 1 \right).\end{aligned}$$

Let $w^{[k]} \triangleq (w_1, \dots, w_k)$ and

$$\zeta(W^{[t-2]}, \widetilde{W}_{t-1} - U_{t-1}, Z) \triangleq \left\| \Psi(W^{[t-2]}, \widetilde{W}_{t-1} - U_{t-1}, Z) + \mathbb{E} \left[g(\widetilde{W}_{t-1} - U_{t-1}, B_t(Z)) \right] \right\|_2^2,$$

we have

$$\begin{aligned}& \mathbb{E}_{W^{[t-2]}, \widetilde{W}_{t-1}, U_{t-1}, Z} \left[\zeta(W^{[t-2]}, \widetilde{W}_{t-1} - U_{t-1}, Z) \right] \\ &= \int_z \int_{w^{[t-2]}} \int_{\widetilde{w}_{t-1}} \int_u \zeta(w^{[t-2]}, \widetilde{w}_{t-1} - u, z) dP_{W_{t-1}|Z, W_{t-2}}(\widetilde{w}_{t-1} - u) dP_{U_{t-1}}(u) dP_{W^{[t-2]}|Z}(w^{[t-2]}) dP_Z(z) \\ &= \int_z \int_{w^{[t-2]}} \int_{\widetilde{w}_{t-1}} \int_{w_{t-1}} \zeta(w^{[t-2]}, w_{t-1}, z) dP_{W_{t-1}|Z, W_{t-2}}(w_{t-1}) dP_{U_{t-1}}(\widetilde{w}_{t-1} - w_{t-1}) dP_{W^{[t-2]}|Z}(w^{[t-2]}) dP_Z(z) \\ &= \int_z \int_{w^{[t-2]}} \int_{w_{t-1}} \zeta(w^{[t-2]}, w_{t-1}, z) dP_{W_{t-1}|Z, W_{t-2}}(w_{t-1}) dP_{W^{[t-2]}|Z}(w^{[t-2]}) dP_Z(z) \\ &= \int_z \int_{w^{[t-1]}} \zeta(w^{[t-2]}, w_{t-1}, z) dP_{W^{[t-1]}, Z}(w^{[t-1]}, z) \\ &= \mathbb{E}_{W^{[t-1]}, Z} [\zeta(W^{[t-1]}, Z)].\end{aligned} \quad (79)$$

Here the first inequality is due to the convolution formulation, and we use $w_{t-1} \triangleq \widetilde{w}_{t-1} - u$ to obtain the second inequality. Plugging Eqs. (78, 79) into Eq. (74), we have

$$\begin{aligned}& I(Z; \widetilde{W}_T) \\ &\leq \sum_{t=1}^T \frac{d}{2} \log \left(\frac{1}{d\sigma_t^2} \mathbb{E} \left[\left\| \Psi(W^{[t-2]}, W_{t-1}, Z) + \mathbb{E} [g(W_{t-1}, B_t(Z))] \right\|_2^2 \right] + 1 \right) \quad (80) \\ &= \sum_{t=1}^T \frac{d}{2} \log \left(\frac{1}{d\sigma_t^2} \mathbb{E} \left[\left\| \frac{\eta}{\sqrt{v_t(W^{[t-1]})} + \epsilon} \odot g(W_{t-1}, B_t(Z)) - \mathbb{E} [g(W_{t-1}, B_t(Z))] \right\|_2^2 \right] + 1 \right).\end{aligned}$$

Combining Eq. (80) with Eq. (73), we have

$$\begin{aligned}& \mathbb{E}_{Z, W_T} [R_u(W_T, Z) - R_m(W_T, Z)] \\ &\leq \sum_{t=1}^T \frac{d}{2} \log \left(\frac{1}{d\sigma_t^2} \mathbb{E} \left[\left\| \frac{\eta}{\sqrt{v_t(W^{[t-1]})} + \epsilon} \odot g(W_{t-1}, B_t(Z)) - \mathbb{E} [g(W_{t-1}, B_t(Z))] \right\|_2^2 \right] + 1 \right) \quad (81) \\ &\quad + \mathbb{E}_{Z, W_T, U_T} [R_m(W_T + U_T, Z) - R_m(W_T, Z)]\end{aligned}$$

Now we discuss how to extend this result to Adam optimization algorithm. For $t \in [T]$, the update rule of Adam is

$$\begin{aligned}m_t &= \beta_1 m_{t-1} + (1 - \beta_1) g(W_{t-1}, B_t(Z)), \\ v_t &= \beta_2 v_{t-1} + (1 - \beta_2) g(W_{t-1}, Z) \odot g(W_{t-1}, B_t(Z)), \\ \hat{v}_t &= \frac{1}{1 - \beta_2^t} v_t, W_t = W_{t-1} - \frac{\eta m_t}{\sqrt{\hat{v}_t} + \epsilon}.\end{aligned}$$

Note that

$$\begin{aligned} W_t &= W_{t-1} - \frac{\eta m_t}{\sqrt{\hat{v}_t} + \epsilon} = W_{t-1} - \sum_{\tau=0}^{t-1} \beta_1^{t-\tau-1} \frac{(1-\beta_1)\eta}{\sqrt{\hat{v}_t} + \epsilon} \odot g(W_\tau, B_{\tau+1}(Z)) \\ &\triangleq W_{t-1} - \Psi(W^{[t-1]}, Z). \end{aligned}$$

Similarly, we construct the weight process as

$$\widetilde{W}_0 = W_0, \widetilde{W}_t = \widetilde{W}_{t-1} - \Psi(W^{[t-1]}, Z) + N_t.$$

By the same technique, one can find that

$$\begin{aligned} &I(Z; \widetilde{W}_T) \\ &\leq \sum_{t=1}^T \frac{d}{2} \log \left(\frac{1}{d\sigma_t^2} \mathbb{E} \left[\left\| \Psi(W^{[t-2]}, W_{t-1}, Z) + \mathbb{E}[g(W_{t-1}, B_t(Z))] \right\|_2^2 \right] + 1 \right) \\ &= \sum_{t=1}^T \frac{d}{2} \log \left(\frac{1}{d\sigma_t^2} \mathbb{E} \left[\left\| \sum_{\tau=0}^{t-1} \beta_1^{t-\tau-1} \frac{(1-\beta_1)\eta}{\sqrt{\hat{v}_t} + \epsilon} \odot g(W_\tau, B_{\tau+1}(Z)) \right\|_2^2 \right] + 1 \right). \end{aligned} \quad (82)$$

Then the upper bound is given by

$$\begin{aligned} &\mathbb{E}_{Z, W_T} [R_u(W_T, Z) - R_m(W_T, Z)] \\ &\leq \sum_{t=1}^T \frac{d}{2} \log \left(\frac{1}{d\sigma_t^2} \mathbb{E} \left[\left\| \sum_{\tau=0}^{t-1} \beta_1^{t-\tau-1} \frac{(1-\beta_1)\eta}{\sqrt{\hat{v}_t} + \epsilon} \odot g(W_\tau, B_{\tau+1}(Z)) \right\|_2^2 \right] + 1 \right) \\ &\quad + \mathbb{E}_{Z, W_T, U_T} [R_m(W_T + U_T, Z) - R_m(W_T, Z)]. \end{aligned} \quad (83)$$

J Proof of Proposition 2

First we have

$$\begin{aligned} &\mathbb{E}_{D_m, D_u^X, W} \left[\frac{1}{u} \sum_{i=m+1}^{m+u} \ell(W, (X_i, Y_i)) - \frac{1}{m} \sum_{i=1}^m \ell(W, (X_i, Y_i)) \right] \\ &= \mathbb{E}_{D_{m+u}} \mathbb{E}_{Z, W|D_{m+u}} [R_u(W, Z) - R_m(W, Z)]. \end{aligned} \quad (84)$$

By the proof of Theorem 1, the following inequality holds for any realization of $W = w$:

$$\mathbb{E}_{Z|D_{m+u}=d_{m+u}} [\exp \{ \lambda (R_u(w, Z) - R_m(w, Z)) \}] \leq \exp \left\{ \frac{\lambda^2 (m+u) C_{m,u}}{8mu} \right\}. \quad (85)$$

Let Z' be the independent copy of Z . By Donsker-Varadhan's variational formula, for any $\lambda \in \mathbb{R}$:

$$\begin{aligned} &D_{\text{KL}}(P_{Z, W|D_{m+u}=d_{m+u}} \| P_{Z', W|D_{m+u}=d_{m+u}}) \\ &\geq \mathbb{E}_{Z, W|D_{m+u}=d_{m+u}} [\lambda (R_u(W, Z) - R_m(W, Z))] - \log \mathbb{E}_{Z', W|D} [\exp \{ \lambda (R_u(W, Z') - R_m(W, Z')) \}] \\ &\geq \mathbb{E}_{Z, W|D=d} [\lambda (R_u(W, Z) - R_m(W, Z))] - \frac{\lambda^2 C_{m,u}}{8} \left(\frac{1}{m} + \frac{1}{u} \right), \end{aligned} \quad (86)$$

which implies that

$$|\mathbb{E}_{Z, W|D=d} [\lambda (R_u(W, Z) - R_m(W, Z))]| \leq \sqrt{\frac{C_{m,u}}{2} \left(\frac{1}{m} + \frac{1}{u} \right) I(Z, W|D=d)}. \quad (87)$$

Taking expectation on both side and plugging into Eq. (84) we have

$$\begin{aligned} &\left| \mathbb{E}_{S_m, S_u^X, W} \left[\frac{1}{u} \sum_{i=m+1}^{m+u} \ell(W, (X_i, Y_i)) - \frac{1}{m} \sum_{i=1}^m \ell(W, (X_i, Y_i)) \right] \right| \\ &= |\mathbb{E}_{S_{m+u}} \mathbb{E}_{Z, W|S_{m+u}} [R_u(W, Z) - R_m(W, Z)]| \\ &\leq \mathbb{E}_{S_{m+u}} |\mathbb{E}_{Z, W|S_{m+u}} [\lambda (R_u(W, Z) - R_m(W, Z))]| \\ &\leq \mathbb{E}_{S_{m+u}} \sqrt{\frac{C_{m,u}}{2} \left(\frac{1}{m} + \frac{1}{u} \right) I(Z, W|S_{m+u})}. \end{aligned} \quad (88)$$

This finishes the first part. For the second part, define $R_\mu(W) \triangleq \mathbb{E}_{(X,Y) \sim \mu} \ell(W, (X, Y))$, then we have

$$\begin{aligned}
& \mathbb{E}_{S_m, S_u^X, W} R(W, S_{m+u}) = \mathbb{E}_{S_{m+u}, Z, W} R(W, S_{m+u}) \\
&= \mathbb{E}_{S_{m+u}} \mathbb{E}_{Z, W | S_{m+u}} \left[\frac{m}{m+u} R_m(W, Z) + \frac{u}{m+u} R_u(W, Z) \right] \\
&\leq \mathbb{E}_{S_{m+u}} \left[\mathbb{E}_{Z, W | S_{m+u}} [R_m(W, Z)] + \frac{u}{m+u} \sqrt{\frac{C_{m,u}}{2} \left(\frac{1}{m} + \frac{1}{u} \right) I(Z, W | S_{m+u})} \right] \\
&= \frac{1}{m} \sum_{i=1}^m \ell(W, (X_i, Y_i)) + \mathbb{E}_{S_{m+u}} \sqrt{\frac{C_{m,u} I(Z, W | S_{m+u})}{2m(m/u + 1)}}.
\end{aligned} \tag{89}$$

By Donsker-Varadhan's variational formula,

$$\begin{aligned}
& \text{D}_{\text{KL}}(P, ||P_{Z', W | D=d}) \\
&\geq \mathbb{E}_{S_{m+u}, Z, W} [\lambda(R_\mu(W) - R(W, S_{m+u}))] - \log \mathbb{E}_{S'_{m+u}, Z', W} [\exp \{\lambda(R_\mu(W) - R(W, S_{m+u}))\}] \\
&= \mathbb{E}_{S_m, S_u^X, W} [\lambda(R_\mu(W) - R(W, S_{m+u}))] - \log \mathbb{E}_{S'_{m+u}, Z', W} [\exp \{\lambda(R_\mu(W) - R(W, S_{m+u}))\}] \\
&\geq \mathbb{E}_{S_m, S_u^X, W} [\lambda(R_\mu(W) - R(W, S_{m+u}))] - \frac{\lambda^2 B^2}{2(m+u)}.
\end{aligned} \tag{90}$$

Therefore, we have

$$\begin{aligned}
& \mathbb{E}_{S_m, S_u^X, W} R_\mu(W) \\
&\leq \mathbb{E}_{S_m, S_u^X, W} R(W, S_{m+u}) + \sqrt{\frac{2I(S_{m+u}, Z; W)}{m+u}} \\
&\leq \frac{1}{m} \sum_{i=1}^m \ell(W, (X_i, Y_i)) + \mathbb{E}_{S_{m+u}} \sqrt{\frac{C_{m,u} I(Z, W | S_{m+u})}{2m(m/u + 1)}} + \sqrt{\frac{2I(S_{m+u}, Z; W)}{m+u}}.
\end{aligned} \tag{91}$$

Also, we have

$$\begin{aligned}
& \left| \mathbb{E}_{S_m, S_u^X, W} \left[\frac{1}{u} \sum_{i=m+1}^{m+u} \ell(W, (X_i, Y_i)) - \frac{1}{m} \sum_{i=1}^m \ell(W, (X_i, Y_i)) \right] \right| \\
&= \left| \mathbb{E}_{S_{m+u}} \mathbb{E}_{Z, W | S_{m+u}} [R_u(W, Z) - R_m(W, Z)] \right| \\
&= \left| \mathbb{E}_{S_{m+u}} \mathbb{E}_{\tilde{Z}, U, W | S_{m+u}} [R_u(W, Z) - R_m(W, Z)] \right| \\
&\leq \frac{1}{m} \sum_{i=1}^m \mathbb{E}_{S_{m+u}} \mathbb{E}_{\tilde{Z} | S_{m+u}} \left| \mathbb{E}_{U, W | \tilde{Z}, S_{m+u}} [\ell(W, \tilde{Z}_{i, S_i}) - \ell(W, \tilde{Z}_{i, 1-S_i})] \right| \\
&= \frac{1}{m} \sum_{i=1}^m \mathbb{E}_{S_{m+u}} \mathbb{E}_{\tilde{Z} | S_{m+u}} \left| \mathbb{E}_{F_i, U_i | \tilde{Z}, S_{m+u}} [g(F_i, S_i, \tilde{Y}_i)] \right|.
\end{aligned} \tag{92}$$