
Can Variance-Based Regularization Improve Domain Generalization?

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Without prior information, domain generalization with only access to multi-domain
2 training data relies on guessing what the test data is. In this work, we consider mild
3 assumptions that there is a distribution over domains and the out-of-distribution
4 data is generated by the shift of the domain distribution. We study a domain-level
5 variance-based regularizer. We show that the variance-regularized method locally
6 approximates the group distributionally robust optimization and embeds the local
7 information into the objective function as a weighting scheme. By taking the
8 empirical domain distribution as an anchor of the location, we propose a weighting
9 correction scheme and provide guarantees of in-distribution generalization. Com-
10 pared to the Empirical Risk Minimization, we prove the potential benefits of our
11 proposed method but do not observe consistent improvements in general.

12 1 Introduction

13 Domain generalization [12, 28] is an out-of-distribution (OOD) generalization problem and has drawn
14 much attention recently [39, 44, 35]. Some recent works consider an ambitious goal that generalizes
15 to "absolutely" unseen domain by learning domain-invariant features. From the perspective of theory,
16 the price of such invariant learning methods is the requirement for harsh assumptions or strong prior
17 information, which is necessary to guarantee that the invariance exists and is identifiable. In this
18 work, we assume that there exists a distribution of domains and the OOD test data is generated by
19 the shift of the domain distribution. Then domain generalization is formulated into a distributionally
20 robust optimization problem (DRO, [9, 11, 10]).

21 Let $\mathbf{z} = (\mathbf{x}, \mathbf{y})$ be a data point consisting of an input vector $\mathbf{x} \in \mathcal{X}$ and the target label $\mathbf{y} \in \mathcal{Y}$.
22 Suppose the training data is structured with respect to a latent domain label:

$$\mathcal{D}_{tr} = \{\mathbf{z}_l, 1 \leq l \leq m\} = \{\{\mathbf{z}_{i,j}, 1 \leq j \leq m_i\}, 1 \leq i \leq n\}, \quad (1)$$

23 where m is the total sample size, m_i is the sample size of the i -th domain and n is the number of
24 domains. We assume that the training domains are randomly drawn from possible domains with
25 a domain distribution Q , i.e. $\mathcal{E}_{tr} = \{e_1, e_2, \dots, e_n\} \subseteq \mathcal{E}$ with $e_i \sim Q$ and the data points under
26 domain e is sampled from the distribution P_e . Let \mathcal{H} be the hypothetical space and $h \in \mathcal{H}$ be a
27 model that maps $\mathbf{x} \in \mathcal{X}$ to $h(\mathbf{x}) \in \mathcal{Y}$. The loss function $\ell(\hat{\mathbf{y}}, \mathbf{y}) : \mathcal{Y} \times \mathcal{Y} \rightarrow [0, M]$ measures
28 how poorly the output $\hat{\mathbf{y}} = h(\mathbf{x})$ predicts the target \mathbf{y} . Denote \mathcal{F} as the collection of the functions
29 $f = \ell(h(\cdot), \cdot) : \mathcal{Z} \rightarrow [0, M]$ with $h \in \mathcal{H}$. The in-domain expected risk and its sample average
30 approximation ([34]) are denoted by

$$R(f|e_i) = \mathbb{E}_{\mathbf{z} \sim P_{e_i}}[f(\mathbf{z})] \quad \text{and} \quad \hat{R}(f|e_i) = \frac{1}{m_i} \sum_{j=1}^{m_i} f(\mathbf{z}_{i,j}) \quad (2)$$

31 respectively. The distribution shift between training and test data is characterized by the change of Q ,
32 while the data distributions $P_e, e \in \mathcal{E}$ are fixed.

33 We study the group distributionally robust optimization problem (group DRO, [21, 30, 33]):

$$\min_{f \in \mathcal{F}} \max_Q \mathbb{E}_{\mathbf{z} \sim P} [f(\mathbf{z})], \quad s.t. \quad P = \int P_e Q(de), \quad D_\phi(Q \| Q_0) \leq \rho, \quad (3)$$

34 where Q_0 is a selected domain distribution, $D_\phi(\cdot \| \cdot)$ stands for the ϕ -divergence ([3, 14]) and the
 35 tuning parameter ρ modulates the distribution shift. Throughout this paper, $D_\phi(\cdot \| \cdot)$ is the χ^2 -
 36 divergence, i.e., $\phi(t) = \frac{1}{2}(t - 1)^2$. Sagawa* et al. [33] consider the empirical optimization problem,

$$\min_{f \in \mathcal{F}} \max_{\mathbf{q} \in \Delta_n} \sum_{i=1}^n q_i \hat{R}(f | e_i) \quad \text{with} \quad \Delta_n = \left\{ (q_1, \dots, q_n) : q_i \geq 0, \sum_{i=1}^n q_i = 1 \right\}.$$

37 Here Δ_n is the $(n - 1)$ -dimensional probability simplex. In this case, the parameter ρ is fixed
 38 and sufficiently large. For more ambitious goals, Krueger et al. [25] propose the minimax risk
 39 extrapolation (MM-REx) that extends the uncertainty region Δ_n into

$$\tilde{\Delta}_n(\alpha) = \left\{ \mathbf{q} = (q_1, \dots, q_n) : q_i \geq \alpha, \sum_{i=1}^n q_i = 1 \right\},$$

40 where the parameter $\alpha \in (-\infty, 1/n]$ modulates the uncertainty region. The negative value of α
 41 extrapolates risks and encourages robustness to large distribution shifts.

42 At the sample level, the DRO loss can be asymptotically approximated by the sum of the ERM loss
 43 [38] and a variance-based regularizer [17], where the negligible error term converges to zero almost
 44 surely. Section 7 in [17] gives general results when the DRO objective is a Hadamard differentiable
 45 functional to P and \mathcal{F} is a P_0 -Donsker class. From the perspective of generalization, the upper
 46 bound of the prediction risk may also have a variance-based regularization term that trades between
 47 approximation error and estimation error [5, 6, 13, 24]. Sample variance penalization [27] replaces
 48 the variance-based regularization with its empirical estimator and gives theoretical guarantees on
 49 the prediction performance. To address the computationally intractable problem caused by the non-
 50 convexity of the regularizer, Namkoong and Duchi [29] and Duchi and Namkoong [16] investigate
 51 the robustly regularized risk, that provides a convex surrogate for variance-regularized loss, and
 52 prove finite-sample and asymptotic results characterizing prediction performance. Back to domain
 53 generalization problem, Krueger et al. [25] develop a variance-regularized empirical loss (V-REx):
 54 $\tilde{R}(f) + \lambda \tilde{V}_{out}(f)$, where

$$\tilde{R}(f) = \frac{1}{n} \sum_{i=1}^n \hat{R}(f | e_i) \quad \text{and} \quad \tilde{V}_{out}(f) = \frac{1}{n} \sum_{i=1}^n \left(\hat{R}(f | e_i) - \tilde{R}(f) \right)^2.$$

55 Xie et al. [41] prove that with high probability, optimizing the regularized loss $\tilde{R}(f) + \lambda \sqrt{\tilde{V}_{out}(f)}$
 56 is equivalent to solve a MM-REx problem.

57 In this work, we refine $\tilde{R}(f)$ and $\tilde{V}_{out}(f)$ based on the intuitive understanding of generalization
 58 and distribution estimation. Recall the problem in (3). In general, Q_0 is the ground-truth domain
 59 distribution and Q belongs to a neighborhood of Q_0 . Therefore, the empirical version of (3) should
 60 replace Q_0 with its empirical approximation over \mathcal{E}_{tr} , i.e.,

$$\hat{\mathbf{q}} = (\hat{q}_1, \hat{q}_2, \dots, \hat{q}_n) = \left(\frac{m_1}{m}, \dots, \frac{m_n}{m} \right).$$

61 However, the existing variance-regularized methods directly replace Q_0 with a discrete uniform
 62 distribution (the center of $\tilde{\Delta}_n(\alpha)$) without considering a consistent and efficient estimator $\hat{\mathbf{q}}$. In the
 63 sample variance penalization, this problem does not exist because the discrete uniform distribution
 64 on sample points (no tie), i.e. the empirical distribution, is a consistent estimator of the ground-truth
 65 data distribution. Consider a new uncertainty region:

$$\mathcal{Q}_{\alpha, \rho}(\hat{\mathbf{q}}) = \tilde{\Delta}_n(\alpha) \cap \{ \mathbf{q} : D_\phi(\mathbf{q} \| \hat{\mathbf{q}}) \leq \rho \}.$$

66 Specifically, any $\mathbf{q} = (q_1, \dots, q_n) \in \mathcal{Q}_{\alpha, \rho}(\hat{\mathbf{q}})$ satisfies

$$q_i \geq \alpha, \quad \sum_{i=1}^n q_i = 1, \quad \sum_{i=1}^n \frac{1}{2} \left(\frac{q_i}{\hat{q}_i} - 1 \right)^2 \hat{q}_i \leq \rho.$$

67 In Section 3.2, we prove that with high probability, the MM-REx problem on $\mathcal{Q}_{\alpha,\rho}(\hat{\mathbf{q}})$ can be
 68 uniformly equivalent to minimize the variance-regularized empirical loss $\hat{R}(f) + \lambda\sqrt{\hat{V}_{out}(f)}$ where

$$\hat{R}(f) = \sum_{i=1}^n \hat{q}_i \hat{R}(f|e_i) \quad \text{and} \quad \hat{V}_{out}(f) = \sum_{i=1}^n \hat{q}_i (\hat{R}(f|e_i) - \hat{R}(f))^2. \quad (4)$$

69 Comparing to $\tilde{R}(f)$ and $\tilde{V}_{out}(f)$, the two terms $\hat{R}(f)$ and $\hat{V}_{out}(f)$ just introduce a weighting scheme
 70 derived from the empirical domain distribution $\hat{\mathbf{q}}$. In addition, the term $\hat{R}(f)$ is the exact ERM loss
 71 [38]. In Section 3.1, we investigate the generalization guarantee of the variance-regularized estimator,
 72

$$\hat{f} = \arg \min_{f \in \mathcal{F}} \hat{R}(f) + \lambda\sqrt{\hat{V}_{out}(f)}, \quad (5)$$

73 via the covering number of the function class \mathcal{F} . Appendix C also provides a version of the general-
 74 ization guarantee with localized Rademacher complexities, which may provide tighter generalization
 75 bounds in some cases. In Section 4, we consider a general uncertainty region $\mathcal{Q}_{\alpha,\rho}(\mathbf{q}_0)$, where the
 76 choice of \mathbf{q}_0 represents a kind of prior knowledge. Similar to the arguments in Section 3, we can also
 77 write \mathbf{q}_0 as a weight assignment and embed it into the variance-regularized loss function. We present
 78 a general form of the proposed method and prove that the optimization equivalence in Section 3.2
 79 still holds when we replace $\mathcal{Q}_{\alpha,\rho}(\hat{\mathbf{q}})$ with $\mathcal{Q}_{\alpha,\rho}(\mathbf{q}_0)$.

80 Our results clearly show that

- 81 • From the perspective of generalization, we propose a weighting correction scheme for
 82 variance-regularized domain generalization methods. The proposed method can outperform
 83 ERM under some cases, which shows the potential competitive edge of the proposed
 84 weighting correction method.
- 85 • We do not observe that our method consistently improves ERM under general cases.
- 86 • The proposed method is robust to the change of the domain distribution Q . From an
 87 optimization perspective, it is equivalent to solve a group DRO problem.

88 2 Preliminaries

89 In this section, we present the rationale for using variance-based regularization to improve the
 90 robustness of generalization. Section 2.1 gives two domain adaptation examples that the test data
 91 is known. We prove that the standard deviation of risk can bound the generalization gap between
 92 training and test data. In Section 2.2, we formulate an invariant learning principle as a hypothesis
 93 testing problem. We point out that penalizing the risk variance can protect the null hypothesis: the
 94 model is invariant across domains.

95 2.1 Risk variance bounds generalization gap

96 We present two simple examples to show that penalizing the standard deviation of risk is a natural
 97 strategy to improve robustness to the domain distribution shift.

98 **Risk Interpolation.** In the first example, we assume the test distribution belongs to the convex
 99 hull of training domains. This is a typical risk interpolation case. Let P^* be the test distribution.
 100 Suppose there exists $\mathbf{q}^* = (q_1^*, \dots, q_n^*) \in \Delta_n$ such that $P^* = \sum_{i=1}^n q_i^* P_{e_i}$, where $P_{e_i}, 1 \leq i \leq n$
 101 are training domains. Then the generalization gap between the training and test data is

$$\text{err}_f = \sum_{i=1}^n q_i^* R(f|e_i) - \sum_{i=1}^n q_i R(f|e_i) = \sum_{i=1}^n (q_i^* - q_i) \left(R(f|e_i) - \sum_{i=1}^n q_i R(f|e_i) \right),$$

102 where $q_i = Q(de_i)/Q(d\mathcal{E}_{tr})$ is the proportion of the training domain e_i in the training data. We
 103 write $\mathbf{q} = (q_1, \dots, q_n)$. By the Cauchy–Schwarz inequality, we have

$$\text{err}_f \leq \sqrt{2D_\phi(\mathbf{q}^* \|\mathbf{q})} \times \sqrt{V_{out}(f)}, \quad (6)$$

104 where $V_{out}(f)$ is the between-domain risk variance over the training domains:

$$V_{out}(f) = \sum_{i=1}^n q_i \left(R(f|e_i) - \sum_{i=1}^n q_i R(f|e_i) \right)^2.$$

105 Notice that $V_{out}(f)$ only depends the training data. Therefore, it is natural to penalize $\sqrt{V_{out}(f)}$ to
 106 obtain a tight upper bound of the test error. The principle here is that if for $\forall e_i \in \mathcal{E}_{tr}$, $R(f|e_i)$ is
 107 a constant that only depends on f , i.e. $V_{out}(f) = 0$, then changes from \mathbf{q} to \mathbf{q}^* cannot cause any
 108 generalization gap.

109 **Sub-population Shift.** Recall that the training data in (1) is structured with respect to a latent domain
 110 label. In this example, the domain label is the class label \mathbf{y} . Therefore, the marginal distribution of
 111 \mathbf{y} is different in the training and test data, and the conditional distribution $P(\mathbf{x}|\mathbf{y})$ is the same. Let
 112 $\mathcal{Y} = \{1, 2, \dots, K\}$. Then the generalization gap between the training and test data is

$$\begin{aligned} \text{err}_f &= \sum_{k=1}^K \mathbb{E}[f(\mathbf{z})|\mathbf{y} = k] \times (P_{e'}(\mathbf{y} = k) - P_e(\mathbf{y} = k)) \\ &\leq \sqrt{2D_\phi(P_{e'}(\mathbf{y})\|P_e(\mathbf{y}))} \times \sqrt{V_{out}(f)}, \end{aligned}$$

113 where

$$V_{out}(f) = \sum_{k=1}^K P_e(\mathbf{y} = k) \left(\mathbb{E}[f(\mathbf{z})|\mathbf{y} = k] - \frac{1}{K} \sum_{k=1}^K \mathbb{E}[f(\mathbf{z})|\mathbf{y} = k] \right)^2$$

114 is the between-class risk variance over the training data. Therefore, the generalization gap is also
 115 bounded above by the between-domain risk variance. If the in-class risks are equal, i.e., $\mathbb{E}[f(\mathbf{z})|\mathbf{y} =$
 116 $k] = \mathbb{E}[f(\mathbf{z})|\mathbf{y} = k']$, $\forall k, k' \in \mathcal{Y}$, then the sub-population shift cannot cause generalization gap.

117 2.2 Penalizing risk variance protects invariant models

118 In this section, we heuristically discuss the relationship between variance-based regularization and
 119 invariant learning. The REX principle [25] presents two training goals: **Reducing training risks** and
 120 **Increasing the similarity of training risks**. Krueger et al. [25] heuristically explain the utility of
 121 V-REx as enforcing the equality of training risks in the limit case $\lambda \rightarrow +\infty$. In some experiments,
 122 V-REx with small λ also shows robust generalization and may outperform ERM. Here we understand
 123 this phenomenon by extending the REX principle to the population level:

- 124 (i) Minimizing the expected risk $R(f)$;
- 125 (ii) Cannot reject the null hypothesis of the test:

$$H_0 : R(f|e) = R(f|e'), \forall e, e' \in \mathcal{E} \quad \text{vs} \quad H_1 : R(f|e) \neq R(f|e'), \exists e, e' \in \mathcal{E}. \quad (7)$$

126 In general, Principle (i) is achieved by minimizing the ERM loss. Next we show that variance-based
 127 regularization is related to the hypothesis testing problem in Principle (ii). Under regular assumptions,
 128 one can use the one-way ANOVA F-test to check the hypothesis testing in (7). The F-test statistic is
 129 the ratio of the between-domain variance to the in-domain variance, i.e.,

$$F = \frac{\hat{V}_{out}(f)}{\hat{V}(f) - \hat{V}_{out}(f)} \quad \text{with} \quad \hat{V}(f) = \frac{1}{m} \sum_{i=1}^n \sum_{j=1}^{m_i} (f(\mathbf{z}_{ij}) - \hat{R}(f))^2.$$

130 Here $\hat{V}_{out}(f)$ and $\hat{R}(f)$ are defined in (4). If F is larger than a threshold, e.g. the $(1 - 5\%)$ -quantile
 131 of a F distribution, one should reject the null hypothesis. Here 5% is the significance level. If
 132 the in-domain variance of a well-trained model is approximately stable, then *penalizing $\hat{V}_{out}(f)$ is*
 133 *equivalent to a constraint that H_0 cannot be rejected.* Therefore our proposed method that penalizes
 134 $\hat{V}_{out}(f)$ is consistent with the REX principle and the regularization term $\hat{V}_{out}(f)$ is a generalized
 135 version of V-REx.

136 3 Variance-Based Regularization

137 Motivated by Section 2, we study a variance-based regularization method for domain generalization,
 138 which minimizes the following empirical loss function:

$$\hat{R}(f) + \lambda \sqrt{\hat{V}_{out}(f)}, \quad (8)$$

139 where λ is a tuning parameter and $\hat{V}_{out}(f)$ is an empirical estimator of the between-domain risk
 140 variance. The proposed loss (8) directly optimizes the ERM principle $\hat{R}(f)$, which is different to
 141 the recent invariant learning methods that minimize $\tilde{R}(f)$, e.g. Invariant Risk Minimization [1].
 142 The regularization term is slightly different to V-REx: (i) The square-root operator is derived from
 143 generalization gap; (ii) *Different to the empirical variance of $R(f|e)$, we penalize the between-domain*
 144 *variance of $f(\mathbf{z})$.*

145 We consider Q_0 in (3) as the training domain distribution and denote the training distribution as
 146 $P_0 = \int P_e Q_0(de)$. To proceed further, we denote more notations as follows:

$$\begin{aligned} R(f) &= \mathbb{E}_{e \sim Q_0}[R(f|e)] = \mathbb{E}_{\mathbf{z} \sim P_0}[f(\mathbf{z})], & V(f) &= \mathbb{E}_{\mathbf{z} \sim P_0}[(f(\mathbf{z}) - R(f))^2], \\ V_{in}(f|e) &= \mathbb{E}_{\mathbf{z} \sim P_e}[(f(\mathbf{z}) - R(f|e))^2], & V_{out}(f) &= \mathbb{E}_{e \sim Q_0}[(R(f|e) - R(f))^2], \end{aligned}$$

147 where $R(f|e)$ is defined in (2). Here $V_{out}(f)$ is the between-domain variance and $V_{in}(f|e)$ is the
 148 in-domain variance of the domain $e \in \mathcal{E}$. According to the decomposition of the total variance, we
 149 have

$$V(f) = \text{Var}(\mathbb{E}[f(\mathbf{z})|e]) + \mathbb{E}[\text{Var}(f|e)] = V_{out}(f) + \mathbb{E}_{e \sim Q_0}[V_{in}(f|e)].$$

150 When Q_0 and P_e are replaced by the corresponding empirical distributions, we rewrite $V(f)$,
 151 $V_{in}(f|e)$ and $V_{out}(f)$ as $\hat{V}(f)$, $\hat{V}_{in}(f|e)$ and $\hat{V}_{out}(f)$ respectively. In the finite-sample setup, the
 152 decomposition of the total variance also holds:

$$\hat{V}(f) = \hat{V}_{out}(f) + \sum_{i=1}^n \frac{m_i}{m} \hat{V}_{in}(f|e).$$

153 3.1 Generalization

154 Since the empirical loss (8) is derived from generalization bounds, we present two versions of the
 155 generalization guarantee. The first result depends on the covering number of the function class \mathcal{F} .
 156 In the Appendix, we also derive a version of the generalization bound with localized Rademacher
 157 complexities, which can provide more refined uniform generalization bounds in some cases.

158 We start with the definition of the covering number. Let \mathcal{F} be a collection of bounded functions
 159 $f: \mathcal{X} \times \mathcal{Y} \rightarrow [0, M]$. Suppose \mathcal{F} is a subset of a metric space with a norm $\|\cdot\|$. We say a collection
 160 $\{f^1, \dots, f^N\} \subseteq \mathcal{F}$ is an ϵ -cover of \mathcal{F} if for each $f \in \mathcal{F}$, there exists f^i such that $\|f - f^i\| \leq \epsilon$. The
 161 covering number of \mathcal{F} is

$$\begin{aligned} N(\mathcal{F}, \epsilon, \|\cdot\|) &:= \inf \left\{ N \in \mathbb{N} : \text{there exists a collection } \{f^1, \dots, f^N\} \right. \\ &\quad \left. \text{which is an } \epsilon\text{-cover of } \mathcal{F} \text{ with respect to } \|\cdot\| \right\}. \end{aligned}$$

162 In the following, we use the ℓ^∞ norm: $\|f - g\|_\infty = \sup_{z \in \mathcal{X} \times \mathcal{Y}} |f(z) - g(z)|$. Now we are ready to
 163 present the following theorem:

164 **Theorem 1** *Let $n \geq 2$ and $\{\mathbf{z}_{i,j}, 1 \leq i \leq n, 1 \leq j \leq m_i\}$ is an i.i.d sample drawn from P_0 .
 165 Suppose $f(z) \in [0, M]$ for any $f \in \mathcal{F}$ and $z \in \mathcal{X} \times \mathcal{Y}$ and the function class \mathcal{F} has the over number:
 166 $N_\epsilon = N(\mathcal{F}, \epsilon, \|\cdot\|_{L^\infty(\mathcal{X} \times \mathcal{Y})})$. Let $0 < \delta < 1$ and*

$$t = \log \frac{(n+2)N_\epsilon}{\delta}, \quad \lambda = \sqrt{\frac{2t}{m-1}}.$$

167 Then we have, with probability at least $1 - \delta$,

$$\begin{aligned}
R(f) &\leq \hat{R}(f) + \lambda \sqrt{\hat{V}_{out}(f)} + \sum_{i=1}^n \lambda \sqrt{\frac{(m_i - 1)V_{in}(f|e_i)}{m}} \\
&\quad + \sum_{i=1}^n \frac{\sqrt{(m-1)m_i}M\lambda^2}{\sqrt{m(m_i-1)}} + \frac{(4m-1)M\lambda^2}{3m} \\
&\quad + \left(2 + \lambda + \sum_{i=1}^n \lambda \sqrt{\frac{(m_i-1)}{m}}\right)\epsilon,
\end{aligned}$$

168 holds for every $f \in \mathcal{F}$.

169 The proof of Theorem 1 is presented in the Appendix B. In some cases, the covering number-
170 based analysis cannot provide a tight generalization bound [4, 5, 36]. Therefore, we also use the
171 local Rademacher complexity [5] to present the generalization of the proposed variance-based
172 regularization. The details and proof are postponed into Appendix C.

173 **Why we study In-Distribution generalization?** Theorem 1 provides the generalization guarantee
174 for the in-distribution (ID) generalization rather than the OOD generalization. But its result gives
175 important insights into the OOD generalization. First, the ID error provides a lower bound for
176 the worst-case OOD error since \mathcal{E}_{tr} is a subset of \mathcal{E} . Second, some empirical studies of OOD
177 generalization have observed a linear relationship between the ID and OOD test error [31, 20, 22].
178 Third, some OOD generalization bounds are derived from a domain adaptation framework [8, 7, 2,
179 42, 43], e.g.,

$$\text{OOD error} \leq \text{ID error} + \text{error gap} + O(\cdot), \quad (9)$$

180 which starts from the ID error and then depicts the error gap. *Most recent works focus on minimising*
181 *the error gap and ignore how their robust (or invariant) methods increase the ID test error.* Fourth,
182 our assumptions are mild and general. We do not impose strong constraint on the test data, e.g.
183 structured generative mechanism, and only assume the domain distribution shift. Therefore, we
184 analyze the ID error of the proposed robust method under mild assumptions.

185 We denote f^* as the optimal function and let \hat{f} be a solution:

$$\hat{f} \in \arg \min_{f \in \mathcal{F}} \hat{R}(f) + \lambda \sqrt{\hat{V}_{out}(f)}.$$

186 Next we study the excess risk of \hat{f} . According to Theorem 1, we obtain the following result.

187 **Corollary 2** Suppose the assumptions in Theorem 1 hold. Let $0 < \delta < 1$ and

$$t = \log \frac{2N_\epsilon + 2}{\delta}, \quad \lambda = \sqrt{\frac{2t}{m-1}}.$$

188 Then, with probability at least $1 - \delta$,

$$\begin{aligned}
R(\hat{f}) - R(f^*) &\leq 2\lambda \sqrt{\frac{(m-1)V(f^*)}{m}} + \sum_{i=1}^n \lambda \sqrt{\frac{m_i \hat{V}_{in}(\hat{f}|e_i)}{m}} \\
&\quad + \left(2 + \lambda + \sum_{i=1}^n \lambda \sqrt{\frac{m_i}{m}}\right)\epsilon + \lambda^2 \frac{4(4m-1)M}{3m}.
\end{aligned}$$

189 **Parametric Example.** Suppose the hypothetical space \mathcal{F} is a class of parametric functions:

$$\mathcal{F} = \{f_\theta(z) : z \in \mathcal{X} \times \mathcal{Y}, \theta \in \Theta \subseteq \mathbb{R}^d\},$$

190 where the parameter set Θ is bounded. Further, for any data point \mathbf{z} , $f_\theta(\mathbf{z})$ is a L -Lipschitz function
191 of θ with respect to ℓ^2 norm on Θ . Then the covering number is bounded above:

$$N_\epsilon \leq \left(1 + \text{diam}(\Theta) \cdot L \cdot \frac{1}{\epsilon}\right)^d, \quad \text{with} \quad \text{diam}(\Theta) = \sup_{\theta, \theta' \in \Theta} \|\theta - \theta'\|_2.$$

192 Then we take

$$\epsilon = \frac{1}{m}, \quad \log N_\epsilon = O(\log m), \quad \lambda = O\left(\sqrt{\frac{\log m}{m}}\right).$$

193 Therefore, by Corollary 2, with probability at least $1 - \delta$,

$$R(\hat{f}) - R(f^*) \leq 2\lambda\sqrt{\frac{(m-1)V(f^*)}{m}} + \sum_{i=1}^n \lambda\sqrt{\frac{m_i\hat{V}_{in}(\hat{f}|e_i)}{m}} + O\left(\frac{\log m}{m}\right). \quad (10)$$

194 **Potential competitive edge.** The second term on the RHS of (10) contains the empirical in-domain
 195 variance $\hat{V}_{in}(\hat{f}|e_i)$. For over-parameterized model, the empirical in-domain variance of \hat{f} can be
 196 close to zero. If there exists an optimal function $f^* \in \arg \min_f R(f)$ such that $V(f^*) = 0$, then the
 197 term $O(\log m/m)$ dominates the convergence rate of the excess risk.¹ For ERM, the convergence
 198 rate of the excess risk is $1/\sqrt{m}$, which is slower than $\log m/m$. *Due to the fast convergence rate, our*
 199 *proposed method can outperform ERM when the sample size m is large enough.*

200 **Cannot consistently outperform ERM.** If there is no optimal function $f^* \in \arg \min_f R(f)$ satisfies
 201 $V(f^*) = 0$, then the first term on the RHS (10) can dominate the excess risk. In this case, the
 202 convergence rate of the the excess risk of our method is $\sqrt{\log m/m}$, which is slower than ERM. This
 203 implies that *if $V(f^*) > 0$ for $\forall f^* \in \arg \min_f R(f)$, ERM can outperform our method when m is*
 204 *large enough.*

205 **OOD generalization.** According to Eq. (9), OOD error can be rewritten as a sum of ID error and error
 206 gap. The distance between the training and test domain distribution can determine the error gap term
 207 under our setup. Furthermore, we only assume that the training and test domain distributions are close
 208 but different, and do not impose any structured generative models, such as structural equation models
 209 [40] or probabilistic graphical models [23]. In other words, we do not introduce prior information and
 210 use mild and general assumptions. Due to the uncertainty of the test data, the error gap should be the
 211 worst-case error gap for the domain distribution shift and hypothetical space, which is independent of
 212 the estimator. This implies that without prior information, the ID error is a reliable metric to infer the
 213 OOD error.

214 **Non-convexity.** Similar to the Sample Variance Penalization [27], the proposed objective function
 215 (8) is in general non-convex and computationally intractable. The proposed regularization term
 216 is non-convex even if the loss function is convex. It is still unclear how to actually minimize the
 217 variance-regularized objective function. Krueger et al. [25] use a penalty annealing scheme to obtain
 218 a good pre-train model. In the Appendix A, we empirically show that our method can use random
 219 initialization without dropping generalization performance.

220 3.2 Optimization

221 In this section, we show that minimizing (8) is equivalent to solving a group DRO problem con-
 222 cerning a local neighbourhood of the empirical domain distribution. Let $\mathbf{q} = (q_1, q_2, \dots, q_n)$ be a
 223 discrete distributions defined on the domain set $\mathcal{E}_{tr} = \{e_1, e_2, \dots, e_n\}$. We consider the following
 224 optimization problem that minimizes

$$\max_{\mathbf{q} \in \mathcal{Q}_{\alpha, \rho}(\hat{\mathbf{q}})} \sum_{i=1}^n q_i \hat{R}(f|e_i), \quad (11)$$

225 which is slightly different to group DRO problem because $\mathcal{Q}_{\alpha}(\hat{\mathbf{q}}, \rho)$ is not centered at the uniform
 226 discrete distribution. We denote $\lambda = \sqrt{2\rho}$ and rewrite the empirical loss in (8) as

$$\mathcal{L}(f; \rho) = \hat{R}(f) + \sqrt{2\rho\hat{V}_{out}(f)}. \quad (12)$$

227 The following theorem shows that the objective (11) is bounded by two variance-regularized functions
 228 in the form of (12).

¹The factor $\log m$ comes from the covering number N_ϵ . If the hypothetical space \mathcal{F} only contains finite models, N_ϵ is a constant and is independent to m . Then the convergence rate of the excess risk is $1/m$.

229 **Theorem 3** Suppose the training dataset \mathcal{D} and a function $f \in \mathcal{F}$ are given. Let ρ_+ be the largest
 230 distance between $\hat{\mathbf{q}}$ and $\mathbf{q} \in \mathcal{Q}_{\alpha, +\infty}(\hat{\mathbf{q}})$ and

$$\rho_- = \frac{\min_i (\alpha/\hat{q}_i - 1)^2 \hat{V}_{out}(f)}{2(\min_i \hat{R}(f|e_i) - \hat{R}(f))^2},$$

231 then we have

$$\mathcal{L}(f; \rho_-) \leq \max_{\mathbf{q} \in \mathcal{Q}_{\alpha, +\infty}(\hat{\mathbf{q}})} \sum_{i=1}^n q_i \hat{R}(f|e_i) \leq \mathcal{L}(f; \rho_+). \quad (13)$$

232 This second inequality in (13) implies that the optimization problem (11) with $\rho = +\infty$ is always
 233 bounded above by the variance-regularized loss with the tuning parameter ρ_+ . On the other hand, we
 234 can also derive a tuning parameter ρ_- depends on the training data and a given model f , and then
 235 prove that $\mathcal{L}(f; \rho_-)$ is a lower boundary of (11). According to the proof of Theorem 3, one can find
 236 that the equality holds:

$$\max_{\mathbf{q} \in \mathcal{Q}_{\alpha, \rho}(\hat{\mathbf{q}})} \sum_{i=1}^n q_i \hat{R}(f|e_i) = \hat{R}(f) + \sqrt{2\rho \hat{V}_{out}(f)},$$

237 when the radius ρ satisfies $\rho \leq \rho_-$. If $\hat{V}_{out}(f)$ is nonzero and ρ is given, the equality holds if and
 238 only if $\forall e_i \in \mathcal{E}_{tr}$,

$$\alpha \leq \hat{q}_i \left(\sqrt{\frac{2\rho}{\hat{V}_{out}(f)}} (\hat{R}(f|e_i) - \hat{R}(f)) + 1 \right). \quad (14)$$

239 Therefore, the parameter α and the radius ρ govern each other.

240 **Sketch of Proof:** We start with a preliminary result: for any α and ρ ,

$$\max_{\mathbf{q} \in \mathcal{Q}_{\alpha, \rho}(\hat{\mathbf{q}})} \sum_{i=1}^n q_i \hat{R}(f|e_i) \leq \mathcal{L}(f; \rho),$$

241 which is directly derived from the Cauchy-Schwarz inequality. By checking the conditions for the
 242 equality, we obtain the constraints in (14). Note that $\mathcal{Q}_{\alpha, +\infty}(\hat{\mathbf{q}}) \subseteq \mathcal{Q}_{+\infty, \rho_+}(\hat{\mathbf{q}})$ since ρ_+ is the largest
 243 distance between $\hat{\mathbf{q}}$ and $\mathbf{q} \in \mathcal{Q}_{\alpha, +\infty}(\hat{\mathbf{q}})$. Hence the second inequality in (13) is trivial. Let \mathbf{q}^* be

$$\mathbf{q}_-^* = \arg \max_{\mathbf{q} \in \mathcal{Q}_{+\infty, \rho_-}(\hat{\mathbf{q}})} \sum_{i=1}^n q_i \hat{R}(f|e_i).$$

244 Furthermore, $\mathbf{q}_-^* \in \mathcal{Q}_{\alpha, \rho_-}(\hat{\mathbf{q}}) \subseteq \mathcal{Q}_{\alpha, +\infty}(\hat{\mathbf{q}})$. Hence the first inequality in (13) holds. \square

245 Theorem 3 shows the equivalence between group DRO and the variance-based regularization in
 246 (12). However, the lower bound $\mathcal{L}(f; \rho_-)$ still depends on the training data and model. Next we use
 247 concentration inequalities and the covering numbers of \mathcal{F} to derive the uniform results.

248 **Theorem 4** Suppose that α is a non-positive scalar and $V'_{out}(f) = V(f) - \sum_i q_i V_{in}(f|e_i) > 0$.
 249 For each training domain, both \hat{q}_i and q_i are larger than $\delta > 0$. We write

$$\rho' = \frac{V'_{out}(f)}{16M^2} \left(\frac{\alpha}{1 - (n-1)\delta} - 1 \right)^2.$$

250 Let $\tau > 0$ and $0 < \eta < 1$ be two constants. Define

$$\mathcal{F}_{\tau, \eta} = \left\{ f \in \mathcal{F} : V(f) \geq \tau, \text{ and } \frac{V_{in}(f|e_i)}{V(f)} \leq \eta, \forall e_i \in \mathcal{E}_{tr} \right\}.$$

251 For any $f \in \mathcal{F}_{\tau, \eta}$, the following expansion uniformly holds:

$$\max_{\mathbf{q} \in \mathcal{Q}_{\alpha, \rho'}(\hat{\mathbf{q}})} \sum_{i=1}^n q_i \hat{R}(f|e_i) = \mathcal{L}(f; \rho'),$$

252 with probability at least $1 - N_{\tau,\eta} \times p$, where $N_{\tau,\eta} = N\left(\mathcal{F}_{\tau,\eta}, \sqrt{\frac{1}{10}(1-\eta)\tau}, \|\cdot\|_{L^\infty(\mathcal{X} \times \mathcal{Y})}\right)$ is the
 253 covering number of $\mathcal{F}_{\tau,\eta}$ and

$$p = \exp\left(-\frac{m(1-\eta)^2\tau}{32M^2} + \frac{1}{16}\right) + \binom{m+n-1}{n-1} \exp\left(-\frac{m(1-\eta)^2\tau^2}{M^4}\right) \\ + \sum_{i=1}^n \exp\left\{-\frac{1}{2M^2m_i} \left(\frac{m_i(1-\eta) + 4\eta}{1+3\eta}\right)^2\right\}.$$

254 4 General Version

255 Recall the uncertainty region in (3): $\{Q : D_\phi(Q\|Q_0) \leq \rho\}$. In Section 3, Q_0 is the ground-truth
 256 domain distribution. In fact, it can be a selected anchor distribution closed to the target test domain.
 257 The choice of Q_0 can be regarded as a kind of prior knowledge and the hyperparameter ρ represents
 258 how strong is the confidence in the prior. We formulate the finite-sample optimization problem as

$$\max_{\mathbf{q} \in \mathcal{Q}_{\alpha,\rho}(\mathbf{q}_0)} \sum_{i=1}^n q_i \hat{R}(f|e_i), \quad (15)$$

259 where \mathbf{q}_0 is the conditional distribution of e given $e \in \mathcal{E}_{tr}$, which is derived from Q_0 . In this problem,
 260 the uncertainty region $\mathcal{Q}_{\alpha,\rho}(\mathbf{q}_0)$ is centered at a discrete distribution \mathbf{q}_0 rather than the uniform
 261 distribution or the empirical distribution $\hat{\mathbf{q}}$. Therefore, we can manually select \mathbf{q}_0 to introduce the
 262 prior information.

263 According to the proof of Theorem 3 in the , the optimization equivalence in Section 3.2 also
 264 holds when we replace $\mathcal{Q}_{\alpha,\rho}(\hat{\mathbf{q}})$ with $\mathcal{Q}_{\alpha,\rho}(\mathbf{q}_0)$. To proceed further, we rewrite $\mathcal{L}(f; \rho, \mathbf{q}_0) =$
 265 $\hat{R}(f, \mathbf{q}_0) + \sqrt{2\rho \hat{V}_{out}(f, \mathbf{q}_0)}$ with

$$\hat{R}(f, \mathbf{q}_0) = \sum_{i=1}^n q_{0,i} \hat{R}(f|e_i) \quad \text{and} \quad \hat{V}_{out}(f, \mathbf{q}_0) = \sum_{i=1}^n q_{0,i} (\hat{R}(f|e_i) - \hat{R}(f, \mathbf{q}_0))^2.$$

266 Then we restate Theorem 3 as the following general version.

267 **Theorem 5** Given the training dataset and a function $f \in \mathcal{F}$, then for any distribution \mathbf{q}_0 , the
 268 inequality always holds:

$$\max_{\mathbf{q} \in \mathcal{Q}_{\alpha,\rho}(\mathbf{q}_0)} \sum_{i=1}^n q_i \hat{R}(f|e_i) \leq \mathcal{L}(f; \rho, \mathbf{q}_0).$$

269 If the between-domain variance $\hat{V}_{out}(f, \mathbf{q}_0)$ is non-zero, the equality holds if and only if $\forall e_i \in \mathcal{E}_{tr}$,

$$\alpha \leq q_{0,i} \left(\sqrt{\frac{2\rho}{\hat{V}_{out}(f, \mathbf{q}_0)}} (\hat{R}(f|e_i) - \hat{R}(f, \mathbf{q}_0)) + 1 \right).$$

270 On the other hand, if α is fixed, the equality holds when the radius of $\mathcal{Q}_{\alpha,\rho}(\mathbf{q}_0)$ satisfies

$$\rho \leq \frac{\min_i (\alpha/q_{0,i} - 1)^2 \hat{V}_{out}(f, \mathbf{q}_0)}{2(\min_i \hat{R}(f|e_i) - \hat{R}(f, \mathbf{q}_0))^2}.$$

271 This result shows the equivalence between the optimization problem (15) and the variance-regularized
 272 loss $\mathcal{L}(f; \rho, \mathbf{q}_0)$. Therefore, for unbalanced domains and any given prior \mathbf{q}_0 , we can still use the
 273 variance-based regularization to approximate the DRO problem. Please refer to the Appendix for the
 274 complete proof of Theorem 5.

275 5 Conclusion

276 In this work, we study a variance-based regularization method for domain generalization. We prove
 277 the guarantees for in-distribution generalization and figure out the potential benefits of our proposed
 278 method compared to ERM. Our proposed objective function is non-convex and the optimization
 279 procedure is computationally intractable. The learnt model can be highly dependent on initialization
 280 or pretraining. In future work, we will consider combining generalization bounds with specific
 281 optimization algorithms to seek fine-grained generalization guarantees.

282 **References**

- 283 [1] Kartik Ahuja, Karthikeyan Shanmugam, Kush Varshney, and Amit Dhurandhar. Invariant risk minimization
284 games. In *International Conference on Machine Learning*, pages 145–155. PMLR, 2020.
- 285 [2] Isabela Albuquerque, João Monteiro, Mohammad Darvishi, Tiago H Falk, and Ioannis Mitliagkas. Adver-
286 sarial target-invariant representation learning for domain generalization. 2020.
- 287 [3] Syed Mumtaz Ali and Samuel D Silvey. A general class of coefficients of divergence of one distribution
288 from another. *Journal of the Royal Statistical Society: Series B (Methodological)*, 28(1):131–142, 1966.
- 289 [4] Peter L Bartlett and Shahar Mendelson. Rademacher and gaussian complexities: Risk bounds and structural
290 results. *Journal of Machine Learning Research*, 3(Nov):463–482, 2002.
- 291 [5] Peter L Bartlett, Olivier Bousquet, Shahar Mendelson, et al. Local rademacher complexities. *The Annals*
292 *of Statistics*, 33(4):1497–1537, 2005.
- 293 [6] Peter L Bartlett, Michael I Jordan, and Jon D McAuliffe. Convexity, classification, and risk bounds.
294 *Journal of the American Statistical Association*, 101(473):138–156, 2006.
- 295 [7] Shai Ben-David, John Blitzer, Koby Crammer, Fernando Pereira, et al. Analysis of representations for
296 domain adaptation. *Advances in neural information processing systems*, 19:137, 2007.
- 297 [8] Shai Ben-David, John Blitzer, Koby Crammer, Alex Kulesza, Fernando Pereira, and Jennifer Wortman
298 Vaughan. A theory of learning from different domains. *Machine learning*, 79(1):151–175, 2010.
- 299 [9] Aharon Ben-Tal, Laurent El Ghaoui, and Arkadi Nemirovski. *Robust optimization*, volume 28. Princeton
300 University Press, 2009.
- 301 [10] Aharon Ben-Tal, Elad Hazan, Tomer Koren, and Shie Mannor. Oracle-based robust optimization via online
302 learning. *Operations Research*, 63(3):628–638, 2015.
- 303 [11] Dimitris Bertsimas, Vishal Gupta, and Nathan Kallus. Data-driven robust optimization. *Mathematical*
304 *Programming*, 167(2):235–292, 2018.
- 305 [12] Gilles Blanchard, Gyemin Lee, and Clayton Scott. Generalizing from several related classification tasks to
306 a new unlabeled sample. *Advances in neural information processing systems*, 24:2178–2186, 2011.
- 307 [13] Stéphane Boucheron, Olivier Bousquet, and Gábor Lugosi. Theory of classification: A survey of some
308 recent advances. *ESAIM: probability and statistics*, 9:323–375, 2005.
- 309 [14] Imre Csiszár. Information-type measures of difference of probability distributions and indirect observation.
310 *studia scientiarum Mathematicarum Hungarica*, 2:229–318, 1967.
- 311 [15] Imre Csiszár. The method of types [information theory]. *IEEE Transactions on Information Theory*, 44(6):
312 2505–2523, 1998.
- 313 [16] John Duchi and Hongseok Namkoong. Variance-based regularization with convex objectives. *The Journal*
314 *of Machine Learning Research*, 20(1):2450–2504, 2019.
- 315 [17] John C Duchi, Peter W Glynn, and Hongseok Namkoong. Statistics of robust optimization: A generalized
316 empirical likelihood approach. *Mathematics of Operations Research*, 2021.
- 317 [18] Ishaan Gulrajani and David Lopez-Paz. In search of lost domain generalization. *arXiv preprint*
318 *arXiv:2007.01434*, 2020.
- 319 [19] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition.
320 In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2016.
- 321 [20] Dan Hendrycks, Steven Basart, Norman Mu, Saurav Kadavath, Frank Wang, Evan Dorundo, Rahul Desai,
322 Tyler Zhu, Samyak Parajuli, Mike Guo, et al. The many faces of robustness: A critical analysis of
323 out-of-distribution generalization. In *Proceedings of the IEEE/CVF International Conference on Computer*
324 *Vision*, pages 8340–8349, 2021.
- 325 [21] Weihua Hu, Gang Niu, Issei Sato, and Masashi Sugiyama. Does distributionally robust supervised learning
326 give robust classifiers? In *International Conference on Machine Learning*, pages 2029–2037. PMLR,
327 2018.

- 328 [22] Pang Wei Koh, Shiori Sagawa, Henrik Marklund, Sang Michael Xie, Marvin Zhang, Akshay Balsubramani,
329 Weihua Hu, Michihiro Yasunaga, Richard Lanus Phillips, Irena Gao, et al. Wilds: A benchmark of in-
330 the-wild distribution shifts. In *International Conference on Machine Learning*, pages 5637–5664. PMLR,
331 2021.
- 332 [23] Daphne Koller and Nir Friedman. *Probabilistic graphical models: principles and techniques*. MIT press,
333 2009.
- 334 [24] Vladimir Koltchinskii et al. Local rademacher complexities and oracle inequalities in risk minimization.
335 *The Annals of Statistics*, 34(6):2593–2656, 2006.
- 336 [25] David Krueger, Ethan Caballero, Joern-Henrik Jacobsen, Amy Zhang, Jonathan Binas, Dinghui Zhang,
337 Remi Le Priol, and Aaron Courville. Out-of-distribution generalization via risk extrapolation (rex). In
338 Marina Meila and Tong Zhang, editors, *Proceedings of the 38th International Conference on Machine*
339 *Learning*, volume 139 of *Proceedings of Machine Learning Research*, pages 5815–5826. PMLR, 18–24
340 Jul 2021. URL <https://proceedings.mlr.press/v139/krueger21a.html>.
- 341 [26] Da Li, Yongxin Yang, Yi-Zhe Song, and Timothy M Hospedales. Deeper, broader and artier domain
342 generalization. In *Proceedings of the IEEE international conference on computer vision*, pages 5542–5550,
343 2017.
- 344 [27] Andreas Maurer and Massimiliano Pontil. Empirical bernstein bounds and sample variance penalization.
345 *arXiv preprint arXiv:0907.3740*, 2009.
- 346 [28] Krikamol Muandet, David Balduzzi, and Bernhard Schölkopf. Domain generalization via invariant feature
347 representation. In *International Conference on Machine Learning*, pages 10–18, 2013.
- 348 [29] Hongseok Namkoong and John C Duchi. Variance-based regularization with convex objec-
349 tives. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan,
350 and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 30.
351 Curran Associates, Inc., 2017. URL [https://proceedings.neurips.cc/paper/2017/file/](https://proceedings.neurips.cc/paper/2017/file/5a142a55461d5fef016acfb927fee0bd-Paper.pdf)
352 [5a142a55461d5fef016acfb927fee0bd-Paper.pdf](https://proceedings.neurips.cc/paper/2017/file/5a142a55461d5fef016acfb927fee0bd-Paper.pdf).
- 353 [30] Yonatan Oren, Shiori Sagawa, Tatsunori Hashimoto, and Percy Liang. Distributionally robust language
354 modeling. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing*
355 *and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages
356 4218–4228, 2019.
- 357 [31] Benjamin Recht, Rebecca Roelofs, Ludwig Schmidt, and Vaishaal Shankar. Do imagenet classifiers
358 generalize to imagenet? In *International Conference on Machine Learning*, pages 5389–5400. PMLR,
359 2019.
- 360 [32] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang,
361 Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. Imagenet large scale visual recognition challenge.
362 *International journal of computer vision*, 115(3):211–252, 2015.
- 363 [33] Shiori Sagawa*, Pang Wei Koh*, Tatsunori B. Hashimoto, and Percy Liang. Distributionally robust neural
364 networks. In *International Conference on Learning Representations*, 2020. URL [https://openreview.](https://openreview.net/forum?id=ryxGuJrFvS)
365 [net/forum?id=ryxGuJrFvS](https://openreview.net/forum?id=ryxGuJrFvS).
- 366 [34] Alexander Shapiro, Darinka Dentcheva, and Andrzej Ruszczyński. *Lectures on stochastic programming:*
367 *modeling and theory*. SIAM, 2014.
- 368 [35] Zheyang Shen, Jiashuo Liu, Yue He, Xingxuan Zhang, Renzhe Xu, Han Yu, and Peng Cui. Towards
369 out-of-distribution generalization: A survey. *arXiv preprint arXiv:2108.13624*, 2021.
- 370 [36] Nathan Srebro, Karthik Sridharan, and Ambuj Tewari. Smoothness, low noise and fast rates. In *Advances*
371 *in neural information processing systems*, pages 2199–2207, 2010.
- 372 [37] Antonio Torralba and Alexei A Efros. Unbiased look at dataset bias. In *CVPR 2011*, pages 1521–1528.
373 IEEE, 2011.
- 374 [38] Vladimir Vapnik. Principles of risk minimization for learning theory. In *Advances in neural information*
375 *processing systems*, pages 831–838, 1992.
- 376 [39] Jindong Wang, Cuiling Lan, Chang Liu, Yidong Ouyang, and Tao Qin. Generalizing to unseen domains:
377 A survey on domain generalization. In Zhi-Hua Zhou, editor, *Proceedings of the Thirtieth International*
378 *Joint Conference on Artificial Intelligence, IJCAI-21*, pages 4627–4635. International Joint Conferences on
379 Artificial Intelligence Organization, 8 2021. doi: 10.24963/ijcai.2021/628. URL [https://doi.org/10.](https://doi.org/10.24963/ijcai.2021/628)
380 [24963/ijcai.2021/628](https://doi.org/10.24963/ijcai.2021/628). Survey Track.

- 381 [40] Sewall Wright. Correlation and causation. 1921.
- 382 [41] Chuanlong Xie, Haotian Ye, Fei Chen, Yue Liu, Rui Sun, and Zhenguo Li. Risk variance penalization.
383 *arXiv preprint arXiv:2006.07544*, 2020.
- 384 [42] Haotian Ye, Chuanlong Xie, Tianle Cai, Ruichen Li, Zhenguo Li, and Liwei Wang. Towards a theoretical
385 framework of out-of-distribution generalization. In A. Beygelzimer, Y. Dauphin, P. Liang, and J. Wortman
386 Vaughan, editors, *Advances in Neural Information Processing Systems*, 2021.
- 387 [43] Guojun Zhang, Han Zhao, Yaoliang Yu, and Pascal Poupart. Quantifying and improving transferability
388 in domain generalization. In A. Beygelzimer, Y. Dauphin, P. Liang, and J. Wortman Vaughan, editors,
389 *Advances in Neural Information Processing Systems*, 2021.
- 390 [44] Kaiyang Zhou, Ziwei Liu, Yu Qiao, Tao Xiang, and Chen Change Loy. Domain generalization: A survey.
391 *arXiv preprint arXiv:2103.02503*, 2021.

392 A Experiments

393 In this section, we present empirical evidence to verify our theoretical results that under mild and
394 general assumptions, *the proposed weighting correction scheme in (5) has better ID generalization*
395 *guarantees than the existing variance-regularized domain generalization methods*. By Theorem 4,
396 we reformulate the objective function in (5) into a batch version. Then we consider the following two
397 settings:

- 398 • **Balanced batch.** This is the standard operation in DomainBed [18] and is commonly used
399 by the existing variance-based regularization methods. In each iteration, the same number
400 of data points are randomly drawn from each training domain to form a batch.
- 401 • **Unbalanced batch.** (Our method) In this setting, we randomly draw data points from each
402 training domain with equal proportions, such that the proportion of each domain in one
403 batch is the same as the proportion of the domains in the entire training data.

404 We consider three variance-regularized domain generalization methods.

- 405 • **Variance.** The V-REx regularization [25] penalizes the domain-level variance of risk without
406 considering the empirical domain distribution. The V-REx estimator is

$$\hat{f} = \arg \min_{f \in \mathcal{F}} \tilde{R}(f) + \lambda \tilde{V}_{out}(f),$$

407 where

$$\tilde{R}(f) = \frac{1}{n} \sum_{i=1}^n \hat{R}(f|e_i) \quad \text{and} \quad \tilde{V}_{out}(f) = \frac{1}{n} \sum_{i=1}^n \left(\hat{R}(f|e_i) - \tilde{R}(f) \right)^2.$$

- 408 • **Standard Deviation.** We also consider the RVP Regularization [41], which slightly changes
409 the penalty term of V-REx into the domain-level standard deviation of risk. Xie et al. [41]
410 provides an understanding of generalization from the perspective of quantile regression and
411 shows that RVP locally approximates DRO. We *remove the scheme of penalty annealing*
412 since it is not involved in group DRO. The RVP estimator is

$$\hat{f} = \arg \min_{f \in \mathcal{F}} \tilde{R}(f) + \lambda \sqrt{\tilde{V}_{out}(f)}.$$

- 413 • **Weighting Correction.** Our proposed method introduces the empirical domain distribution
414 as a weighting scheme into the objective function. We also *remove the scheme of penalty*
415 *annealing*. Our proposed estimator is

$$\hat{f} = \arg \min_{f \in \mathcal{F}} \hat{R}(f) + \lambda \sqrt{\hat{V}_{out}(f)},$$

416 where

$$\hat{R}(f) = \sum_{i=1}^n \hat{q}_i \hat{R}(f|e_i) \quad \text{and} \quad \hat{V}_{out}(f) = \sum_{i=1}^n \hat{q}_i \left(\hat{R}(f|e_i) - \hat{R}(f) \right)^2.$$

417 **Implementation Details.** We consider two datasets: PACS [26] and VLCS [37] We use ResNet50
418 as neural network architecture [19] and start from a pretrained model on ImageNet [32]. In order
419 to fairly evaluate the different regularization, we follow the DomainBed Benchmark to randomly
420 select 20 groups of hyperparameter combinations and repeated the experiment three times for each
421 hyperparameter group. The model is selected according to the training domain validation accuracy,
422 that is the in-distribution validation accuracy. The hyper-parameters includes batch size, learning
423 rate, weight decay, iterations of penalty annealing, and the regularization parameter λ . The other
424 experimental settings are the same as those in Gulrajani and Lopez-Paz [18].

425 The results are reported in Table 1. One can find that the improvement of the weighting correction
426 scheme is statistically significant compared to the original VREx method.

Table 1: The ID prediction accuracy on PACS and VLCS.

PACS						
Balance	Regularization Method	A	C	P	S	Avg
×	Weighting Correction	96.9 ± 0.1	97.0 ± 0.2	96.4 ± 0.1	97.2 ± 0.2	96.9 ± 0.1
×	Standard Deviation	67.4 ± 0.3	54.9 ± 1.0	51.0 ± 1.4	82.4 ± 0.4	63.9 ± 0.2
×	Variance	20.5 ± 0.5	19.6 ± 0.2	20.5 ± 0.5	36.7 ± 5.9	24.3 ± 1.7
✓	Standard Deviation	82.0 ± 0.6	81.3 ± 0.1	77.2 ± 0.6	86.1 ± 0.5	81.7 ± 0.4
✓	Variance	97.0 ± 0.1	96.6 ± 0.1	96.2 ± 0.2	97.0 ± 0.2	96.7 ± 0.1
VLCS						
Balance	Regularization Method	C	L	S	V	Avg
×	Weighting Correction	81.9 ± 0.2	87.6 ± 0.1	85.7 ± 0.3	83.4 ± 0.0	84.7 ± 0.0
×	Standard Deviation	43.7 ± 0.0	45.0 ± 0.2	48.2 ± 0.7	46.1 ± 0.5	45.8 ± 0.3
×	Variance	53.8 ± 1.5	49.1 ± 0.9	53.3 ± 0.9	52.7 ± 1.7	52.2 ± 0.5
✓	Standard Deviation	48.0 ± 2.3	47.5 ± 0.7	48.9 ± 0.5	53.6 ± 0.4	49.5 ± 0.7
✓	Variance	81.3 ± 0.2	87.3 ± 0.2	85.4 ± 0.5	83.1 ± 0.2	84.3 ± 0.2

427 B Proof of Theorem 1

428 **Theorem.** Let $n \geq 2$ and $\{\mathbf{z}_{i,j}, 1 \leq i \leq n, 1 \leq j \leq m_i\}$ is an i.i.d sample drawn from P_0 . Suppose
429 the function set \mathcal{F} has cover numbers

$$N_\epsilon = N(\mathcal{F}, \epsilon, \|\cdot\|_{L^\infty(\mathcal{X} \times \mathcal{Y})}),$$

430 and for any $f \in \mathcal{F}$ and $z \in \mathcal{X} \times \mathcal{Y}$, $f(z) \in [0, M]$. Then we have for every $f \in \mathcal{F}$,

$$\begin{aligned} R(f) &\leq \hat{R}(f) + \sqrt{\frac{2\hat{V}_{out}(f)t}{m-1}} + \sum_{i=1}^n \sqrt{\frac{2(m_i-1)V_{in}(f|e_i)t}{m(m-1)}} \\ &\quad + \sum_{i=1}^n \frac{2\sqrt{m_i}Mt}{\sqrt{m(m_i-1)(m-1)}} + \frac{2(4m-1)Mt}{3m(m-1)} \\ &\quad + \left(2 + \sqrt{\frac{2t}{m-1}} + \sum_{i=1}^n \sqrt{\frac{2(m_i-1)t}{m(m-1)}}\right)\epsilon, \end{aligned}$$

431 with probability at least $1 - (n+2)N_\epsilon \exp(-t)$.

432 **Proof:** By the Bernstein inequality, with probability at least $1 - \exp(-t)$,

$$R(f) \leq \hat{R}(f) + \sqrt{\frac{2V(f)t}{m}} + \frac{2Mt}{3m},$$

433 holds for any given function f . Furthermore, by Theorem 10 in [27],

$$\sqrt{V(f)} \leq \sqrt{\frac{m}{m-1}\hat{V}(f)} + \frac{\sqrt{2mM^2t}}{m-1},$$

434 holds with probability larger than $1 - \exp(-t)$. Then,

$$\begin{aligned} R(f) &\leq \hat{R}(f) + \sqrt{\frac{2\hat{V}(f)t}{m-1}} + \frac{2Mt}{m-1} + \frac{2Mt}{3m}, \\ &= \hat{R}(f) + \sqrt{\frac{2\hat{V}(f)t}{m-1}} + \frac{2(4m-1)Mt}{3m(m-1)}, \end{aligned}$$

435 holds with probability larger than $1 - 2\exp(-t)$. According to the decomposition of the total variance,
436 we know that

$$\begin{aligned} \hat{V}(f) &= \frac{1}{m} \sum_{i=1}^n \sum_{j=1}^{m_i} (f(\mathbf{z}_{i,j}) - \hat{R}(f|e_i) + \hat{R}(f|e_i) - \hat{R}(f))^2 \\ &= \frac{1}{m} \sum_{i=1}^n \sum_{j=1}^{m_i} (f(\mathbf{z}_{i,j}) - \hat{R}(f|e_i))^2 + (\hat{R}(f|e_i) - \hat{R}(f))^2 \\ &= \sum_{i=1}^n \frac{m_i}{m} \hat{V}_{in}(f|e_i) + \sum_{i=1}^n \hat{q}_i (\hat{R}(f|e_i) - \hat{R}(f))^2 \\ &= \sum_{i=1}^n \frac{m_i}{m} \hat{V}_{in}(f|e_i) + \hat{V}_{out}(f). \end{aligned}$$

437 Hence, with probability larger than $1 - 2\exp(-t)$,

$$R(f) \leq \hat{R}(f) + \sqrt{\frac{2\hat{V}_{out}(f)t}{m-1}} + \sum_{i=1}^n \sqrt{\frac{2m_i \hat{V}_{in}(f|e_i)t}{m(m-1)}} + \frac{2(4m-1)Mt}{3m(m-1)}.$$

438 By applying Theorem 10 in [27],

$$\sqrt{V_{in}(f|e_i)} \geq \sqrt{\frac{m_i}{m_i-1} \hat{V}_{in}(f|e_i)} - \frac{\sqrt{2m_i M^2 t}}{m_i-1},$$

439 holds with probability smaller than $\exp(-t)$. Hence we have,

$$\begin{aligned} R(f) &\leq \hat{R}(f) + \sqrt{\frac{2\hat{V}_{out}(f)t}{m-1}} + \sum_{i=1}^n \sqrt{\frac{2(m_i-1)V_{in}(f|e_i)t}{m(m-1)}} \\ &\quad + \sum_{i=1}^n \frac{2\sqrt{m_i}Mt}{\sqrt{m(m_i-1)(m-1)}} + \frac{2(4m-1)Mt}{3m(m-1)}, \end{aligned}$$

440 holds with probability larger than $1 - (2+n)\exp(-t)$.

441 Next, we consider a set of functions $\{f^1, \dots, f^{N_\epsilon}\}$, which is a minimal ϵ -cover of the function space
442 \mathcal{F} of size

$$N_\epsilon = N(\mathcal{F}, \epsilon, \|\cdot\|_{L^\infty(\mathcal{X} \times \mathcal{Y})}).$$

443 Then, for any $f \in \mathcal{F}$, there exists f^j , $1 \leq j \leq N_\epsilon$ such that $\|f - f^j\|_{L^\infty(\mathcal{X} \times \mathcal{Y})} \leq \epsilon$. Therefore,

$$\begin{aligned} R(f) &\leq R(f^j) + \epsilon \\ &\leq \hat{R}(f^j) + \sqrt{\frac{2\hat{V}_{out}(f^j)t}{m-1}} + \sum_{i=1}^n \sqrt{\frac{2(m_i-1)V_{in}(f^j|e_i)t}{m(m-1)}} \\ &\quad + \sum_{i=1}^n \frac{2\sqrt{m_i}Mt}{\sqrt{m(m_i-1)(m-1)}} + \frac{2(4m-1)Mt}{3m(m-1)} + \epsilon, \end{aligned}$$

444 with probability larger than $1 - (2+n)\exp(-t)$. Notice that $\hat{R}(f^j) \leq \hat{R}(f) + \epsilon$ and

$$\begin{aligned} \sqrt{\hat{V}_{out}(f^j)} &\leq \sqrt{\hat{V}_{out}(f)} + \sqrt{\hat{V}_{out}(f^j - f)} \leq \sqrt{\hat{V}_{out}(f)} + \epsilon, \\ \sqrt{\hat{V}_{in}(f^j|e_i)} &\leq \sqrt{\hat{V}_{in}(f|e_i)} + \sqrt{\hat{V}_{in}(f^j - f|e_i)} \leq \sqrt{\hat{V}_{in}(f|e_i)} + \epsilon. \end{aligned}$$

445 Therefore, for every $f \in \mathcal{F}$,

$$\begin{aligned}
R(f) &\leq \hat{R}(f) + \sqrt{\frac{2\hat{V}_{out}(f)t}{m-1}} + \sum_{i=1}^n \sqrt{\frac{2(m_i-1)V_{in}(f|e_i)t}{m(m-1)}} \\
&\quad + \sum_{i=1}^n \frac{2\sqrt{m_i}Mt}{\sqrt{m(m_i-1)(m-1)}} + \frac{2(4m-1)Mt}{3m(m-1)} \\
&\quad + \left(2 + \sqrt{\frac{2t}{m-1}} + \sum_{i=1}^n \sqrt{\frac{2(m_i-1)t}{m(m-1)}}\right)\epsilon,
\end{aligned}$$

446 holds with probability at least $1 - (n+2)N_\epsilon \exp(-t)$.

447

□

448 Let

$$t = \log \frac{(n+2)N_\epsilon}{\delta}, \quad \text{and} \quad \lambda = \sqrt{\frac{2t}{m-1}}.$$

449 Then we have, with probability at least $1 - \delta$,

$$\begin{aligned}
R(f) &\leq \hat{R}(f) + \lambda\sqrt{\hat{V}_{out}(f)} + \sum_{i=1}^n \lambda\sqrt{\frac{(m_i-1)V_{in}(f|e_i)}{m}} \\
&\quad + \sum_{i=1}^n \frac{\sqrt{(m-1)m_i}M\lambda^2}{\sqrt{m(m_i-1)}} + \frac{(4m-1)M\lambda^2}{3m} \\
&\quad + \left(2 + \sqrt{\frac{2t}{m-1}} + \sum_{i=1}^n \sqrt{\frac{2(m_i-1)t}{m(m-1)}}\right)\epsilon,
\end{aligned}$$

450 for every $f \in \mathcal{F}$. Hence Theorem 1 is proved.

451

□

452 C More results for Theorem 1

453 In this section, we use the local Rademacher complexity [5] to present the generalization of the
454 proposed variance-based regularization. We start with the definition of the sub-root function, which
455 can be used to bound the local Rademacher complexity.

456 **Definition 6 ([5], Definition 3.1)** A function $\psi : [0, \infty) \rightarrow [0, \infty)$ is sub-root if it is non-negative,
457 non-decreasing and if $r \mapsto \psi(r)/\sqrt{r}$ is non-increasing for $r > 0$.

458 For any nontrivial sub-root function ψ , i.e., not the constant function $\psi \equiv 0$, it is continuous and has
459 a unique positive fix point $r^* = \psi(r^*)$. In addition, for all $r > 0$, $r \geq \psi(r)$ if and only if $r^* \leq r$.
460 We consider the local Rademacher complexity:

$$\mathbb{E}[\mathcal{R}_n(\{cf(z) : f \in \mathcal{F}, c \in [0, 1], \text{ and } \mathbb{E}_{\mathbf{z} \sim P}[c^2 f^2(\mathbf{z})] \leq r\})] \quad (16)$$

461 which is also used by [16]. Here the notation \mathbb{E} in (16) takes the expectations with respect to the
462 Rademacher random variables. We denote a sub-root function $\psi_m(r)$ as an upper bound of the
463 localized Rademacher complexity:

$$\psi_m(r) \geq \mathbb{E}[\mathcal{R}_n(\{cf(z) : f \in \mathcal{F}, c \in [0, 1], \text{ and } \mathbb{E}_{\mathbf{z} \sim P}[c^2 f^2(\mathbf{z})] \leq r\})]. \quad (17)$$

464 The solution of $\psi_m(r) = r$ is denoted as r_m^* . When the distribution P is replaced by P_{e_i} , the upper
465 bound sub-root function and the corresponding fixed point are written as $\psi_{m,i}$ and $r_{m,i}^*$.

466 **Theorem 7** Let \mathcal{F} be a collection of bounded functions $f : \mathcal{X} \times \mathcal{Y} \rightarrow [0, M]$ satisfying the localiza-
467 tion inequality (17) for some sub-root function $\psi_m(r)$ ($\psi_{m,i}(r)$) with root r_m^* ($r_{m,i}^*$). Let

$$B_m = \frac{1}{m}(t + \log \lceil \log \frac{m}{t} \rceil) \quad \text{and} \quad C_m = 2((2e + 84M)B_m + 36r_m^*),$$

468 where $\lceil \cdot \rceil$ stands for the ceiling function. Then, for every $f \in \mathcal{F}$,

$$\begin{aligned}
R(f) &\leq (1 + \sqrt{2C_m}) \left(\hat{R}(f) + \frac{\sqrt{2C_m}}{1 + \sqrt{2C_m}} \sqrt{\hat{V}_{out}(f)} \right) \\
&\quad + \sqrt{\frac{3}{2} C_m \sum_{i=1}^n \frac{m_i}{m} \mathbb{E}[f^2(\mathbf{z}) | e_i]} \\
&\quad + \sqrt{144C_m M^2 \sum_{i=1}^n \frac{m_i}{m} r_{m,i}^* + C_m \frac{nMt}{m} (4 + \frac{7}{3}M)} \\
&\quad + \sqrt{144C_m M^2 r_m^* + \frac{C_m Mt}{m} (4 + \frac{7}{3}M) + 6r_m^* + 14MB_m},
\end{aligned}$$

469 with probability at least $1 - (2 + n) \exp(-t)$.

470 **Proof:** Before proving the theorem, we state a useful lemma that provides a version of uniform
471 Bernstein's inequality by measuring the complexity of the localized functions that near the optimum
472 of an empirical risk.

473 **Lemma 8** [[16], Lemma 17 and Lemma 18] *Let \mathcal{F} be a collection of bounded functions $f : \mathcal{X} \times \mathcal{Y} \rightarrow$
474 $[0, M]$ satisfying the localization inequality (17) for some sub-root function $\psi_m(r)$ with root r_m^* . Let*

$$B_m = \frac{1}{m} (t + \log \lceil \log \frac{m}{t} \rceil) \quad \text{and} \quad \eta > 0.$$

475 Then with probability at least $1 - \exp(-t)$, for every $f \in \mathcal{F}$,

$$|R(f) - \hat{R}(f)| \leq (\sqrt{2eB_m} + 6\sqrt{r_m^* + \frac{7}{3}MB_m}) \sqrt{\mathbb{E}[f^2]} + 6r_m^* + 14MB_m, \quad (18)$$

476

$$\mathbb{E}[f^2] \leq \hat{\mathbb{E}}[f^2] + \frac{1}{\eta} \hat{\mathbb{E}}[f^2] + 72M^2(1 + \eta)r_m^* + \frac{Mt}{m} (4 + \frac{7}{3}M), \quad (19)$$

477 and

$$\hat{\mathbb{E}}[f^2] \leq \mathbb{E}[f^2] + \frac{\eta}{\eta + 1} \mathbb{E}[f^2] + 72M^2(1 + \eta)r_m^* + \frac{Mt}{m} (4 + \frac{7}{3}M), \quad (20)$$

478 where $\hat{\mathbb{E}}[f^2(\mathbf{z})] = \frac{1}{m} \sum_{i=1}^n \sum_{j=1}^{m_i} f^2(\mathbf{z}_{i,j})$.

479 **Proof:** Please refer to Appendix D.1 and D.2 of [16] for a complete proof of the lemma.

480

□

481 Now we turn to the proof of Theorem 7. We denote $C_m = 2((2e + 84M)B_m + 36r_m^*)$. It is easy to
482 see that

$$\begin{aligned}
&(\sqrt{2eB_m} + 6\sqrt{r_m^* + 7MB_m/3})^2 \\
&= 2eB_m + (36r_m^* + 84MB_m) + 2\sqrt{2eB_m} \sqrt{36r_m^* + 84MB_m} \\
&\leq 2(2eB_m + 36r_m^* + 84MB_m) = 2C_m.
\end{aligned}$$

483 Then the inequality (18) implies that for every $f \in \mathcal{F}$,

$$R(f) \leq \hat{R}(f) + \sqrt{C_m \mathbb{E}[f^2(\mathbf{z})]} + 6r_m^* + 14MB_m,$$

484 holds with probability at least $1 - \exp(-t)$. By (19) with $\eta = 1$, we have for all $f \in \mathcal{F}$,

$$R(f) \leq \hat{R}(f) + \sqrt{2C_m \hat{\mathbb{E}}[f^2(\mathbf{z})]} + \sqrt{144C_m M^2 r_m^* + \frac{C_m Mt}{m} (4 + \frac{7}{3}M) + 6r_m^* + 14MB_m},$$

485 with probability at least $1 - 2 \exp(-t)$. Notice that

$$\begin{aligned}\hat{\mathbb{E}}[f^2(\mathbf{z})] &= \frac{1}{m} \sum_{i=1}^n \sum_{j=1}^{m_i} (f(\mathbf{z}_{i,j}) - \hat{R}(f|e_i) + \hat{R}(f|e_i) - \hat{R}(f) + \hat{R}(f))^2 \\ &= \frac{1}{m} \sum_{i=1}^n \sum_{j=1}^{m_i} (f(\mathbf{z}_{i,j}) - \hat{R}(f|e_i))^2 + (\hat{R}(f|e_i) - \hat{R}(f))^2 + \hat{R}(f)^2 \\ &= \hat{R}(f)^2 + \hat{V}_{out}(f) + \sum_{i=1}^n \frac{m_i}{m} \hat{V}_{in}(f|e_i).\end{aligned}$$

486 Then we have

$$\begin{aligned}R(f) &\leq (1 + \sqrt{2C_m})\hat{R}(f) + \sqrt{2C_m\hat{V}_{out}(f)} + \sqrt{C_m \sum_{i=1}^n \frac{m_i}{m} \hat{V}_{in}(f|e_i)} \\ &\quad + \sqrt{144C_m M^2 r_m^* + \frac{C_m M t}{m} (4 + \frac{7}{3}M) + 6r_m^* + 14MB_m}.\end{aligned}$$

Next we deal with the upper bound of $\hat{V}_{in}(f|e_i)$. We denote

$$\hat{\mathbb{E}}[f^2|e_i] = \frac{1}{m_i} \sum_{j=1}^{m_i} f^2(\mathbf{z}_{i,j}).$$

487 According to (20), for every $f \in \mathcal{F}$,

$$\hat{\mathbb{E}}[f^2(\mathbf{z})|e_i] \leq \frac{2\eta + 1}{\eta + 1} \mathbb{E}[f^2(\mathbf{z})|e_i] + 72M^2(1 + \eta)r_{m,i}^* + \frac{Mt}{m_i} (4 + \frac{7}{3}M),$$

488 holds with probability at least $1 - \exp(-t)$. Let $\eta = 1$. Then, for every $f \in \mathcal{F}$,

$$\begin{aligned}R(f) &\leq (1 + \sqrt{2C_m})\hat{R}(f) + \sqrt{2C_m\hat{V}_{out}(f)} + \sqrt{\frac{3}{2}C_m \sum_{i=1}^n \frac{m_i}{m} \mathbb{E}[f^2(\mathbf{z})|e_i]} \\ &\quad + \sqrt{144C_m M^2 \sum_{i=1}^n \frac{m_i}{m} r_{m,i}^* + C_m \frac{nMt}{m} (4 + \frac{7}{3}M)} \\ &\quad + \sqrt{144C_m M^2 r_m^* + \frac{C_m M t}{m} (4 + \frac{7}{3}M) + 6r_m^* + 14MB_m},\end{aligned}$$

489 with probability at least $1 - (2 + n) \exp(-t)$.

490

□

491 D Proof of Corollary 2

492 **Corollary.** Let $n \geq 2$ and $\{\mathbf{z}_{i,j}, 1 \leq i \leq n, 1 \leq j \leq m_i\}$ is an i.i.d sample. Suppose the function set
493 \mathcal{F} has cover numbers

$$N_\epsilon = N(\mathcal{F}, \epsilon, \|\cdot\|_{L^\infty(\mathcal{X} \times \mathcal{Y})}),$$

494 and for any $f \in \mathcal{F}$ and $z \in \mathcal{X} \times \mathcal{Y}$, $f(z) \in [0, M]$. Let

$$t = \log \frac{2N_\epsilon + 2}{\delta}, \quad \text{and} \quad \lambda = \sqrt{\frac{2t}{m-1}}.$$

495 Then, with probability at least $1 - \delta$,

$$\begin{aligned}R(\hat{f}) - R(f^*) &\leq 2\lambda \sqrt{\frac{(m-1)V(f^*)}{m}} + \sum_{i=1}^n \lambda \sqrt{\frac{m_i \hat{V}_{in}(\hat{f}|e_i)}{m}} \\ &\quad + \left(2 + \lambda + \sum_{i=1}^n \lambda \sqrt{\frac{m_i}{m}}\right) \epsilon + \lambda^2 \frac{4(4m-1)M}{3m}.\end{aligned}$$

496 **Proof:** According to the proof of Theorem 1,

$$R(f) \leq \hat{R}(f) + \sqrt{\frac{2\hat{V}_{out}(f)t}{m-1}} + \sum_{i=1}^n \sqrt{\frac{2m_i\hat{V}_{in}(f|e_i)t}{m(m-1)}} + \frac{2(4m-1)Mt}{3m(m-1)},$$

497 holds with probability larger than $1 - 2\exp(-t)$. Next, we consider a set of functions $\{f^1, \dots, f^{N_\epsilon}\}$,
 498 which is a minimal ϵ -cover of the function space \mathcal{F} of size $N_\epsilon = N(\mathcal{F}, \epsilon, \|\cdot\|_{L^\infty(\mathcal{X} \times \mathcal{Y})})$. Then, for
 499 any $f \in \mathcal{F}$, there exists f^j , $1 \leq j \leq N_\epsilon$ such that $\|f - f^j\|_{L^\infty(\mathcal{X} \times \mathcal{Y})} \leq \epsilon$. Therefore,

$$\begin{aligned} R(f) &\leq R(f^j) + \epsilon \\ &\leq \hat{R}(f^j) + \sqrt{\frac{2\hat{V}_{out}(f^j)t}{m-1}} + \sum_{i=1}^n \sqrt{\frac{2m_i\hat{V}_{in}(f^j|e_i)t}{m(m-1)}} \\ &\quad + \frac{2(4m-1)Mt}{3m(m-1)} + \epsilon, \end{aligned}$$

500 with probability larger than $1 - (2+n)\exp(-t)$. Notice that $\hat{R}(f^j) \leq \hat{R}(f) + \epsilon$ and

$$\begin{aligned} \sqrt{\hat{V}_{out}(f^j)} &\leq \sqrt{\hat{V}_{out}(f)} + \sqrt{\hat{V}_{out}(f^j - f)} \leq \sqrt{\hat{V}_{out}(f)} + \epsilon, \\ \sqrt{\hat{V}_{in}(f^j|e_i)} &\leq \sqrt{\hat{V}_{in}(f|e_i)} + \sqrt{\hat{V}_{in}(f^j - f|e_i)} \leq \sqrt{\hat{V}_{in}(f|e_i)} + \epsilon. \end{aligned}$$

501 Therefore, for every $f \in \mathcal{F}$,

$$\begin{aligned} R(f) &\leq \hat{R}(f) + \sqrt{\frac{2\hat{V}_{out}(f)t}{m-1}} + \sum_{i=1}^n \sqrt{\frac{2m_i\hat{V}_{in}(f|e_i)t}{m(m-1)}} \\ &\quad + \frac{2(4m-1)Mt}{3m(m-1)} + \left(2 + \sqrt{\frac{2t}{m-1}} + \sum_{i=1}^n \sqrt{\frac{2m_it}{m(m-1)}}\right)\epsilon, \end{aligned} \tag{21}$$

502 holds with probability at least $1 - 2N_\epsilon \exp(-t)$. Since

$$\hat{f} = \arg \min_{f \in \mathcal{F}} \hat{R}(f) + \sqrt{\frac{2\hat{V}_{out}(f)t}{m-1}},$$

503 then we have

$$\hat{R}(\hat{f}) + \sqrt{\frac{2\hat{V}_{out}(\hat{f})t}{m-1}} \leq \hat{R}(f^*) + \sqrt{\frac{2\hat{V}_{out}(f^*)t}{m-1}}.$$

504 By the Bernstein inequality, with probability at least $1 - \exp(-t)$,

$$\hat{R}(f^*) \leq R(f^*) + \sqrt{\frac{2V(f^*)t}{m}} + \frac{2Mt}{3m}. \tag{22}$$

505 By Theorem 10 in [27],

$$\sqrt{\hat{V}_{out}(f^*)} \leq \sqrt{\hat{V}(f^*)} \leq \sqrt{\frac{m-1}{m}V(f^*)} + \sqrt{\frac{2M^2t}{m-1}}, \tag{23}$$

506 holds with probability larger than $1 - \exp(-t)$. Combining (21), (22) and (23),

$$\begin{aligned} R(\hat{f}) &\leq R(f^*) + 2\sqrt{\frac{2V(f^*)t}{m}} + \sum_{i=1}^n \sqrt{\frac{2m_i\hat{V}_{in}(\hat{f}|e_i)t}{m(m-1)}} \\ &\quad + \frac{2Mt}{3m} + \frac{2Mt}{m-1} + \frac{2(4m-1)Mt}{3m(m-1)} \\ &\quad + \left(2 + \sqrt{\frac{2t}{m-1}} + \sum_{i=1}^n \sqrt{\frac{2m_it}{m(m-1)}}\right)\epsilon, \end{aligned}$$

507 holds with probability larger than $1 - (2N_\epsilon + 2)\exp(-t)$.

508 □

509 **E Proof of Theorem 3**

510 **Theorem** Suppose the training dataset \mathcal{D} and a function $f \in \mathcal{F}$ are given. Let ρ_+ be the largest
 511 distance between $\mathbf{q} \in \mathcal{Q}_\alpha(\hat{\mathbf{q}}, +\infty)$ and

$$\rho_- = \frac{\min_i (\alpha/\hat{q}_i - 1)^2 \hat{V}_{out}(f)}{2(\min_i \hat{R}(f|e_i) - \hat{R}(f))^2}.$$

512 Then we have

$$\mathcal{L}(f; \rho_-) \leq \max_{\mathbf{q} \in \mathcal{Q}_{\alpha, +\infty}(\hat{\mathbf{q}})} \sum_{i=1}^n q_i \hat{R}(f|e_i) \leq \mathcal{L}(f; \rho_+).$$

513 In this section, we decompose the complete proof into the three steps.

514 **Step 1.** Given a training dataset and a function $f \in \mathcal{F}$, the following inequality always holds:

$$\max_{\mathbf{q} \in \mathcal{Q}_{\alpha, \rho}(\hat{\mathbf{q}})} \sum_{i=1}^n q_i \hat{R}(f|e_i) \leq \hat{R}(f) + \sqrt{2\rho \hat{V}_{out}(f)}. \quad (24)$$

515 **Proof:** Since $\sum_{i=1}^n q_i = 1$ and $\sum_{i=1}^n \hat{q}_i = 1$, we have

$$\begin{aligned} \sum_{i=1}^n q_i \hat{R}(f|e_i) &= \sum_{i=1}^n (q_i - \hat{q}_i) \hat{R}(f|e_i) + \sum_{i=1}^n \hat{q}_i \hat{R}(f|e_i) \\ &= \hat{R}(f) + \sum_{i=1}^n (q_i - \hat{q}_i) \hat{R}(f|e_i). \end{aligned}$$

516 Since $\sum_{i=1}^n (q_i - \hat{q}_i)C = 0$ holds for any constant C ,

$$\sum_{i=1}^n (q_i - \hat{q}_i) \hat{R}(f|e_i) = \sum_{i=1}^n (q_i - \hat{q}_i) (\hat{R}(f|e_i) - \hat{R}(f)).$$

517 Thus the max problem in (24) is equivalent to maximize

$$\max_{\mathbf{q} \in \mathcal{Q}_{\alpha, \rho}(\hat{\mathbf{q}})} \sum_{i=1}^n (q_i - \hat{q}_i) (\hat{R}(f|e_i) - \hat{R}(f)).$$

518 By the Cauchy-Schwarz inequality,

$$\begin{aligned} &\sum_{i=1}^n (q_i - \hat{q}_i) (\hat{R}(f|e_i) - \hat{R}(f)) \\ &= \sum_{i=1}^n \frac{q_i - \hat{q}_i}{\sqrt{\hat{q}_i}} \sqrt{\hat{q}_i} (\hat{R}(f|e_i) - \hat{R}(f)) \\ &\leq \sqrt{\sum_{i=1}^n \frac{(q_i - \hat{q}_i)^2}{\hat{q}_i}} \times \sqrt{\sum_{i=1}^n \hat{q}_i (\hat{R}(f|e_i) - \hat{R}(f))^2} \\ &\leq \sqrt{2\rho} \times \sqrt{\hat{V}_{out}(f)} \end{aligned}$$

519 □

520 **Step 2.** If the between-domain variance $\hat{V}_{out}(f)$ is non-zero, the equality holds if and only if
 521 $\forall e_i \in \mathcal{E}_{tr}$,

$$\alpha \leq \hat{q}_i \left(\sqrt{\frac{2\rho}{\hat{V}_{out}(f)}} (\hat{R}(f|e_i) - \hat{R}(f)) + 1 \right).$$

522 On the other hand, if α is fixed, the equality holds when the radius of $\mathcal{Q}_{\alpha,\rho}(\hat{\mathbf{q}})$ satisfies

$$\rho \leq \frac{(\alpha/\hat{q}_i - 1)^2 \hat{V}_{out}(f)}{2(\min_i \hat{R}(f|e_i) - \hat{R}(f))^2}.$$

523 **Proof:** The equality in (24) is attained if and only if the following requirements hold at the same time:

524 (i) There exists a constant c such that $\forall 1 \leq i \leq n$,

$$\frac{q_i - \hat{q}_i}{\sqrt{\hat{q}_i}} = c\sqrt{\hat{q}_i}(\hat{R}(f|e_i) - \hat{R}(f)).$$

525 (ii) The χ^2 divergence between \mathbf{q} and $\hat{\mathbf{q}}$ achieves ρ , that is

$$\sum_{i=1}^n \frac{(q_i - \hat{q}_i)^2}{\hat{q}_i} = 2\rho.$$

526 It is easy to see

$$c^2 \sum_{i=1}^n \hat{q}_i (\hat{R}(f|e_i) - \hat{R}(f))^2 = 2\rho \quad \Rightarrow \quad c = \sqrt{\frac{2\rho}{\hat{V}_{out}(f)}}.$$

527 Then the discrete distribution \mathbf{q} satisfies (i) is

$$q_i = \sqrt{\frac{2\rho}{\hat{V}_{out}(f)}} \hat{q}_i (\hat{R}(f|e_i) - \hat{R}(f)) + \hat{q}_i.$$

528 Since \mathbf{q} belongs to $\mathcal{Q}_{\alpha,\rho}(\hat{\mathbf{q}})$, the only constraint here is $q_i \geq \alpha$, $\forall e_i$ which holds if and only if
529 $\forall e_i \in \mathcal{E}_{tr}$,

$$\alpha \leq \hat{q}_i \left(\sqrt{\frac{2\rho}{\hat{V}_{out}(f)}} (\hat{R}(f|e_i) - \hat{R}(f)) + 1 \right).$$

530 On the other hand, if α is fixed and non-positive, the constraint implies that the radius ρ of $\mathcal{Q}_{\alpha,\rho}(\hat{\mathbf{q}})$
531 should be sufficiently small:

$$\begin{aligned} \alpha &\leq \hat{q}_i \left(\sqrt{\frac{2\rho}{\hat{V}_{out}(f)}} (\hat{R}(f|e_i) - \hat{R}(f)) + 1 \right), \\ \Leftrightarrow \frac{\alpha}{\hat{q}_i} - 1 &\leq \sqrt{\frac{2\rho}{\hat{V}_{out}(f)}} (\hat{R}(f|e_i) - \hat{R}(f)), \\ \Leftrightarrow \frac{\alpha/\hat{q}_i - 1}{\min_i (\hat{R}(f|e_i) - \hat{R}(f))} &\geq \sqrt{\frac{2\rho}{\hat{V}_{out}(f)}}, \\ \Leftrightarrow \rho &\leq \frac{\min_i (\alpha/\hat{q}_i - 1)^2 \hat{V}_{out}(f)}{2(\min_i \hat{R}(f|e_i) - \hat{R}(f))^2}. \end{aligned}$$

532 Hence the proof is finished.

533

□

534 **Step 3.** Proof of (13) in Theorem 3.

535 **Proof:** First, ρ_+ is the largest distance between $\hat{\mathbf{q}}$ and $\mathbf{q} \in \mathcal{Q}_{\alpha,+\infty}(\hat{\mathbf{q}})$. Therefore,

$$\mathcal{Q}_{\alpha,+\infty}(\hat{\mathbf{q}}) = \mathcal{Q}_{\alpha,\rho_+}(\hat{\mathbf{q}}) \subseteq \mathcal{Q}_{+\infty,\rho_+}(\hat{\mathbf{q}}).$$

536 Hence the second inequality in (13) is trivial. Let \mathbf{q}^* be the solution:

$$\mathbf{q}^* = \arg \max_{\mathbf{q} \in \mathcal{Q}_{+\infty,\rho_-}(\hat{\mathbf{q}})} \sum_{i=1}^n q_i \hat{R}(f|e_i),$$

537 According to **Step 2**,

$$\max_{\mathbf{q} \in \mathcal{Q}_{+\infty, \rho_-}(\hat{\mathbf{q}})} \sum_{i=1}^n q_i \hat{R}(f|e_i) = \mathcal{L}(f; \rho_-),$$

538 and

$$\mathbf{q}_-^* \in \mathcal{Q}_{+\infty, \rho_-}(\hat{\mathbf{q}}) = \mathcal{Q}_{\alpha, \rho_-}(\hat{\mathbf{q}}) \subseteq \mathcal{Q}_{\alpha, +\infty}(\hat{\mathbf{q}}).$$

539 Hence the first inequality in (13) is proved. □

540

541 **F Proof of Theorem 4**

542 In this section, we start with a preliminary result.

543 **Theorem 9** *Suppose that α is a non-positive scalar and*

$$V'_{out}(f) = V(f) - \sum_{i=1}^n q_i V_{in}(f|e_i) > 0.$$

544 *For each training domain, both \hat{q}_i and q_i are larger than $\delta > 0$. Let*

$$\rho' = \frac{V'_{out}(f)}{8M^2} \left(\frac{\alpha}{1 - (n-1)\delta} - 1 \right)^2.$$

545 *If $n > 2$ and m is sufficiently large such that*

$$\frac{m}{4} V'_{out}(f) > V(f),$$

546 *then, given $f \in \mathcal{F}$, the following expansion uniformly holds:*

$$\max_{\mathbf{q} \in \mathcal{Q}_{\alpha, \rho'}(\hat{\mathbf{q}})} \sum_{i=1}^n q_i \hat{R}(f|e_i) = \mathcal{L}(f; \rho'),$$

547 *with probability at least*

$$1 - \exp\left(-\frac{\left(\frac{m}{4}V'_{out}(f) - V(f)\right)^2}{2M^2(m-1)V(f)}\right) - \binom{m+n-1}{n-1} \exp\left(-\frac{mV'_{out}(f)^2}{M^4}\right) \\ - \sum_{i=1}^n \exp\left(-\frac{(V_{in}(f|e_i) + \frac{m_i}{4}V'_{out}(f))^2}{2M^2m_i(V_{in}(f|e_i) + \frac{1}{4}V'_{out}(f))^2}\right).$$

548 **Proof:** Note that $\hat{R}(f|e) \in [0, M]$ for any $f \in \mathcal{F}$ and $e \in \mathcal{E}_{tr}$. Thus for any training data \mathbb{D} ,

$$\left(\min_i \hat{R}(f|e_i) - \hat{R}(f)\right)^2 \leq M^2.$$

549 In addition, since $\alpha \leq 0$ and $q_i, \hat{q}_i \geq \delta$,

$$\min_i \left(\frac{\alpha}{\hat{q}_i} - 1\right)^2 \leq \left(\frac{\alpha}{1 - (n-1)\delta} - 1\right)^2.$$

550 Hence, to satisfying

$$\rho' \leq \frac{\min_i (\alpha/\hat{q}_i - 1)^2 \hat{V}_{out}(f)}{2(\min_i \hat{R}(f|e_i) - \hat{R}(f))^2},$$

551 it suffices to show that

$$\hat{V}_{out}(f) \geq \frac{1}{4} V'_{out}(f).$$

552 Notice that

$$\hat{V}_{out}(f) = \hat{V}(f) - \sum_{i=1}^n \hat{q}_i \hat{V}_{in}(f|e_i) = V'_{out}(f) + I_1 + I_2 + I_3,$$

553 where

$$\begin{aligned} I_1 &= \hat{V}(f) - V(f), \\ I_2 &= - \sum_{i=1}^n (\hat{q}_i - q_i) \hat{V}_{in}(f|e_i), \\ I_3 &= - \sum_{i=1}^n q_i (\hat{V}_{in}(f|e_i) - V_{in}(f|e_i)). \end{aligned}$$

554 By Theorem 10 in [27], for any $\delta > 0$,

$$P\left(\frac{m}{m-1} \hat{V}(f) < V(f) - \delta\right) \leq \exp\left(-\frac{(m-1)\delta^2}{2M^2V(f)}\right).$$

555 We take

$$\delta = \frac{m}{4(m-1)} V'_{out}(f) - \frac{1}{m-1} V(f).$$

556 Therefore,

$$\begin{aligned} P\left(I_1 < -\frac{1}{4} V'_{out}(f)\right) &= P\left(\hat{V}(f) < V(f) - \frac{1}{4} V'_{out}(f)\right) \\ &\leq \exp\left(-\frac{\left(\frac{m}{4} V'_{out}(f) - V(f)\right)^2}{2M^2(m-1)V(f)}\right). \end{aligned}$$

557 For the term I_2 ,

$$|I_2| \leq \sum_{i=1}^n |\hat{q}_i - q_i| \hat{V}_{in}(f|e_i) \leq \|\hat{\mathbf{q}} - \mathbf{q}\|_1 \frac{M^2}{4}.$$

558 According to the Pinsker's inequality and the method of types [15], for any $\delta > 0$,

$$\|\hat{\mathbf{q}} - \mathbf{q}\|_1 \leq \sqrt{2D_{KL}(\hat{\mathbf{q}}\|\mathbf{q})} < \delta$$

559 with probability at least

$$1 - \binom{m+n-1}{n-1} \exp(-m\delta^2).$$

560 We take $\delta = V'_{out}(f)/M^2$. Then we know

$$P\left(I_2 < -\frac{1}{4} V'_{out}(f)\right) \leq \binom{m+n-1}{n-1} \exp\left(-\frac{mV'_{out}(f)^2}{M^4}\right).$$

561 Next, we deal with I_3 . By Theorem 10 in [27], for any $\delta > 0$ and any $e_i \in \mathcal{E}_{tr}$,

$$P\left(\frac{m}{m-1} \hat{V}_{in}(f|e_i) > V_{in}(f|e_i) + \delta\right) \leq \exp\left(-\frac{(m_i-1)\delta^2}{2M^2(V_{in}(f|e_i) + \delta)}\right).$$

562 Let

$$\delta = \frac{1}{m_i-1} V_{in}(f|e_i) + \frac{m_i}{4(m_i-1)} V'_{out}(f).$$

563 Then we have

$$\begin{aligned} P\left(\hat{V}_{in}(f|e_i) > V_{in}(f|e_i) + \frac{1}{4} V'_{out}(f)\right) &= P\left(\frac{m}{m-1} \hat{V}_{in}(f|e_i) > V_{in}(f|e_i) + \delta\right) \\ &\leq \exp\left(-\frac{(V_{in}(f|e_i) + \frac{m_i}{4} V'_{out}(f))^2}{2M^2 m_i (V_{in}(f|e_i) + \frac{1}{4} V'_{out}(f))^2}\right). \end{aligned}$$

564 Therefore $I_3 < -\frac{1}{4}V'_{out}(f)$ with probability smaller than

$$\sum_{i=1}^n \exp\left(-\frac{(V_{in}(f|e_i) + \frac{m_i}{4}V'_{out}(f))^2}{2M^2m_i(V_{in}(f|e_i) + \frac{1}{4}V'_{out}(f))^2}\right)$$

565 Combining the results of I_1 , I_2 and I_3 , we have

$$\hat{V}_{out}(f) \geq \frac{1}{4}V'_{out}(f),$$

566 with probability larger than

$$1 - \exp\left(-\frac{(\frac{m}{4}V'_{out}(f) - V(f))^2}{2M^2(m-1)V(f)}\right) - \binom{m+n-1}{n-1} \exp\left(-\frac{mV'_{out}(f)^2}{M^4}\right) \\ - \sum_{i=1}^n \exp\left(-\frac{(V_{in}(f|e_i) + \frac{m_i}{4}V'_{out}(f))^2}{2M^2m_i(V_{in}(f|e_i) + \frac{1}{4}V'_{out}(f))^2}\right).$$

567 Hence the proof is finished.

568

□

569 Next, we extend Theorem 9 to a more general variant with respect to the family of functions \mathcal{F} .

570 **Theorem.** Suppose that α is a non-positive scalar and $V'_{out}(f) = V(f) - \sum_i q_i V_{in}(f|e_i) > 0$. For
571 each training domain, both \hat{q}_i and q_i are larger than $\delta > 0$. We write

$$\rho' = \frac{V'_{out}(f)}{16M^2} \left(\frac{\alpha}{1 - (n-1)\delta} - 1 \right)^2.$$

572 Let $\tau > 0$ and $0 < \eta < 1$ be two constants. Define

$$\mathcal{F}_{\tau,\eta} = \{f \in \mathcal{F} : V(f) \geq \tau, \text{ and } \frac{V_{in}(f|e_i)}{V(f)} \leq \eta, \forall e_i \in \mathcal{E}_{tr}\}.$$

573 For any $f \in \mathcal{F}_{\tau,\eta}$, the following expansion uniformly holds:

$$\max_{\mathbf{q} \in \mathcal{Q}_{\alpha,\rho'}(\hat{\mathbf{q}})} \sum_{i=1}^n q_i \hat{R}(f|e_i) = \mathcal{L}(f; \rho'),$$

574 with probability at least $1 - N_{\tau,\eta} \times p$, where $N_{\tau,\eta} = N\left(\mathcal{F}_{\tau,\eta}, \sqrt{\frac{1}{10}(1-\eta)\tau}, \|\cdot\|_{L^\infty(\mathcal{X} \times \mathcal{Y})}\right)$ is the
575 covering number of $\mathcal{F}_{\tau,\eta}$ and

$$p = \exp\left(-\frac{m(1-\eta)^2\tau}{32M^2} + \frac{1}{16}\right) + \binom{m+n-1}{n-1} \exp\left(-\frac{m(1-\eta)^2\tau^2}{M^4}\right) \\ + \sum_{i=1}^n \exp\left\{-\frac{1}{2M^2m_i} \left(\frac{m_i(1-\eta) + 4\eta}{1+3\eta}\right)^2\right\}.$$

576 **Proof:** We consider a set of functions $\{f^1, \dots, f^N\}$, which is a minimal ϵ -cover of $\mathcal{F}_{\tau,\eta}$ of size

$$N_{\tau,\eta} = N(\mathcal{F}_{\tau,\eta}, \epsilon, \|\cdot\|_{L^\infty(\mathcal{X} \times \mathcal{Y})}).$$

577 Define the event

$$\mathfrak{E} = \left\{ \hat{V}_{out}(f^i) \geq \frac{1}{4}V'_{out}(f^i), \text{ for } i = 1, \dots, N \right\}.$$

578 Recall the proof of Theorem 3 and the terms I_1 to I_3 . For any $f \in \mathcal{F}_{\tau,\eta}$,

$$P\left(I_1 < -\frac{1}{4}V'_{out}(f)\right) \leq \exp\left(-\frac{(\frac{m}{4}V'_{out}(f) - V(f))^2}{2M^2(m-1)V(f)}\right) \\ = \exp\left(-\frac{m^2V'_{out}(f)^2}{32M^2(m-1)V(f)} + \frac{mV'_{out}(f)}{4M^2(m-1)} - \frac{V(f)}{2M^2(m-1)}\right) \\ \leq \exp\left(-\frac{m^2V'_{out}(f)^2}{32M^2(m-1)V(f)} + \frac{1}{16}\right) \\ \leq \exp\left(-\frac{m(1-\eta)^2\tau}{32M^2} + \frac{1}{16}\right).$$

579 For the term I_2 ,

$$\begin{aligned} P\left(I_2 < -\frac{1}{4}V'_{out}(f)\right) &\leq \binom{m+n-1}{n-1} \exp\left(-\frac{mV'_{out}(f)^2}{M^4}\right) \\ &\leq \binom{m+n-1}{n-1} \exp\left(-\frac{m(1-\eta)^2\tau^2}{M^4}\right). \end{aligned}$$

580 For any $e_i \in \mathcal{E}_{tr}$,

$$\begin{aligned} -\frac{(V_{in}(f|e_i) + \frac{m_i}{4}V'_{out}(f))^2}{2M^2m_i(V_{in}(f|e_i) + \frac{1}{4}V'_{out}(f))^2} &= -\frac{1}{2M^2m_i} \left(\frac{\frac{m_i-1}{4}V'_{out}(f)}{V_{in}(f|e_i) + \frac{1}{4}V'_{out}(f)} + 1 \right)^2 \\ &= -\frac{1}{2M^2m_i} \left(\frac{m_i-1}{4\frac{V_{in}(f|e_i)}{V'_{out}(f)} + 1} + 1 \right)^2 \\ &\leq -\frac{1}{2M^2m_i} \left(\frac{m_i-1}{4\frac{\eta}{1-\eta} + 1} + 1 \right)^2 \\ &= -\frac{1}{2M^2m_i} \left(\frac{(m_i-1)(1-\eta)}{1+3\eta} + 1 \right)^2 \\ &= -\frac{1}{2M^2m_i} \left(\frac{m_i(1-\eta) + 4\eta}{1+3\eta} \right)^2 \end{aligned}$$

581 Therefore, for the term I_3 ,

$$P\left(I_3 < -\frac{1}{4}V'_{out}(f)\right) \leq \sum_{i=1}^n \exp\left\{-\frac{1}{2M^2m_i} \left(\frac{m_i(1-\eta) + 4\eta}{1+3\eta} \right)^2\right\}.$$

582 Furthermore, we have for any $f \in \mathcal{F}_{\tau,\eta}$,

$$\hat{V}_{out}(f) \geq \frac{1}{4}V'_{out}(f),$$

583 with probability larger than $1-p$, where

$$\begin{aligned} p &= \exp\left(-\frac{m(1-\eta)^2\tau}{32M^2} + \frac{1}{16}\right) + \binom{m+n-1}{n-1} \exp\left(-\frac{m(1-\eta)^2\tau^2}{M^4}\right) \\ &\quad + \sum_{i=1}^n \exp\left\{-\frac{1}{2M^2m_i} \left(\frac{m_i(1-\eta) + 4\eta}{1+3\eta} \right)^2\right\}. \end{aligned}$$

584 Then, combining the results of $\{f^j, j = 1, \dots, N\}$,

$$P(\mathfrak{E}) \geq 1 - N(\mathcal{F}_{\tau,\eta}, \epsilon, \|\cdot\|_{L^\infty(\mathcal{X} \times \mathcal{Y})}) \times p.$$

Given the event \mathfrak{E} , for any $f \in \mathcal{F}$, there exists f^j such that

$$\|f - f^j\|_{L^\infty(\mathcal{X} \times \mathcal{Y})} \leq \epsilon.$$

585 Hence,

$$\begin{aligned} \hat{V}_{out}(f) &\geq \hat{V}_{out}(f^j) - \epsilon^2 \\ &\geq \frac{1}{4}V'_{out}(f^j) - \epsilon^2 \geq \frac{1}{4}V'_{out}(f) - \frac{5}{4}\epsilon^2. \end{aligned}$$

586 We take

$$\epsilon = \sqrt{\frac{1}{10}(1-\eta)\tau}.$$

587 Then,

$$\hat{V}_{out}(f) \geq \frac{1}{4}V'_{out}(f) - \frac{5}{4}\epsilon^2 \geq \frac{1}{8}V'_{out}(f).$$

588 Hence the proof is finished.

589

□

590 **G Proof of Theorem 5**

591 **Theorem.** Given the training dataset and a function $f \in \mathcal{F}$, then for any anchor distribution \mathbf{q}_0 , the
 592 inequality always holds:

$$\max_{\mathbf{q} \in \mathcal{Q}_\alpha(\mathbf{q}_0, \rho)} \sum_{i=1}^n q_i \hat{R}(f|e_i) \leq \hat{R}(f, \mathbf{q}_0) + \sqrt{2\rho \hat{V}_{out}(f, \mathbf{q}_0)},$$

593 where

$$\begin{aligned} \hat{R}(f, \mathbf{q}_0) &= \frac{1}{n} \sum_{i=1}^n q_{0,i} \hat{R}(f|e_i), \\ \hat{V}_{out}(f, \mathbf{q}_0) &= \sum_{i=1}^n q_{0,i} (\hat{R}(f|e_i) - \hat{R}(f))^2. \end{aligned}$$

594 If the between-domain variance $\hat{V}_{out}(f, \mathbf{q}_0)$ is non-zero, the equality holds if and only if $\forall e_i \in \mathcal{E}_{tr}$,

$$\alpha \leq q_{0,i} \left(\sqrt{\frac{2\rho}{\hat{V}_{out}(f, \mathbf{q}_0)}} (\hat{R}(f|e_i) - \hat{R}(f, \mathbf{q}_0)) + 1 \right).$$

595 On the other hand, if α is fixed, the equality holds when the radius of $\mathcal{Q}_\alpha(\mathbf{q}_0, \rho)$ satisfies,

$$\rho \leq \frac{\min_i (\alpha/q_{0,i} - 1)^2 \hat{V}_{out}(f, \mathbf{q}_0)}{2(\min_i \hat{R}(f|e_i) - \hat{R}(f, \mathbf{q}_0))^2}.$$

596 The proof here is similar to that of Theorem 3.

597 **Proof:** Since $\sum_{i=1}^n q_i = 1$ and $\sum_{i=1}^n q_{0,i} = 1$, we have

$$\begin{aligned} \sum_{i=1}^n q_i \hat{R}(f|e_i) &= \sum_{i=1}^n q_{0,i} \hat{R}(f|e_i) + \sum_{i=1}^n (q_i - q_{0,i}) \hat{R}(f|e_i) \\ &= \hat{R}(f, \mathbf{q}_0) + \sum_{i=1}^n (q_i - q_{0,i}) \hat{R}(f|e_i). \end{aligned}$$

598 Since $\sum_{i=1}^n (q_i - \hat{q}_i) = 0$, then we have

$$\sum_{i=1}^n (q_i - q_{0,i}) \hat{R}(f|e_i) = \sum_{i=1}^n (q_i - q_{0,i}) (\hat{R}(f|e_i) - \hat{R}(f, \mathbf{q}_0)).$$

599 Thus the max problem with respect to \mathbf{q} is equivalent to maximize

$$\max_{\mathbf{q} \in \mathcal{Q}_\alpha(\mathbf{q}_0, \rho)} \sum_{i=1}^n (q_i - q_{0,i}) (\hat{R}(f|e_i) - \hat{R}(f, \mathbf{q}_0)).$$

600 By the Cauchy-Schwarz inequality,

$$\begin{aligned} &\sum_{i=1}^n (q_i - q_{0,i}) (\hat{R}(f|e_i) - \hat{R}(f, \mathbf{q}_0)) \\ &= \sum_{i=1}^n \frac{q_i - q_{0,i}}{\sqrt{q_{0,i}}} \sqrt{q_{0,i}} (\hat{R}(f|e_i) - \hat{R}(f, \mathbf{q}_0)) \\ &\leq \sqrt{\sum_{i=1}^n \frac{(q_i - q_{0,i})^2}{q_{0,i}}} \times \sqrt{\sum_{i=1}^n q_{0,i} (\hat{R}(f|e_i) - \hat{R}(f, \mathbf{q}_0))^2} \\ &\leq \sqrt{2\rho} \times \sqrt{\hat{V}_{out}(f, \mathbf{q}_0)} \end{aligned}$$

601 The equality is attained if and only if the following requirements hold at the same time:

602 (i) There exists a constant c such that $\forall 1 \leq i \leq n$,

$$\frac{q_i - q_{0,i}}{\sqrt{q_{0,i}}} = c\sqrt{q_{0,i}}(\hat{R}(f|e_i) - \hat{R}(f, \mathbf{q}_0)).$$

603 (ii) The χ^2 divergence between \mathbf{q} and \mathbf{q}_0 achieves ρ :

$$\sum_{i=1}^n \frac{(q_i - q_{0,i})^2}{q_{0,i}} = 2\rho.$$

604 It is easy to see

$$c^2 \sum_{i=1}^n q_{0,i} (\hat{R}(f|e_i) - \hat{R}(f, \mathbf{q}_0))^2 = 2\rho \Rightarrow c = \sqrt{\frac{2\rho}{\hat{V}_{out}(f, \mathbf{q}_0)}}.$$

605 Then the discrete distribution \mathbf{q} satisfies (i) is

$$q_i = \sqrt{\frac{2\rho}{\hat{V}_{out}(f, \mathbf{q}_0)}} q_{0,i} (\hat{R}(f|e_i) - \hat{R}(f, \mathbf{q}_0)) + q_{0,i}.$$

606 Since \mathbf{q} belongs to $\mathcal{Q}_\alpha(\mathbf{q}_0, \rho)$, the constraint here is $q_i \geq \alpha, \forall e_i$ which holds if and only if

$$\alpha \leq q_{0,i} \left(\sqrt{\frac{2\rho}{\hat{V}_{out}(f)}} (\hat{R}(f|e_i) - \hat{R}(f)) + 1 \right), \forall e_i \in \mathcal{E}_{tr}.$$

607 On the other hand, if α is fixed and non-positive, the constraint implies that the radius ρ of $\mathcal{Q}_\alpha(\mathbf{q}_0, \rho)$
608 should be sufficiently small:

$$\begin{aligned} \alpha &\leq q_{0,i} \left(\sqrt{\frac{2\rho}{\hat{V}_{out}(f, \mathbf{q}_0)}} (\hat{R}(f|e_i) - \hat{R}(f, \mathbf{q}_0)) + 1 \right) \\ \Leftrightarrow \frac{\alpha}{q_{0,i}} - 1 &\leq \sqrt{\frac{2\rho}{\hat{V}_{out}(f, \mathbf{q}_0)}} (\hat{R}(f|e_i) - \hat{R}(f, \mathbf{q}_0)), \\ \Leftrightarrow \frac{\alpha/q_{0,i} - 1}{\min_i (\hat{R}(f|e_i) - \hat{R}(f, \mathbf{q}_0))} &\geq \sqrt{\frac{2\rho}{\hat{V}_{out}(f, \mathbf{q}_0)}}, \\ \Leftrightarrow \rho &\leq \frac{\min_i (\alpha/q_{0,i} - 1)^2 \hat{V}_{out}(f, \mathbf{q}_0)}{2(\min_i \hat{R}(f|e_i) - \hat{R}(f, \mathbf{q}_0))^2}. \end{aligned}$$

609 Hence the proof is finished.

610

□