# Dropout: Explicit Forms and Capacity Control

Raman Arora\* Peter Bartlett† Poorya Mianjy‡ Nathan Srebro§

March 10, 2020

#### Abstract

We investigate the capacity control provided by dropout in various machine learning problems. First, we study dropout for matrix completion, where it induces a data-dependent regularizer that, in expectation, equals the weighted trace-norm of the product of the factors. In deep learning, we show that the data-dependent regularizer due to dropout directly controls the Rademacher complexity of the underlying class of deep neural networks. These developments enable us to give concrete generalization error bounds for the dropout algorithm in both matrix completion as well as training deep neural networks. We evaluate our theoretical findings on real-world datasets, including MovieLens, MNIST, and Fashion-MNIST.

## 1 Introduction

Dropout is a popular algorithmic regularization technique for training deep neural networks that aims at "breaking co-adaptation" among neurons by randomly dropping them at training time (Hinton et al., 2012). Dropout has been shown effective across a wide range of machine learning tasks, from classification (Srivastava et al., 2014; Szegedy et al., 2015) to regression (Toshev & Szegedy, 2014). Notably, dropout is considered an essential component in the design of AlexNet (Krizhevsky et al., 2012), which won the prominent ImageNet challenge in 2012 with a significant margin and helped transform the field of computer vision.

Dropout regularizes the empirical risk by randomly perturbing the model parameters during training. A natural first step toward understanding generalization due to dropout, therefore, is to instantiate the explicit form of the regularizer due to dropout. In linear regression, with dropout applied to the input layer (i.e., on the input features), the explicit regularizer was shown to be akin to a data-dependent ridge penalty Srivastava et al. (2014); Wager et al. (2013); Baldi & Sadowski (2013); Wang & Manning (2013). In factored models dropout yields more exotic forms of regularization. For instance, dropout induces regularizer that behaves similar to nuclear norm regularization in matrix factorization Cavazza et al. (2018), in single hidden-layer linear networks Mianjy et al. (2018), and in deep linear networks Mianjy & Arora (2019). However, none of the works above discuss how the induced regularizer provides capacity control, or equivalently, help us establish generalization bounds for dropout.

In this paper, we provide an answer to this question. We give *explicit forms* of the regularizers induced by dropout for the matrix sensing problem and two-layer neural networks with ReLU activations. Further, we establish *capacity control* due to dropout and give precise generalization bounds. Our key contributions are as follows.

1. In Section 2, we study dropout for matrix completion, wherein, the matrix factors are dropped randomly during training. We show that this algorithmic procedure induces a data-dependent regularizer that behaves similar to the weighted trace-norm which has been shown to yield strong generalization guarantees for matrix completion (Foygel et al., 2011).

<sup>\*</sup>Johns Hopkins University, email: arora@cs.jhu.edu

<sup>†</sup>University of California, Berkeley, email: bartlett@cs.berkeley.edu

<sup>&</sup>lt;sup>‡</sup>Johns Hopkins University, email: mianjy@jhu.edu

<sup>§</sup>TTI Chicago, email: nati@ttic.edu

- 2. In Section 3, we study dropout in two-layer ReLU networks. We show that the regularizer induced by dropout is a data-dependent measure that behaves as  $\ell_2$ -path norm Neyshabur et al. (2015a), and establish data-dependent generalization bounds.
- 3. In Section 5, we present empirical evaluations that confirm our theoretical findings for matrix completion and deep regression on real world datasets including the MovieLens data, as well as the MNIST and Fashion MNIST datasets.

#### 1.1 Related Work

Dropout was first introduced by Hinton et al. (2012) as an effective heuristic for algorithmic regularization, yielding lower test errors on the MNIST and TIMIT datasets. In a subsequent work, Srivastava et al. (2014) reported similar improvements over several tasks in computer vision (on CIFAR-10/100 and ImageNet datasets), speech recognition, text classification and genetics.

Thenceforth, dropout has been widely used in training state-of-the-art systems for several tasks including large-scale visual recognition Szegedy et al. (2015), large vocabulary continuous speech recognition Dahl et al. (2013), image question answering Yang et al. (2016), handwriting recognition Pham et al. (2014), sentiment prediction and question classification Kalchbrenner et al. (2014), dependency parsing Chen & Manning (2014), and brain tumor segmentation Havaei et al. (2017).

Following the empirical success of dropout, there have been several studies in recent years aimed at establishing theoretical underpinnings of why and how dropout helps with generalization. Early work of Baldi & Sadowski (2013) showed that for a single linear unit (and a single sigmoid unit, approximately), dropout amounts to weight decay regularization on the weights. A similar result was shown by McAllester (2013) in a PAC-Bayes setting. For generalized linear models, Wager et al. (2013) established that dropout performs an adaptive regularization which is equivalent to a data-dependent scaling of the weight decay penalty. In their follow-up work, Wager et al. (2014) show that for linear classification, under a generative assumption on the data, dropout improves the convergence rate of the generalization error. In this paper, we focus on predictors represented in a factored form and give generalization bounds for matrix learning problems and single hidden layer ReLU networks.

In a related line of work, Helmbold & Long (2015) study the structural properties of the dropout regularizer in the context of linear classification. They characterize the landscape of the dropout criterion in terms of unique minimizers and establish non-monotonic and non-convex nature of the regularizer. In a follow up work, Helmbold & Long (2017) extend their analysis to dropout in deep ReLU networks and surprisingly find that the nature of regularizer is different from that in linear classification. In particular, they show that unlike weight decay, dropout regularizer in deep networks can grow exponentially with depth and remains invariant to rescaling of inputs, outputs, and network weights. We confirm some of these findings in our theoretical analysis. However, counter to the claims of Helmbold & Long (2017), we argue that dropout does indeed prevent co-adaptation.

In a closely related approach as ours, the works of Zhai & Wang (2018), Gao & Zhou (2016), and Wan et al. (2013) bound the Rademacher complexity of deep neural networks trained using dropout. In particular, Gao & Zhou (2016) show that the Rademacher complexity of the target class decreases polynomially or exponentially, for shallow and deep networks, respectively, albeit they assume additional norm bounds on the weight vectors. Similarly, the works of Wan et al. (2013) and Zhai & Wang (2018) assume that certain norms of the weights are bounded, and show that the Rademacher complexity of the target class decreases with dropout rates. We argue in this paper that dropout alone does not directly control the norms of the weight vectors; therefore, each of the works above fail to capture the practice. We emphasize that none of the previous works provide a generalization guarantee, i.e., a bound on the gap between the population risk and the empirical risk, merely in terms of the value of the explicit regularizer due to dropout. We give a first such result for dropout in the context of matrix completion and for a single hidden layer ReLU network.

There are a bunch of other works that do not fall into any of the categories above, and, in fact, are somewhat unrelated to the focus in this paper. Nonetheless, we discuss them here for completeness. For instance, Gal & Ghahramani (2016) study dropout as Bayesian approximation, Bank & Giryes (2018)

draw insights from frame theory to connect the notion of equiangular tight frames with dropout training in auto-encoders. Also, some recent works have considered variants of dropout. For instance, Mou et al. (2018) consider a variant of dropout, which they call "truthful" dropout, that ensures that the output of the randomly perturbed network is unbiased. However, rather than bound generalization error, Mou et al. (2018) bound the gap between the population risk and the dropout objective, i.e., the empirical risk plus the explicit regularizer. Li et al. (2016) study a yet another variant based on multinomoal sampling (different nodes are dropped with different rates), and establish sub-optimality bounds for stochastic optimization of linear models (for convex Lipschitz loss functions).

Matrix Factorization with Dropout. Our study of dropout is motivated in part by recent works of Cavazza et al. (2018), Mianjy et al. (2018), and Mianjy & Arora (2019). This line of work was initiated by Cavazza et al. (2018), who studied dropout for low-rank matrix factorization without constraining the rank of the factors or adding an explicit regularizer to the objective. They show that dropout in the context of matrix factorization yields an explicit regularizer whose convex envelope is given by nuclear norm. This result is further strengthened by Mianjy et al. (2018) who show that induced regularizer is indeed nuclear norm.

While matrix factorization is not a learning problem per se (for instance, what is training versus test data), in follow-up works by Mianjy et al. (2018) and Mianjy & Arora (2019), the authors show that training deep linear networks with  $\ell_2$ -loss using dropout reduces to the matrix factorization problem if the marginal distribution of the input feature vectors is assumed to be isotropic, i.e.,  $\mathbb{E}[\mathbf{x}\mathbf{x}^{\top}] = \mathbf{I}$ . We note that this is a strong assumption. If we do not assume isotropy, we show that dropout induces a data-dependent regularizer which amounts to a simple scaling of the parameters and, therefore, does not control capacity in any meaningful way. We revisit this discussion in Section 4.

To summarize, while we are motivated by Cavazza et al. (2018), the problem setup, the nature of statements in this paper, and the tools we use are different from that in Cavazza et al. (2018). Our proofs are simple and quickly verified. We do build closely on the prior work of Mianjy et al. (2018).

However, different from Mianjy et al. (2018), we rigorously argue for dropout in matrix completion by 1) showing that the induced regularizer is equal to weighted trace-norm, which as far as we know, is a novel result, 2) giving strong generalization bounds, and 3) providing extensive experimental evidence that dropout provides state of the art performance on one of the largest datasets in recommendation systems research. Beyond that we rigorously extend our results to two layer ReLU networks, describe the explicit regularizer, bound the Rademacher complexity of the hypothesis class controlled by dropout, show precise generalization bounds, and support them with empirical results.

#### 1.2 Notation and Preliminaries

We are primarily interested in understanding how dropout controls the capacity of the hypothesis class when using dropout for training. To that end, we consider Rademacher complexity, a sample dependent measure of complexity of a hypothesis class that can directly bound the generalization gap (Bartlett & Mendelson, 2002). Formally, let  $S = \{(x_1, y_1), \dots, (x_n, y_n)\}$  be a sample of size n. Then, the empirical Rademacher complexity of a function class  $\mathcal{F}$  with respect to S, and the expected Rademacher complexity

are defined, respectively, as

$$\mathfrak{R}_{\mathcal{S}}(\mathcal{F}) = \mathbb{E}_{\sigma} \sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^{n} \sigma_{i} f(\mathbf{x}_{i}), \quad \mathfrak{R}_{n}(\mathcal{F}) = \mathbb{E}_{\mathbf{x}}[\mathfrak{R}_{\mathcal{S}}(\mathcal{F})],$$

where  $\sigma_i$  are i.i.d. Rademacher random variables.

# 2 Matrix Sensing

We begin with understanding dropout for matrix sensing, a problem which arguably is an important instance of a matrix learning problem with lots of applications, and is well understood from a theoretical perspective. The problem setup is the following. Let  $M_* \in \mathbb{R}^{d_2 \times d_0}$  be a matrix with rank  $r_* := \operatorname{Rank}(M_*)$ . Let  $A^{(1)}, \ldots, A^{(n)}$  be a set of measurement matrices of the same size as  $M_*$ . The goal of matrix sensing is to recover the matrix  $M_*$  from n observations of the form  $y_i = \langle M_*, A^{(i)} \rangle$  such that  $n \ll d_2 d_0$ .

A natural approach is to represent the matrix in terms of factors and solve the following *empirical risk* minimization problem:

$$\min_{\mathbf{U},\mathbf{V}} \widehat{L}(\mathbf{U},\mathbf{V}) := \widehat{\mathbb{E}}_i (y_i - \langle \mathbf{U}\mathbf{V}^\top, \mathbf{A}^{(i)} \rangle)^2$$
 (1)

where  $U = [u_1, \ldots, u_{d_1}] \in \mathbb{R}^{d_2 \times d_1}$ ,  $V = [v_1, \ldots, v_{d_1}] \in \mathbb{R}^{d_0 \times d_1}$ . When the number of factors is unconstrained, i.e., when  $d_1 \gg r_*$ , there exist many "bad" empirical minimizers, i.e., those with a large true risk  $L(U, V) := \mathbb{E}(y - \langle UV^\top, A \rangle)^2$ . Interestingly, Li et al. (2018) showed recently that under a restricted isometry property (RIP), despite the existence of such poor ERM solutions, gradient descent with proper initialization is implicitly biased towards finding solutions with minimum nuclear norm – this is an important result which was first conjectured and empirically verified by Gunasekar et al. (2017). We do not make an RIP assumption here. Further, we argue that for the most part, modern machine learning systems employ explicit regularization techniques. In fact, as we show in the experimental section, the implicit bias due to (stochastic) gradient descent does not prevent it from blatant overfitting in the matrix completion problem.

We propose solving the ERM problem (1) using dropout, where at training time, corresponding columns of U and V are dropped uniformly at random. As opposed to an *implicit* effect of gradient descent, dropout *explicitly* regularizes the empirical objective. It is then natural to ask, in the case of matrix sensing, if dropout also biases the ERM towards certain low norm solutions. To answer this question, we begin with the observation that dropout can be viewed as an instance of SGD on the following objective Cavazza et al. (2018); Mianjy et al. (2018):

$$\widehat{L}_{\text{drop}}(\mathbf{U}, \mathbf{V}) = \widehat{\mathbb{E}}_j \mathbb{E}_{\mathbf{B}}(y_j - \langle \mathbf{U} \mathbf{B} \mathbf{V}^\top, \mathbf{A}^{(j)} \rangle)^2, \tag{2}$$

where  $B \in \mathbb{R}^{d_1 \times d_1}$  is a diagonal matrix whose diagonal elements are Bernoulli random variables distributed as  $B_{ii} \sim \frac{1}{1-p} \text{Ber}(1-p)$ . It is easy to show that for  $p \in [0,1)$ :

$$\widehat{L}_{\text{drop}}(\mathbf{U}, \mathbf{V}) = \widehat{L}(\mathbf{U}, \mathbf{V}) + \frac{p}{1-p} \widehat{R}(\mathbf{U}, \mathbf{V}), \tag{3}$$

where  $\widehat{R}(\mathbf{U}, \mathbf{V}) := \sum_{i=1}^{d_1} \widehat{\mathbb{E}}_j(\mathbf{u}_i^{\top} \mathbf{A}^{(j)} \mathbf{v}_i)^2$  is a data-dependent term that captures the *explicit* regularizer due to dropout. A similar result was shown by Cavazza et al. (2018) and Mianjy et al. (2018), but we provide a proof for completeness (see Proposition 2 in the Appendix).

We show that the *explicit* regularizer concentrates around its expected value w.r.t. the data distribution (see Lemma 2 in the Appendix). Furthermore, given that we seek a minimum of  $\hat{L}_{\text{drop}}$ , it suffices to consider the factors with the minimal value of the regularizer among all that yield the same empirical loss. This motivates studying the the following distribution-dependent *induced* regularizer:

$$\Theta(\mathbf{M}) := \min_{\mathbf{U} \mathbf{V}^{\top} = \mathbf{M}} R(\mathbf{U}, \mathbf{V}), \quad \text{where } R(\mathbf{U}, \mathbf{V}) := \mathbb{E}_{\mathbf{A}}[\widehat{R}(\mathbf{U}, \mathbf{V})].$$

For a wide range of random measurements,  $\Theta(\cdot)$  turns out to be a "suitable" regularizer. Here, we instantiate two important examples (see Proposition 3 in the Appendix).

Gaussian Measurements. For all  $j \in [n]$ , let  $A^{(j)}$  be standard Gaussian matrices. In this case, it is easy to see that  $L(U, V) = \|M_* - UV^\top\|_F^2$  and we recover the matrix factorization problem. Furthermore, we know from Cavazza et al. (2018); Mianjy & Arora (2019) that dropout regularizer acts as trace-norm regularization, i.e.,  $\Theta(M) = \frac{1}{d_*} \|M\|_*^2$ .

**Matrix Completion.** For all  $j \in [n]$ , let  $A^{(j)}$  be an indicator matrix whose (i, k)-th element is selected randomly with probability p(i)q(k), where p(i) and q(k) denote the probability of choosing the i-th row and the k-th column, respectively. Then

$$\Theta(\mathbf{M}) = \frac{1}{d_1} \| \operatorname{diag}(\sqrt{p}) \mathbf{U} \mathbf{V}^{\top} \operatorname{diag}(\sqrt{q}) \|_{*}^{2}$$

is the weighted trace-norm studied by Srebro & Salakhutdinov (2010) and Foygel et al. (2011).

These observations are specifically important because they connect dropout, an algorithmic heuristic in deep learning, to strong complexity measures that are empirically effective as well as theoretically well understood. To illustrate, here we give a generalization bound for matrix completion using dropout in terms of the value of the *explicit* regularizer at the minimizer.

**Theorem 1.** Assume that  $d_2 \geq d_0$  and  $\|\mathbf{M}_*\| \leq 1$ . Furthermore, assume that  $\min_{i,k} p(i)q(k) \geq \frac{\log(d_2)}{n\sqrt{d_2d_0}}$ . Let  $(\mathbf{U}, \mathbf{V})$  be a minimizer of the dropout ERM objective in equation (3). Let  $\alpha$  be such that  $R(\mathbf{U}, \mathbf{V}) \leq \alpha/d_1$ . Then, for any  $\delta \in (0,1)$ , the following generalization bounds holds with probability at least  $1-\delta$  over a sample of size n:

$$L(g(\mathbf{U}\mathbf{V}^{\top})) \le \widehat{L}(\mathbf{U}, \mathbf{V}) + 8\sqrt{\frac{2\alpha d_2 \log(d_2) + \frac{1}{4}\log(2/\delta)}{n}}$$

where g(M) thresholds M at  $\pm 1$ , i.e.  $g(M)(i,j) = \max\{-1, \min\{1, M(i,j)\}\}$ , and  $L(g(UV^\top)) := \mathbb{E}(y - \langle g(UV^\top), A \rangle)^2$  is the  $true\ risk$  of  $g(UV^\top)$ .

The proof of Theorem 1 follows from standard generalization bounds for  $\ell_2$  loss (Mohri et al., 2018) based on the Rademacher complexity (Bartlett & Mendelson, 2002) of the class of functions with weighted trace-norm bounded by  $\sqrt{\alpha}$ , i.e.  $\mathcal{M}_{\alpha} := \{M : \|\operatorname{diag}(\sqrt{p})M\operatorname{diag}(\sqrt{q})\|_*^2 \le \alpha\}$ . The non-degeneracy condition  $\min_{i,j} p(i)q(j) \ge \frac{\log(d_2)}{n\sqrt{d_2}d_0}$  is required to obtain a bound on the Rademacher complexity of  $\mathcal{M}_{\alpha}$ , as established by Foygel et al. (2011).

We note that for large enough sample size,  $\widehat{R}(U, V) \approx R(U, V) \approx \Theta(UV^{\top}) = \frac{1}{d_1} \|\operatorname{diag}(\sqrt{p})UV^{\top}\operatorname{diag}(\sqrt{q})\|_*^2$ , where the second approximation is due the fact that the pair (U, V) is a minimizer. That is, compared to the weighted trace-norm, the value of the explicit regularizer at the minimizer roughly scales as  $1/d_1$ . Hence the assumption  $\widehat{R}(U, V) \leq \alpha/d_1$  in the statement of the corollary.

In practice, for models that are trained with dropout, the training error  $\widehat{L}(U,V)$  is negligible (see Figure 1 for experiments on the MovieLens dataset). Moreover, given that the sample size is large enough, the third term can be made arbitrarily small. Having said that, the second term, which is  $O(\sqrt{\alpha d_2/n})$ , dominates the right hand side of generalization error bound in Theorem 9. In Appendix, we also give optimistic generalization bounds that decay as  $O(ad_2/n)$ .

Finally, the required sample size heavily depends on the value of the explicit regularizer (i.e.,  $\alpha/d_1$ ) at a minimizer, and hence, on the dropout rate p. In particular, increasing the dropout rate increases the regularization parameter  $\lambda := \frac{p}{1-p}$ , thereby intensifying the penalty due to the explicit regularizer. Intuitively, a larger dropout rate p results in a smaller  $\alpha$ , thereby a tighter generalization gap can be guaranteed. We show through experiments that that is indeed the case in practice.

## 3 Non-linear Networks

Next, we focus on neural networks with a single hidden layer. Let  $\mathcal{X} \subseteq \mathbb{R}^{d_0}$  and  $\mathcal{Y} \subseteq [-1,1]^{d_2}$  denote the input and output spaces, respectively. Let  $\mathcal{D}$  denote the joint probability distribution on  $\mathcal{X} \times \mathcal{Y}$ . Given n

examples  $\{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^n \sim \mathcal{D}^n$  drawn i.i.d. from the joint distribution and a loss function  $\ell: \mathcal{Y} \times \mathcal{Y} \to \mathbb{R}$ , the goal of learning is to find a hypothesis  $f_{\mathbf{w}}: \mathcal{X} \to \mathcal{Y}$ , parameterized by w, that has a small population risk  $L(f_{\mathbf{w}}) := \mathbb{E}_{\mathcal{D}}[\ell(f_{\mathbf{w}}(\mathbf{x}), \mathbf{y})]$ .

We focus on the squared  $\ell_2$  loss, i.e.,  $\ell(y,y') = \|y-y'\|^2$ , and study the generalization properties of the dropout algorithm for minimizing the *empirical risk*  $\widehat{L}(f_w) := \widehat{\mathbb{E}}_i[\|y_i - f_w(x_i)\|^2]$ . We consider the hypothesis class associated with feed-forward neural networks with 2 layers, i.e., functions of the form  $f_w(x) = U\sigma(V^\top x)$ , where  $U = [u_1, \dots, u_{d_1}] \in \mathbb{R}^{d_2 \times d_1}$ ,  $V = [v_1, \dots, v_{d_1}] \in \mathbb{R}^{d_0 \times d_1}$  are the weight matrices. The parameter w is the collection of weight matrices  $\{U, V\}$  and  $\sigma : \mathbb{R} \to \mathbb{R}$  is the ReLU activation function applied entrywise to an input vector.

As in Section 2, we view dropout as an instance of stochastic gradient descent on the following *dropout* objective:

 $\widehat{L}_{\text{drop}}(\mathbf{w}) := \widehat{\mathbb{E}}_i \mathbb{E}_{\mathbf{B}} \| \mathbf{y}_i - \mathbf{U} \mathbf{B} \sigma(\mathbf{V}^\top \mathbf{x}_i) \|^2, \tag{4}$ 

where B is a diagonal random matrix with diagonal elements distributed identically and independently as  $B_{ii} \sim \frac{1}{1-p} Bern(1-p)$ ,  $i \in [d_1]$ , for some dropout rate p. We seek to understand the explicit regularizer due to dropout:

$$\widehat{R}(\mathbf{w}) := \widehat{L}_{\text{drop}}(\mathbf{w}) - \widehat{L}(\mathbf{w}). \tag{5}$$

We denote the output of the *i*-th hidden node on an input vector x by  $a_i(x) \in \mathbb{R}$ ; for example,  $a_2(x) = \sigma(v_2^\top x)$ . Similarly, the vector  $\mathbf{a}(x) \in \mathbb{R}^{d_1}$  denotes the activation of the hidden layer on input x. Using this notation, we can rewrite the objective in (4) as  $\widehat{L}_{drop}(\mathbf{w}) := \mathbb{E}_i \mathbb{E}_{\mathbf{B}} \|\mathbf{y}_i - \mathbf{UBa}(\mathbf{x}_i)\|^2$ . It is then easy to show that the regularizer due to dropout in (5) is given as (see Proposition 4 in Appendix):

$$\widehat{R}(\mathbf{w}) = \frac{p}{1 - p} \sum_{j=1}^{d_1} \|\mathbf{u}_j\|^2 \widehat{a}_j^2, \text{ where } \widehat{a}_j = \sqrt{\widehat{\mathbb{E}}_i a_j(\mathbf{x}_i)^2}.$$

The explicit regularizer  $\widehat{R}(\mathbf{w})$  is a summation over hidden nodes, of the product of the squared norm of the outgoing weights with the empirical second moment of the output of the corresponding neuron. We should view it as a data-dependent variant of the  $\ell_2$  path-norm of the network, studied recently by Neyshabur et al. (2015b) and shown to yield capacity control in deep learning. Indeed, if we consider ReLU activations and input distributions that are symmetric and isotropic Mianjy et al. (2018), the expected regularizer is equal to the sum over all paths from input to output of the product of the squares of weights along the paths, i.e.,

$$R(\mathbf{w}) := \mathbb{E}[\widehat{R}(\mathbf{w})] = \frac{1}{2} \sum_{i_0, i_1, i_2 = 1}^{d_0, d_1, d_2} \mathbf{U}(i_2, i_1)^2 \mathbf{V}(i_0, i_1)^2,$$

which is precisely the squared  $\ell_2$  path-norm of the network. We refer the reader to Proposition 5 in the Appendix for a formal statement and proof.

#### 3.1 Generalization Bounds

To understand the generalization properties of dropout, we focus on the following distribution-dependent class

$$\mathcal{F}_{\alpha} := \{ f_{\mathbf{w}} : \mathbf{x} \mapsto \mathbf{u}^{\top} \sigma(\mathbf{V}^{\top} \mathbf{x}), \sum_{i=1}^{d_1} |u_i| a_i \le \alpha \}.$$

where  $u \in \mathbb{R}^{d_1}$  is the top layer weight vector,  $u_i$  denotes the *i*-th entry of u, and  $a_i^2 := \mathbb{E}_{\mathbf{x}}[\hat{a}_i^2] = \mathbb{E}_{\mathbf{x}}[a_i(\mathbf{x})^2]$  is the expected squared neural activation for the *i*-th hidden node. For simplicity, we focus on networks with one output neuron  $d_2 = 1$ ; extension to multiple output neurons is rather straightforward.

We argue that networks that are trained with dropout belong to the class  $\mathcal{F}_{\alpha}$ , for a small value of  $\alpha$ . In particular, by Cauchy-Schwartz inequality, it is easy to to see that  $\sum_{i=1}^{d_1} |u_i| a_i \leq \sqrt{d_1 R(\mathbf{w})}$ . Thus, for a fixed width, dropout implicitly controls the function class  $\mathcal{F}_{\alpha}$ . More importantly, this inequality is loose

if a small subset of hidden nodes  $\mathcal{J} \subset [d_1]$  "co-adapt" in a way that for all  $j \in [d_1] \setminus \mathcal{J}$ , the other hidden nodes are almost inactive, i.e.  $u_j a_j \approx 0$ . In other words, by minimizing the expected regularizer, dropout is biased towards networks where the gap between  $R(\mathbf{w})$  and  $(\sum_{i=1}^{d_1} |u_i|a_i)^2/d_1$  is small, which in turn happens if  $|u_i|a_i \approx |u_j|a_j, \forall i,j \in [d_1]$ . In this sense, dropout breaks "co-adaptation" between neurons by promoting solutions with nearly equal contribution from hidden neurons.

As we mentioned in the introduction, a bound on the dropout regularizer is not sufficient to guarantee a bound on a norm-based complexity measures that are common in the deep learning literature (see, e.g. Golowich et al. (2018) and the references therein), whereas a norm bound on the weight vector would imply a bound on the explicit regularizer due to dropout. Formally, we show the following.

**Proposition 1.** For any C > 0, there exists a distribution on the unit Euclidean sphere, and a network  $f_w : x \mapsto \sigma(w^\top x)$ , such that  $R(w) = \sqrt{\mathbb{E}\sigma(w^\top x)^2} \le 1$ , while ||w|| > C.

In other words, even though we connect the dropout regularizer to path-norm, the data-dependent nature of the regularizer prevents us from leveraging that connection in data-independent manner (i.e., for all distributions). At the same time, making strong distributional assumptions (as in Proposition 5) would be impractical. Instead, we argue for the following milder condition on the input distribution which we show as sufficient to ensure generalization.

**Assumption 1** ( $\beta$ -retentive). The marginal input distribution is  $\beta$ -retentive for some  $\beta \in (0, 1/2]$ , if for any non-zero vector  $\mathbf{v} \in \mathbb{R}^d$ , it holds that  $\mathbb{E}\sigma(\mathbf{v}^\top \mathbf{x})^2 \geq \beta \mathbb{E}(\mathbf{v}^\top \mathbf{x})^2$ .

Intuitively, what the assumption implies is that the variance (aka, the information or signal in the data) in the pre-activation at any node in the network is not quashed considerably due to the non-linearity. In fact, no reasonable training algorithm should learn weights where  $\beta$  is small. However, we steer clear from algorithmic aspects of dropout training, and make the assumption above for every weight vector as we will need it when carrying out a union bound.

We now present the first main result of this section, which bounds the Rademacher complexity of  $\mathcal{F}_{\alpha}$  in terms of  $\alpha$ , the retentiveness coefficient  $\beta$ , and the Mahalanobis norm of the data with respect to the pseudo-inverse of the second moment, i.e.  $\|\mathbf{X}\|_{\mathbf{C}^{\dagger}}^2 = \sum_{i=1}^n \mathbf{x}_i^{\top} \mathbf{C}^{\dagger} \mathbf{x}_i$ .

**Theorem 2.** For any sample  $S = \{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^n$  of size n,

$$\mathfrak{R}_{\mathcal{S}}(\mathcal{F}_{\alpha}) \le \frac{2\alpha \|\mathbf{X}\|_{\mathbf{C}^{\dagger}}}{n\sqrt{\beta}}.$$

Furthermore, it holds for the expected Rademacher complexity that  $\mathfrak{R}_n(\mathcal{F}_\alpha) \leq 2\alpha \sqrt{\frac{\operatorname{Rank}(C)}{\beta n}}$ .

First, note that the bound depends on the quantity  $\|X\|_{C^{\dagger}}$  which can be in the same order as  $\|X\|_{F}$  with both scaling as  $\approx \sqrt{nd_0}$ ; the latter is more common in the literature Neyshabur et al. (2018); Bartlett et al. (2017); Neyshabur et al. (2017); Golowich et al. (2018); Neyshabur et al. (2015b). This is unfortunately unavoidable, unless one makes stronger distributional assumptions.

Second, as we discussed earlier, the dropout regularizer directly controls the value of  $\alpha$ , thereby controlling the Rademacher complexity in Theorem 2. This bound also gives us a bound on the Rademacher complexity of the networks trained using dropout. To see that, consider the following class of networks with bounded explicit regularizer, i.e.,  $\mathcal{H}_r := \{h_{\mathbf{w}} : \mathbf{x} \mapsto \mathbf{u}^\top \sigma(\mathbf{V}^\top \mathbf{x}), \ R(\mathbf{u}, \mathbf{V}) \leq r\}$ . Then, Theorem 2 yields  $\mathfrak{R}_{\mathcal{S}}(\mathcal{H}_r) \leq \frac{2\sqrt{d_1 r} \|\mathbf{X}\|_{\mathbf{C}^\dagger}}{n\sqrt{\beta}}$ . In fact, we can show that this bound is tight up to  $1/\sqrt{\beta}$  by a reduction to the linear case. Formally, we show the following.

**Theorem 3** (Lowerbound on  $\mathfrak{R}_{\mathcal{S}}$ ). There is a constant c such that for any scalar r > 0,  $\mathfrak{R}_{\mathcal{S}}(\mathcal{H}_r) \geq \frac{c\sqrt{d_1r}\|X\|_{C^{\frac{1}{r}}}}{2}$ .

Moreover, it is easy to give a generalization bound based on Theorem 2 that depends only on the distribution dependent quantities  $\alpha$  and  $\beta$ . Let  $g_{\mathbf{w}}(\cdot) := \max\{-1, \min\{1, f_{\mathbf{w}}(\cdot)\}\}$  project the network output  $f_{\mathbf{w}}$  onto the range [-1, 1]. We have the following generalization gurantees for  $g_{\mathbf{w}}$ .

	plain SGD		dropout			
$\mathbf{width}$	last iterate	best iterate	p = 0.1	p = 0.2	p = 0.3	p = 0.4
$d_1 = 30$	0.8041	0.7938	0.7805	0.785	0.7991	0.8186
$d_1 = 70$	0.8315	0.7897	0.7899	0.7771	0.7763	0.7833
$d_1 = 110$	0.8431	0.7873	0.7988	0.7813	0.7742	0.7743
$d_1 = 150$	0.8472	0.7858	0.8042	0.7852	0.7756	0.7722
$d_1 = 190$	0.8473	0.7844	0.8069	0.7879	0.7772	0.772

Table 1: MovieLens dataset: Test RMSE of plain SGD as well as the dropout algorithm with various dropout rates for various factorization sizes. The grey cells shows the best performance(s) in each row.

Corollary 1. For any  $w \in \mathcal{F}_{\alpha}$ , for any  $\delta \in (0,1)$ , the following generalization bound holds with probability at least  $1 - \delta$  over a sample S of size n

$$L_{\mathcal{D}}(g_{\mathbf{w}}) \le \widehat{L}_{\mathcal{S}}(g_{\mathbf{w}}) + \frac{16\alpha \|\mathbf{X}\|_{\mathbf{C}^{\dagger}}}{\sqrt{\beta}n} + 12\sqrt{\frac{\log(2/\delta)}{2n}}$$

 $\beta$ -independent Bounds. Geometrically,  $\beta$ -retentiveness requires that for any hyperplane passing through the origin, both halfspaces contribute significantly to the second moment of the data in the direction of the normal vector. It is not clear, however, if  $\beta$  can be estimated efficiently, given a dataset. Nonetheless, when  $\mathcal{X} \subseteq \mathbb{R}^{d_0}_+$ , which is the case for image datasets, a simple symmetrization technique, described below, allows us to give bounds that are  $\beta$ -independent; note that the bound still depend on the sample as  $\|\mathbf{X}\|_{\mathbb{C}^1}$ . Here is the randomized symmetrization we propose. Given a training sample  $\mathcal{S} = \{(\mathbf{x}_i, y_i), i \in [n]\}$ , consider the following randomized perturbation,  $\mathcal{S}' = \{(\zeta_i \mathbf{x}_i, y_i), i \in [n]\}$ , where  $\zeta_i$ 's are i.i.d. Rademacher random variables. We give a generalization bound (w.r.t. the original data distribution) for the hypothesis class with bounded regularizer w.r.t. the perturbed data distribution.

Corollary 2. Given an i.i.d. sample  $S = \{(x_i, y_i)\}_{i=1}^n$ , let

$$\mathcal{F}'_{\alpha} := \{ f_{\mathbf{w}} : \mathbf{x} \mapsto u^{\top} \sigma(V^{\top} \mathbf{x}), \sum_{i=1}^{d_1} |u_i| a'_i \le \alpha \},$$

where  $a_i'^2 := \mathbb{E}_{\mathbf{x},\zeta}[a_i(\zeta\mathbf{x})^2]$ . For any  $\mathbf{w} \in \mathcal{F}'_{\alpha}$ , for any  $\delta \in (0,1)$ , the following generalization bound holds with probability at least  $1-\delta$  over a sample of size n and the randomization in symmetrization

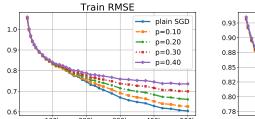
$$L_{\mathcal{D}}(g_{\mathbf{w}}) \le 2\widehat{L}_{\mathcal{S}'}(g_{\mathbf{w}}) + \frac{46\alpha \|\mathbf{X}\|_{\mathbf{C}^{\dagger}}}{n} + 24\sqrt{\frac{\log(2/\delta)}{2n}}.$$

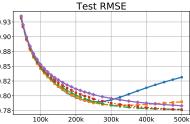
Note that the population risk of the clipped predictor  $g_{\mathbf{w}}(\cdot) := \max\{-1, \min\{1, f_{\mathbf{w}}(\cdot)\}\}$  is bounded in terms of the empirical risk on  $\mathcal{S}'$ . Finally, we verify in Section 5 that symmetrization of the training set, on MNIST and FashionMNIST datasets, does not have an effect on the performance of the trained models.

# 4 Role of Parametrization

In this section, we argue that parametrization plays an important role in determining the nature of the inductive bias.

We begin by considering matrix sensing in non-factorized form, which entails minimizing  $\widehat{L}(M) := \widehat{\mathbb{E}}_i(y_i - \langle \operatorname{vec}(M), \operatorname{vec}(A^{(i)}) \rangle)^2$ , where  $\operatorname{vec}(M)$  denotes the column vectorization of M. Then, the expected explicit regularizer due to dropout equals  $R(M) = \frac{p}{1-p} \| \operatorname{vec}(M) \|_{\operatorname{diag}(C)}^2$ , where  $C = \mathbb{E}[\operatorname{vec}(A) \operatorname{vec}(A)^{\top}]$  is the second moment of the measurement matrices. For instance, with Gaussian measurements, the second moment equals the identity matrix, in which case, the regularizer reduces to the Frobenius norm of the parameters





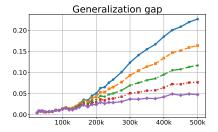


Figure 1: MovieLens dataset: the training error (**left**), the test error (**middle**), and the generalization gap (**right**) for plain SGD and dropout with  $p \in \{0.1, 0.2, 0.3, 0.4\}$  as a function of the number of iterations. The factorization size is  $d_1 = 70$ .

 $R(M) = \frac{p}{1-p} ||M||_F^2$ . While such a ridge penalty yields a useful inductive bias in linear regression, it is not "rich" enough to capture the kind of inductive bias that provides rank control in matrix sensing.

However, simply representing the hypotheses in a factored form alone is not sufficient in terms of imparting a rich inductive bias to the learning problem. Recall that in linear regression, dropout, when applied on the input features, yields ridge regularization. However, if we were to represent the linear predictor in terms of a deep linear network, then we argue that the effect of dropout is markedly different. Consider a deep linear network,  $f_w: x \mapsto W_k \cdots W_1 x$  with a single output neuron. In this case, Mianjy & Arora (2019) show that  $\nu \|f\|_{\widehat{C}}^2 = \min_{f_w = f} \widehat{R}(w)$ , where  $\nu$  is a regularization parameter independent of the parameters w. Consequently, in deep linear networks with a single output neuron, dropout reduces to solving

$$\min_{\mathbf{u} \in \mathbb{R}^{d_0}} \widehat{\mathbb{E}}_i (y_i - \mathbf{u}^\top \mathbf{x}_i)^2 + \nu \|\mathbf{u}\|_{\widehat{\mathbf{C}}}^2.$$

All the minimizers of the above problem are solutions to the system of linear equations  $(1 + \frac{\nu}{n})XX^{\top}u = Xy$ , where  $X = [x_1, \dots, x_n] \in \mathbb{R}^{d_0 \times n}, y = [y_1; \dots; y_n] \in \mathbb{R}^n$  are the design matrix and the response vector, respectively. In other words, the dropout regularizer manifests itself merely as a scaling of the parameters which does not seem to be useful in any meaningful way.

What we argue above may at first seem to contradict the results of Section 2 on matrix sensing, which is arguably an instance of regression with a two-layer linear network. Note though that casting matrix sensing in a factored form as a linear regression problem requires us to use a convolutional structure. This is easy to check since

$$\begin{split} \langle UV^\top, A \rangle &= \langle \operatorname{vec} \left( U^\top \right), \operatorname{vec} \left( V^\top A^\top \right) \rangle \\ &= \langle \operatorname{vec} \left( U^\top \right), \left( I_{d_2} \otimes V^\top \right) \operatorname{vec} \left( A^\top \right) \rangle, \end{split}$$

where  $\otimes$  is the Kronecker product, and we used the fact that  $\text{vec}(AB) = (I \otimes A) \text{ vec}(B)$  for any pair of matrices A, B. The expression  $(I \otimes V^{\top})$  represents a fully connected convolutional layer with  $d_1$  filters specified by columns of V. The convolutional structure in addition to dropout is what imparts the problem of matrix sensing the nuclear norm regularization. For nonlinear networks, however, a simple feed-forward structure suffices as we saw in Section 3.

# 5 Experimental Results

In this section, we report our empirical findings on real world datasets. All results are averaged over 50 independent runs with random initialization.

#### 5.1 Matrix Completion

We evaluate dropout on the MovieLens dataset Harper & Konstan (2016), a publicly available collaborative filtering dataset that contains 10M ratings for 11K movies by 72K users of the online movie recommender ser-

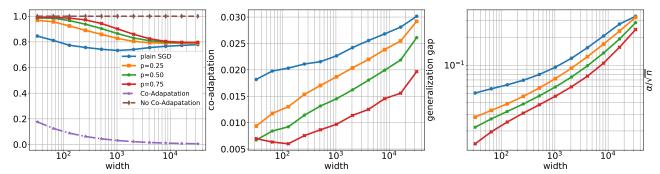


Figure 2: (left) "co-adaptation"; (middle) generalization gap; and (right)  $\alpha/\sqrt{n}$  as a function of the width of networks trained with dropout on MNIST. In left figure, the dashed brown and dotted purple lines represent minimal and maximal co-adaptations, respectively.

vice MovieLens.

We initialize the factors using the standard He initialization scheme. We train the model for 100 epochs over the training data, where we use a fixed learning rate of 1r = 1, and a batch size of 2000. We report the results for plain SGD (p = 0.0) as well as the dropout algorithm with  $p \in \{0.1, 0.2, 0.3, 0.4\}$ .

Figure 1 shows the progress in terms of the training and test error as well as the gap between them as a function of the number of iterations for  $d_1 = 70$ . It can be seen that plain SGD is the fastest in minimizing the empirical risk. The dropout rate clearly determines the trade-off between the approximation error and the estimation error: as the dropout rate p increases, the algorithm favors less complex solutions that suffer larger empirical error (left figure) but enjoy smaller generalization gap (right figure). The best trade-off here seems to be achieved by a moderate dropout rate of p = 0.3. We observe similar behaviour for different factorization sizes; please see the Appendix for additional plots with factorization sizes  $d_1 \in \{30, 110, 150, 190\}$ .

It is remarkable, how even in the "simple" problem of matrix completion, plain SGD lacks a proper inductive bias. As seen in the middle plot, without *explicit* regularization, in particular, without early stopping or dropout, SGD starts overfitting. We further illustrate this in Table 1, where we compare the test root-mean-squared-error (RMSE) of plain SGD with the dropout algorithm, for various factorization sizes. To show the superiority of dropout over SGD with early stopping, we give SGD the advantage of having access to the *test set* (and not a separate validation set), and report the best iterate in the third column. Even with this impractical privilege, dropout performs better (> 0.01 difference in test RMSE).

#### 5.2 Neural Networks

We train 2-layer neural networks with and without dropout, on MNIST dataset of handwritten digits and Fashion MNIST dataset of Zalando's article images, each of which contains 60K training examples and 10K test examples, where each example is a  $28 \times 28$  grayscale image associated with a label from 10 classes. We extract two classes  $\{4,7\}$  and label them as  $\{-1,+1\}^{-1}$ . The learning rate in all experiments is set to  $\mathbf{lr} = 1e - 3$ . We train the models for 30 epochs over the training set. We run the experiments both with and without symmetrization. Here we only report the results with symmetrization, and on the MNIST dataset. For experiments without symmetrization, and experiments on FashionMNIST, please see the Appendix. We remark that under the above experimental setting, the trained networks achieve 100% training accuracy.

For any node  $i \in [d_1]$ , we define its flow as  $\psi_i := |u_i|a_i$  (respectively  $\psi_i := |u_i|a_i'$  for symmetrized data), which measures the overall contribution of a node to the output of the network. Co-adaptation occurs when a small subset of nodes dominate the overall function of the network. We argue that  $\phi(w) = \frac{\|\psi\|_1}{\sqrt{d_1\|\psi\|_2}}$  is a suitable measure of co-adaptation (or lack thereof) in a network parameterized by w. In case of high co-adaptation, only a few nodes have a high flow, which implies  $\phi(w) \approx \frac{1}{\sqrt{d_1}}$ . At the other end of the

<sup>&</sup>lt;sup>1</sup>We observe similar results across other choices of target classes.

spectrum, all nodes are equally active, in which case  $\phi(w) \approx 1$ . Figure 2 (left) illustrates this measure as a function of the network width for several dropout rates  $p \in \{0, 0.25, 0.5, 0.75\}$ . In particular, we observe that a higher dropout rate corresponds to less co-adapation. More interestingly, even plain SGD is *implicitly* biased towards networks with less co-adapation. Moreover, for a fixed dropout rate, the regularization effect due to dropout decreases as we increase the width. Thus, it is natural to expect more co-adaptation as the network becomes wider, which is what we observe in the plots.

The generalization gap is plotted in Figure 2 (middle). As expected, increasing dropout rate decreases the generalization gap, uniformly for all widths. In our experiments, the generalization gap increases with the width of the network. The figure on the right shows the quantity  $\alpha/\sqrt{n}$  that shows up in the Rademacher complexity bounds in Section 3. We note that, the bound on the Rademacher complexity is predictive of the generalization gap, in the sense that a smaller bound corresponds to a curve with smaller generalization gap.

## 6 Conclusion

Motivated by the success of dropout in deep learning, we study a dropout algorithm for matrix sensing and show that it enjoys strong generalization guarantees as well as competitive test performance on the MovieLens dataset. We then focus on deep regression under the squared loss and show that the regularizer due to dropout serves as a strong complexity measure for the underlying class of neural networks, using which we give a generalization error bound in terms of the value of the regularizer.

# Acknowledgements

This research was supported, in part, by NSF BIGDATA award IIS-1546482 and NSF CAREER award IIS-1943251. The seeds of this work were sown during the summer 2019 workshop on the Foundations of Deep Learning at the Simons Institute for the Theory of Computing. Raman Arora acknowledges the support provided by the Institute for Advanced Study, Princeton, New Jersey as part of the special year on Optimization, Statistics, and Theoretical Machine Learning.

## References

- Baldi, P. and Sadowski, P. J. Understanding dropout. In Advances in Neural Information Processing Systems (NIPS), pp. 2814–2822, 2013.
- Bank, D. and Giryes, R. On the relationship between dropout and equiangular tight frames. arXiv preprint arXiv:1810.06049, 2018.
- Bartlett, P. L. and Mendelson, S. Rademacher and gaussian complexities: Risk bounds and structural results. Journal of Machine Learning Research, 3(Nov):463–482, 2002.
- Bartlett, P. L., Foster, D. J., and Telgarsky, M. J. Spectrally-normalized margin bounds for neural networks. In *Advances in Neural Information Processing Systems*, pp. 6240–6249, 2017.
- Cavazza, J., Haeffele, B. D., Lane, C., Morerio, P., Murino, V., and Vidal, R. Dropout as a low-rank regularizer for matrix factorization. *Int. Conf. on Artificial Intelligence and Statistics (AISTATS)*, 2018.
- Chen, D. and Manning, C. A fast and accurate dependency parser using neural networks. In *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)*, pp. 740–750, 2014.
- Dahl, G. E., Sainath, T. N., and Hinton, G. E. Improving deep neural networks for lvcsr using rectified linear units and dropout. In 2013 IEEE international conference on acoustics, speech and signal processing, pp. 8609–8613. IEEE, 2013.

- Foygel, R., Shamir, O., Srebro, N., and Salakhutdinov, R. R. Learning with the weighted trace-norm under arbitrary sampling distributions. In *Advances in Neural Information Processing Systems*, pp. 2133–2141, 2011.
- Gal, Y. and Ghahramani, Z. Dropout as a bayesian approximation: Representing model uncertainty in deep learning. In *Int. Conf. Machine Learning (ICML)*, 2016.
- Gao, W. and Zhou, Z.-H. Dropout rademacher complexity of deep neural networks. *Science China Information Sciences*, 59(7):072104, 2016.
- Golowich, N., Rakhlin, A., and Shamir, O. Size-independent sample complexity of neural networks. In Conference On Learning Theory, pp. 297–299, 2018.
- Gunasekar, S., Woodworth, B. E., Bhojanapalli, S., Neyshabur, B., and Srebro, N. Implicit regularization in matrix factorization. In *Advances in Neural Information Processing Systems*, pp. 6151–6159, 2017.
- Harper, F. M. and Konstan, J. A. The movielens datasets: History and context. Acm transactions on interactive intelligent systems (tiis), 5(4):19, 2016.
- Havaei, M., Davy, A., Warde-Farley, D., Biard, A., Courville, A., Bengio, Y., Pal, C., Jodoin, P.-M., and Larochelle, H. Brain tumor segmentation with deep neural networks. *Medical image analysis*, 35:18–31, 2017.
- Helmbold, D. P. and Long, P. M. On the inductive bias of dropout. *Journal of Machine Learning Research* (*JMLR*), 16:3403–3454, 2015.
- Helmbold, D. P. and Long, P. M. Surprising properties of dropout in deep networks. *The Journal of Machine Learning Research*, 18(1):7284–7311, 2017.
- Hinton, G. E., Srivastava, N., Krizhevsky, A., Sutskever, I., and Salakhutdinov, R. R. Improving neural networks by preventing co-adaptation of feature detectors. arXiv preprint arXiv:1207.0580, 2012.
- Kalchbrenner, N., Grefenstette, E., and Blunsom, P. A convolutional neural network for modelling sentences. arXiv preprint arXiv:1404.2188, 2014.
- Krizhevsky, A., Hinton, G., et al. Learning multiple layers of features from tiny images. Technical report, Citeseer, 2009.
- Krizhevsky, A., Sutskever, I., and Hinton, G. E. Imagenet classification with deep convolutional neural networks. In *Advances in neural information processing systems*, pp. 1097–1105, 2012.
- Li, Y., Ma, T., and Zhang, H. Algorithmic regularization in over-parameterized matrix sensing and neural networks with quadratic activations. In *Conference On Learning Theory*, pp. 2–47, 2018.
- Li, Z., Gong, B., and Yang, T. Improved dropout for shallow and deep learning. In *Advances in neural information processing systems*, pp. 2523–2531, 2016.
- McAllester, D. A pac-bayesian tutorial with a dropout bound. arXiv preprint arXiv:1307.2118, 2013.
- Mianjy, P. and Arora, R. On dropout and nuclear norm regularization. In *International Conference on Machine Learning*, 2019.
- Mianjy, P., Arora, R., and Vidal, R. On the implicit bias of dropout. In *International Conference on Machine Learning*, pp. 3537–3545, 2018.
- Mohri, M., Rostamizadeh, A., and Talwalkar, A. Foundations of machine learning. MIT press, 2018.
- Mou, W., Zhou, Y., Gao, J., and Wang, L. Dropout training, data-dependent regularization, and generalization bounds. In *International Conference on Machine Learning*, pp. 3642–3650, 2018.

- Neyshabur, B., Salakhutdinov, R. R., and Srebro, N. Path-sgd: Path-normalized optimization in deep neural networks. In *Advances in Neural Information Processing Systems*, pp. 2422–2430, 2015a.
- Neyshabur, B., Tomioka, R., and Srebro, N. Norm-based capacity control in neural networks. In *Conference on Learning Theory*, pp. 1376–1401, 2015b.
- Neyshabur, B., Bhojanapalli, S., and Srebro, N. A pac-bayesian approach to spectrally-normalized margin bounds for neural networks. arXiv preprint arXiv:1707.09564, 2017.
- Neyshabur, B., Li, Z., Bhojanapalli, S., LeCun, Y., and Srebro, N. Towards understanding the role of over-parametrization in generalization of neural networks. arXiv preprint arXiv:1805.12076, 2018.
- Pham, V., Bluche, T., Kermorvant, C., and Louradour, J. Dropout improves recurrent neural networks for handwriting recognition. In 2014 14th international conference on frontiers in handwriting recognition, pp. 285–290. IEEE, 2014.
- Srebro, N. and Salakhutdinov, R. R. Collaborative filtering in a non-uniform world: Learning with the weighted trace norm. In *Advances in Neural Information Processing Systems*, pp. 2056–2064, 2010.
- Srebro, N., Sridharan, K., and Tewari, A. Optimistic rates for learning with a smooth loss. arXiv preprint arXiv:1009.3896, 2010.
- Srivastava, N., Hinton, G. E., Krizhevsky, A., Sutskever, I., and Salakhutdinov, R. Dropout: a simple way to prevent neural networks from overfitting. *Journal of Machine Learning Research (JMLR)*, 15(1), 2014.
- Szegedy, C., Liu, W., Jia, Y., Sermanet, P., Reed, S., Anguelov, D., Erhan, D., Vanhoucke, V., and Rabinovich, A. Going deeper with convolutions. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 1–9, 2015.
- Toshev, A. and Szegedy, C. Deeppose: Human pose estimation via deep neural networks. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2014.
- Vershynin, R. High-dimensional probability: An introduction with applications in data science, volume 47. Cambridge University Press, 2018.
- Wager, S., Wang, S., and Liang, P. S. Dropout training as adaptive regularization. In *Advances in Neural Information Processing Systems (NIPS)*, 2013.
- Wager, S., Fithian, W., Wang, S., and Liang, P. S. Altitude training: Strong bounds for single-layer dropout. In Adv. Neural Information Processing Systems, 2014.
- Wan, L., Zeiler, M., Zhang, S., Le Cun, Y., and Fergus, R. Regularization of neural networks using dropconnect. In *International conference on machine learning*, pp. 1058–1066, 2013.
- Wang, S. and Manning, C. Fast dropout training. In international conference on machine learning, pp. 118–126, 2013.
- Yang, Z., He, X., Gao, J., Deng, L., and Smola, A. Stacked attention networks for image question answering. In Proceedings of the IEEE conference on computer vision and pattern recognition, pp. 21–29, 2016.
- Zhai, K. and Wang, H. Adaptive dropout with rademacher complexity regularization. In *International Conference on Learning Representations*, 2018. URL https://openreview.net/forum?id=S1uxsyeOZ.

# Supplementary Materials for "Dropout: Explicit Forms and Capacity Control"

# A Auxiliary Results

**Lemma 1** (Khintchine-Kahane inequality). Let  $\{\epsilon_i\}_{i=1}^n$  be i.i.d. Rademacher random variables, and  $\{x\}_{i=1}^n \subset \mathbb{R}^d$ . Then there exist a universal constants c > 0 such that

$$\mathbb{E}\|\sum_{i=1}^n \epsilon_i \mathbf{x}_i\| \ge c \sqrt{\sum_{i=1}^n \|\mathbf{x}_i\|^2}.$$

**Theorem 4** (Hoeffding's inequality: Theorem 2.6.2 Vershynin (2018)). Let  $X_1, \ldots, X_N$  be independent, mean zero, sub-Gaussian random variables. Then, for every  $t \ge 0$ , we have

$$\mathbb{P}\left(\left|\widehat{\mathbb{E}}_{i}X_{i}\right| \geq t\right) \leq 2e^{-\frac{ct^{2}N^{2}}{\sum_{i=1}^{N}\|X_{i}\|_{\psi_{2}}^{2}}}$$

**Theorem 5** (Theorem 3.1 of Mohri et al. (2018)). Let  $\mathcal{G}$  be a family of functions mapping from  $\mathcal{Z}$  to [0,1]. Then, for any  $\delta > 0$ , with probability at least  $1 - \delta$  over a sample  $\mathcal{S} = \{z_1, \ldots, z_n\}$ , the following holds for all  $g \in \mathcal{G}$ 

$$E[g(z)] \le \frac{1}{n} \sum_{i=1}^{n} g(z_i) + 2\Re_n(\mathcal{G}) + \sqrt{\frac{\log(1/\delta)}{2n}}$$

$$E[g(z)] \le \frac{1}{n} \sum_{i=1}^{n} g(z_i) + 2\Re_{\mathcal{S}}(\mathcal{G}) + 3\sqrt{\frac{\log(1/\delta)}{2n}}$$

**Theorem 6** (Theorem 10.3 of Mohri et al. (2018)). Assume that  $||h - f||_{\infty} \leq M$  for all  $h \in \mathcal{H}$ . Then, for any  $\delta > 0$ , with probability at least  $1 - \delta$  over a sample  $\mathcal{S} = \{(x_i, y_i), i \in [n]\}$  of size n, the following inequalities holds uniformly for all  $h \in \mathcal{H}$ .

$$\mathbb{E}[|h(\mathbf{x}) - f(\mathbf{x})|^2] \le \widehat{\mathbb{E}}_i |h(\mathbf{x}_i) - f(\mathbf{x}_i)|^2 + 4M\mathfrak{R}_n(\mathcal{H}) + M^2 \sqrt{\frac{\log(2/\delta)}{2n}}$$

$$\mathbb{E}[|h(\mathbf{x}) - f(\mathbf{x})|^2] \le \widehat{\mathbb{E}}_i |h(\mathbf{x}_i) - f(\mathbf{x}_i)|^2 + 4M\mathfrak{R}_{\mathcal{S}}(\mathcal{H}) + 3M^2 \sqrt{\frac{\log(2/\delta)}{2n}}$$

**Theorem 7** (Based on Theorem 1 in Srebro et al. (2010)). Let  $\mathcal{X}$  and  $\mathcal{Y} = [-1, 1]$  denote the input space and the label space, respectively. Let  $\mathcal{H} \subseteq \{f : \mathcal{X} \to \mathcal{Y}\}$  be the target function class. For any  $f \in \mathcal{H}$ , and any  $(x, y) \in \mathcal{X} \times \mathcal{Y}$ , let  $\ell(f, x, y) := (f(x) - y)^2$  be the squared loss. Let  $L(f) = \mathbb{E}_{\mathcal{D}}[\ell(f, x, y)]$  be the population risk with respect to the joint distribution  $\mathcal{D}$  on  $\mathcal{X} \times \mathcal{Y}$ . For any  $\delta > 0$ , with probability at least  $1 - \delta$  over a sample of size n, we have for any  $f \in \mathcal{H}$ :

$$L(f) \le L_* + K\left(\sqrt{L_*}\left(\sqrt{2}\log(n)^{1.5}\mathfrak{R}_n(\mathcal{H}) + \sqrt{\frac{4\log\frac{1}{\delta}}{n}}\right) + 2\log(n)^3\mathfrak{R}_n^2(\mathcal{H}) + \frac{4\log\frac{1}{\delta}}{n}\right)$$

where  $L_* := \min_{f \in \mathcal{H}} L(f)$ , and K is a numeric constant derived from Srebro et al. (2010).

**Theorem 8** (Theorem 3.3 in Mianjy et al. (2018)). For any pair of matrices  $U \in \mathbb{R}^{d_2 \times d_1}$ ,  $V \in \mathbb{R}^{d_0 \times d_1}$ , there exist a rotation matrix  $Q \in SO(d_1)$  such that rotated matrices  $\tilde{U} := UQ$ ,  $\tilde{V} := VQ$  satisfy  $\|\tilde{u}_i\| \|\tilde{v}_i\| = \frac{1}{d_1} \|UV^{\top}\|_*$ , for all  $i \in [d_1]$ .

**Theorem 9** (Theorem 1 in Foygel et al. (2011)). Assume that  $p(i)q(j) \geq \frac{\log(d_2)}{n\sqrt{d_2d_0}}$  for all  $i \in [d_2], j \in [d_0]$ . For any  $\alpha > 0$ , let  $\mathcal{M}_{\alpha} := \{M \in \mathbb{R}^{d_2 \times d_1} : \|\operatorname{diag}(\sqrt{p})M\operatorname{diag}(\sqrt{q})\|_*^2 \leq \alpha\}$  be the class of linear transformations with weighted trace-norm bounded with  $\sqrt{\alpha}$ . Then the expected Rademacher complexity of  $\mathcal{M}_{\alpha}$  is bounded as follows:

 $\mathfrak{R}_n(\mathcal{M}_{\alpha}) \leq O\left(\sqrt{\frac{\alpha d_2 \log(d_2)}{n}}\right)$ 

# B Matrix Sensing

**Proposition 2** (Dropout regularizer in matrix sensing). The following holds for any  $p \in [0,1)$ :

$$\widehat{L}_{drop}(U, V) = \widehat{L}(U, V) + \lambda \widehat{R}(U, V), \tag{6}$$

where  $\widehat{R}(U, V) = \sum_{i=1}^{d_1} \widehat{\mathbb{E}}_j(u_i^\top A^{(j)} v_i)^2$  and  $\lambda = \frac{p}{1-p}$  is the regularization parameter.

Proof of Proposition 2. Similar statements and proofs can be found in several previous works Srivastava et al. (2014); Wang & Manning (2013); Cavazza et al. (2018); Mianjy et al. (2018). For completeness, we include a proof here. The following equality follows from the definition of variance:

$$\mathbb{E}_{\mathbf{b}}[(y_i - \langle \mathbf{U}\mathbf{B}\mathbf{V}^\top, \mathbf{A}^{(i)}\rangle)^2] = \left(\mathbb{E}_{\mathbf{b}}[y_i - \langle \mathbf{U}\mathbf{B}\mathbf{V}^\top, \mathbf{A}^{(i)}\rangle]\right)^2 + \mathbf{Var}(y_i - \langle \mathbf{U}\mathbf{B}\mathbf{V}^\top, \mathbf{A}^{(i)}\rangle)$$
(7)

Recall that for a Bernoulli random variable  $B_{ii}$ , we have  $\mathbb{E}[B_{ii}] = 1$  and  $Var(B_{ii}) = \frac{p}{1-p}$ . Thus, the first term on right hand side is equal to  $(y_i - \langle UV^\top, A^{(i)} \rangle)^2$ . For the second term we have

$$\operatorname{Var}(y_{i} - \langle \operatorname{UBV}^{\top}, \operatorname{A}^{(i)} \rangle) = \operatorname{Var}(\sum_{j=1}^{d_{1}} \operatorname{B}_{jj} \operatorname{u}_{j}^{\top} \operatorname{A}^{(i)} \operatorname{v}_{j}) = \sum_{j=1}^{d_{1}} (\operatorname{u}_{j}^{\top} \operatorname{A}^{(i)} \operatorname{v}_{j})^{2} \operatorname{Var}(\operatorname{B}_{jj}) = \frac{p}{1-p} \sum_{j=1}^{d_{1}} (\operatorname{u}_{j}^{\top} \operatorname{A}^{(i)} \operatorname{v}_{j})^{2}$$

Plugging the above into Equation (7) and averaging over samples we get

$$\begin{split} \widehat{L}_{\text{drop}}(\mathbf{U}, \mathbf{V}) &= \widehat{\mathbb{E}}_i \mathbb{E}_{\mathbf{b}}[(y_i - \langle \mathbf{U} \mathbf{B} \mathbf{V}^\top, \mathbf{A}^{(i)} \rangle)^2] \\ &= \widehat{\mathbb{E}}_i (y_i - \langle \mathbf{U} \mathbf{V}^\top, \mathbf{A}^{(i)} \rangle)^2 + \widehat{\mathbb{E}}_i \frac{p}{1 - p} \sum_{j=1}^{d_1} (\mathbf{u}_j^\top \mathbf{A}^{(i)} \mathbf{v}_j)^2 \\ &= \widehat{L}(\mathbf{U}, \mathbf{V}) + \frac{p}{1 - p} \widehat{R}(\mathbf{U}, \mathbf{V}). \end{split}$$

which completes the proof.

**Lemma 2** (Concentration in matrix completion). For  $\ell \in [n]$ , let  $A^{(\ell)}$  be an indicator matrix whose (i,j)-th element is selected according to some distribution. Assume U, V is such that  $||U^{\top}||_{2,\infty}||V||_{\infty,\infty} \leq \gamma$ . Then, with probability at least  $1 - \delta$  over a sample of size n, we have that

$$|R(U, V) - \widehat{R}(U, V)| \le \frac{C\gamma^2 \sqrt{\log(2/\delta)}}{\sqrt{n}}.$$

Proof of Lemma 2. Define  $X_{\ell} := \sum_{w=1}^{d_1} (\mathbf{u}_w^{\top} \mathbf{A}^{(\ell)} \mathbf{v}_w)^2$  and observe that

$$\begin{split} X_{\ell} &= \sum_{w=1}^{d_1} \left( \sum_{i,j} \mathbf{U}_{iw} \mathbf{V}_{jw} \mathbf{A}_{ij}^{(\ell)} \right)^2 = \sum_{w=1}^{d_1} \sum_{i,i',j,j'} \mathbf{U}_{iw} \mathbf{U}_{i'w} \mathbf{V}_{jw} \mathbf{A}_{ij}^{(\ell)} \mathbf{A}_{i'j'}^{(\ell)} \\ &= \sum_{w=1}^{d_1} \sum_{i,j} \mathbf{U}_{iw}^2 \mathbf{V}_{jw}^2 \mathbf{A}_{ij}^{(\ell)} \leq \max_{i,j} \sum_{w=1}^{d_1} \mathbf{U}_{iw}^2 \mathbf{V}_{jw}^2 \\ &\leq \max_{i,j} \|\mathbf{U}(i,:)\|^2 \|\mathbf{V}(j,:)\|_{\infty}^2 = \|\mathbf{U}^{\top}\|_{2,\infty}^2 \|\mathbf{V}\|_{\infty,\infty}^2 \leq \gamma^2 \end{split}$$

where the third equality follows because for an indicator matrix  $A^{(\ell)}$ , it holds that  $A^{(\ell)}_{ij}A^{(\ell)}_{i'j'}=0$  if  $(i,j)\neq (i',j')$ . Thus,  $X_{w,\ell}$  is a sub-Gaussian (more strongly, bounded) random variable with mean  $\mathbb{E}[X_\ell]=R(U,V)$  and sub-Gaussian norm  $\|X_\ell\|_{\psi_2}\leq \gamma^2/\ln(2)$ . Furthermore,  $\|X_\ell-R(U,V)\|_{\psi_2}\leq C'\|X_\ell\|_{\psi_2}\leq C\gamma^2$ , for some absolute constants C', C (Lemma 2.6.8 of Vershynin (2018)). Using Theorem 4, for  $t=Cd_1\sqrt{\frac{\log 2/\delta}{n}}$  we get that:

$$\mathbb{P}\left(\left|\widehat{R}(\mathbf{U},\mathbf{V}) - R(\mathbf{U},\mathbf{V})\right| \ge t\right) = \mathbb{P}\left(\left|\frac{1}{n}\sum_{\ell=1}^{n}X_{\ell} - R(\mathbf{U},\mathbf{V})\right| \ge C\gamma^{2}\sqrt{\frac{\log 2/\delta}{n}}\right) \le \delta$$

which completes the proof.

**Proposition 3.** [Induced regularizer] For  $j \in [n]$ , let  $A^{(j)}$  be an indicator matrix whose (i,k)-th element is selected randomly with probability p(i)q(k), where p(i) and q(k) denote the probability of choosing the i-th row and the k-th column. Then  $\Theta(M) = \frac{1}{d_1} \|\operatorname{diag}(\sqrt{p})UV^{\top}\operatorname{diag}(\sqrt{q})\|_*^2$ .

Proof of Proposition 3. For any pair of factors (U, V) it holds that

$$R(\mathbf{U}, \mathbf{V}) = \sum_{i=1}^{d_1} \mathbb{E}(\mathbf{u}_i^{\top} \mathbf{A} \mathbf{v}_i)^2 = \sum_{i=1}^{d_1} \sum_{j=1}^{d_2} \sum_{k=1}^{d_0} p(j) q(k) (\mathbf{u}_i^{\top} \mathbf{e}_j \mathbf{e}_k^{\top} \mathbf{v}_i)^2$$

$$= \sum_{i=1}^{d_1} \sum_{j=1}^{d_2} \sum_{k=1}^{d_0} p(j) q(k) \mathbf{U}(j, i)^2 \mathbf{V}(k, i)^2 = \sum_{i=1}^{d_1} \|\operatorname{diag}(\sqrt{p}) \mathbf{u}_i\|^2 \|\operatorname{diag}(\sqrt{q}) \mathbf{v}_i\|^2$$

We can now lower bound the right hand side above as follows:

$$\begin{split} R(\mathbf{U}, \mathbf{V}) &\geq \frac{1}{d_1} \left( \sum_{i=1}^{d_1} \|\operatorname{diag}(\sqrt{p}) \mathbf{u}_i\| \|\operatorname{diag}(\sqrt{q}) \mathbf{v}_i\| \right)^2 \\ &= \frac{1}{d_1} \left( \sum_{i=1}^{d_1} \|\operatorname{diag}(\sqrt{p}) \mathbf{u}_i \mathbf{v}_i^\top \operatorname{diag}(\sqrt{q}) \|_* \right)^2 \\ &\geq \frac{1}{d_1} \left( \|\operatorname{diag}(\sqrt{p}) \sum_{i=1}^{d_1} \mathbf{u}_i \mathbf{v}_i^\top \operatorname{diag}(\sqrt{q}) \|_* \right)^2 = \frac{1}{d_1} \|\operatorname{diag}(\sqrt{p}) \mathbf{U} \mathbf{V}^\top \operatorname{diag}(\sqrt{q}) \|_*^2 \end{split}$$

where the first inequality is due to Cauchy-Schwartz and the second inequality follows from the triangle inequality. The equality right after the first inequality follows from the fact that for any two vectors  $a, b, \|ab^\top\|_* = \|ab^\top\| = \|a\|\|b\|$ . Since the inequalities hold for any U, V, it implies that

$$\Theta(\mathbf{U}\mathbf{V}^{\top}) \ge \frac{1}{d_1} \|\operatorname{diag}(\sqrt{p})\mathbf{U}\mathbf{V}^{\top}\operatorname{diag}(\sqrt{q})\|_*^2.$$

Applying Theorem 8 on  $(\operatorname{diag}(\sqrt{p})U, \operatorname{diag}(\sqrt{p})V)$ , there exist a rotation matrix Q such that

$$\|\operatorname{diag}(\sqrt{p})\operatorname{Uq}_i\|\|\operatorname{diag}(\sqrt{q})\operatorname{Vq}_i\| = \frac{1}{d_1}\|\operatorname{diag}(\sqrt{p})\operatorname{UV}^\top\operatorname{diag}(\sqrt{q})\|_*$$

We evaluate the expected dropout regularizer at UQ, VQ:

$$\begin{split} R(\mathbf{U}\mathbf{Q},\mathbf{V}\mathbf{Q}) &= \sum_{i=1}^{d_1} \|\operatorname{diag}(\sqrt{p})\mathbf{U}\mathbf{q}_i\|^2 \|\operatorname{diag}(\sqrt{q})\mathbf{V}\mathbf{q}_i\|^2 \\ &= \sum_{i=1}^{d_1} \frac{1}{d_1^2} \|\operatorname{diag}(\sqrt{p})\mathbf{U}\mathbf{V}^{\top}\operatorname{diag}(\sqrt{q})\|_*^2 = \frac{1}{d_1} \|\operatorname{diag}(\sqrt{p})\mathbf{U}\mathbf{V}^{\top}\operatorname{diag}(\sqrt{q})\|_*^2 \leq \Theta(\mathbf{U}\mathbf{V}^{\top}) \end{split}$$

which completes the proof of the first part.

*Proof of Theorem 1.* We use Theorem 6 to bound the population risk in terms of the Rademacher complexity of the target class. Define the class of predictors with weighted trace-norm bounded by  $\sqrt{\alpha}$ , i.e.

$$\mathcal{M}_{\alpha} = \{ M : \| \operatorname{diag}(\sqrt{p}) M \operatorname{diag}(\sqrt{q}) \|_{*}^{2} \le \alpha \}.$$

In particular dropout empirical risk minimizers U, V belong to this class:

$$\|\operatorname{diag}(\sqrt{p})\operatorname{UV}^{\top}\operatorname{diag}(\sqrt{q})\|_{*}^{2} = d_{1}\Theta(\operatorname{UV}^{\top}) \leq d_{1}R(\operatorname{U},\operatorname{V}) \leq \alpha$$

where the first inequality holds by definition of the induced regularizer, and the second inequality follows from the assumption of the theorem. Since g is a contraction, by Talagrand's lemma and Theorem 9, we have that  $\Re_n(g \circ \mathcal{M}_\alpha) \leq \Re_n(\mathcal{M}_\alpha) \leq \sqrt{\frac{\alpha d_2 \log(d_2)}{n}}$ . To obtain the maximum deviation parameter M in Theorem 6, we note that the assumption  $\|\mathbf{M}_*\| \leq 1$  implies that  $|\mathbf{M}_*(i,j)| \leq 1$  for all i,j, so that  $g(\mathbf{M}_*) = \mathbf{M}_*$ . We have that:

$$\max_{\mathbf{A}} |\langle \mathbf{M}_* - g(\mathbf{U}\mathbf{V}^\top), \mathbf{A} \rangle| = \max_{i,j} |\langle \mathbf{M}_* - g(\mathbf{U}\mathbf{V}^\top), \mathbf{e}_i \mathbf{e}_j^\top \rangle| \leq \max_{i,j} |\mathbf{M}_*(i,j)| + \max_{i,j} |\langle \mathbf{U}\mathbf{V}^\top, \mathbf{e}_i \mathbf{e}_j^\top \rangle| \leq \|\mathbf{M}_*\| + 1 \leq 2$$

Let  $L(g(\mathbf{U}\mathbf{V}^{\top})) := \mathbb{E}(y - \langle g(\mathbf{U}\mathbf{V}^{\top}), \mathbf{A} \rangle)^2$  and  $\widehat{L}(g(\mathbf{U}\mathbf{V}^{\top})) := \widehat{\mathbb{E}}_i(y_i - \langle g(\mathbf{U}\mathbf{V}^{\top}), \mathbf{A}^{(i)} \rangle)^2$  denote the true risk and the empirical risk of  $g(\mathbf{U}\mathbf{V}^{\top})$ , respectively. Plugging the above results in Theorem 6, we get

$$L(g(\mathbf{U}, \mathbf{V})) \le \widehat{L}(g(\mathbf{U}, \mathbf{V})) + 8\mathfrak{R}_n(g \circ \mathcal{M}_\alpha) + 4\sqrt{\frac{\log(2/\delta)}{2n}}$$
$$\le \widehat{L}(\mathbf{U}, \mathbf{V}) + 8\sqrt{\frac{\alpha d_2 \log(d_2)}{n}} + 4\sqrt{\frac{\log(2/\delta)}{2n}}$$
$$\le \widehat{L}(\mathbf{U}, \mathbf{V}) + 8\sqrt{\frac{2\alpha d_2 \log(d_2) + \frac{1}{4}\log(2/\delta)}{n}}$$

where the second inequality holds since  $\widehat{L}(g(U, V)) \leq \widehat{L}(U, V)$ .

## **B.1** Optimistic Rates

As we discussed in the main text, under additional assumptions on the value of  $\alpha$ , it is possible to give optimistic generalization bounds that decay as  $\tilde{O}(\alpha d_2/n)$ . This result is given as the following theorem.

**Theorem 10.** Assume that  $d_2 \ge d_0$  and  $||\mathbf{M}_*|| \le 1$ . Furthermore, assume that  $\min_{i,k} p(i)q(k) \ge \frac{\log(d_2)}{n\sqrt{d_2d_0}}$ . Let  $(\mathbf{U}, \mathbf{V})$  be a minimizer of the dropout ERM objective in equation (3). Let  $\alpha$  be such that  $\max\{R(\mathbf{U}, \mathbf{V}), \Theta(\mathbf{M}_*)\} \le 1$ 

 $\alpha/d_1$ . Then, for any  $\delta \in (0,1)$ , the following generalization bounds holds with probability at least  $1-\delta$  over a sample of size n:

$$L(g(\mathbf{U}\mathbf{V}^{\top})) \le \frac{2K\log(n)^3 \alpha d_2 \log(d_2) + 4K\log(1/\delta)}{n}$$

where K is an absolute constant Srebro et al. (2010), g(M) thresholds M between [-1,1], and  $L(g(UV^{\top})) := \mathbb{E}(y - \langle g(UV^{\top}), A \rangle)^2$  is the true risk of  $g(UV^{\top})$ .

Proof of Theorem 10. We use Theorem 7 to bound the population risk in terms of the Rademacher complexity of the target class. Define the class of predictors with weighted trace-norm bounded by  $\sqrt{\alpha}$ , i.e.

$$\mathcal{M}_{\alpha} = \{ M : \| \operatorname{diag}(\sqrt{p}) \operatorname{M} \operatorname{diag}(\sqrt{q}) \|_{*}^{2} \le \alpha \}.$$

In particular dropout empirical risk minimizers U, V belong to this class:

$$\|\operatorname{diag}(\sqrt{p})UV^{\top}\operatorname{diag}(\sqrt{q})\|_{*}^{2} = d_{1}\Theta(UV^{\top}) \leq d_{1}R(U,V) \leq \alpha$$

where the first inequality holds by definition of the induced regularizer, and the second inequality follows from the assumption of the theorem. Moreover, by assumption  $\Theta(M_*) \leq \alpha$ , we have that  $M_* \in \mathcal{M}_{\alpha}$ . With this, we get that

$$L_* := \min_{M \in g \circ \mathcal{M}_{\alpha}} L(M) \le L(g(M_*)) = L(g(M_*)) = 0.$$

Since g is a contraction, by Talagrand's lemma and Theorem 9, we have that  $\mathfrak{R}_n(g \circ \mathcal{M}_\alpha) \leq \mathfrak{R}_n(\mathcal{M}_\alpha) \leq \sqrt{\frac{\alpha d_2 \log(d_2)}{n}}$ . Plugging the above in Theorem 6, we get

$$L(g(\mathbf{U}, \mathbf{V})) \le 2K\log(n)^3 \mathfrak{R}_n^2(g \circ \mathcal{M}_\alpha) + \frac{4K\log\frac{1}{\delta}}{n}$$
$$\le \frac{2K\log(n)^3 \alpha d_2 \log(d_2) + 4K\log\frac{1}{\delta}}{n}$$

## C Non-linear Neural Networks

Proposition 4 (Dropout regularizer in deep regression).

$$\widehat{L}_{drop}(\mathbf{w}) = \widehat{L}(\mathbf{w}) + \widehat{R}(\mathbf{w}), \quad where \quad \widehat{R}(\mathbf{w}) = \lambda \sum_{j=1}^{d_1} \|u_j\|^2 \widehat{a}_j^2.$$

where  $\hat{a}_j = \sqrt{\hat{\mathbb{E}}_i a_j(\mathbf{x}_i)^2}$  and  $\lambda = \frac{p}{1-p}$  is the regularization parameter.

Proof of Proposition 4. Similar statements and proofs can be found in several previous works Srivastava et al. (2014); Wang & Manning (2013); Cavazza et al. (2018); Mianjy et al. (2018). Here we include a proof for completeness. Recall that  $\mathbb{E}[B_{ii}] = 1$  and  $\text{Var}(B_{ii}) = \frac{p}{1-p}$ . Conditioned on x, y in the current mini-batch, we have that:

$$\mathbb{E}_{\mathbf{B}} \|\mathbf{y} - \mathbf{U}^{\top} \mathbf{B} \mathbf{a}(\mathbf{x}) \|^{2} = \sum_{i=1}^{d_{2}} (\mathbb{E}_{\mathbf{B}} [y_{i} - \mathbf{u}_{i}^{\top} \mathbf{B} \mathbf{a}(\mathbf{x})])^{2} + \sum_{i=1}^{d_{2}} \mathbf{Var}(y_{i} - \mathbf{u}_{i}^{\top} \mathbf{B} \mathbf{a}(\mathbf{x}))$$

Since  $\mathbb{E}[B] = I$ , the first term on right hand side is equal to  $\|y - U^{\top}a(x)\|^2$ . For the second term we have

$$\sum_{i=1}^{d_2} \operatorname{Var}(y_i - \mathbf{u}_i^{\top} \operatorname{Ba}(\mathbf{x})) = \sum_{i=1}^{d_2} \operatorname{Var}(\sum_{j=1}^{d_1} \operatorname{U}_{j,i} \operatorname{B}_{jj} a_j(\mathbf{x})) = \sum_{i=1}^{d_2} \sum_{j=1}^{d_1} (\operatorname{U}_{j,i} a_j(\mathbf{x}))^2 \operatorname{Var}(\operatorname{B}_{jj}) = \frac{p}{1-p} \sum_{j=1}^{d_1} \|\mathbf{u}_j\|^2 a_j(\mathbf{x})^2$$

Thus, conditioned on the sample (x, y), we have that

$$\mathbb{E}_{\mathbf{B}}[\|\mathbf{y} - \mathbf{U}^{\top} \mathbf{B} \mathbf{a}(\mathbf{x})\|^{2}] = \|\mathbf{y} - \mathbf{U}^{\top} \mathbf{a}(\mathbf{x})\|^{2} + \frac{p}{1 - p} \sum_{i=1}^{d_{1}} \|\mathbf{u}_{i}\|^{2} a_{j}(\mathbf{x})^{2}$$

Now taking the empirical average with respect to x, y, we get

$$\widehat{L}_{\text{drop}}(\mathbf{w}) = \widehat{L}(\mathbf{w}) + \frac{p}{1-p} \sum_{j=1}^{d_1} \|\mathbf{u}_j\|^2 \widehat{a}_j^2 = \widehat{L}(\mathbf{w}) + \widehat{R}(\mathbf{w})$$

which completes the proof.

**Proposition 5.** Consider a two layer neural network  $f_w(\cdot)$  with ReLU activation functions in the hidden layer. Furthermore, assume that the marginal input distribution  $\mathbb{P}_{\mathcal{X}}(\mathbf{x})$  is symmetric and isotropic, i.e.,  $\mathbb{P}_{\mathcal{X}}(\mathbf{x}) = \mathbb{P}_{\mathcal{X}}(-\mathbf{x})$  and  $\mathbb{E}[\mathbf{x}\mathbf{x}^{\top}] = I$ . Then the following holds for the expected 5 due to dropout:

$$R(\mathbf{w}) := \mathbb{E}[\widehat{R}(\mathbf{w})] = \frac{\lambda}{2} \sum_{i_0, i_1, i_2 = 1}^{d_0, d_1, d_2} U(i_1, i_2)^2 V(i_1, i_0)^2, \tag{8}$$

Proof of Proposition 5. Using Proposition 4, we have that:

$$R(\mathbf{w}) = \mathbb{E}[\widehat{R}(\mathbf{w})] = \lambda \sum_{j=1}^{d_1} \|\mathbf{u}_j\|^2 \mathbb{E}[\sigma(\mathbf{V}(j,:)^{\top}\mathbf{x})^2]$$

It remains to calculate the quantity  $\mathbb{E}_{\mathbf{x}}[\sigma(\mathbf{V}(j,:)^{\top}\mathbf{x})^{2}]$ . By symmetry assumption, we have that  $\mathbb{P}_{\mathcal{X}}(\mathbf{x}) = \mathbb{P}_{\mathcal{X}}(-\mathbf{x})$ . As a result, for any  $\mathbf{v} \in \mathbb{R}^{d_{0}}$ , we have that  $\mathbb{P}(\mathbf{v}^{\top}\mathbf{x}) = \mathbb{P}(-\mathbf{v}^{\top}\mathbf{x})$  as well. That is, the random variable  $z_{j} := \mathbf{W}_{1}(j,:)^{\top}\mathbf{x}$  is also symmetric about the origin. It is easy to see that  $\mathbb{E}_{z}[\sigma(z)^{2}] = \frac{1}{2}\mathbb{E}_{z}[z^{2}]$ .

$$\mathbb{E}_{z}[\sigma(z)^{2}] = \int_{-\infty}^{\infty} \sigma(z)^{2} d\mu(z) = \int_{0}^{\infty} \sigma(z)^{2} d\mu(z) = \int_{0}^{\infty} z^{2} d\mu(z) = \frac{1}{2} \int_{\infty}^{\infty} z^{2} d\mu(z) = \frac{1}{2} \mathbb{E}_{z}[z^{2}].$$

Plugging back the above identity in the expression of  $R(\mathbf{w})$ , we get that

$$R(\mathbf{w}) = \lambda \sum_{j=1}^{d_1} \|\mathbf{u}_j\|^2 \mathbb{E}[(\mathbf{V}(j,:)^\top \mathbf{x})^2] = \frac{\lambda}{2} \sum_{j=1}^{d_1} \|\mathbf{u}_j\|^2 \|\mathbf{V}(j,:)\|^2$$

where the second equality follows from the assumption that the distribution is isotropic.

*Proof of Proposition 1.* For  $\delta \in (0, \frac{1}{2})$ , consider the following random variable:

$$\mathbf{x} = \begin{cases} [1;0] & \text{with probability } \delta \\ [\frac{-\delta}{1-\delta}; \frac{\sqrt{1-2\delta}}{1-\delta}] & \text{with probability } \frac{1-\delta}{2} \\ [\frac{-\delta}{1-\delta}; -\frac{\sqrt{1-2\delta}}{1-\delta}] & \text{with probability } \frac{1-\delta}{2} \end{cases}$$

It is easy to check that the x has zero mean and is supported on the unit sphere. Consider the vector  $\mathbf{w} = [\frac{1}{\sqrt{\delta}}; 0]$ . It is easy to check that x satisfies  $R(\mathbf{w}) = \sqrt{\mathbb{E}\sigma(\mathbf{w}^{\top}\mathbf{x})^2} = 1$ ; however, for any given C, it holds that  $\|\mathbf{w}\| \geq C$  as long as we let  $\delta = C^2$ .

Proof of Theorem 2. For any  $j \in [h]$ , let  $a_j^2 := \mathbb{E}[\sigma(\mathbf{v}_j^\top \mathbf{x})^2]$  denote the average squared activation of the j-th node with respect to the input distribution. Given n i.i.d. samples  $\mathcal{S} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ , the empirical Rademahcer complexity is bounded as follows:

$$\mathfrak{R}_{\mathcal{S}}(\mathcal{F}_{\alpha}) = \mathbb{E}_{\zeta} \sup_{f_{\{\mathbf{u},\mathbf{v}\}} \in \mathcal{F}_{\alpha}} \frac{1}{n} \sum_{j=1}^{h} u_{j} a_{j} \sum_{i=1}^{n} \zeta_{i} \frac{\sigma(\mathbf{v}_{j}^{\top} \mathbf{x}_{i})}{a_{j}}$$

$$\leq \mathbb{E}_{\zeta} \sup_{f_{\{\mathbf{u},\mathbf{v}\}} \in \mathcal{F}_{\alpha}} \frac{1}{n} \sum_{j=1}^{h} |u_{j} a_{j}| \left| \sum_{i=1}^{n} \zeta_{i} \frac{\sigma(\mathbf{v}_{j}^{\top} \mathbf{x}_{i})}{a_{j}} \right|$$

$$\leq \mathbb{E}_{\zeta} \left[ \left( \sup_{f_{\{\mathbf{u},\mathbf{v}\}} \in \mathcal{F}_{\alpha}} \sum_{j=1}^{h} |u_{j} a_{j}| \right) \left( \sup_{\mathbf{v}} \max_{j \in [h]} \left| \frac{1}{n} \sum_{i=1}^{n} \zeta_{i} \frac{\sigma(\mathbf{v}_{j}^{\top} \mathbf{x}_{i})}{a_{j}} \right| \right) \right]$$

where we used the fact that the supremum of product of positive functions is upperbounded by the product of the supremums. By definition of  $\mathcal{F}_{\alpha}$ , the first term on the right hand side is bounded by  $\alpha$ . To bound the second term in the right hand side, we note that the maximum over rows of  $V^{\top}$  can be absorbed into the supremum.

$$\frac{1}{n} \mathbb{E}_{\zeta} \sup_{\mathbf{v}} \left| \sum_{i=1}^{n} \zeta_{i} \frac{\sigma(\mathbf{v}^{\top} \mathbf{x}_{i})}{\sqrt{\mathbb{E}[\sigma(\mathbf{v}^{\top} \mathbf{x})^{2}]}} \right| = \frac{1}{n} \mathbb{E}_{\zeta} \sup_{\mathbb{E}[\sigma(\mathbf{v}^{\top} \mathbf{x})^{2}] \leq 1} \left| \sum_{i=1}^{n} \zeta_{i} \sigma(\mathbf{v}^{\top} \mathbf{x}_{i}) \right| \\
\leq \frac{2}{n} \mathbb{E}_{\zeta} \sup_{\mathbb{E}[\sigma(\mathbf{v}^{\top} \mathbf{x})^{2}] \leq 1} \sum_{i=1}^{n} \zeta_{i} \sigma(\mathbf{v}^{\top} \mathbf{x}_{i}) \\
\leq \frac{2}{n} \mathbb{E}_{\zeta} \sup_{\beta \mathbb{E}(\mathbf{v}^{\top} \mathbf{x})^{2} \leq 1} \sum_{i=1}^{n} \zeta_{i} \sigma(\mathbf{v}^{\top} \mathbf{x}_{i}) \tag{\beta-retentiveness}$$

Let  $C^{\dagger}$  be the pseudo-inverse of C. We perform the following change the variable:  $w \leftarrow C^{-\dagger/2}v$ .

$$R.H.S. \leq \frac{2}{n} \mathbb{E}_{\zeta} \sup_{\mathbb{E}[(\mathbf{w}^{\top} \mathbf{C}^{\dagger/2} \mathbf{x})^{2}] \leq 1/\beta} \sum_{i=1}^{n} \zeta_{i} \mathbf{w}^{\top} \mathbf{C}^{\dagger/2} \mathbf{x}_{i}$$

$$= \frac{2}{n} \mathbb{E}_{\zeta} \sup_{\|\mathbf{w}\|^{2} \leq 1/\beta} \langle \mathbf{w}, \sum_{i=1}^{n} \zeta_{i} \mathbf{C}^{\dagger/2} \mathbf{x}_{i} \rangle$$

$$= \frac{2}{n\sqrt{\beta}} \mathbb{E}_{\zeta} \| \sum_{i=1}^{n} \zeta_{i} \mathbf{C}^{\dagger/2} \mathbf{x}_{i} \|$$

$$\leq \frac{2}{n\sqrt{\beta}} \sqrt{\mathbb{E}_{\zeta} \| \sum_{i=1}^{n} \zeta_{i} \mathbf{C}^{\dagger/2} \mathbf{x}_{i} \|^{2}} = \frac{2}{n\sqrt{\beta}} \sqrt{\sum_{i=1}^{n} \mathbf{x}_{i}^{\top} \mathbf{C}^{\dagger} \mathbf{x}_{i}}$$

where the last inequality holds due to Jensen's inequality. To bound the expected Rademacher complexity, we take the expected value of both sides with respected to sample S, which gives the following:

$$\mathfrak{R}_n(\mathcal{F}_{\alpha}) = \mathbb{E}_{\mathbf{x}}[\mathfrak{R}_{\mathcal{S}}(\mathcal{F}_{\alpha})] \leq \frac{2}{n\sqrt{\beta}} \mathbb{E}_{\mathcal{S}} \sqrt{\sum_{i=1}^{n} \mathbf{x}_i^{\top} \mathbf{C}^{\dagger} \mathbf{x}_i} \leq \frac{2}{n\sqrt{\beta}} \sqrt{\sum_{i=1}^{n} \mathbb{E}_{\mathbf{x}_i}[\mathbf{x}_i^{\top} \mathbf{C}^{\dagger} \mathbf{x}_i]},$$

where the last inequality holds again due to Jensen's inequality. Finally, we have that  $\mathbb{E}_{\mathbf{x}_i}\mathbf{x}_i^{\mathsf{T}}\mathbf{C}^{\dagger}\mathbf{x}_i = \mathbb{E}_{\mathbf{x}_i}\langle\mathbf{x}_i\mathbf{x}_i^{\mathsf{T}},\mathbf{C}^{\dagger}\rangle = \langle\mathbf{C},\mathbf{C}^{\dagger}\rangle = \mathrm{Rank}(\mathbf{C})$ , which completes the proof of the Theorem.

*Proof of Theorem 3.* For simplicity, assume that the width of the hidden layer is even. Consider the linear function class:

$$\mathcal{G}_r := \{g_{\mathbf{w}} : \mathbf{x} \mapsto \mathbf{w}^\top \mathbf{x}, \ \mathbb{E}(\mathbf{w}^\top \mathbf{x})^2 \le d_1 r / 2\}.$$

Recall that  $\mathcal{H}_r := \{h_w : x \mapsto u^\top \sigma(V^\top x), \ R(u, V) \leq r\}$ . First, we argue that  $\mathcal{G}_r \subset \mathcal{H}_r$ . Let  $g_w \in \mathcal{G}_r$ ; we show that there exist u, V such that  $g_w = f_{u,V}$  and  $f_{u,V} \in \mathcal{H}_r$ . Define  $u := \frac{2}{d_1}[1; -1; \cdots 1; -1] \in \mathbb{R}^{d_1}$ , and let  $V = w(e_1 - e_2 + e_3 - e_4 + \cdots + e_{d_1-1} - e_{d_1})^\top$ , where  $e_i \in \mathbb{R}^{d_1}$  is the *i*-th standard basis vector. It's easy to see that

$$f_{\mathbf{u},\mathbf{V}}(\mathbf{x}) = \mathbf{u}^{\top} \sigma(\mathbf{V}^{\top} \mathbf{x}) = \sum_{i=1}^{d_1} u_i \sigma(\mathbf{v}_i^{\top} \mathbf{x})$$

$$= \sum_{i=1}^{d_1} \frac{2}{d_1} (-1)^{i-1} \sigma(\mathbf{v}_i^{\top} \mathbf{x})$$

$$= \sum_{i=1}^{d_1/2} \frac{2}{d_1} (\sigma(\mathbf{v}_{2i-1}^{\top} \mathbf{x}) - \sigma(\mathbf{v}_{2i}^{\top} \mathbf{x}))$$

$$= \sum_{i=1}^{d_1/2} \frac{2}{d_1} (\sigma(\mathbf{w}^{\top} \mathbf{x}) - \sigma(-\mathbf{w}^{\top} \mathbf{x})) = \mathbf{w}^{\top} \mathbf{x} = g_{\mathbf{w}}.$$

Furthermore, it holds for the explicit regularizer that

$$\begin{split} R(\mathbf{u}, \mathbf{V}) &= \sum_{i=1}^{d_1} u_i^2 \mathbb{E} \sigma(\mathbf{v}_i^\top \mathbf{x})^2 = \sum_{i=1}^{d_1/2} \frac{4}{d_1^2} \left( \mathbb{E} \sigma(\mathbf{v}_{2i-1}^\top \mathbf{x})^2 + \mathbb{E} \sigma(\mathbf{v}_{2i}^\top \mathbf{x})^2 \right) \\ &= \sum_{i=1}^{d_1/2} \frac{4}{d_1^2} \mathbb{E} [\sigma(\mathbf{w}^\top \mathbf{x})^2 + \sigma(-\mathbf{w}^\top \mathbf{x})^2] \\ &= \frac{2}{d_1} \mathbb{E}(\mathbf{w}^\top \mathbf{x})^2 \leq r \end{split}$$

Thus, we have that  $\mathcal{G}_r \subset \mathcal{H}_r$ , and the following inequalities follow.

$$\mathfrak{R}_{\mathcal{S}}(\mathcal{H}_r) \geq \mathfrak{R}_{\mathcal{S}}(\mathcal{G}_r) = \mathbb{E}_{\epsilon_i} \sup_{g_{\mathbf{w}} \in \mathcal{G}_r} \frac{1}{n} \sum_{i=1}^n \epsilon_i g_{\mathbf{w}}(\mathbf{x}_i)$$

$$= \mathbb{E}_{\epsilon_i} \sup_{\mathbb{E}(\mathbf{w}^{\top} \mathbf{x})^2 \leq d_1 r/2} \frac{1}{n} \sum_{i=1}^n \epsilon_i \mathbf{w}^{\top} \mathbf{x}_i$$

$$= \mathbb{E}_{\epsilon_i} \sup_{\mathbf{w}^{\top} \mathbf{C} \mathbf{w} \leq d_1 r/2} \frac{1}{n} \langle \mathbf{w}, \sum_{i=1}^n \epsilon_i \mathbf{x}_i \rangle$$

$$= \mathbb{E}_{\epsilon_i} \sup_{\|\mathbf{C}^{1/2} \mathbf{w}\|^2 \leq d_1 r/2} \frac{1}{n} \langle \mathbf{C}^{1/2} \mathbf{w}, \sum_{i=1}^n \epsilon_i \mathbf{C}^{-\dagger/2} \mathbf{x}_i \rangle$$

$$= \frac{\sqrt{d_1 r}}{\sqrt{2n}} \mathbb{E}_{\epsilon_i} \| \sum_{i=1}^n \epsilon_i \mathbf{C}^{\dagger/2} \mathbf{x}_i \|$$

$$\geq \frac{c\sqrt{d_1 r}}{\sqrt{2n}} \sqrt{\sum_{i=1}^n \|\mathbf{C}^{\dagger/2} \mathbf{x}_i\|^2} = \frac{c\sqrt{d_1 r} \|\mathbf{X}\|_{\mathbf{C}^{\dagger}}}{\sqrt{2n}}$$

where the last inequality follows from Khintchine-Kahane inequality in Lemma 1.

Next, we define some function classes that will be used frequently in the proofs.

**Definition 1.** For any closed subset  $[a,b] \subset \mathbb{R}$ , let  $\Pi_{[a,b]}(y) := \max\{a, \min\{b,y\}\}$ . For z := (x,y) and  $f : \mathcal{X} \to \mathcal{Y}$ , define the squared loss  $\ell_2(f,z) := (1-yf(x))^2$ . For a given value  $\alpha > 0$ , consider the following classes

$$\mathcal{W}_{\alpha} := \{ \mathbf{w} = (u, V) \in \mathbb{R}^{d_1} \times \mathbb{R}^{d_0 \times d_1}, \sum_{i=1}^{d_1} |u_i| \sqrt{\mathbb{E}\sigma(v_i^\top \mathbf{x})^2} \le \alpha \}$$

$$\mathcal{F}_{\alpha} := \{ f_{\mathbf{w}} : \mathbf{x} \mapsto u^\top \sigma(V^\top \mathbf{x}), \ \mathbf{w} \in \mathcal{W}_{\alpha} \},$$

$$\mathcal{G}_{\alpha} := \Pi_{[-1,1]} \circ \mathcal{F}_{\alpha} = \{ g_{\mathbf{w}} = \Pi_{[-1,-1]} \circ f_{\mathbf{w}}, \ f_{\mathbf{w}} \in \mathcal{F}_{\alpha} \}$$

$$\mathcal{L}_{\alpha} := \{ \ell_2 : (g_{\mathbf{w}}, z) \mapsto (y - g_{\mathbf{w}}(\mathbf{x}))^2, \ g_{\mathbf{w}} \in \mathcal{G}_{\alpha} \}$$

**Lemma 3.** Let  $W_{\alpha}$ ,  $\mathcal{F}_{\alpha}$ ,  $\mathcal{G}_{\alpha}$ ,  $\mathcal{L}_{\alpha}$  be as defined in Definition 1. Then the following holds true:

- 1.  $\Re_{\mathcal{S}}(\mathcal{G}_{\alpha}) \leq \Re_{\mathcal{S}}(\mathcal{F}_{\alpha})$ .
- 2. If  $\mathcal{Y} = \{-1, +1\}$  (binary classification), then it holds that  $\mathfrak{R}_{\mathcal{S}}(\mathcal{L}_{\alpha}) \leq 2\mathfrak{R}_{\mathcal{S}}(\mathcal{G}_{\alpha})$ .

*Proof.* Since  $\Pi_{[-1,-1]}(\cdot)$  is 1-Lipschitz, by Talagrand's contraction lemma, we have that  $\mathfrak{R}_{\mathcal{S}}(\mathcal{G}_{\alpha}) \leq \mathfrak{R}_{\mathcal{S}}(\mathcal{F}_{\alpha})$ . The second claim follows from

$$\mathfrak{R}_{\mathcal{S}}(\mathcal{L}_{\alpha}) = \mathbb{E}_{\zeta} \sup_{\mathbf{w} \in \mathcal{W}} \frac{1}{n} \sum_{i=1}^{n} \zeta_{i} (y_{i} - g_{\mathbf{w}}(\mathbf{x}_{i}))^{2}$$

$$= \mathbb{E}_{\zeta} \sup_{\mathbf{w} \in \mathcal{W}} \frac{1}{n} \sum_{i=1}^{n} \zeta_{i} (1 - y_{i} g_{\mathbf{w}}(\mathbf{x}_{i}))^{2} \qquad (y_{i} \in \{-1, +1\})$$

$$\leq 2\mathbb{E}_{\zeta} \sup_{\mathbf{w} \in \mathcal{W}} \frac{1}{n} \sum_{i=1}^{n} \zeta_{i} y_{i} g_{\mathbf{w}}(\mathbf{x}_{i})$$

$$= 2\mathbb{E}_{\zeta} \sup_{\mathbf{w} \in \mathcal{W}} \frac{1}{n} \sum_{i=1}^{n} \zeta_{i} g_{\mathbf{w}}(\mathbf{x}_{i}) = 2\mathfrak{R}_{\mathcal{S}}(\mathcal{G}_{\alpha})$$

where the first inequality follows from Talagrand's contraction lemma due to the fact that  $h(z) = (1-z)^2$  is 2-Lipschitz for  $z \in [-1,1]$ , and the penultimate holds true since for any fixed  $(y_i)_{i=1}^n \in \{-1,+1\}^n$ , the distribution of  $(\zeta_1y_1,\ldots,\zeta_ny_n)$  is the same as that of  $(\zeta_1,\ldots,\zeta_n)$ .

*Proof of Corollary 1.* We use the standard generalization bound in Theorem 6 for class  $\mathcal{G}_{\alpha}$ :

$$L_{\mathcal{D}}(g_{\mathbf{w}}) \leq \widehat{L}_{\mathcal{S}}(g_{\mathbf{w}}) + 4M\mathfrak{R}_{\mathcal{S}}(\mathcal{G}_{\alpha}) + 3M^{2}\sqrt{\frac{\log(2/\delta)}{2n}}$$

$$\leq \widehat{L}_{\mathcal{S}}(g_{\mathbf{w}}) + 8\mathfrak{R}_{\mathcal{S}}(\mathcal{F}_{\alpha}) + 12\sqrt{\frac{\log(2/\delta)}{2n}} \qquad \text{(Lemma 3)}$$

$$\leq \widehat{L}_{\mathcal{S}}(g_{\mathbf{w}}) + \frac{16\alpha\|\mathbf{X}\|_{\mathbf{C}^{\dagger}}}{\sqrt{\beta}n} + 12\sqrt{\frac{\log(2/\delta)}{2n}} \qquad \text{(Theorem 2)}$$

where second inequality follows because the maximum deviation parameter M in Theorem 6 is bounded as

$$M = \sup_{\mathbf{w} \in \mathcal{W}} \sup_{(\mathbf{x}, y) \in \mathcal{X} \times \mathcal{Y}} |y - g_{\mathbf{w}}(\mathbf{x})| \le \sup_{\mathbf{w} \in \mathcal{W}} \sup_{(\mathbf{x}, y) \in \mathcal{X} \times \mathcal{Y}} |y| + |g_{\mathbf{w}}(\mathbf{x})| \le 2.$$

Proof of Corollary 2. Recall that the input is jointly distributed as  $(x, y) \sim \mathcal{D}$ . For  $\mathcal{X} \subseteq \mathbb{R}^{d_0}_+$ , let  $\mathcal{X}' = \mathcal{X} \cup -\mathcal{X}$  be the symmetrized input domain. Let  $\zeta$  be a Rademacher random variable. Denote the symmetrized input by  $x' = \zeta x$ , and the joint distribution of (x', y) by  $\mathcal{D}'$ . By construction,  $\mathcal{D}'$  is centrally symmetric w.r.t. x',

i.e., it holds for all  $(x, y) \in \mathcal{X} \times \mathcal{Y}$  that  $\mathcal{D}'(x, y) = \mathcal{D}'(-x, y) = \frac{1}{2}\mathcal{D}(x, y)$ . As a result, population risk with respect to the original distribution  $\mathcal{D}$  can be bounded in terms of the population risk with respect to the symmetrized distribution  $\mathcal{D}'$  as follows:

$$L_{\mathcal{D}}(f) := \mathbb{E}_{\mathcal{D}}[\ell(f(\mathbf{x}), y)]$$

$$\leq \mathbb{E}_{\mathcal{D}}[\ell(f(\mathbf{x}), y) + \ell(f(-\mathbf{x}), y)]$$

$$= 2\mathbb{E}_{\mathcal{D}}[\frac{1}{2}\ell(f(\mathbf{x}), y) + \frac{1}{2}\ell(f(-\mathbf{x}), y)]$$

$$= 2\mathbb{E}_{\mathcal{D}}\mathbb{E}_{\zeta}[\ell(f(\zeta\mathbf{x}), y) \mid \mathbf{x}, y]$$

$$= 2\mathbb{E}_{\mathcal{D}'}[\ell(f(\mathbf{x}'), y)] = 2L_{\mathcal{D}'}(f)$$
(9)

Moreover, since  $\mathcal{D}'$  is centrally symmetric, Assumption 1 holds with  $\beta = \frac{1}{2}$ . The proof of Corollary 2 follows by doubling the right hand side of inequalities in Corollary 1, and substituting  $\beta = \frac{1}{2}$ .

#### C.1 Classification

Although in the main text we only focus on the task of regression with squared loss, it is not hard to extend the results to binary classification. In particular, the following two Corollaries bound the miss-classification error in terms of the training error and the Rademacher complexity of the target class, with and without symmetrization.

Corollary 3. Consider a binary classification setting where  $\mathcal{Y} = \{-1, +1\}$ . For any  $w \in \mathcal{F}_{\alpha}$ , for any  $\delta \in (0, 1)$ , the following generalization bound holds with probability at least  $1 - \delta$  over  $\mathcal{S} = \{(x_i, y_i)\}_{i=1}^n \sim \mathcal{D}^n$ :

$$\mathbb{P}\{yf_{\mathbf{w}}(\mathbf{x}) < 0\} \le \widehat{L}_{\mathcal{S}}(g_{\mathbf{w}}) + \frac{8\alpha \|\mathbf{X}\|_{\mathbf{C}^{\dagger}}}{\sqrt{\beta}n} + 4\sqrt{\frac{\log(1/\delta)}{2n}}$$

where  $g_{\mathbf{w}}(\cdot) = \max\{-1, \min\{1, f_{\mathbf{w}}(\cdot)\}\}$  projects the network output onto the range [-1, 1].

Proof of Corollary 3. We use the standard generalization bound in Theorem 5. Recall that  $g_w = \Pi_{[-1,1]}(f_w)$ , where  $\Pi_{[-1,1]}(y) = \max\{-1, \max\{1, y\}\}$  projects onto the range [-1,1]. It is easy to bound the classification error of  $f_w$  in terms of the  $\ell_2$ -loss of  $g_w$ :

$$\mathbb{P}\{\operatorname{sgn}(f_{\mathbf{w}}(\mathbf{x})) \neq y\} = \mathbb{P}\{yf_{\mathbf{w}}(\mathbf{x}) < 0\} = \mathbb{E}[1_{yf_{\mathbf{w}}(\mathbf{x}) < 0}] = \mathbb{E}[1_{yg_{\mathbf{w}}(\mathbf{x}) < 0}] \leq \mathbb{E}(1 - yg_{\mathbf{w}}(\mathbf{x}))^2 = L_{\mathcal{D}}(g_{\mathbf{w}}). \tag{10}$$

We use Theorem 5 for class  $\frac{1}{4}\mathcal{L}_{\alpha}$  to get the generalization bound as follows:

$$\frac{1}{4}L_{\mathcal{D}}(g_{\mathbf{w}}) \leq \frac{1}{4}\widehat{L}_{\mathcal{S}}(g_{\mathbf{w}}) + 2\mathfrak{R}_{\mathcal{S}}(\frac{1}{4}\mathcal{L}_{\alpha}) + \sqrt{\frac{\log(1/\delta)}{2n}}$$

$$\Rightarrow L_{\mathcal{D}}(g_{\mathbf{w}}) \leq \widehat{L}_{\mathcal{S}}(g_{\mathbf{w}}) + 4\mathfrak{R}_{\mathcal{S}}(\mathcal{F}_{\alpha}) + 4\sqrt{\frac{\log(1/\delta)}{2n}} \qquad \text{(by Lemma 3)}$$

$$\Rightarrow L_{\mathcal{D}}(g_{\mathbf{w}}) \leq \widehat{L}_{\mathcal{S}}(g_{\mathbf{w}}) + \frac{8\alpha \|\mathbf{X}\|_{\mathbf{C}^{\dagger}}}{\sqrt{\beta}n} + 4\sqrt{\frac{\log(1/\delta)}{2n}} \qquad \text{(by Theorem 2)}$$

Corollary 4. Consider a binary classification setting where  $\mathcal{Y} = \{-1, +1\}$ . For any  $w \in \mathcal{F}'_{\alpha}$ , for any  $\delta \in (0, 1)$ , the following generalization bound holds with probability at least  $1 - \delta$  over a sample of size n and the randomization in symmetrization

$$\mathbb{P}\{yg_{\mathbf{w}}(\mathbf{x}) < 0\} \le 2\widehat{L}_{\mathcal{S}'}(g_{\mathbf{w}}) + \frac{23\alpha' \|\mathbf{X}\|_{\mathbf{C}^{\dagger}}}{n} + 8\sqrt{\frac{\log(1/\delta)}{2n}}$$

Proof of Corollary 4. Akin to proof of Corollary 2, we have that  $L_{\mathcal{D}}(f) \leq 2L_{\mathcal{D}'}(f)$ , and the marginal distribution is  $\frac{1}{2}$ -retentive. Proof of Corollary 4 follows by doubling the right hand side of inequalities in Corollary 3, and substituting  $\beta = \frac{1}{2}$ .

# D Additional Experiments

In this section, we include additional plots which was not reported in the main paper due to the space limitations.

## D.1 Matrix Completion

Figure 1 in the main paper shows comparisons between plain SGD and the dropout algorithm on the MovieLens dataset for a factorization size of  $d_1 = 70$ . The observation that we make with regard to those plots is not at all limited to the specific choice of the factorization size. In Figure 1 here, we report similar experiments with factorization sizes  $d_1 \in \{30, 110, 150, 190\}$ . It can be seen that the overall behaviour of plain SGD and dropout are very similar in all experiments. In particular, plain SGD always achieves the best training error but it has the largest generalization gap. Furthermore, increasing the dropout rate increases the training error but results in a tighter generalization gap.

It can be seen that an appropriate choice of the dropout rate always perform better than the plain SGD in terms of the test error. For instance, a dropout rate of p = 0.2 seems to always outperform plain SGD. Moreover, as the factorization size increases, the function class becomes more complex, and a larger value of the dropout rate is more helpful. For example, when  $d_1 = 30$ , the dropout with rates p = 0.3, 0.4 fail to achieve a good test performance, where as for larger factorization sizes  $(d_1 \in \{110, 150, 190\})$ , they consistently outperform plain SGD as well as other dropout rates.

#### D.2 Shallow Neural Networks

In Figure 2, we plot the co-adaptation measure, the generalization gap, as well as the complexity measure  $\alpha/\sqrt{n}$  as a function of width of the network, for FashionMNIST with symmetrization, and for MNIST without symmetrization.

The co-adaptation plot is very similar to Figure 2 in the main text. In particular, 1) increasing the dropout rate results in less co-adaptation; 2) even plain SGD is biased towards networks with less co-adaptation; and 3) as the networks becomes wider, the co-adaptation curves corresponding to plain SGD converge to those of dropout. We also make similar observations for the generalization gap as well as the complexity term  $\alpha/\sqrt{n}$ . In particular, 1) a higher dropout rate corresponds to a lower generalization gap, uniformly for all widths; 2) the generalization gap is higher for wider networks; and 3) curves with smaller complexity terms in the right plot correspond to curves with smaller generalization gaps in the middle plot.

### D.3 Deep Neural Networks

In Section 3, we derived generalization bounds that scale with the explicit regularizer as  $O(\sqrt{\frac{\text{width }R(w)}{n}})$ . Although our theoretical analysis is limited to two-layer networks; empirically, we show in Figure 3 that the generalization gap correlates well with this measure even for deep neural networks. In particular, we train deep convolutional neural networks with a dropout layer on top of the feature extractor, i.e. the top hidden layer. Let feature, denote the *i*-th hidden node in the top hidden layer. Akin to the derivation presented in Proposition 4, it is easy to see that the (expected) explicit regularizer is given by  $R(w) = \frac{p}{1-p} \sum_{i=1}^{width} \|\mathbf{u}_i\|^2 a_i^2$ , where width is the width of the top hidden layer, U denotes the top layer weight matrix, and  $a_i^2 = \mathbb{E}_{\mathbf{x}}[\text{feature}_i(\mathbf{x})^2]$  is the second moment of the *i*-th node in the top hidden layer.

We train convolutional neural networks with and without dropout, on MNIST, Fashion MNIST, and CIFAR-10. The CIFAR-10 dataset consists of 60 K  $32 \times 32$  color images in 10 classes, with 6 k images per class, divided into a training set and a test set of sizes 50 K and 10 K respectively Krizhevsky et al. (2009). We do not perform symmetrization in these experiments. In contrast with the experiments in the previous section, here we run the experiments on full datasets, representing each of the ten classes as a one-hot target vector.

For MNIST and Fashion MNIST datasets, we use a convolutional neural network with one convolutional layer and two fully connected layers. The convolutional layer has 16 convolutional filters, padding and stride

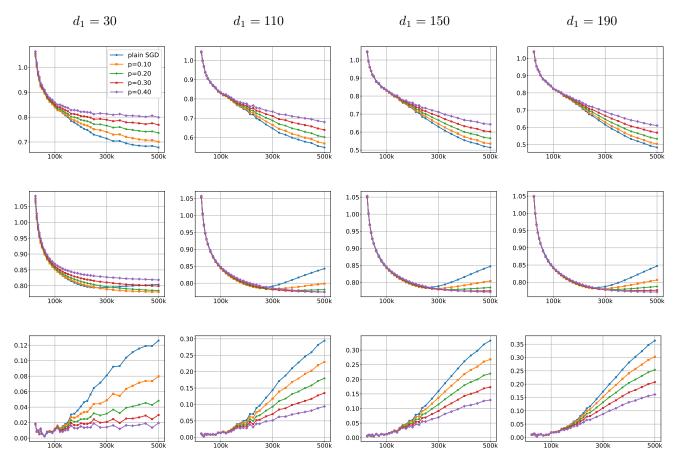


Figure 1: MovieLens dataset: the training error (**top**), the test error (**middle**), and the generalization gap for plain SGD as well as dropout with  $p \in \{0.25, 0.50, 0.75\}$  as a function of the number of iterates, for different factorization sizes  $d_1 = 30$  (first column),  $d_1 = 110$  (second column),  $d_1 = 150$  (third column), and  $d_1 = 190$  (forth column).

of 2, and kernel size of 5. We report experiments on networks with the width of the top hidden layer chosen from width  $\in \{2^6, 2^7, 2^8, 2^9, 2^{10}, 2^{11}\}$ . In all the experiments, a fixed learning rate lr = 0.5 and a mini-batch of size 256 is used to perform the updates. We train the models for 30 epochs over the whole training set.

For CIFAR-10, we use an AlexNet Krizhevsky et al. (2012), where the layers are modified accordingly to match the dataset. The only difference here is that we apply dropout to the top hidden layer, whereas in Krizhevsky et al. (2012), dropout is used on top of the second and the third hidden layers from the top. We report experiments on networks with the width of the top hidden layer chosen from width  $\in \{2^5, 2^6, 2^7, 2^8, 2^9, 2^{10}, 2^{11}, 2^{12}\}$ . In all the experiments, an initial learning rate lr = 5 and a mini-batch of size 256 is used to perform the updates. We train the models for 100 epochs over the whole training set. We decay the learning rate by a factor of 10 every 30 epochs.

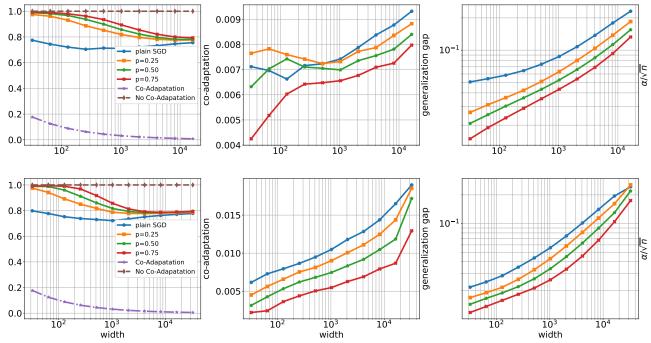


Figure 2: (left) "co-adaptation"; (middle) generalization gap; and (right)  $\alpha/\sqrt{n}$  (top) with symmetrization on FashionMNIST; and (bottom) without symmetrization on MNIST. In left column, the dashed brown and dotted purple lines represent minimal and maximal co-adaptations, respectively.

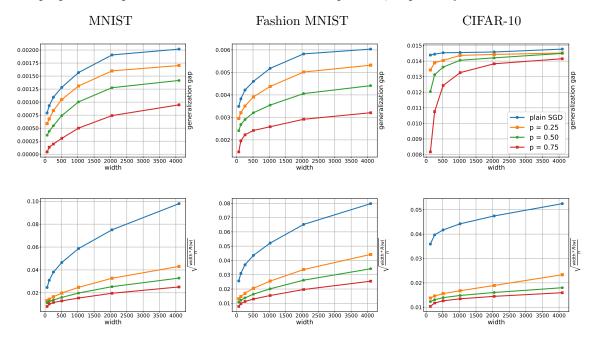


Figure 3: (**top**) generalization gap and (**bottom**) the complexity measure  $(\sqrt{\frac{\text{width} \cdot R(w)}{n}})$  as a function of the width of the top hidden layer on (**left**) MNIST, (**middle**) Fashion MNIST, and (**right**) CIFAR-10.