

---

# Learning Neural Networks with Adaptive Regularization

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1       Feed-forward neural networks can be understood as a combination of an interme-  
2       diate representation and a linear hypothesis. While most previous works aim to  
3       diversify the representations, we explore the complementary direction by perform-  
4       ing an adaptive and data-dependent regularization motivated by the empirical Bayes  
5       method. Specifically, we propose to construct a matrix-variate normal prior (on  
6       weights) whose covariance matrix has a Kronecker product structure. This structure  
7       is designed to capture the correlations in neurons through backpropagation. Under  
8       the assumption of this Kronecker factorization, the prior encourages neurons to  
9       borrow statistical strength from one another. Hence, it leads to an adaptive and  
10      data-dependent regularization when training networks on small datasets. To opti-  
11      mize the model, we present an efficient block coordinate descent algorithm with  
12      analytical solutions. Empirically, we demonstrate that the proposed method helps  
13      networks converge to local optima with smaller stable ranks and spectral norms.  
14      These properties suggest better generalizations and we present empirical results  
15      to support this expectation. We also verify the effectiveness of the approach on  
16      multiclass classification and multitask regression problems with various network  
17      structures.

## 18   1 Introduction

19   Although deep neural networks have been widely applied in various domains [19, 25, 29], usually its  
20   parameters are learned via the principle of maximum likelihood, hence its success crucially hinges  
21   on the availability of large scale datasets. When training rich models on small datasets, explicit  
22   regularization techniques are crucial to alleviate overfitting. Previous works have explored various  
23   regularization [42] and data augmentation [19, 41] techniques to learn diversified representations.  
24   In this paper, we look into an alternative direction by proposing an adaptive and data-dependent  
25   regularization method to encourage neurons of the same layer to share statistical strength. The goal of  
26   our method is to prevent overfitting when training (large) networks on small dataset. Our key insight  
27   stems from the famous argument by Efron [8] in the literature of the empirical Bayes method: *It is*  
28   *beneficial to learn from the experience of others*. From an algorithmic perspective, we argue that the  
29   connection weights of neurons in the same layer (row/column vectors of the weight matrix) will be  
30   correlated with each other through the backpropagation learning. Hence, by learning the correlations  
31   of the weight matrix, a neuron can “borrow statistical strength” from other neurons in the same layer.  
32   As an illustrating example, consider a simple setting where the input  $\mathbf{x} \in \mathbb{R}^d$  is fully connected to a  
33   hidden layer  $\mathbf{h} \in \mathbb{R}^p$ , which is further fully connected to the single output  $\hat{y} \in \mathbb{R}$ . Let  $\sigma(\cdot)$  be the  
34   nonlinear activation function, e.g., ReLU [36],  $W \in \mathbb{R}^{p \times d}$  be the connection matrix between the  
35   input layer and the hidden layer, and  $\mathbf{a} \in \mathbb{R}^p$  be the vector connecting the output and the hidden layer.  
36   Without loss of generality, ignoring the bias term in each layer, we have:  $\hat{y} = \mathbf{a}^T \mathbf{h}$ ,  $\mathbf{h} = \sigma(W\mathbf{x})$ .  
37   Consider using the usual  $\ell_2$  loss function  $\ell(\hat{y}, y) = \frac{1}{2}|\hat{y} - y|^2$  and take the derivative of  $\ell(\hat{y}, y)$  w.r.t.

38  $W$ . We obtain the update formula in backpropagation as  $W \leftarrow W - \alpha(\hat{y} - y)(\mathbf{a} \circ \mathbf{h}') \mathbf{x}^T$ , where  
39  $\mathbf{h}'$  is the componentwise derivative of  $\mathbf{h}$  w.r.t. its input argument, and  $\alpha > 0$  is the learning rate.  
40 Realize that  $(\mathbf{a} \circ \mathbf{h}') \mathbf{x}^T$  is a rank 1 matrix, and the component of  $\mathbf{h}'$  is either 0 or 1. Hence, the  
41 update for each row vector of  $W$  is linearly proportional to  $\mathbf{x}$ . Note that the observation holds for any  
42 input pair  $(\mathbf{x}, y)$ , so the update formula implies that the row vectors of  $W$  are correlated with each  
43 other. Although in this example we only discuss a one-hidden-layer network, it is straightforward to  
44 verify that the gradient update formula for general feed-forward networks admits the same rank one  
45 structure. The above observation leads us to the following question:

46 *Can we define a prior distribution over  $W$  that captures the correlations through*  
47 *the learning process for better generalization?*

48 **Our Contributions.** To answer the above question, we develop an adaptive regularization method for  
49 neural nets inspired by the empirical Bayes method. Motivated by the example above, we propose a  
50 matrix-variate normal prior whose covariance matrix admits a Kronecker product structure to capture  
51 the correlations between different neurons. Using tools from convex analysis, we present an efficient  
52 block coordinate descent algorithm with analytical solutions to optimize the model. Empirically, we  
53 show the proposed method helps the network converge to local optima with smaller stable ranks and  
54 spectral norms, and we verify the effectiveness of the approach on both multiclass classification and  
55 multitask regression problems with various network structures.

## 56 2 Preliminary

57 **Notation and Setup** We use lowercase letter to represent scalar and lowercase bold letter to denote  
58 vector. Capital letter, e.g.,  $X$ , is reserved for matrix. Calligraphic letter, such as  $\mathcal{D}$ , is used to denote  
59 set. We write  $\text{Tr}(A)$  as the trace of a matrix  $A$ ,  $\det(A)$  as the determinant of  $A$  and  $\text{vec}(A)$  as  
60  $A$ 's vectorization by column.  $[n]$  is used to represent the set  $\{1, \dots, n\}$  for any integer  $n$ . Other  
61 notations will be introduced whenever needed. Suppose we have access to a training set  $\mathcal{D}$  of  $n$  pairs  
62 of data instances  $(\mathbf{x}_i, y_i), i \in [n]$ . We consider the supervised learning setting where  $\mathbf{x}_i \in \mathcal{X} \subseteq \mathbb{R}^d$   
63 and  $y_i \in \mathcal{Y}$ . Let  $p(y | \mathbf{x}, \mathbf{w})$  be the conditional distribution of  $y$  given  $\mathbf{x}$  with parameter  $\mathbf{w}$ . The  
64 parametric form of the conditional distribution is assumed be known. In this paper, we assume the  
65 model parameter  $\mathbf{w}$  is sampled from a prior distribution  $p(\mathbf{w} | \theta)$  with hyperparameter  $\theta$ . On the  
66 other hand, given  $\mathcal{D}$ , the posterior distribution of  $\mathbf{w}$  is denoted by  $p(\mathbf{w} | \mathcal{D}, \theta)$ .

67 **The Empirical Bayes Method** To compute the predictive distribution, we need access to the value  
68 of the hyperparameter  $\theta$ . However, complete information about the hyperparameter  $\theta$  is usually not  
69 available in practice. To this end, empirical Bayes method [1, 9, 10, 12, 39] proposes to estimate  $\theta$   
70 from the data directly using the marginal distribution:

$$\hat{\theta} = \arg \max_{\theta} p(\mathcal{D} | \theta) = \arg \max_{\theta} \int p(\mathcal{D} | \mathbf{w}) \cdot p(\mathbf{w} | \theta) d\mathbf{w}. \quad (1)$$

71 Under specific choice of the likelihood function  $p(\mathbf{x}, y | \mathbf{w})$  and the prior distribution  $p(\mathbf{w} | \theta)$ , e.g.,  
72 conjugate pairs, we can solve the above integral in closed form. In certain cases we can even obtain  
73 an analytic solution of  $\hat{\theta}$ , which can then be plugged into the prior distribution. At a high level, by  
74 learning the hyperparameter  $\theta$  in the prior distribution directly from data, the empirical Bayes method  
75 provides us a principled and data-dependent way to obtain an estimator of  $\mathbf{w}$ . In fact, when both the  
76 prior and the likelihood functions are normal, it has been formally shown that the empirical Bayes  
77 estimators, e.g., the James-Stein estimator [23] and the Efron-Morris estimator [11], dominate the  
78 classic maximum likelihood estimator (MLE) in terms of quadratic loss for every choice of the model  
79 parameter  $\mathbf{w}$ . At a colloquial level, the success of the empirical Bayes method can be attributed to  
80 the effect of “*borrowing statistical strength*” [8], which also makes it a powerful tool in multitask  
81 learning [30, 46] and meta-learning [15].

## 82 3 Learning with Adaptive Regularization

83 In this section we first propose an adaptive regularization (AdaReg) method, which is inspired by the  
84 empirical Bayes method, for learning neural networks. We then combine our observation in Sec. 1  
85 to develop an efficient adaptive learning algorithm with matrix-variate normal prior. Through our  
86 derivation, we provide several connections and interpretations with other learning paradigms.

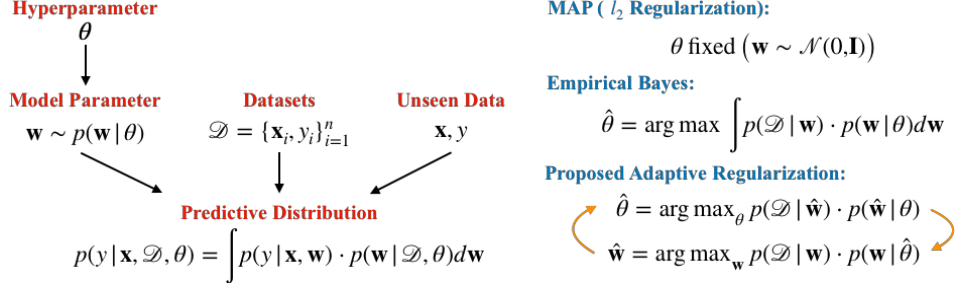


Figure 1: Illustration for Bayes/ Empirical Bayes, and our proposed adaptive regularization.

### 3.1 The Proposed Adaptive Regularization

When the likelihood function  $p(\mathcal{D} | \mathbf{w})$  is implemented as a neural network, the marginalization in (1) over model parameter  $\mathbf{w}$  cannot be computed exactly. Nevertheless, instead of performing expensive Monte-Carlo simulation, we propose to estimate both the model parameter  $\mathbf{w}$  and the hyperparameter  $\theta$  in the prior simultaneously from the joint distribution  $p(\mathcal{D}, \mathbf{w} | \theta) = p(\mathcal{D} | \mathbf{w}) \cdot p(\mathbf{w} | \theta)$ . Specifically, given an estimate  $\hat{\mathbf{w}}$  of the model parameter, by maximizing the joint distribution w.r.t.  $\theta$ , we can obtain  $\hat{\theta}$  as an approximation of the maximum marginal likelihood estimator. As a result, we can use  $\hat{\theta}$  to further refine the estimate  $\hat{\mathbf{w}}$  by maximizing the posterior distribution as follows:

$$\hat{\mathbf{w}} \leftarrow \max_{\mathbf{w}} p(\mathbf{w} | \mathcal{D}) = \max_{\mathbf{w}} p(\mathcal{D} | \mathbf{w}) \cdot p(\mathbf{w} | \hat{\theta}). \quad (2)$$

The maximizer of (2) can in turn be used in an updated joint distribution. Formally, we can define the following optimization problem that characterizes our Adaptive Regularization (AdaReg) framework:

$$\max_{\mathbf{w}} \max_{\theta} \log p(\mathcal{D} | \mathbf{w}) + \log p(\mathbf{w} | \theta). \quad (3)$$

It is worth connecting the optimization problem (3) to the classic maximum a posteriori (MAP) inference and also discuss their difference. If we drop the inner optimization over the hyperparameter  $\theta$  in the prior distribution. Then for any fixed value  $\hat{\theta}$ , (3) reduces to MAP with the prior defined by the specific choice of  $\hat{\theta}$ , and the maximizer  $\hat{\mathbf{w}}$  corresponds to the mode of the posterior distribution given by  $\hat{\theta}$ . From this perspective, the optimization problem in (3) actually defines a series of MAP inference problems, and the sequence  $\{\hat{\mathbf{w}}_j(\hat{\theta}_j)\}_j$  defines a solution path towards the final model parameter. On the algorithmic side, the optimization problem (3) also suggests a natural block coordinate descent algorithm where we alternatively optimize over  $\mathbf{w}$  and  $\theta$  until the convergence of the objective function. An illustration of the framework is shown in Fig. 1.

### 3.2 Neural Network with Matrix-Normal Prior

Inspired by the observation from Sec. 1, we propose to define a matrix-variate normal distribution [16] over the connection weight matrix  $W$ :  $W \sim \mathcal{MN}(0_{p \times d}, \Sigma_r, \Sigma_c)$ , where  $\Sigma_r \in \mathbb{S}_{++}^p$  and  $\Sigma_c \in \mathbb{S}_{++}^d$  are the row and column covariance matrices, respectively.<sup>1</sup> Equivalently, one can understand the matrix-variate normal distribution over  $W$  as a multivariate normal distribution with a Kronecker product covariance structure over  $\text{vec}(W)$ :  $\text{vec}(W) \sim \mathcal{N}(0_{p \times d}, \Sigma_c \otimes \Sigma_r)$ . It is then easy to check that the marginal prior distributions over the row and column vectors of  $W$  are given by:

$$W_{i:} \sim \mathcal{N}(\mathbf{0}_d, [\Sigma_r]_{ii} \cdot \Sigma_c), \quad W_{:j} \sim \mathcal{N}(\mathbf{0}_p, [\Sigma_c]_{jj} \cdot \Sigma_r).$$

We point out that the Kronecker product structure of the covariance matrix exactly captures our prior about the connection matrix  $W$ : the fan-in/fan-out of neurons in the same layer (row/column vectors of  $W$ ) are correlated with the same correlation matrix in the prior, and they only differ at the scales.

For illustration purpose, let us consider the simple feed-forward network discussed in Sec. 1. Consider a reparametrization of the model by defining  $\Omega_r := \Sigma_r^{-1}$  and  $\Omega_c := \Sigma_c^{-1}$  to be the corresponding precision matrices and plug in the prior distribution into the our AdaReg framework (see (3)). After

<sup>1</sup>The probability density function is given by  $p(W | \Sigma_r, \Sigma_c) = \frac{\exp(-\text{Tr}(\Sigma_r^{-1} W \Sigma_c^{-1} W^T)/2)}{(2\pi)^{pd/2} \det(\Sigma_r)^{d/2} \det(\Sigma_c)^{p/2}}$ .

120 routine algebraic simplifications, we reach the following concrete optimization problem:

$$\begin{aligned} \min_{W, \mathbf{a}} \min_{\Omega_r, \Omega_c} \quad & \frac{1}{2n} \sum_{i \in [n]} (\hat{y}(\mathbf{x}_i; W, \mathbf{a}) - y_i)^2 + \lambda \|\Omega_r^{1/2} W \Omega_c^{1/2}\|_F^2 - \lambda (d \log \det(\Omega_r) + p \log \det(\Omega_c)) \\ \text{subject to} \quad & uI_p \preceq \Omega_r \preceq vI_p, \quad uI_d \preceq \Omega_c \preceq vI_d \end{aligned} \quad (4)$$

121 where  $\lambda$  is a constant that only depends on  $p$  and  $d$ ,  $0 < u \leq v$  and  $uv = 1$ . Note that the constraint  
122 is necessary to guarantee the feasible set to be compact so that the optimization problem is well  
123 formulated and a minimum is attainable.<sup>2</sup> It is not hard to show that in general the optimization  
124 problem (4) is not jointly convex in terms of  $\{\mathbf{a}, W, \Omega_r, \Omega_c\}$ , and this holds even if the activation  
125 function is linear. However, as we will show later, for any fixed  $\mathbf{a}, W$ , the reparametrization makes  
126 the partial optimization over  $\Omega_r$  and  $\Omega_c$  bi-convex. More importantly, we can derive an efficient  
127 algorithm that finds the optimal  $\Omega_r(\Omega_c)$  for any fixed  $\mathbf{a}, W, \Omega_c(\Omega_r)$  in  $O(\max\{d^3, p^3\})$  time with  
128 closed form solutions. This allows us to apply our algorithm to networks of large sizes, where  
129 a typical hidden layer can contain thousands of nodes. Note that this is in contrast to solving a  
130 general semi-definite programming (SDP) problem using black-box algorithm, e.g., the interior-point  
131 method [35], which is computationally intensive and hard to scale to networks with moderate sizes.  
132 Before we delve into the details on solving (4), it is instructive to discuss some of its connections and  
133 differences to other learning paradigms.

134 **Maximum-A-Posteriori Estimation.** Essentially, for model parameter  $W$ , (4) defines a sequence of  
135 MAP problems where each MAP is indexed by the pair of precision matrices  $(\Omega_r^{(t)}, \Omega_c^{(t)})$  at iteration  $t$ .  
136 Equivalently, at each stage of the optimization, we can interpret (4) as placing a matrix variate normal  
137 prior on  $W$  where the precision matrix in the prior is given by  $\Omega_r^{(t)} \otimes \Omega_c^{(t)}$ . From this perspective, if  
138 we fix  $\Omega_r^{(t)} = I_p$  and  $\Omega_c^{(t)} = I_d$ ,  $\forall t$ , then (4) naturally reduces to learning with  $\ell_2$  regularization [26].  
139 More generally, for non-diagonal precision matrices, the regularization term for  $W$  becomes:

$$\|\Omega_r^{1/2} W \Omega_c^{1/2}\|_F^2 = \|\text{vec}(\Omega_r^{1/2} W \Omega_c^{1/2})\|_2^2 = \|(\Omega_c^{1/2} \otimes \Omega_r^{1/2}) \text{vec}(W)\|_2^2,$$

140 and this is exactly the Tikhonov regularization [13] imposed on  $W$  where the Tikhonov matrix  $\Gamma$  is  
141 given by  $\Gamma := \Omega_c^{1/2} \otimes \Omega_r^{1/2}$ . But instead of manually designing the regularization matrix  $\Gamma$  to improve  
142 the conditioning of the estimation problem, we propose to also learn both precision matrices (so  $\Gamma$  as  
143 well) from data. From an algorithmic perspective,  $\Gamma^T \Gamma = \Omega_c \otimes \Omega_r$  serves as a preconditioning matrix  
144 w.r.t. model parameter  $W$  to reshape the gradient according to the geometry of the data [7, 17, 18].

145 **Volume Minimization.** Let us consider the  $\log \det(\cdot)$  function over the positive definite cone. It  
146 is well known that the log-determinant function is concave [3]. Hence for any pair of matrices  
147  $A_1, A_2 \in \mathbb{S}_{++}^m$ , the following inequality holds:

$$\log \det(A_1) \leq \log \det(A_2) + \langle \nabla \log \det(A_2), A_1 - A_2 \rangle = \log \det(A_2) + \text{Tr}(A_2^{-1} A_1) - m. \quad (5)$$

148 Applying the above inequality twice by fixing  $A_1 = W \Omega_c W^T / 2d$ ,  $A_2 = \Sigma_r$  and  $A_1 =$   
149  $W^T \Omega_r W / 2p$ ,  $A_2 = \Sigma_c$  respectively leads to the following inequalities:

$$\begin{aligned} d \log \det(W \Omega_c W^T / 2d) &\leq -d \log \det(\Omega_r) + \frac{1}{2} \text{Tr}(\Omega_r W \Omega_c W^T) - dp, \\ p \log \det(W^T \Omega_r W / 2p) &\leq -p \log \det(\Omega_c) + \frac{1}{2} \text{Tr}(\Omega_r W \Omega_c W^T) - dp. \end{aligned}$$

150 Realize  $\text{Tr}(\Omega_r W \Omega_c W^T) = \|\Omega_r^{1/2} W \Omega_c^{1/2}\|_F^2$ . Summing the above two inequalities leads to:

$$d \log \det(W \Omega_c W^T) + p \log \det(W^T \Omega_r W) \leq \|\Omega_r^{1/2} W \Omega_c^{1/2}\|_F^2 - (d \log \det(\Omega_r) + p \log \det(\Omega_c)) + c, \quad (6)$$

151 where  $c$  is a constant that only depends on  $d$  and  $p$ . Recall that  $|\det(A^T A)|$  computes the squared  
152 volume of the parallelepiped spanned by the column vectors of  $A$ . Hence (6) gives us a natural  
153 interpretation of the objective function in (4): the regularizer essentially upper bounds the log-volume  
154 of the two parallelepipeds spanned by the row and column vectors of  $W$ . But instead of measuring the  
155 volume using standard Euclidean inner product, it also takes into account the local curvatures defined  
156 by  $\Sigma_r$  and  $\Sigma_c$ , respectively. For vectors with fixed lengths, the volume of the parallelepiped spanned  
157 by them becomes smaller when they are more linearly correlated, either positively or negatively. At a  
158 colloquial level, this means that the regularizer in (4) forces fan-in/fan-out of neurons at the same  
159 layer to be either positively or negatively correlated with each other, and this corresponds exactly to  
160 the effect of sharing statistical strengths.

<sup>2</sup>The constraint  $uv = 1$  is only for the ease of presentation in the following part and can be readily removed.

---

**Algorithm 1** Block Coordinate Descent for Adaptive Regularization
 

---

**Input:** Initial value  $\phi^{(0)} := \{\mathbf{a}^{(0)}, W^{(0)}\}$ ,  $\Omega_r^{(0)} \in \mathbb{S}_{++}^p$  and  $\Omega_c^{(0)} \in \mathbb{S}_{++}^d$ , first-order optimization algorithm  $\mathfrak{A}$ .

- 1: **for**  $t = 1, \dots, \infty$  until convergence **do**
- 2:   Fix  $\Omega_r^{(t-1)}, \Omega_c^{(t-1)}$ , optimize  $\phi^{(t)}$  by backpropagation and algorithm  $\mathfrak{A}$
- 3:    $\Omega_r^{(t)} \leftarrow \text{INVTHRESHOLD}(W^{(t)}\Omega_c^{(t-1)}W^{(t)T}, d, u, v)$
- 4:    $\Omega_c^{(t)} \leftarrow \text{INVTHRESHOLD}(W^{(t)T}\Omega_r^{(t)}W^{(t)}, p, u, v)$
- 5: **end for**
- 6: **procedure**  $\text{INVTHRESHOLD}(\Delta, m, u, v)$
- 7:   Compute SVD:  $Q\text{diag}(\mathbf{r})Q^T = \text{SVD}(\Delta)$
- 8:   Hard thresholding  $\mathbf{r}' \leftarrow \mathbb{T}_{[u,v]}(m/\mathbf{r})$
- 9:   **return**  $Q\text{diag}(\mathbf{r}')Q^T$
- 10: **end procedure**

---

### 161 3.3 The Algorithm

162 In this section we describe a block coordinate descent algorithm to optimize the objective function  
 163 in (4) and detail how to efficiently solve the matrix optimization subproblems in closed form using  
 164 tools from convex analysis. Due to space limit, we defer proofs and detailed derivation to appendix.  
 165 Given a pair of constants  $0 < u \leq v$ , we define the following thresholding function  $\mathbb{T}_{[u,v]}(x)$ :

$$\mathbb{T}_{[u,v]}(x) := \max\{u, \min\{v, x\}\}. \quad (7)$$

166 We summarize our block coordinate descent algorithm to solve (4) in Alg. 1. In each iteration, Alg. 1  
 167 takes a first-order algorithm  $\mathfrak{A}$ , e.g., the stochastic gradient descent, to optimize the parameters of the  
 168 neural network by backpropagation. It then proceeds to compute the optimal solutions for  $\Omega_r$  and  $\Omega_c$   
 169 using  $\text{INVTHRESHOLD}$  as a sub-procedure. Alg. 1 terminates when a stationary point is found.

170 We now proceed to show that the procedure  $\text{INVTHRESHOLD}$  finds the optimal solution given all the  
 171 other variables fixed. Due to the symmetry between  $\Omega_r$  and  $\Omega_c$  in (4), we will only prove this for  $\Omega_r$ ,  
 172 and similar arguments can be applied to  $\Omega_c$  as well. Fix both  $W$ ,  $\Omega_c$  and ignore all the terms that do  
 173 not depend on  $\Omega_r$ , the sub-problem on optimizing  $\Omega_r$  becomes:

$$\min_{\Omega_r} \quad \text{Tr}(\Omega_r W \Omega_c W^T) - d \log \det(\Omega_r), \quad \text{subject to} \quad uI_p \preceq \Omega_r \preceq vI_p. \quad (8)$$

174 It is not hard to show that the optimization problem (8) is convex. Define the constraint set  
 175  $\mathcal{C} := \{A \in \mathbb{S}_{++}^p \mid uI_p \preceq A \preceq vI_p\}$  and the indicator function  $\mathbb{I}_{\mathcal{C}}(A) = 0$  iff  $A \in \mathcal{C}$  else  
 176  $\infty$ . Given the convexity of (8), we can use the indicator function to first transform (8) into an  
 177 unconstrained one and use the first-order optimality condition to characterize the optimal solu-  
 178 tion:  $0 \in \partial(\frac{1}{d} \text{Tr}(\Omega_r W \Omega_c W^T) - \log \det(\Omega_r) + \mathbb{I}_{\mathcal{C}}(\Omega_r)) = W \Omega_c W^T / d - \Omega_r^{-1} + \mathcal{N}_{\mathcal{C}}(\Omega_r)$ , where  
 179  $\mathcal{N}_{\mathcal{C}}(A) := \{B \in \mathbb{S}^p \mid \text{Tr}(B^T(Z - A)) \leq 0, \forall Z \in \mathcal{C}\}$  is the normal cone w.r.t.  $\mathcal{C}$  at  $A$ . With the help  
 180 of Lemma 1 in appendix, equivalently, we have  $\Omega_r^{-1} - W \Omega_c W^T / d \in \mathcal{N}_{\mathcal{C}}(\Omega_r)$ . Geometrically, this  
 181 means that the optimum  $\Omega_r^{-1}$  is the Euclidean projection of  $W \Omega_c W^T / d$  onto  $\mathcal{C}$ . Hence in order to  
 182 solve (8), it suffices if we can solve the following Euclidean projection problem efficiently, where  
 183  $\widetilde{\Omega}_r \in \mathbb{S}^p$  is a real symmetric matrix:

$$\min_{\Omega_r} \quad \|\Omega_r - \widetilde{\Omega}_r\|_F^2, \quad \text{subject to} \quad uI_p \preceq \Omega_r \preceq vI_p. \quad (9)$$

184 The following theorem characterizes the optimal solution to the above Euclidean projection problem:

185 **Theorem 1.** Let  $\widetilde{\Omega}_r \in \mathbb{S}^p$  with eigendecomposition as  $\widetilde{\Omega}_r = Q\Lambda Q^T$  and  $\text{Proj}_{\mathcal{C}}(\cdot)$  be the Euclidean  
 186 projection operator onto  $\mathcal{C}$ , then  $\text{Proj}_{\mathcal{C}}(\widetilde{\Omega}_r) = Q\mathbb{T}_{[u,v]}(\Lambda)Q^T$ .

187 **Corollary 1.** Let  $W \Omega_c W^T$  be eigendecomposed as  $Q\text{diag}(\mathbf{r})Q^T$ , then the optimal solution to (8) is  
 188 given by  $Q\mathbb{T}_{[u,v]}(d/\mathbf{r})Q^T$ .

189 Similar arguments can be made to derive the solution for  $\Omega_c$  in (4). The final algorithm is very  
 190 simple as it only contains one SVD, hence its time complexity is  $O(\max\{d^3, p^3\})$ . Note that the total  
 191 number of parameters in the network is at least  $\Omega(dp)$ , hence the algorithm is efficient as it scales  
 192 sub-quadratically in terms of number of parameters in the network.

Table 1: Explained variance of different methods on 7 regression tasks from the SARCOS dataset.

Method	1st	2nd	3rd	4th	5th	6th	7th
<b>MTL</b>	0.4418	0.3472	0.5222	0.5036	0.6024	0.4727	0.5298
<b>MTL-Dropout</b>	0.4413	0.3271	0.5202	0.5063	0.6036	0.4711	0.5345
<b>MTL-BN</b>	0.4768	0.3770	0.5396	0.5216	0.6117	0.4936	0.5479
<b>MTL-DeCoV</b>	0.4027	0.3137	0.4703	0.4515	0.5229	0.4224	0.4716
<b>MTL-AdaReg</b>	<b>0.4769</b>	<b>0.3969</b>	<b>0.5485</b>	<b>0.5308</b>	<b>0.6202</b>	<b>0.5085</b>	<b>0.5561</b>

## 193 4 Experiments

194 In this section we demonstrate the effectiveness of AdaReg in learning practical deep neural networks  
 195 on real-world datasets. We report generalization, optimization as well as stability results.

### 196 4.1 Experimental Setup

197 **Multiclass Classification (MNIST & CIFAR10):** In this experiment, we show that AdaReg provides  
 198 an effective regularization on the network parameters. To this end, we use a convolutional neural  
 199 network as our baseline model. To show the effect of regularization, we gradually increase the  
 200 training set size. In MNIST we use the step from 60 to 60,000 (11 different experiments) and in  
 201 CIFAR10 we consider the step from 5,000 to 50,000 (10 different experiments). For each training  
 202 set size, we repeat the experiments for 10 times. The mean along with its standard deviation are  
 203 shown as the statistics. Moreover, since both the optimization and generalization of neural networks  
 204 are sensitive to the size of minibatches [14, 24], we study two minibatch settings for 256 and 2048,  
 205 respectively. In our method, we place a matrix-variate normal prior over the weight matrix of the last  
 206 softmax layer, and we use Alg. 1 to optimize both the model weights and two covariance matrices.

207 **Multitask Regression (SARCOS):** SARCOS relates to an inverse dynamics problem for a seven  
 208 degree-of-freedom (DOF) SARCOS anthropomorphic robot arm [44]. The goal of this task is to  
 209 map from a 21-dimensional input space (7 joint positions, 7 joint velocities, 7 joint accelerations) to  
 210 the corresponding 7 joint torques. Hence there are 7 tasks and the inputs are shared among all the  
 211 tasks. The training set and test set contain 44,484 and 4,449 examples, respectively. Again, we apply  
 212 AdaReg on the last layer weight matrix, where each row corresponds to a separate task vector.

213 We compare AdaReg with classic regularization methods in the literature, including weight decay,  
 214 dropout [42], batch normalization (BN) [22] and the DeCov method [6]. We also note that we  
 215 fix all the hyperparameters such as learning rate to be the same for all the methods. We report  
 216 evaluation metrics on test set as a measure of generalization. To understand how the proposed  
 217 adaptive regularization helps in optimization, we visualize the trajectory of the loss function during  
 218 training. Lastly, we also present the inferred correlation of the weight matrix for qualitative study.

### 219 4.2 Results and Analysis

220 **Multiclass Classification (MNIST & CIFAR10):** Results on the multiclass classification for dif-  
 221 ferent training sizes are show in Fig. 2. For both MNIST and CIFAR10, we find AdaReg, Weight  
 222 Decay, and Dropout are the effective regularization methods, while Batch Normalization and DeCov  
 223 vary in different settings. Batch Normalization suffers from large batch size in CIFAR10 (comparing  
 224 Fig. 2 (c) and (d)) but is not sensitive to batch size in MNIST (comparing Fig. 2 (a) and (b)). The  
 225 performance deterioration in large batch size of Batch Normalization is also observed by [21]. DeCov,  
 226 on the other hand, improves the generalization in MNIST with batch size 256 (see Fig. 2 (a)), while  
 227 it demonstrates only comparable or even worse performance in other settings. To conclude, as  
 228 training set size grows, AdaReg consistently performs better generalization as comparing to other  
 229 regularization methods. We also note that AdaReg is not sensitive to the size of minibatches while  
 230 most of the methods suffer from large minibatches. In appendix, we show the combination of AdaReg  
 231 with other generalization methods can usually lead to even better results.

232 **Multitask Regression (SARCOS):** In this experiment we are interested in investigating whether  
 233 AdaReg can lead to better generalization for multiple related regression problems. To do so, we  
 234 report the explained variance as a normalized metric, e.g., one minus the ratio between mean squared  
 235 error and the variance of different methods in Table 1. The larger the explained variance, the better

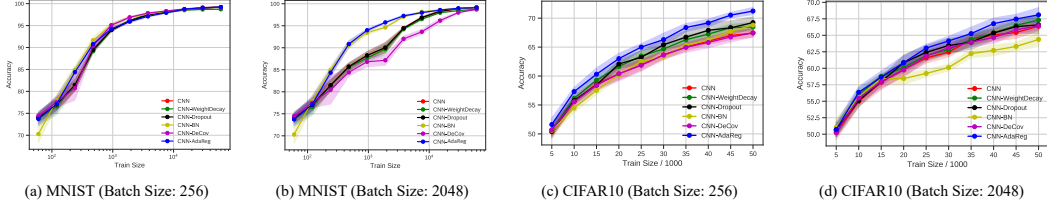


Figure 2: Generalization performance on MNIST and CIFAR10. AdaReg improves generalization under both minibatch settings.

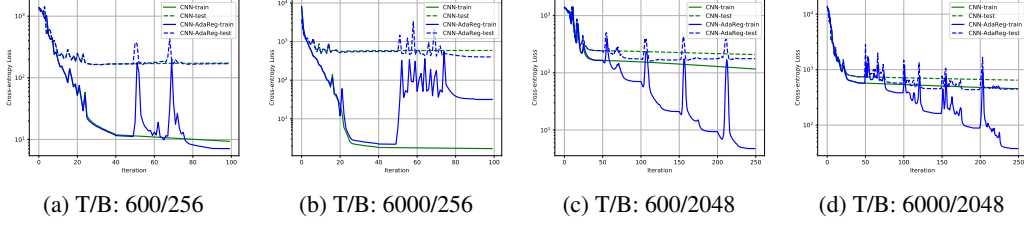


Figure 3: Optimization trajectory of AdaReg on MNIST with training size/batch size on training and test sets. AdaReg helps to converge to better local optima. Note the log-scale on  $y$ -axis.

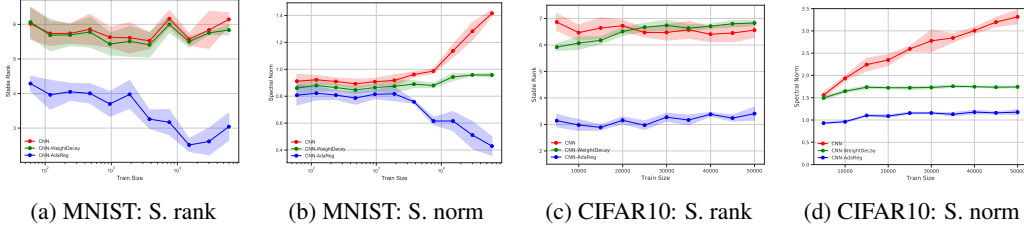
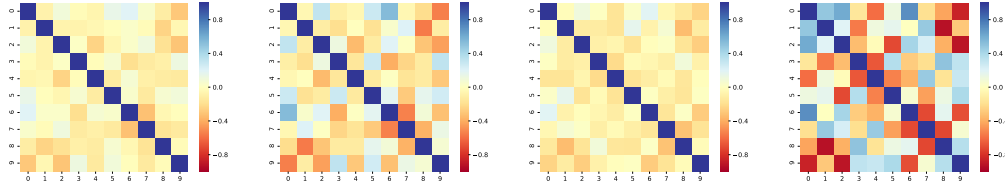


Figure 4: Comparisons of stable ranks (S. rank) and spectral norms (S. norm) from different methods on MNIST and CIFAR10.  $x$ -axis corresponds to the training size.

the predictive performance. In this case we observe a consistent improvement of AdaReg over other competitors on all the 7 regression tasks. We would like to emphasize that all the experiments share exactly the same experimental protocol, including network structure, optimization algorithm, training iteration, etc, so that the performance differences can only be explained by different ways of regularizations. For better visualization, we also plot the result in appendix.

**Optimization:** It has recently been empirically shown that BN helps optimization not by reducing internal covariate shift, but instead by smoothing the landscape of the loss function [40]. To understand how AdaReg improves generalization, in Fig. 3, we plot the values of the cross entropy loss function on both the training and test sets during optimization using Alg. 1. The experiment is performed in MNIST with batch size 256/2048. In this experiment, we fix the number of outer loop to be 2/5 and each block optimization over network weights contains 50 epochs. Because of the stochastic optimization over model weights, we can see several unstable peaks in function value around iteration 50 when trained with AdaReg, which corresponds to the transition phase between two consecutive outer loops with different row/column covariance matrices. In all the cases AdaReg converges to better local optima of the loss landscape, which lead to better generalization on the test set as well because they have smaller loss values on the test set when compared with training without AdaReg.

**Stable rank and spectral norm:** Given a matrix  $W$ , the stable rank of  $W$ , denoted as  $\text{srnk}(W)$ , is defined as  $\text{srnk}(W) := \|W\|_F^2 / \|W\|_2^2$ . As its name suggests, the stable rank is more stable than the rank because it is largely unaffected by tiny singular values. It has recently been shown [37, Theorem 1] that the generalization error of neural networks crucially depends on both the stable ranks and the spectral norms of connection matrices in the network. Specifically, it can be shown that the generalization error is upper bounded by  $O(\sqrt{\prod_{j=1}^L \|W_j\|_2^2 \sum_{j=1}^L \text{srnk}(W_j) / n})$ , where  $L$  is the number of layers in the network. Essentially, this upper bound suggests that smaller spectral norm (smoother function mapping) and stable rank (skewed spectrum) leads to better generalization.



(a) CNN, Acc: 89.34    (b) AdaReg, Acc: 92.50    (c) CNN, Acc: 98.99    (d) AdaReg, Acc: 99.19

Figure 5: Correlation matrix of the weight matrix in the softmax layer. The left two correspond to dataset with training size 600 and the right two with size 60,000. Acc means the test set accuracy.

To understand why AdaReg improves generalization, in Fig. 4, we plot both the stable rank and the spectral norm of the weight matrix in the last layer of the CNNs used in our MNIST and CIFAR10 experiments. We compare 3 methods: CNN without any regularization, CNN trained with weight decay and CNN with AdaReg. For each setting we repeat the experiments for 5 times, and we plot the mean along with its standard deviation. From Fig. 4a and Fig. 4c it is clear that AdaReg leads to a significant reduction in terms of the stable rank when compared with weight decay, and this effect is consistent in all the experiments with different training size. Similarly, in Fig. 4b and Fig. 4d we plot the spectral norm of the weight matrix. Again, both weight decay and AdaReg help reduce the spectral norm in all settings, but AdaReg plays a more significant role than the usual weight decay. Combining the experiments with the generalization upper bound introduced above, we can see that training with AdaReg leads to an estimator of  $W$  that has lower stable rank and smaller spectral norm, which explains why it achieves a better generalization performance. Furthermore, this observation holds on the SARCOS datasets as well, and we show the results in the appendix.

**Correlation Matrix:** To verify that AdaReg imposes the effect of “sharing statistical strength” during training, we visualize the weight matrix of the softmax layer by computing the corresponding correlation matrix, as shown in Fig. 5. In Fig. 5, darker color means stronger correlation. We conduct two experiments with training size 600 and 60,000 respectively. As we can observe, training with AdaReg leads to weight matrix with stronger correlations, and this effect is more evident when the training set is large. This is consistent with our analysis of sharing statistical strengths. As a sanity check, from Fig. 5 we can also see that similar digits, e.g., 1 and 7, share a positive correlation while dissimilar ones, e.g., 1 and 8, share a negative correlation.

## 5 Related Work

Despite the name, empirical Bayes method is in fact a frequentist approach to obtain estimator with favorable properties. On the other hand, truly Bayesian inference would instead put a posterior distribution over model weights to characterize the uncertainty during training [2, 20, 32]. However, due to the complexity of nonlinear neural networks, analytic posterior is not available, hence strong independent assumptions over model weight have to be made in order to achieve computationally tractable variational solution. Typically, both the prior and the variational posterior are assumed to fully factorize over model weights. As an exception, Louizos and Welling [31], Sun et al. [43] seek to learn Bayesian neural nets where they approximate the intractable posterior distribution using matrix-variate Gaussian distribution. The prior for weights are still assumed to be known and fixed. As a comparison, we use matrix-variate Gaussian as the prior distribution and we learn the hyperparameter in the prior from data. Hence our method does not belong to Bayesian neural nets: we instead use the empirical Bayes principle to derive adaptive regularization method in order to have better generalization, as done in [4, 38].

## 6 Conclusion

Inspired by empirical Bayes method, we propose an adaptive regularization (AdaReg) with matrix-variate normal prior for model parameters in deep neural networks. The prior encourages neurons to borrow statistical strength from other neurons during the learning process, and it provides an effective regularization when training networks on small datasets. To optimize the model, we design an efficient block coordinate descent algorithm to learn both model weights and the covariance structures. Empirically, on three datasets we demonstrate that AdaReg improves generalization by finding better local optima with smaller spectral norms and stable ranks.

## References

- [1] José M Bernardo and Adrian FM Smith. Bayesian theory, 2001.
- [2] Charles Blundell, Julien Cornebise, Koray Kavukcuoglu, and Daan Wierstra. Weight uncertainty in neural networks. *arXiv preprint arXiv:1505.05424*, 2015.
- [3] Stephen Boyd and Lieven Vandenberghe. *Convex optimization*. Cambridge university press, 2004.
- [4] Philip J Brown, James V Zidek, et al. Adaptive multivariate ridge regression. *The Annals of Statistics*, 8(1):64–74, 1980.
- [5] Rich Caruana, Steve Lawrence, and C Lee Giles. Overfitting in neural nets: Backpropagation, conjugate gradient, and early stopping. In *Advances in neural information processing systems*, pages 402–408, 2001.
- [6] Michael Cogswell, Faruk Ahmed, Ross Girshick, Larry Zitnick, and Dhruv Batra. Reducing overfitting in deep networks by decorrelating representations. *arXiv preprint arXiv:1511.06068*, 2015.
- [7] John Duchi, Elad Hazan, and Yoram Singer. Adaptive subgradient methods for online learning and stochastic optimization. *Journal of Machine Learning Research*, 12(Jul):2121–2159, 2011.
- [8] Bradley Efron. *Large-scale inference: empirical Bayes methods for estimation, testing, and prediction*, volume 1. Cambridge University Press, 2012.
- [9] Bradley Efron and Trevor Hastie. *Computer age statistical inference*, volume 5. Cambridge University Press, 2016.
- [10] Bradley Efron and Carl Morris. Stein’s estimation rule and its competitors—an empirical Bayes approach. *Journal of the American Statistical Association*, 68(341):117–130, 1973.
- [11] Bradley Efron and Carl Morris. Stein’s paradox in statistics. *Scientific American*, 236(5): 119–127, 1977.
- [12] Andrew Gelman, John B Carlin, Hal S Stern, David B Dunson, Aki Vehtari, and Donald B Rubin. *Bayesian data analysis*. CRC press, 2013.
- [13] Gene H Golub, Michael Heath, and Grace Wahba. Generalized cross-validation as a method for choosing a good ridge parameter. *Technometrics*, 21(2):215–223, 1979.
- [14] Priya Goyal, Piotr Dollár, Ross Girshick, Pieter Noordhuis, Lukasz Wesolowski, Aapo Kyrola, Andrew Tulloch, Yangqing Jia, and Kaiming He. Accurate, large minibatch sgd: training imagenet in 1 hour. *arXiv preprint arXiv:1706.02677*, 2017.
- [15] Erin Grant, Chelsea Finn, Sergey Levine, Trevor Darrell, and Thomas Griffiths. Recasting gradient-based meta-learning as hierarchical bayes. *arXiv preprint arXiv:1801.08930*, 2018.
- [16] Arjun K Gupta and Daya K Nagar. *Matrix variate distributions*. Chapman and Hall/CRC, 2018.
- [17] Vineet Gupta, Tomer Koren, and Yoram Singer. A unified approach to adaptive regularization in online and stochastic optimization. *arXiv preprint arXiv:1706.06569*, 2017.
- [18] Elad Hazan, Amit Agarwal, and Satyen Kale. Logarithmic regret algorithms for online convex optimization. *Machine Learning*, 69(2-3):169–192, 2007.
- [19] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- [20] José Miguel Hernández-Lobato and Ryan Adams. Probabilistic backpropagation for scalable learning of bayesian neural networks. In *International Conference on Machine Learning*, pages 1861–1869, 2015.

- [21] Elad Hoffer, Itay Hubara, and Daniel Soudry. Train longer, generalize better: closing the generalization gap in large batch training of neural networks. In *Advances in Neural Information Processing Systems*, pages 1731–1741, 2017.
- [22] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. *arXiv preprint arXiv:1502.03167*, 2015.
- [23] William James and Charles Stein. Estimation with quadratic loss. In *Proceedings of the fourth Berkeley symposium on mathematical statistics and probability*, volume 1, pages 361–379, 1961.
- [24] Nitish Shirish Keskar, Dheevatsa Mudigere, Jorge Nocedal, Mikhail Smelyanskiy, and Ping Tak Peter Tang. On large-batch training for deep learning: Generalization gap and sharp minima. *arXiv preprint arXiv:1609.04836*, 2016.
- [25] Alex Krizhevsky and Geoffrey Hinton. Learning multiple layers of features from tiny images. 2009.
- [26] Anders Krogh and John A Hertz. A simple weight decay can improve generalization. In *Advances in neural information processing systems*, pages 950–957, 1992.
- [27] Alex Kulesza and Ben Taskar. Learning determinantal point processes. In *Proceedings of the Twenty-Seventh Conference on Uncertainty in Artificial Intelligence*, pages 419–427. AUAI Press, 2011.
- [28] Alex Kulesza, Ben Taskar, et al. Determinantal point processes for machine learning. *Foundations and Trends® in Machine Learning*, 5(2–3):123–286, 2012.
- [29] Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. Deep learning. *nature*, 521(7553):436, 2015.
- [30] Mingsheng Long, Zhangjie Cao, Jianmin Wang, and S Yu Philip. Learning multiple tasks with multilinear relationship networks. In *Advances in Neural Information Processing Systems*, pages 1594–1603, 2017.
- [31] Christos Louizos and Max Welling. Structured and efficient variational deep learning with matrix gaussian posteriors. In *International Conference on Machine Learning*, pages 1708–1716, 2016.
- [32] David JC MacKay. A practical bayesian framework for backpropagation networks. *Neural computation*, 4(3):448–472, 1992.
- [33] Zelda Mariet and Suvrit Sra. Diversity networks: Neural network compression using determinantal point processes. *arXiv preprint arXiv:1511.05077*, 2015.
- [34] James Martens and Roger Grosse. Optimizing neural networks with kronecker-factored approximate curvature. In *International conference on machine learning*, pages 2408–2417, 2015.
- [35] Sanjay Mehrotra. On the implementation of a primal-dual interior point method. *SIAM Journal on optimization*, 2(4):575–601, 1992.
- [36] Vinod Nair and Geoffrey E Hinton. Rectified linear units improve restricted boltzmann machines. In *Proceedings of the 27th international conference on machine learning (ICML-10)*, pages 807–814, 2010.
- [37] Behnam Neyshabur, Srinadh Bhojanapalli, David McAllester, and Nathan Srebro. A pac-bayesian approach to spectrally-normalized margin bounds for neural networks. *arXiv preprint arXiv:1707.09564*, 2017.
- [38] Samuel D Oman. A different empirical bayes interpretation of ridge and stein estimators. *Journal of the Royal Statistical Society: Series B (Methodological)*, 46(3):544–557, 1984.
- [39] Herbert Robbins. An empirical bayes approach to statistics. Technical report, Columbia University, New York City, United States, 1956.

- 394 [40] Shibani Santurkar, Dimitris Tsipras, Andrew Ilyas, and Aleksander Madry. How does batch  
395 normalization help optimization?(no, it is not about internal covariate shift). *arXiv preprint*  
396 *arXiv:1805.11604*, 2018.
- 397 [41] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale  
398 image recognition. *arXiv preprint arXiv:1409.1556*, 2014.
- 399 [42] Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov.  
400 Dropout: A simple way to prevent neural networks from overfitting. *The Journal of Machine*  
401 *Learning Research*, 15(1):1929–1958, 2014.
- 402 [43] Shengyang Sun, Changyou Chen, and Lawrence Carin. Learning structured weight uncertainty  
403 in bayesian neural networks. In *Artificial Intelligence and Statistics*, pages 1283–1292, 2017.
- 404 [44] Sethu Vijayakumar and Stefan Schaal. Locally weighted projection regression: Incremental  
405 real time learning in high dimensional space. In *Proceedings of the Seventeenth International*  
406 *Conference on Machine Learning*, pages 1079–1086. Morgan Kaufmann Publishers Inc., 2000.
- 407 [45] Guodong Zhang, Shengyang Sun, David Duvenaud, and Roger Grosse. Noisy natural gradient  
408 as variational inference. *arXiv preprint arXiv:1712.02390*, 2017.
- 409 [46] Han Zhao, Otilia Stretcu, Alex Smola, and Geoff Gordon. Efficient multitask feature and  
410 relationship learning. In *Proceedings of the Thirty-Fifth Conference on Uncertainty in Artificial*  
411 *Intelligence*. AUAI Press, 2019.

412 In this appendix we first describe more related work, and then present missing proofs in the main  
 413 section. We also provide detailed description of our experiments.

## 414 A More Related Work

415 Different kinds of regularization approaches have been studied and designed for neural networks,  
 416 e.g., weight decay [26], early stopping [5], Dropout [42] and the more recent DeCov [6] method.  
 417 BN was proposed to reduce the internal covariate shift during training, but recently it has been  
 418 empirically shown to actually smooth the landscape of the loss function [40]. As a comparison, we  
 419 propose AdaReg as an adaptive regularization method, with the aim to reduce overfitting by allowing  
 420 neurons to share statistical strengths. From the optimization perspective, learning the row and column  
 421 covariance matrices help to converge to better local optimum that also generalizes better.

422 The Kronecker factorization assumption has also been applied in the literature of neural networks  
 423 to approximate the Fisher information matrix in second-order optimization methods [34, 45]. The  
 424 main idea here is to approximate the curvature of the loss function’s landscape, in order to achieve  
 425 better convergence speed compared with first-order method while maintaining the tractability of such  
 426 computation.

427 Determinantal point process (DPP) has been previously applied to compress neural networks [33].  
 428 Specifically, a DPP kernel is placed over the activations of neurons from the same layer, and then  
 429 neurons with similar activations over a fixed dataset are merged into a single one. However, it is  
 430 well known that DPPs can capture only negative correlations [27, 28], and as a result they do not  
 431 stimulate neurons to learn from the experience of other neurons. As a comparison, by explicitly  
 432 learning both precision (covariance) matrices, our framework can account for both positive and  
 433 negative correlations among fan-in/fan-out of neurons from the same layer.

## 434 B Detailed Derivation and Proofs of Our Algorithm

435 We first show that the optimization problem (8) is convex:

436 **Proposition 1.** The optimization problem (8) is convex.

437 *Proof.* It is clear that the objective function is convex: the trace term is linear in  $\Omega_r$  and it is well-  
 438 known that the  $\log \det(\cdot)$  is concave in the positive definite cone [3], hence it trivially follows that  
 439  $\text{Tr}(\Omega_r W \Omega_c W^T) - d \log \det(\Omega_r)$  is convex in  $\Omega_r$ .

440 It remains to show that the constraint set is also convex. Let  $\Omega_1, \Omega_2$  be any feasible points, i.e.,  
 441  $uI_p \preceq \Omega_1 \preceq vI_p$  and  $uI_p \preceq \Omega_2 \preceq vI_p$ . Let  $\forall t \in (0, 1)$ , we have:

$$\|t\Omega_1 + (1-t)\Omega_2\|_2 \leq t\|\Omega_1\|_2 + (1-t)\|\Omega_2\|_2 \leq tv + (1-t)v = v,$$

442 where we use  $\|\cdot\|_2$  to denote the spectral norm of a matrix. Now since both  $\Omega_1$  and  $\Omega_2$  are positive  
 443 definite, the spectral norm is also the largest eigenvalue, hence this shows that  $t\Omega_1 + (1-t)\Omega_2 \preceq vI_p$ .

444 To show the other direction, we use the Courant-Fischer characterization of eigenvalues. Let  $\lambda_{\min}(A)$   
 445 denote the minimum eigenvalue of a real symmetric matrix  $A$ , then by the Courant-Fischer min-max  
 446 theorem, we have:

$$\lambda_{\min}(A) := \min_{\mathbf{x} \neq 0, \|\mathbf{x}\|_2=1} \|A\mathbf{x}\|_2.$$

447 For the matrix  $t\Omega_1 + (1-t)\Omega_2$ , let  $\mathbf{x}^*$  be the vector corresponding to the minimum eigenvalue, hence  
 448 we have:

$$\begin{aligned} \lambda_{\min}(t\Omega_1 + (1-t)\Omega_2) &= \min_{\mathbf{x} \neq 0, \|\mathbf{x}\|_2=1} \|(t\Omega_1 + (1-t)\Omega_2)\mathbf{x}\|_2 \\ &= (t\Omega_1 + (1-t)\Omega_2)\mathbf{x}^* \\ &\geq t\lambda_{\min}(\Omega_1) + (1-t)\lambda_{\min}(\Omega_2) \\ &\geq tu + (1-t)u \\ &= u, \end{aligned}$$

449 which also means that  $t\Omega_1 + (1-t)\Omega_2 \succeq uI_p$ , and this completes the proof. ■

450 The following key lemma characterizes the structure of the normal cone:

451 **Lemma 1.** Let  $\Omega_r \in \mathcal{C}$ , then  $\mathcal{N}_{\mathcal{C}}(\Omega_r) = -\mathcal{N}_{\mathcal{C}}(\Omega_r^{-1})$ .

452 *Proof.* Let  $S \in \mathcal{N}_{\mathcal{C}}(\Omega_r)$ . We want to show  $-S \in \mathcal{N}_{\mathcal{C}}(\Omega_r^{-1})$ . By definition of the normal cone, since  
453  $S \in \mathcal{N}_{\mathcal{C}}(\Omega_r)$ , we have:

$$\text{Tr}(SZ) \leq \text{Tr}(S\Omega_r), \quad \forall Z \in \mathcal{C}$$

454 Now realize that  $\Omega_r \in \mathcal{C}$  and  $\mathcal{C}$  is a compact set, it follows  $\Omega_r$  is the solution of the following linear  
455 program:

$$\max \quad \text{Tr}(SZ), \quad \text{subject to} \quad Z \in \mathcal{C}$$

456 Since both  $S$  and  $Z$  are real symmetric matrix, we can decompose them as  $Z := Q_Z \Lambda_Z Q_Z^T$  and  
457  $S := Q_S \Lambda_S Q_S^T$ , where both  $Q_Z, Q_S$  are orthogonal matrices and  $\Lambda_Z, \Lambda_S$  are diagonal matrices with  
458 the corresponding eigenvalues in decreasing order. Plug them into the objective function, we have:

$$\text{Tr}(SZ) = \text{Tr}(Q_S \Lambda_S Q_S^T Q_Z \Lambda_Z Q_Z^T) = \text{Tr}(\Lambda_S Q_S^T Q_Z \Lambda_Z Q_Z^T Q_S).$$

459 Define  $K := Q_S^T Q_Z$  and  $D = K \circ K$ , where we use  $\circ$  to denote the Hadamard product between two  
460 matrices. Since both  $Q_S$  and  $Q_Z$  are orthogonal matrices, we know that  $K$  is also orthogonal, which  
461 implies:

$$\sum_{j=1}^p D_{ij} = 1, \forall i \in [p], \quad \text{and} \quad \sum_{i=1}^p D_{ij} = 1, \forall j \in [p].$$

462 As a result,  $D$  is a doubly stochastic matrix and we can further simplify the objective function as:

$$\text{Tr}(\Lambda_S Q_S^T Q_Z \Lambda_Z Q_Z^T Q_S) = \text{Tr}(\Lambda_S K \Lambda_Z K^T) = \lambda_S^T D \lambda_Z = \sum_{i,j=1}^p \lambda_{S,i} D_{ij} \lambda_{Z,j},$$

463 where  $\lambda_S$  and  $\lambda_Z$  are  $p$  dimensional vectors that contain the eigenvalues of  $S$  and  $Z$  in decreasing  
464 order, respectively. Now for any  $\lambda_S$  and  $\lambda_Z$  in decreasing order, we have:

$$u \sum_{i=1}^p \lambda_{S,i} \leq \sum_{i=1}^p \lambda_{S,i} \lambda_{Z,1+p-i} \leq \sum_{i,j=1}^p \lambda_{S,i} D_{ij} \lambda_{Z,j} \leq \sum_{i=1}^p \lambda_{S,i} \lambda_{Z,i} \leq v \sum_{i=1}^p \lambda_{S,i} \quad (10)$$

465 From (10), in order for  $\Omega_r$  to maximize the linear program, it must hold that  $D = K = I_p$  and all the  
466 eigenvalues of  $\Omega_r$  are  $v$ . But due to the assumption that  $uv = 1$ , in this case we also know that all  
467 the eigenvalues of  $\Omega_r^{-1}$  are  $1/v = u$ , hence  $\Omega_r^{-1}$  also minimizes the above linear program, which  
468 implies:

$$\text{Tr}(S\Omega_r^{-1}) \leq \text{Tr}(SZ), \quad \forall Z \in \mathcal{C} \Leftrightarrow \text{Tr}(-S(Z - \Omega_r^{-1})) \leq 0 \quad \forall Z \in \mathcal{C}.$$

469 In other words, we have  $-S \in \mathcal{N}_{\mathcal{C}}(\Omega_r^{-1})$ . Using exactly the same arguments it is clear to see that the  
470 other direction also holds, hence we have  $\mathcal{N}_{\mathcal{C}}(\Omega_r) = -\mathcal{N}_{\mathcal{C}}(\Omega_r^{-1})$ . ■

471 Based on the previous first-order optimality condition, it is clear to see that Lemma 1 implies  
472  $W\Omega_c W^T/d - \Omega_r^{-1} \in \mathcal{N}_{\mathcal{C}}(\Omega_r^{-1})$ . Geometrically, this means that the optimum  $\Omega_r^{-1}$  is the Euclidean  
473 projection of  $W\Omega_c W^T/d$  onto  $\mathcal{C}$ . Hence we proceed to derive the projection operator:

474 **Theorem 1.** Let  $\widetilde{\Omega}_r \in \mathbb{S}^p$  with eigendecomposition as  $\widetilde{\Omega}_r = Q\Lambda Q^T$  and  $\text{Proj}_{\mathcal{C}}(\cdot)$  be the Euclidean  
475 projection operator onto  $\mathcal{C}$ , then  $\text{Proj}_{\mathcal{C}}(\widetilde{\Omega}_r) = Q\mathbb{T}_{[u,v]}(\Lambda)Q^T$ .

476 *Proof.* Since  $\Omega_r \in \mathcal{C}$  is real and symmetric, we can reparametrize  $\Omega_r$  as  $\Omega_r := U\Lambda_{\Omega_r}U^T$  where  $U$   
477 is an orthogonal matrix and  $\Lambda_{\Omega_r}$  is a diagonal matrix whose entries corresponds to the eigenvalues of  
478  $\Omega_r$ . Recall that  $U$  corresponds to a rigid transformation that preserves length, so we have:

$$\|\Omega_r - \widetilde{\Omega}_r\|_F^2 = \|U\Lambda_{\Omega_r}U^T - UU^T\widetilde{\Omega}_rUU^T\|_F^2 = \|\Lambda_{\Omega_r} - U^T\widetilde{\Omega}_rU\|_F^2 \quad (11)$$

479 Define  $B := U^T\widetilde{\Omega}_rU$ . Now by using the fact that  $\widetilde{\Omega}_r$  can be eigendecomposed as  $\widetilde{\Omega}_r = Q\Lambda Q^T$ , we  
480 can further simplify (11) as:

$$\|\Lambda_{\Omega_r} - U^T\widetilde{\Omega}_rU\|_F^2 = \sum_{i \in [p]} (\Lambda_{\Omega_r,ii} - B_{ii})^2 + \sum_{i \neq j} B_{ij}^2 \geq \sum_{i \in [p]} (\Lambda_{\Omega_r,ii} - B_{ii})^2 \geq \sum_{i \in [p]} (\mathbb{T}_{[u,v]}(B_{ii}) - B_{ii})^2,$$

481 where the last inequality holds because  $u \leq \Lambda_{\Omega_r,ii} \leq v, \forall i \in [p]$ . In order to achieve the first  
482 equality,  $B = U^T\widetilde{\Omega}_rU$  should be a diagonal matrix, which means  $U^TQ = I_p \Leftrightarrow U = Q$ . In this  
483 case,  $\text{diag}(B) = \Lambda$ . To achieve the second equality, simply let  $\Lambda_{\Omega_r} = \mathbb{T}_{[u,v]}(\text{diag}(B)) = \mathbb{T}_{[u,v]}(\Lambda)$ ,  
484 which completes the proof. ■

Table 2: Stable rank and spectral norm on SARCOS.

	Stable Rank	Spectral Norm
MTL	4.48	0.96
MTL-WeightDecay	4.83	0.92
MTL-AdaReg	<b>2.88</b>	<b>0.70</b>

## C More Experiments

In this section we first describe the network structures used in our main experiments and present more experimental results.

### C.1 Network Structures

**Multiclass Classification (MNIST & CIFAR10):** We use a convolutional neural network as our baseline model. The network used in the experiment has the following structure:  $\text{CONV}_{5 \times 5 \times 1 \times 10} - \text{CONV}_{5 \times 5 \times 10 \times 20} - \text{FC}_{320 \times 50} - \text{FC}_{50 \times 10}$ . The notation  $\text{CONV}_{5 \times 5 \times 1 \times 10}$  denotes a convolutional layer with kernel size  $5 \times 5$  from depth 1 to 10; the notation  $\text{FC}_{320 \times 50}$  denotes a fully connected layer with size  $320 \times 50$ . Similarly, CIFAR10 considers the structure:  $\text{CONV}_{5 \times 5 \times 3 \times 10} - \text{CONV}_{5 \times 5 \times 10 \times 20} - \text{FC}_{500 \times 500} - \text{FC}_{500 \times 500} - \text{FC}_{500 \times 10}$ .

**Multitask Regression (SARCOS):** The network structure is given by  $\text{FC}_{21 \times 256} - \text{FC}_{256 \times 100} - \text{FC}_{100 \times 7}$ .

### C.2 Stable Rank and Spectral Norm on SARCOS

We also show the experimental results of stable ranks and spectral norms on the SARCOS dataset. For the SARCOS dataset, the weight matrix being regularized is of dimension  $100 \times 7$ . Again, we compare the results using three methods: MTL, MTL-WeightDecay and MTL-AdaReg. As can be observed from Table 2, compared with the weight decay regularization, AdaReg greatly reduces both the stable rank and the spectral norm of learned weight matrix, which also helps to explain why MTL-AdaReg generalizes better compared with MTL and MTL-WeightDecay.

### C.3 Combination

As discussed in the main text, combining the proposed AdaReg with BN can further improve the generalization performance, due to the complementary effects between these two approaches: BN helps smoothing the landscape of the loss function while AdaReg also changes the curvature via the row and column covariance matrices (see Fig. 6).

On the other hand, we do not observe significant difference when combining AdaReg with Dropout on this dataset. While we are not clear what is the exact reason for this effect, we conjecture this is due to the fact that Dropout works as a regularizer that prevents coadaptation while AdaReg instead encourages neurons to learn from each other.

### C.4 Ablations

In all the experiments, the AdaReg algorithm is performed on the softmax layer. Here, we study the effects of applying AdaReg algorithm in all CONV/FC layers, all CONV layers, all FC layers, and the last FC layer (i.e., softmax layer). We first discuss how we handle the convolutions in our AdaReg algorithm. Consider a convolutional layer with {input channel, output channel, kernel width, kernel height} being  $\{a, b, k_w, k_h\}$ , we vectorize the original 4-D tensor to be a 2-D matrix of size  $ak_wk_h \times b$ . The AdaReg algorithm can therefore be directly applied on this transformed matrix. Next, we perform the experiment on MNIST with batch size 2048 in Fig. 7. The training set size here is chosen as  $\{128, 256, 512, 1024, 2048, 4096, 8192, 16384, 32768, 60000\}$ .

We find that simply applying the AdaReg algorithm in the softmax layer reaches best generalization as comparing to applying AdaReg on more layers. The improvement is more obvious when the training set size is small. We argue that neural networks can be realized as a combination of a complex

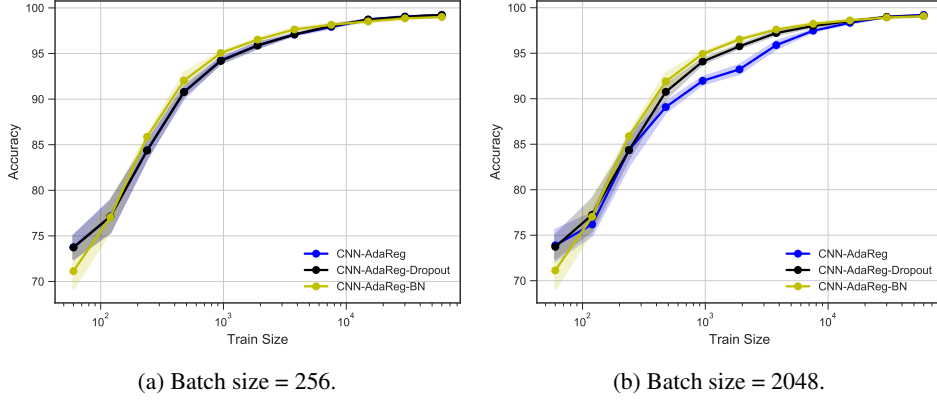


Figure 6: Combine AdaReg with BN and Dropout on MNIST.

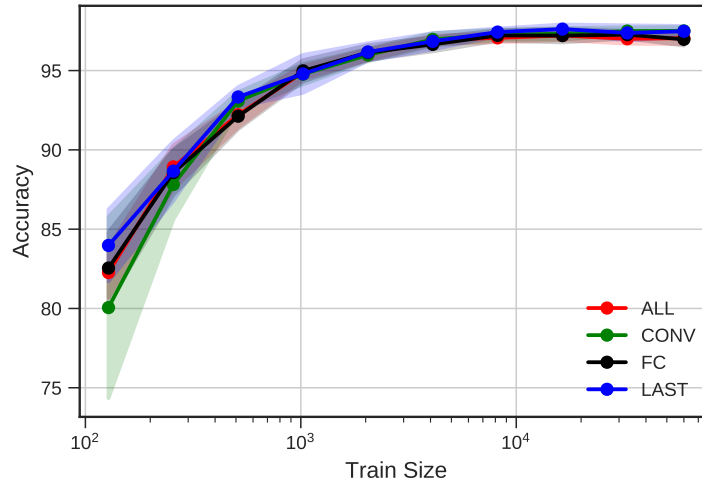


Figure 7: Applying AdaReg on different layers in neural networks for MNIST with batch size 2048.

nonlinear transformation (i.e., feature extraction) and a linear model (i.e., softmax layer). Since AdaReg represents a correlation learning in the weight matrix, it implies that implicit correlations of neurons can also be discovered. In the real world setting, different tasks should be correlated. Therefore, applying AdaReg in the linear model shall improve the model performance by discovering these task correlations. On the contrary, the nonlinear features should be decorrelated for the purpose of generalization. Hence, applying AdaReg in previous layers may lead to adversarial effect.

### C.5 Covariance matrices in the prior

One byproduct that AdaReg brings to us is the learned row and column covariance matrices, which can be used in exploratory data analysis to understand the correlations between learned features and different output tasks. To this end, we visualize both the row and column covariance matrices in Fig. 8. The two covariance matrices on the first row correspond to the ones learned on a training set with 600 instances while the two on the second row are trained with the full dataset on MNIST.

From Fig. 8 we can make the following observations: the structure of both covariance matrices become more evident when trained with larger dataset, and this is consistent with the Bayesian principle because more data provide more evidence. Second, we observe in our experiments that the variances of both matrices are small. In fact, the variance of the row covariance matrix  $\Sigma_r$  achieves the lower bound limit  $u$  at convergence. Lastly, comparing the row covariance matrix  $\Sigma_r$  in Fig. 8 with the one computed from model weights in Fig. 5, we can see that both matrices exhibit the same correlation patterns, except that the one obtained from model weights are more evident, which is

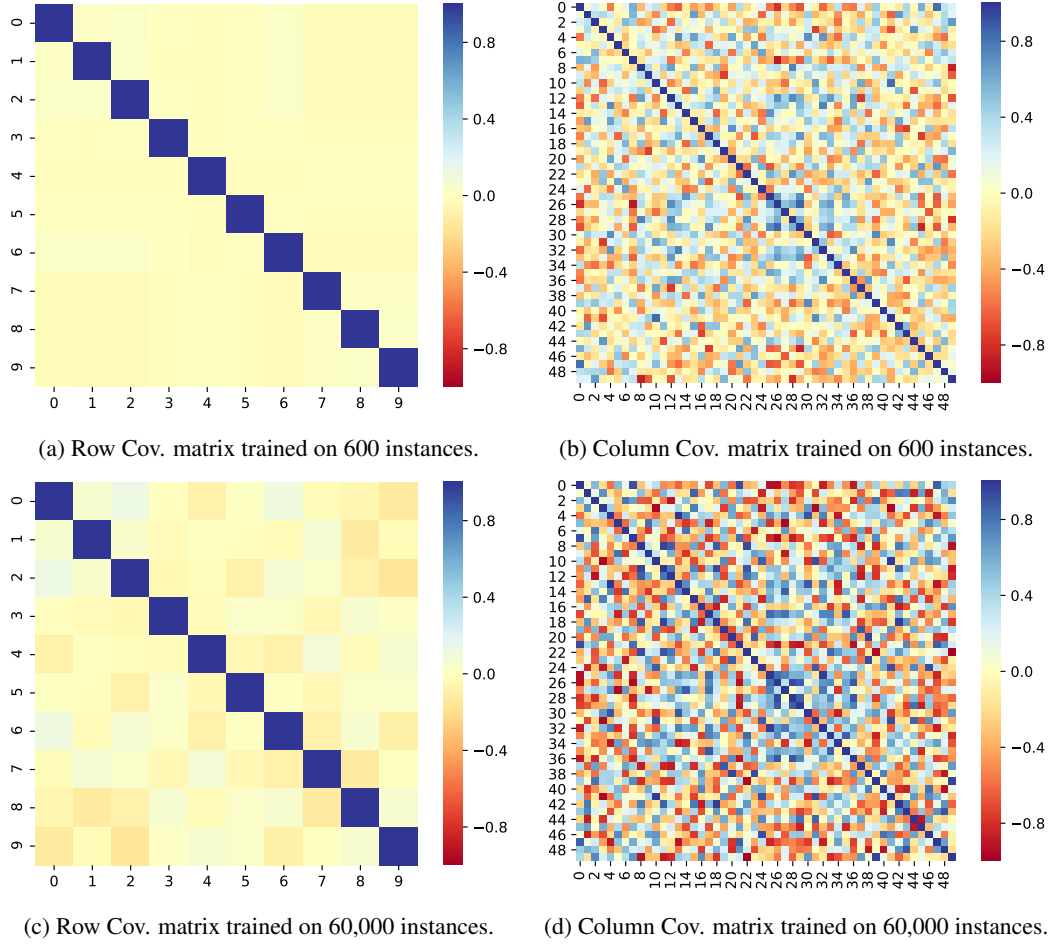


Figure 8: Recovered row covariance matrix  $\Sigma_r$  and column covariance matrix  $\Sigma_c$  in the prior distribution on MNIST.

544 due to the fact that model weights are closer to data evidence than the row covariance matrix in the  
 545 Bayesian hierarchy.

546 On the other hand, the column covariance matrix in Fig. 8 also exhibit rich correlations between the  
 547 learned features, e.g., the neurons in the penultimate layer. Again, with more data, these patterns  
 548 become more evident.

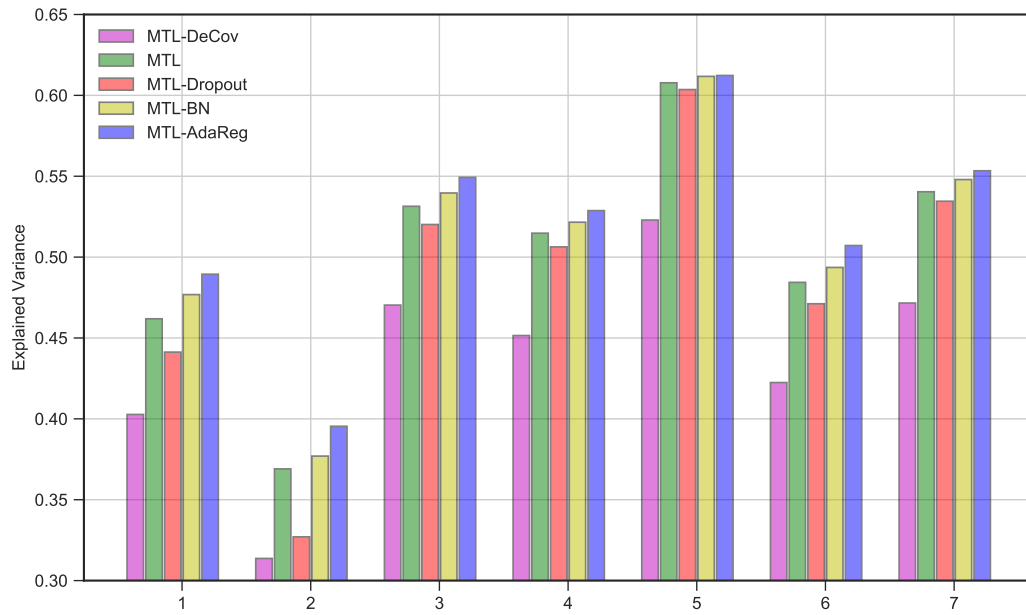


Figure 9: Explained variance of different methods on 7 regression tasks from the SARCOS dataset.