

---

# Statistical Mechanical Analysis of Neural Network Pruning

---

Rupam Acharyya<sup>1</sup>

Ankani Chattoraj<sup>\*2</sup>

Boyu Zhang<sup>\*1</sup>

Shouman Das<sup>3</sup>

Daniel Štefankovič<sup>1</sup>

<sup>1</sup>Computer Science Dept., University of Rochester, Rochester, New York, USA

<sup>1</sup>Brain and Cognitive Science Dept., University of Rochester, Rochester, New York, USA

<sup>1</sup>Mathematics Dept., University of Rochester, Rochester, New York, USA

## Abstract

Deep learning architectures with a huge number of parameters are often compressed using *pruning* techniques to ensure computational efficiency of inference during deployment. Despite multitude of empirical advances, there is a lack of theoretical understanding of the effectiveness of different pruning methods. We inspect different pruning techniques under the statistical mechanics formulation of a teacher-student framework and derive their generalization error (GE) bounds. It has been shown that *Determinantal Point Process* (DPP) based *node* pruning method is notably superior to competing approaches when tested on real datasets. Using GE bounds in the aforementioned setup we provide theoretical guarantees for their empirical observations. Another consistent finding in literature is that sparse neural networks (*edge pruned*) generalize better than dense neural networks (*node pruned*) for a fixed number of parameters. We use our theoretical setup to prove this finding and show that even the baseline *random edge pruning* method performs better than the *DPP node pruning* method. We also validate this empirically on real datasets.

## 1 INTRODUCTION

Deep neural networks have achieved impressive results in a wide variety of applications such as classification [23, 31], image processing [30, 4], natural language processing [8, 7, 42], etc. Most of these networks use millions and sometimes even billions of parameters which makes inference computationally expensive and memory intensive [8]. To address this, researchers explore pruning techniques with the primary goal of comparing performance on real

datasets. The broad scientific paradigm explored by most pruning techniques is to empirically and heuristically determine either how to prune a network or what to prune in a network (sometimes both). In this work, we take a step towards theoretical understanding of these two prime aspects of pruning methods.

We compare the quality of different pruning methods for feedforward neural networks under the *teacher-student* framework [37, 38, 39, 13] in the thermodynamic limit (input dimension goes to infinity) using *generalization error bounds* (GE), a theoretical measure of performance of machine learning models on unseen test data [46].

A fairly recent work by [34] empirically investigates a node pruning technique where a diverse subset of nodes are preserved in a given layer using Determinantal Point Process (DPP) [32, 24]. We provide theoretical guarantees for their empirical observations thereby showing that DPP based node pruning outperforms two standard paradigms of pruning (magnitude based node pruning and random node pruning). Thus, in the first part of this paper, we take a step towards theoretical understanding of the question: how to prune?

For the second part of this work we focus our attention to the study by [6]. This study reviewed multiple papers across decade on various pruning methods and closely analyzed their empirical results to conclude that sparse models obtained after edge/connection (used interchangeably) pruning outperforms dense ones obtained after node pruning for a fixed number of parameters. We extend our theoretical setup and compare node and edge pruning techniques which are within the scope of our investigation, to provide a theoretical justification of their empirical observation driven claim, thereby addressing the question: what to prune?

Our work has multiple contributions with regard to theoretical advancements in the domain of pruning:

- We use GE bounds on the teacher-student framework to compare different pruning methods within a class,

---

\*equal contribution

which to the best of our knowledge, is the first theoretical advance in comparing pruning methods.

- We prove that DPP node pruning outperforms random and importance node pruning methods, previously shown by [34] empirically.
- We also theoretically show and validate on real datasets (MNIST and CIFAR10), that baseline random edge pruning performs better than DPP node pruning (superior in the node pruning regime explored in this paper) which is consistent with empirical observations from pruning literature that sparse models outperform dense models [6].

## 2 RELATED WORK

**Pruning Methods:** Studies under node pruning regime remove entire neurons/nodes (used interchangeably henceforth) keeping the networks dense [17, 29, 18]. Our work is closely related to [34], where a DPP sampling technique is used to select a set of diverse neurons/nodes to be preserved during pruning. The authors also introduce a *reweighting* procedure to compensate contributions of the pruned neurons in the network. Finally, they compare DIVNET (DPP node pruning with reweighting as in [34]) with random and importance node pruning [17] on real datasets. Seminal studies on edge pruning [26, 16] remove unimportant network weights based on the Hessian of the network’s error function. Among others, alternative approaches include low rank matrix factorization of the final weight layers [40] or pruning the unimportant connections below a threshold [15]. Though dense networks can benefit from modern hardware, sparse models outperform dense ones for a fixed number of parameters across domains [27, 21, 14]. In a recent review this is highlighted based on observations from investigating 81 studies on pruning techniques [6].

The various existing methods can be broadly subsumed into a couple of categories [6]. These categories are mainly governed by the principles of pruning heuristics. First category is the magnitude-based approaches which have been extensively studied both globally and layerwise [15, 11]. As per [6], magnitude-based approaches are not only good and common baselines in the literature but they also give comparable performance to other methods such as the gradient-based methods [27, 47]. Another category is the random pruning which serves as an useful baseline for showing superior performance of any other pruning technique. We hence show all our theoretical results w.r.t these two categories, random pruning and importance pruning (same in concept as magnitude based pruning). We do not focus on any specific algorithm within these categories but explore the general concept for theoretical results. There are recent advances in pruning techniques which are complementary to these approaches, such as, being data independent [5, 43], single shot [28, 45] etc. However, these are beyond the scope of

our investigation.

**Theoretical Advances Towards Understanding Neural Networks:** Despite promising performance in empirical data, providing theoretical guarantees for neural networks remains a known challenge. Researchers have explained the training dynamics of neural network from the information theoretic perspective [44, 41]. In another direction of work the learning dynamics of neural networks with infinitely wide hidden layers are explored [19, 9, 2, 48]. Pioneering work by [37, 38, 39] analyzes the generalization dynamics from the statistical mechanics perspective on *teacher-student* framework [12] to understand the performance of neural networks on unseen test data. All our theoretical analyses throughout this work closely follow [1, 13], who analyzed results for the case where the student networks are over parameterized, i.e., it has more number of hidden nodes than the teacher network.

## 3 PRELIMINARIES

**Determinantal Point Process (DPP):** DPP [32] is a probability distribution over power set of a ground set  $\mathcal{G}$ , here finite. DPP is a special case of negatively associated distributions [20] which assigns higher probability mass on diverse subsets. Formally, a DPP with a marginal kernel  $L$  ( $\in \mathbb{R}^{|\mathcal{G}| \times |\mathcal{G}|}$ ) is:  $\mathbb{P}[\mathbf{Y} = Y] = \frac{\det(L_Y)}{\det(L+I)}$ , where  $Y \subseteq \mathcal{G}$  and  $L_Y$  is the principal submatrix defined by the indices of  $Y$ . We use  $k$ -DPP to denote the probability distribution over subsets of fixed size  $k$ .

**DPP Node Pruning:** [34] uses DPP to propose a novel node pruning method for feedforward neural network. They define information at node  $i$  of layer  $l$  as  $\mathbf{a}_i^l (= (a_{i1}^l, \dots, a_{in}^l))$ , where  $a_{ij}^l$  is the activity of node  $i$  of layer  $l$  on  $j^{\text{th}}$  input. Here  $\mathbf{a}_i^l = g(\mathbf{b}_i^l)$ , where  $\mathbf{b}_i^l = \sum_{j=1}^{n_{l-1}} w_{ji}^{l-1} \mathbf{a}_j^{l-1}$  is the information at node  $i$  of layer  $l$  before activation. A layer is pruned by choosing a subset of hidden nodes using a DPP kernel:  $\mathbf{L}$  ( $= \mathbf{L}' + \varepsilon \mathbf{I}$ ), where,  $\mathbf{L}'_{st} = \exp(-\beta \|\mathbf{a}_s^l - \mathbf{a}_t^l\|^2)$  and  $\beta$  is a bandwidth parameter. The matrix  $\mathbf{L}$  is of dimension  $n_l \times n_l$ , as total number of nodes in layer  $l$  is  $n_l$ . By the property of DPP, this procedure will keep a diverse subset of nodes for each layer w.r.t. information obtained from the training data. A *reweighting* technique (see Section 2.2 of [34]) is then applied to outgoing edges of retained nodes to compensate for information lost in that layer due to node removal.

**Remark:** DIVNET denotes DPP node pruning with reweighting as in [34].

**Online Learning in Teacher-Student Setup [13]:** We use a two-layer perceptron which has  $N$  input units,  $M$  hidden units and 1 output unit as the *teacher network* to generate labels for i.i.d Gaussian input,  $\mathbf{x}^t = (x_1^t, \dots, x_N^t)$  where  $x_i^t \sim \mathcal{N}(0, 1) \forall i \in \{1, \dots, N\}$ . Let  $\theta^* = \{\mathbf{w}^* (\in \mathbb{R}^{M \times N}), \mathbf{v}^* \in \mathbb{R}^M\}$  denote the fixed parameters of the teacher network.

Table 1: Notations used in Theorems

Notations	Explanations	Notations	Explanations	Notations	Explanations
$n$	number of inputs	$N$	dimension of the input	$n_l$	number of nodes in layer $l$
$v_i^l$	$i^{th}$ node in layer $l$ ( $1 \leq i \leq n_l$ )	$a_{ij}^l$	activation of $v_i^l$ on $j^{th}$ input	$M$	number of teacher hidden nodes
$e_{ij}^l$	edge from $v_i^l$ to $v_j^{l+1}$ ( $1 \leq i \leq n_l$ and $1 \leq j \leq n_{l+1}$ )	$w_{ij}^l$	weight of $e_{ij}^l$ ( $1 \leq i \leq n_l$ and $1 \leq j \leq n_{l+1}$ )	$K$	number of student hidden nodes
$k_n$	number of student hidden nodes kept after node pruning	$k_e$	number of incoming edges of a hidden node kept after edge pruning	$v^*$	second layer weight of teacher network

The label  $y^t$  of the input  $\mathbf{x}^t$  ( $t = 1, 2, \dots$ ) is given as,

$$y^t = \sum_{m=1}^M v_m^* g\left(\frac{w_m^* \mathbf{x}^t}{\sqrt{N}}\right) + \sigma \zeta^t, \quad (1)$$

where  $\zeta^t \sim \mathcal{N}(0, 1)$  is the output noise, and  $g$  is the sigmoid activation function. The input and teacher generated labels ( $\{\mathbf{x}^1, y^1, \dots\}$ ) are used to train a two-layer *student network* with  $N$  input units,  $K$  hidden units ( $K \geq M$ ) and 1 output unit using online SGD learning method. We consider the quadratic training loss, i.e.,

$$L(\theta) = \frac{1}{2} \left[ \sum_{k=1}^K v_k g\left(\frac{w_k \mathbf{x}^t}{\sqrt{N}}\right) - y^t \right]^2, \quad (2)$$

where  $\theta = \{\mathbf{w}, \mathbf{v}\}$  denotes the parameter of the student network. [13] showed that GE  $\varepsilon(f)$  (expected error on the unseen data, for details see S31 of [13]) for the student network is a function of the following *order parameters*,

$$Q_{ik} = \frac{w_i^T w_k}{N}, \quad R_{in} = \frac{w_i^T w_n^*}{N}, \quad R_{mn} = \frac{w_m^* T w_n^*}{N}. \quad (3)$$

Intuitively, these order parameters measure the similarities between and within the hidden nodes of teacher and student networks. Our theoretical results assume [13]:

- (A1) If  $\mathbf{x} = (x_1, \dots, x_N)$  is an input then  $x_i \in \mathcal{N}(0, 1)$ . Also,  $N \rightarrow \infty$ .
- (A2) Both the teacher and the student networks have only one hidden layer.
- (A3)  $K \geq M$  and  $K = Z \cdot M$  where  $Z \in \mathbb{Z}^+$ .
- (A4) The activation in the hidden layer is sigmoidal for both teacher and student network.
- (A5) The output  $\in \mathbb{R}$  (i.e., regression problem).
- (A6) The order parameters (see section 3) satisfy the ansatz as in (S58) - (S60) of [13].
- (A7) No noise is added to the labels generated by the teacher network, i.e.,  $\sigma = 0$  in (1).

## 4 GE OF PRUNED NETWORK IN TEACHER-STUDENT SETUP

We compare the performance of student networks pruned using different techniques as in Table 2 by analyzing their

GE (see Figure 1). For node and edge pruning comparison, we choose the parameters  $k_n$  and  $k_e$  (see Table 1) such that the total number of parameters of the networks remain same, i.e., they satisfy,

$$\frac{k_n}{K} = \lim_{N \rightarrow \infty} \frac{k_e}{N} = c, \quad (4)$$

where  $c \in [0, 1]$  is a constant. It is important to note that since we assume that the number of student nodes is more than the number of teacher nodes, which means multiple student nodes learn the same teacher node (see Figure 3 of [13]; also in Figure 1: two student hidden nodes learn one teacher hidden node, shown in same color). From [13], we know that, in noiseless case ( $\sigma = 0$  in (1)), the student network learns the teacher network completely when trained till convergence, i.e., the GE becomes 0. When we prune the student network, this GE increases, which we then analyze for different types of pruning under certain assumptions (see Section 3 (A1)-(A7)).

### 4.1 COMPARING NODE PRUNING METHODS

We theoretically show that the increment in GE due to DIVNET is less than that for random and importance node pruning methods, justifying the empirical findings of [34]. The proof proceeds with the following steps: (1) Theorem 1 provides a closed form expression of the GE after DPP node pruning. (2) Theorem 2 shows that: (a) GE of random node pruning is greater than GE of DPP node pruning (b) GE of random node pruning with reweighting is greater than GE of DIVNET (c) GE of importance node pruning is greater than GE of DIVNET.

**Theorem 1.** Assume (A1) – (A7). Let  $k_n \leq M$  nodes are selected by the DPP Node pruning method,

$$\varepsilon_{k_n}^{DPPNode}(f) = (v^*)^2 \left[ \frac{k_n}{6} \left(1 - \frac{1}{Z}\right)^2 + \frac{M - k_n}{6} \right] \quad (5)$$

and

$$\hat{\varepsilon}_{k_n}^{DPPNode}(f) = (M - k_n) \times \frac{(v^*)^2}{6}. \quad (6)$$

**Proof Idea of Theorem 1:** Proof of the above theorem (details in the appendix C) is based on two factors: (1) Results

Table 2: Different pruning methods and notations for their GE. Here  $f$  denotes the pruned student network. u.a.r. and w.p. stand for *uniformly at random* and *with probability* respectively.

Pruning Method	Procedure	Retained Parameters	GE without reweighting	GE with reweighting
Random Node	Keep $k_n$ nodes u.a.r.	$k_n$ hidden nodes	$\mathcal{E}_{k_n}^{RandNode}(f)$	$\hat{\mathcal{E}}_{k_n}^{RandNode}(f)$
Importance Node	[17]	$k_n$ hidden nodes	$\mathcal{E}_{k_n}^{ImpNode}(f)$	$\hat{\mathcal{E}}_{k_n}^{ImpNode}(f)$
DPP Node	see Section 3	$k_n$ hidden nodes	$\mathcal{E}_{k_n}^{DPPNode}(f)$	$\hat{\mathcal{E}}_{k_n}^{DPPNode}(f)$
Random Edge	Keep an edge w.p. $c$ for each hidden node	$k_e$ incoming edges per hidden node	$\mathcal{E}_{k_e}^{RandEdge}(f)$	$\hat{\mathcal{E}}_{k_e}^{RandEdge}(f)$

from [13] assure that analyzing the *order parameters* is enough to obtain closed form of GE. (2) We exploit the observation that the expected kernel of the DPP node pruning is same as the order parameter  $Q$  (see appendix B for proof and Figure 2 E) which, following [13], is a block diagonal matrix with  $M$  blocks. By property of DPP, the pruning method will retain a subset of student hidden nodes with at most 1 hidden node from each block when  $k_n \leq M$  (see Figure 2 G).

**Remark 1.** *As the expected DPP kernel is block-diagonal matrix, the stochasticity in subset selection via DPP does not impact GE when subset size is fixed and it only depends on size of pruned subsets.*

**Remark 2.** *Our theorem uses  $k_n \leq M$ , however, in practice the kernel may have non-zero off-diagonal entries when the assumption (A1) about input data is violated. As a result the probability of sampling a subset of size  $k_n > M$  may be nonzero.*

**Connection to Lottery Ticket Hypothesis:** An interesting direction of research is to find small sub-networks from an overparameterized network with comparable performance. The existence of such networks is hypothesized in Lottery Ticket hypothesis [10]. Interestingly, recent work shows that pruning helps find such networks even without retraining [36, 33] and in our work we explore a sub-network in the teacher student setup.

Note that from Eq (6), when  $M$  student nodes are kept after pruning, i.e.,  $k_n = M$ , then the GE of the DPP node pruned network is 0 which is GE of the original student network. Hence, from the fact that  $K > M$  we can conclude that DPP node pruning can find out the winning ticket, i.e., a small sub-network with much less number of parameters than the original unpruned network but with same performance guarantee.

**Theorem 2.** *Assume (A1) – (A7). Then for  $k_n \leq M$  we have,*

$$\mathbb{E}_f \left[ \mathcal{E}_{k_n}^{RandNode}(f) \right] > \mathcal{E}_{k_n}^{DPPNode}(f') \quad (7)$$

and

$$\mathbb{E}_f \left[ \hat{\mathcal{E}}_{k_n}^{RandNode}(f) \right] > \hat{\mathcal{E}}_{k_n}^{DPPNode}(f') \quad (8)$$

and,

$$\mathcal{E}_{k_n}^{ImpNode}(f') > \hat{\mathcal{E}}_{k_n}^{DPPNode}(f'), \quad (9)$$

*i.e., DPP node pruning outperforms random node pruning in the above setup. Here the expectation is taken over the the subsets of hidden nodes of size  $k_n$  chosen u.a.r.*

**Remark 3.** *Reweighting for DPP/random node pruning follow procedure in Section 2.2 of [34].*

**Proof Idea of Theorem 2:** In random and importance node pruning, two student nodes which learn the same teacher node may both survive after pruning with non-zero probability, unlike DPP node pruning (Figure 1 (B)). Hence, more teacher nodes may remain unexplained by the student network after random or importance node pruning, resulting in increased GE (details in appendix C).

Together, Theorem 1 and 2 gives theoretical guarantees for all empirical results of [34]. Theorem 1 further allows us to show that DIVNET indeed satisfies the stronger version of Lottery Ticket Hypothesis as recently explored in [36, 35]. Importance node pruning with reweighting may be better than DIVNET and was not explored in [34].

## 4.2 COMPARING NODE AND EDGE PRUNING METHODS

In random edge pruning method, for each student hidden node, an incoming edge is kept with probability  $c = \lim_{N \rightarrow \infty} \frac{k_e}{N}$ . Majority of empirical studies throughout literature use random edge or node pruning as a baseline for empirical comparison (see papers in [6]) making it an obvious candidate for our theoretical comparisons as well. It has been shown empirically by [34] and theoretically by us that DPP node pruning is an above baseline node pruning method. In this section we show that baseline random edge pruning outperforms DPP node pruning which is consistent with the empirical observations that sparse models outperform dense models (section 3.2 of [6]). Specifically, here we show that GE after random edge pruning is less than GE after DPP node pruning. Our proof proceeds as follows: (1) Theorem 3 gives a closed form expression for the GE after random

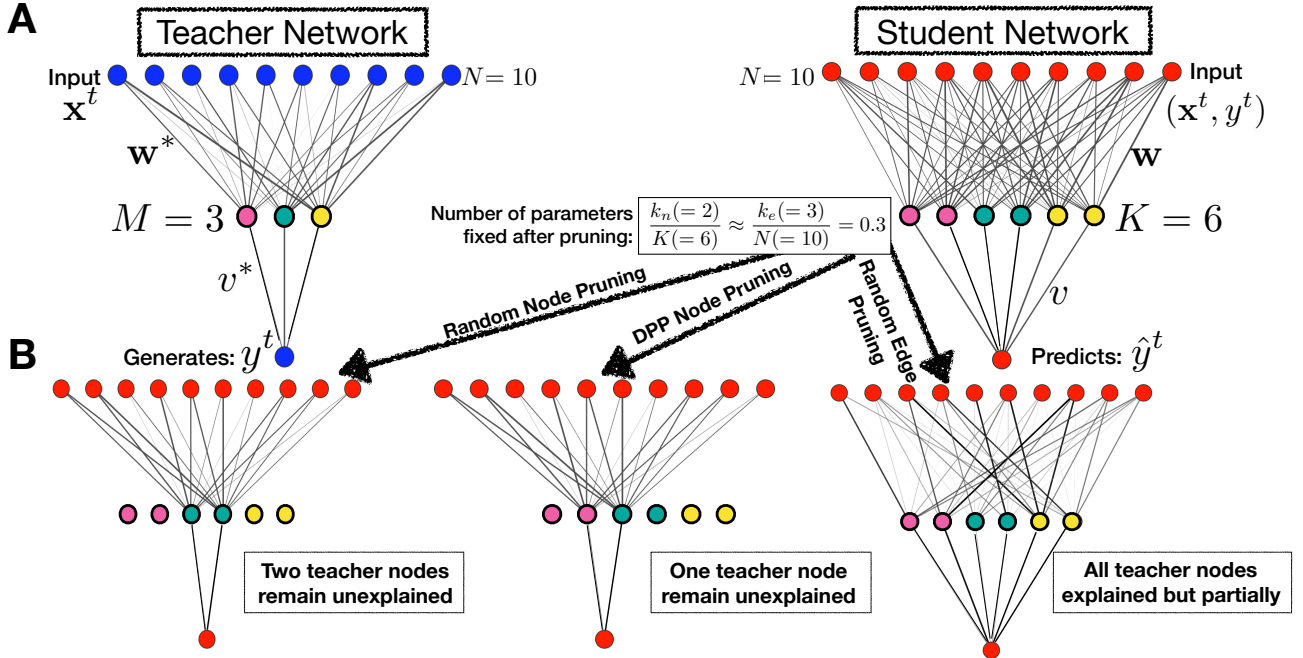


Figure 1: (A) Two layer teacher-student framework: A teacher neural network with 3 hidden nodes (left) and a student network with 6 hidden nodes (right). Input data (i.i.d) along with its label generated by teacher network are fed to student network to predict. (B) Intuitive example for 3 types of pruning on student network. For  $k_n = 2$ , random node pruning might only be able to explain 1 teacher hidden node, whereas DPP node pruning will always retain (partial) information about 2 teacher hidden nodes, hence performs better. Random edge pruning retains sparse information about all 3 teacher nodes which is enough to outperform DPP node pruning. All notations follow Table 1.

edge pruning (2) Theorem 4 then shows that GE of random edge pruning is less than GE of DPP node pruning.

**Theorem 3.** Assume (A1) – (A7). Consider the random edge pruning method with parameter  $\lim_{N \rightarrow \infty} \frac{k_e}{N} = c$  (here  $c$  is a constant between 0 and 1). Then the GE  $\epsilon_c^{\text{RandEdge}}(\mathbb{E}[f])$  is,

$$\frac{M(v^*)^2}{\pi} \left[ \frac{1}{Z} \arcsin \frac{c}{1+c} + \left(1 - \frac{1}{Z}\right) \arcsin \frac{c^2}{1+c} + \frac{\pi}{6} - 2 \arcsin \frac{c}{\sqrt{2(1+c)}} \right]. \quad (10)$$

**Remark 4.** Theorem 3 gives closed form for “GE of the expected network” after pruning instead of the “expected GE of the network” after pruning. However, in the thermodynamic limit ( $N \rightarrow \infty$ ), the order parameters as in Section 3 are highly concentrated near their expected values and the two quantities hence become equal.

**Theorem 4.** Assume (A1) – (A7). Let  $k_n$  and  $c$  satisfy (4), and  $0 \leq c \leq \frac{1}{Z}$  and  $Z \geq 4$ . Then

$$\epsilon_{k_n}^{\text{DPPNode}}(f) \geq \epsilon_c^{\text{RandEdge}}(\mathbb{E}[f]), \quad (11)$$

i.e., Random edge pruning outperforms DPP node pruning in the above setup.

**Proof Idea of Theorem 4:** When  $k_n \leq M$ , node pruned student network leaves at least  $(M - k_n)$  teacher nodes unexplained, whereas after random edge pruning, student network can retain at least partial information about every teacher node (see Figure 1 (B)). After a pruning routine, the sum of partial information about all teacher nodes in an edge pruned student network dominates the sum of information for the explained subset of teacher nodes in a node pruned student network.

**Observations:** From Theorem 2 and 4, we conclude that random edge pruning outperforms random node pruning. Further, using Theorem 2 and the intuition that importance edge pruning is better than random edge pruning, we expect that importance edge pruning will outperform importance node pruning. Figure 2 confirms this empirically in the teacher student setup. These observations leads to the conjecture that for a fixed pruning method, edge pruning outperforms node pruning.

**Conjecture 1.** Assume (A1) – (A7). Let  $k_n$  and  $c$  satisfy (4) and Prune denotes a fixed pruning method (e.g. Rand, Imp) which can be applied to both node and edge. Then,  $\exists c_\epsilon \in (0, 1]$  such that for  $0 \leq c \leq c_\epsilon$ ,

$$\epsilon_{k_n}^{\text{PruneNode}}(f) \geq \epsilon_c^{\text{PruneEdge}}(f). \quad (12)$$

Together, Theorem 3, 4 and Conjecture 1 are consistent with

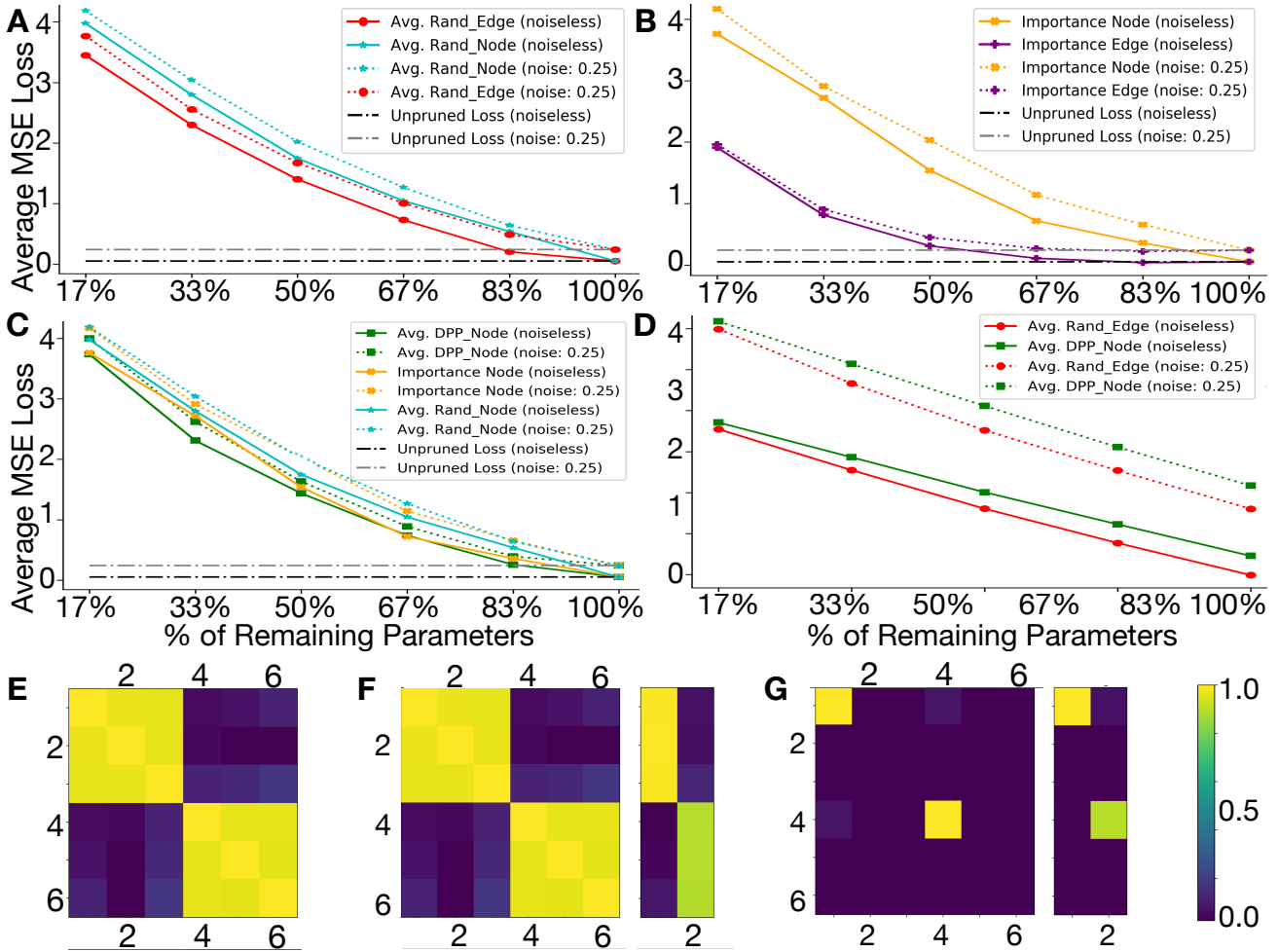


Figure 2: Simulation results in teacher student setup,  $M = 2$  and  $K = 6$  for (A-G). (A-B) Edge pruned networks perform better than node pruned networks in all 3 types of pruning methods (random (A), importance (B)), validating Conjecture 1. (C) DPP Node pruning performs better than importance and random node pruning (Theorem 2) (D) Baseline random edge pruning beats DPP node pruning (Theorem 4). For (D),  $M = 5$  and  $K = 20$ . (E) The kernel of DPP node pruning is same as  $Q$  (F) Order parameters,  $Q$  (same as (E)) and  $R$  of the unpruned student network. (G) When only keeping 2 nodes, DPP node pruned student network keeps one from each block shown in (G).

empirical observations of [6]: sparse networks after edge pruning perform better on the unseen test data than dense networks after node pruning with fixed number of parameters. To the best of our knowledge, [6] based their claims from empirical observations of pruning studies in which the pruned networks were not reweighted. This motivated our choice of comparing GE for DPP node pruning and random edge pruning without any reweighting. However, with reweighting from [34], GE of DPP node pruning will be less than GE of random edge pruning, highlighting the impact of reweighting proposed by [34] (proof and details in appendix C).

We find that GE analysis on teacher-student setup is flexible for various pruning methods and this framework can be extended to theoretically understand other pruning methods which are outside the scope of this work.

## 5 EXPERIMENTS

### 5.1 SIMULATIONS

We run the DPP node, random edge/node, and importance edge/node pruning simulations under the teacher-student setup. For all the simulations, we sampled the 800000 i.i.d input samples from  $\mathcal{N}(0, 1)$  as training data and 80000 as testing data. Following notations from Table 1, we set  $M = 2$ ,  $K = 6$ ,  $N = 500$ , and  $v^* = 4$ . The first layer teacher network weights  $\mathbf{w}^*$  and all the student network parameters  $\theta = \{\mathbf{w}, \mathbf{v}\}$  were drawn independently from  $\mathcal{N}(0, 1)$  as initialization. We choose learning rate  $\eta = 0.50$ , and it is scaled to  $\frac{\eta}{\sqrt{N}}$  for  $\mathbf{w}$  and  $\frac{\eta}{N}$  for  $\mathbf{v}$ . We run the simulations for both noiseless ( $\sigma = 0$  in (1)) and noisy ( $\sigma = 0.25$ ) output labels. For comparisons between node and edge pruning, we

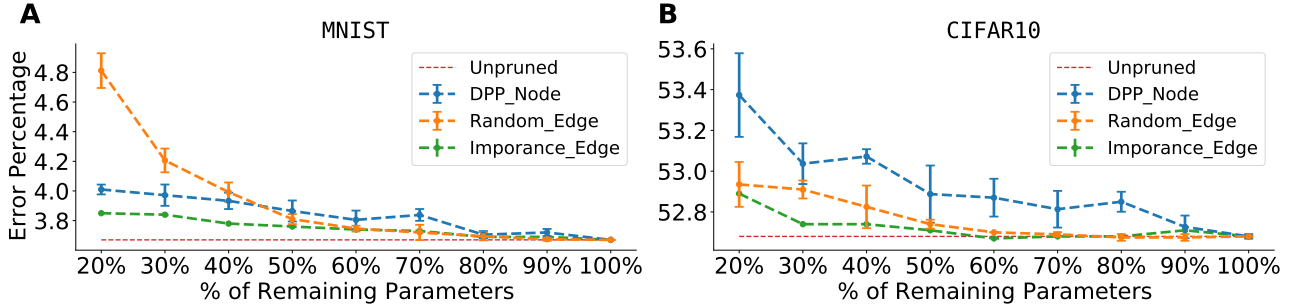


Figure 3: Comparing different edge pruning methods with DPP Node pruning method on the MNIST (A) and CIFAR10 (B) dataset. Horizontal axis represents the percentage of remaining parameters in 1<sup>st</sup> layer after pruning. The vertical axis shows corresponding test error. Both magnitude based edge pruning method (importance pruning) and baseline random edge pruning method outperforms DPP Node pruning which confirms Theorem 4 and the conjecture proposed in [6].

use the node-to-edge ratio [1 : 83, 2 : 166, 3 : 250, 4 : 333, 5 : 417, 6 : 500] to keep the number of parameters the same, given  $N = 500$ ,  $K = 6$ , and  $M = 2$ . In addition, we run the same simulation with  $K = 5$  and  $M = 20$ , see Figure 2D. For other simulation details and results, see appendix. Note that no pruning method undergoes reweighting for reported simulation results which we therefore use to verify and validate our theoretical results without reweighting.

#### Key Observations:

- The expected kernel of the DPP node pruning and the  $Q$  matrix are the same which we exploit for Theorem 1 (Figure 2E,F).
- For  $k_n = 2$  and  $M = 2$ , DPP node pruning chooses exactly one node from each of the diagonal block of the kernel (see Figure 2G) which validates Theorem 1.
- DPP node pruning outperforms random and importance node in both noisy and noiseless case (see Figure 2C) which confirms Theorem 2.
- Random edge pruning is better than DPP node pruning for  $c \leq \frac{1}{Z}$  with  $Z = 4$  and  $M = 5$  in both noisy and noiseless cases (see Figure 2D), validating Theorem 4.
- We see Conjecture 1 holds for random, importance edge and node pruning (see Figure 2A,B)

## 5.2 REAL DATA

In this section, we compare DIVNET by [34] with random edge pruning with reweighting, and importance edge pruning with reweighting on the MNIST [25] and CIFAR10 [22] datasets. We used the exact same network architectures as in Table 1 of [34] for MNIST and CIFAR10, respectively. Note that, for the real data we consider network structures with multiple layers. Following [34], we performed all pruning methods on the first layer. We compare the number of parameters as  $k_e = \frac{k_n(d_{\text{input}} + h_2)}{h_1} - h_2$  where  $k_e$  is the number of edges kept for each node in edge pruning, and  $k_n$  is the

number of nodes kept in the hidden layer for node pruning;  $d_{\text{input}}$ ,  $h_1$ , and  $h_2$  represent the dimension of the input, size of the first hidden layer, and size of the second hidden layer, respectively. As in [34],  $h_1 = h_2$ . We trained our model until the training error reaches predefined thresholds (Table 1 in [34]) and then perform the pruning. For hyperparameters and other details, see F.

**Remark:** Note that we have not presented the results comparing different node pruning methods among themselves as they were already discussed in [34].

#### Key observations:

- Baseline random edge pruning method outperforms DIVNET across all percentages of parameters retained in the network for CIFAR10 dataset shown in Fig 3 B. However, for MNIST dataset, DIVNET performs better than random edge initially but if  $> 40\%$  of parameters are retained in the network random edge outperforms DIVNET (see Fig 3 A).
- Importance edge pruning performs better than both DIVNET and the baseline random edge pruning method on both the real data sets which highlighting the potential of magnitude based pruning method (see Fig 3 A and B).

## 6 DISCUSSION AND FUTURE WORK

Our work takes the first step to develop theoretical comparison for empirical observations of pruning methods in feed forward neural networks. We identify the usefulness of teacher-student setup for providing theoretical guarantees of pruning methods. We then use this setup to theoretically show that DIVNET should indeed outperform random and importance node pruning techniques. We further show that random edge pruning outperforms DPP node pruning providing a theoretical proof for the popular empirical observation:

sparse (node) networks perform better than dense (edge) pruned networks for fixed number of parameters. Finally, we also are able to show that DIVNET satisfies a stronger version of the Lottery Ticket Hypothesis. Our work consolidates the understanding of a particular class of node and edge pruning theoretically.

When comparing two neural networks, using the number of parameters may not always be the optimal choice, instead, measuring the *capacity* and *expressiveness* of neural networks [3] can provide new insights. All our theoretical results have been proved on single hidden layer neural networks which gives future scope of extending them to multiple hidden layer networks. However, our empirical results hold for neural networks with multiple hidden layers suggesting the possibility of generalization of our results.

Throughout this work, we focus only on pruning methods in which a feedforward pre-trained neural network is pruned once without retraining. We choose this class for two primary reasons: (1) it is feasible to make theoretical comparisons with closed form solutions of GE, and, (2) with some assumptions, it has been shown by recent studies [36, 33, 35] that every sufficiently over-parameterized network contains a sub network that, even without training, achieves comparable accuracy to the trained large network. This proven conjecture is even stronger than the Lottery Ticket Hypothesis [10]. Hence, comparing performance of pruning methods within the aforementioned class in the teacher-student setup allowed us explore the existence of such a sub network.

We compare our theoretical results with random pruning and importance pruning which subsumes ideas underlying vast majority of pruning techniques and do not focus on any specific algorithm. A more specific algorithm based justification can also be an extension (may not always be trivial however) of this paradigm.

We introduce the teacher-student setup for proving results related to pruning methods which can further be extended to prove other empirical results in the pruning domain. Such theoretical insights can also be used as a means to guide development of theory-motivated new and better pruning algorithms on other neural network architectures like CNNs and RNNs in future work.

## References

- [1] Madhu S Advani and Andrew M Saxe. High-dimensional dynamics of generalization error in neural networks. *arXiv preprint arXiv:1710.03667*, 2017.
- [2] Sanjeev Arora, Simon S Du, Wei Hu, Zhiyuan Li, Ruslan Salakhutdinov, and Ruosong Wang. On exact computation with an infinitely wide neural net. *arXiv preprint arXiv:1904.11955*, 2019.
- [3] Sanjeev Arora, Rong Ge, Yingyu Liang, Tengyu Ma, and Yi Zhang. Generalization and equilibrium in generative adversarial nets (gans). In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pages 224–232. JMLR. org, 2017.
- [4] Vijay Badrinarayanan, Alex Kendall, and Roberto Cipolla. Segnet: A deep convolutional encoder-decoder architecture for image segmentation. *IEEE transactions on pattern analysis and machine intelligence*, 39(12):2481–2495, 2017.
- [5] Mussay Ben, Margarita Osadchy, Vladimir Braverman, Samson Zhou, and Dan Feldman. Data-independent neural pruning via coresets. In *International Conference on Learning Representations (ICLR)*, 2020.
- [6] Davis Blalock, Jose Javier Gonzalez Ortiz, Jonathan Frankle, and John Gutttag. What is the state of neural network pruning? *arXiv preprint arXiv:2003.03033*, 2020.
- [7] Li Deng and Yang Liu. *Deep learning in natural language processing*. Springer, 2018.
- [8] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- [9] Simon S Du, Xiyu Zhai, Barnabas Poczos, and Aarti Singh. Gradient descent provably optimizes over-parameterized neural networks. In *International Conference on Learning Representations*, 2018.
- [10] Jonathan Frankle and Michael Carbin. The lottery ticket hypothesis: Finding sparse, trainable neural networks. In *International Conference on Learning Representations*, 2018.
- [11] Trevor Gale, Erich Elsen, and Sara Hooker. The state of sparsity in deep neural networks. *arXiv preprint arXiv:1902.09574*, 2019.
- [12] Elizabeth Gardner and Bernard Derrida. Three unfinished works on the optimal storage capacity of networks. *Journal of Physics A: Mathematical and General*, 22(12):1983, 1989.
- [13] Sebastian Goldt, Madhu Advani, Andrew M Saxe, Florent Krzakala, and Lenka Zdeborová. Dynamics of stochastic gradient descent for two-layer neural networks in the teacher-student setup. In *Advances in Neural Information Processing Systems*, pages 6979–6989, 2019.
- [14] Scott Gray, Alec Radford, and Diederik P Kingma. Gpu kernels for block-sparse weights. *arXiv preprint arXiv:1711.09224*, 3, 2017.



- [15] Song Han, Huizi Mao, and William J Dally. Deep compression: Compressing deep neural networks with pruning, trained quantization and huffman coding. *arXiv preprint arXiv:1510.00149*, 2015.
- [16] Babak Hassibi and David G Stork. Second order derivatives for network pruning: Optimal brain surgeon. In *Advances in neural information processing systems*, pages 164–171, 1993.
- [17] Tianxing He, Yuchen Fan, Yanmin Qian, Tian Tan, and Kai Yu. Reshaping deep neural network for fast decoding by node-pruning. In *2014 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 245–249. IEEE, 2014.
- [18] Yihui He, Xiangyu Zhang, and Jian Sun. Channel pruning for accelerating very deep neural networks. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 1389–1397, 2017.
- [19] Arthur Jacot, Franck Gabriel, and Clément Hongler. Neural tangent kernel: convergence and generalization in neural networks. In *Proceedings of the 32nd International Conference on Neural Information Processing Systems*, pages 8580–8589, 2018.
- [20] Kumar Joag-Dev, Frank Proschan, et al. Negative association of random variables with applications. *The Annals of Statistics*, 11(1):286–295, 1983.
- [21] Nal Kalchbrenner, Erich Elsen, Karen Simonyan, Seb Noury, Norman Casagrande, Edward Lockhart, Florian Stimberg, Aaron van den Oord, Sander Dieleman, and Koray Kavukcuoglu. Efficient neural audio synthesis. *arXiv preprint arXiv:1802.08435*, 2018.
- [22] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009.
- [23] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. In *Advances in neural information processing systems*, pages 1097–1105, 2012.
- [24] Alex Kulesza, Ben Taskar, et al. Determinantal point processes for machine learning. *Foundations and Trends® in Machine Learning*, 5(2–3):123–286, 2012.
- [25] Yann LeCun, Corinna Cortes, and Chris Burges. Mnist handwritten digit database. 2010. URL <http://yann.lecun.com/exdb/mnist>, 3(1), 2010.
- [26] Yann LeCun, John S Denker, and Sara A Solla. Optimal brain damage. In *Advances in neural information processing systems*, pages 598–605, 1990.
- [27] Namhoon Lee, Thalaiyasingam Ajanthan, Stephen Gould, and Philip HS Torr. A signal propagation perspective for pruning neural networks at initialization. *arXiv preprint arXiv:1906.06307*, 2019.
- [28] Namhoon Lee, Thalaiyasingam Ajanthan, and Philip Torr. Snip: Single-shot network pruning based on connection sensitivity. In *International Conference on Learning Representations*, 2018.
- [29] Hao Li, Asim Kadav, Igor Durdanovic, Hanan Samet, and Hans Peter Graf. Pruning filters for efficient convnets. *arXiv preprint arXiv:1608.08710*, 2016.
- [30] Geert Litjens, Thijs Kooi, Babak Ehteshami Bejnordi, Arnaud Arindra Adiyoso Setio, Francesco Ciompi, Mohsen Ghahfoorian, Jeroen Awm Van Der Laak, Bram Van Ginneken, and Clara I Sánchez. A survey on deep learning in medical image analysis. *Medical image analysis*, 42:60–88, 2017.
- [31] Weibo Liu, Zidong Wang, Xiaohui Liu, Nianyin Zeng, Yurong Liu, and Fuad E Alsaadi. A survey of deep neural network architectures and their applications. *Neurocomputing*, 234:11–26, 2017.
- [32] Odile Macchi. The coincidence approach to stochastic point processes. *Advances in Applied Probability*, 7(1):83–122, 1975.
- [33] Eran Malach, Gilad Yehudai, Shai Shalev-Schwartz, and Ohad Shamir. Proving the lottery ticket hypothesis: Pruning is all you need. In *International Conference on Machine Learning*, pages 6682–6691. PMLR, 2020.
- [34] Zelda Mariet and Suvrit Sra. Diversity networks: Neural network compression using determinantal point processes. In *International Conference on Learning Representations*, 2016.
- [35] Laurent Orseau, Marcus Hutter, and Omar Rivasplata. Logarithmic pruning is all you need. In *Proceedings of the 34th International Conference on Neural Information Processing Systems*, 2020.
- [36] Vivek Ramanujan, Mitchell Wortsman, Aniruddha Kembhavi, Ali Farhadi, and Mohammad Rastegari. What’s hidden in a randomly weighted neural network? In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11893–11902, 2020.
- [37] David Saad and Sara A Solla. Exact solution for on-line learning in multilayer neural networks. *Physical Review Letters*, 74(21):4337, 1995.
- [38] David Saad and Sara A Solla. On-line learning in soft committee machines. *Physical Review E*, 52(4):4225, 1995.
- [39] David Saad and Sara A Solla. Learning with noise and regularizers in multilayer neural networks. In *Advances in Neural Information Processing Systems*, pages 260–266, 1997.

- [40] Tara N Sainath, Brian Kingsbury, Vikas Sindhwani, Ebru Arisoy, and Bhuvana Ramabhadran. Low-rank matrix factorization for deep neural network training with high-dimensional output targets. In *2013 IEEE international conference on acoustics, speech and signal processing*, pages 6655–6659. IEEE, 2013.
- [41] Andrew M Saxe, Yamini Bansal, Joel Dapello, Madhu Advani, Artemy Kolchinsky, Brendan D Tracey, and David D Cox. On the information bottleneck theory of deep learning. *Journal of Statistical Mechanics: Theory and Experiment*, 2019(12):124020, 2019.
- [42] Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D Manning, Andrew Y Ng, and Christopher Potts. Recursive deep models for semantic compositionality over a sentiment treebank. In *Proceedings of the 2013 conference on empirical methods in natural language processing*, pages 1631–1642, 2013.
- [43] Hidenori Tanaka, Daniel Kunin, Daniel LK Yamins, and Surya Ganguli. Pruning neural networks without any data by iteratively conserving synaptic flow. In *Proceedings of the 34th International Conference on Neural Information Processing Systems*, 2020.
- [44] Naftali Tishby and Noga Zaslavsky. Deep learning and the information bottleneck principle. In *2015 IEEE Information Theory Workshop (ITW)*, pages 1–5. IEEE, 2015.
- [45] Joost van Amersfoort, Milad Alizadeh, Sebastian Farquhar, Nicholas Lane, and Yarín Gal. Single shot structured pruning before training. *arXiv preprint arXiv:2007.00389*, 2020.
- [46] Vladimir N Vapnik. An overview of statistical learning theory. *IEEE transactions on neural networks*, 10(5):988–999, 1999.
- [47] Ruichi Yu, Ang Li, Chun-Fu Chen, Jui-Hsin Lai, Vlad I Morariu, Xintong Han, Mingfei Gao, Ching-Yung Lin, and Larry S Davis. Nisp: Pruning networks using neuron importance score propagation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 9194–9203, 2018.
- [48] Difan Zou, Yuan Cao, Dongruo Zhou, and Quanquan Gu. Gradient descent optimizes over-parameterized deep relu networks. *Machine Learning*, 109(3):467–492, 2020.

## A GE IN TWO LAYER NETWORK

For the theoretical analysis we consider the following assumptions from [13]

- (A1) If  $\mathbf{x} = (x_1, \dots, x_N)$  is an input then  $x_i \in \mathcal{N}(0, 1)$ . Also,  $N \rightarrow \infty$ .
- (A2) Both the teacher and the student networks have only one hidden layer.
- (A3)  $M, K$  denotes the number of hidden nodes for the teacher and student network respectively and  $K \geq M$  and  $K = Z \cdot M$  where  $Z \in \mathbb{Z}^+$ .
- (A4) The activation in the hidden layer is sigmoidal for both teacher and student network.
- (A5) The output  $\in \mathbb{R}$  (i.e., regression problem).
- (A6) The order parameters satisfy the ansatz as in (S58) - (S60) of [13].
- (A7) No noise is added to the labels generated by the teacher network, i.e.,  $\sigma = 0$ .

With the above assumptions, authors of [13] gave a closed form of the GE as follows:

$$\varepsilon_g = f_1(Q) + f_2(T) - f_3(R, Q, T) \quad (13)$$

where,

$$f_1(Q) = \frac{1}{\pi} \sum_{i,k} v_i v_k \arcsin \frac{Q_{ik}}{\sqrt{1+Q_{ii}}\sqrt{1+Q_{kk}}} \quad (14)$$

$$f_2(T) = \frac{1}{\pi} \sum_{n,m} v_n^* v_m^* \arcsin \frac{T_{nm}}{\sqrt{1+T_{nn}}\sqrt{1+T_{mm}}} \quad (15)$$

$$f_3(R, Q, T) = \frac{2}{\pi} \sum_{i,n} v_i v_n^* \arcsin \frac{R_{in}}{\sqrt{1+Q_{ii}}\sqrt{1+T_{nn}}} \quad (16)$$

where  $Q, R, T$  are the order parameters as defined in main text. We also have the assumption (4) about the relation between number of edges and nodes kept after pruning.

## B PROPERTIES OF DPP KERNEL

In main text we see that each node in the hidden layer of a student network carries certain amount of information about the training data and it is captured in a vector form. We create an information matrix by accumulating the information vectors of these hidden nodes. For simplicity of theoretical analysis, we have considered the kernel as the inner product of the information matrix. In the thermodynamic limit, the inner product is divided by the input dimension. Formally, if  $\mathbf{h}_i$  and  $\mathbf{h}_j$  are the information at  $i^{\text{th}}$  hidden node and  $j^{\text{th}}$  hidden node respectively, then

$$L_{ij} = \frac{1}{N} \frac{1}{n} \mathbf{h}_i^T \mathbf{h}_j$$

where  $n$  is the total number of training examples. It can be seen that the analysis for the kernel defined in main text is similar. Note that all analyses are for the student network trying learn from the teacher network. Refer to main text for details of notations.

**Lemma 1.** Assume (A1) - (A7). Then the expected kernel of DPP Node for the hidden layer is the order parameter  $Q$ .

*Proof of Lemma 1.* For the two-layer teacher-student setup, the hidden layer gets information  $(\mathbf{h}_1, \dots, \mathbf{h}_K)$  from the input layer, where  $\mathbf{h}_i = (h_{i1}, \dots, h_{in})$  and  $h_{ij} (= t_j^T \mathbf{w}_i)$  is the information at  $i^{\text{th}}$  hidden node on  $j^{\text{th}}$  input data ( $t_j$ ). Hence,

$$\mathbf{h}_i^T \mathbf{h}_j = \sum_{k=1}^n h_{ik} h_{jk} = \sum_{k=1}^n t_k^T \mathbf{w}_i \cdot t_k^T \mathbf{w}_j = \sum_{k=1}^n \mathbf{w}_i^T t_k \cdot t_k^T \mathbf{w}_j = \sum_{k=1}^n \mathbf{w}_i^T (t_k t_k^T) \mathbf{w}_j$$

But for the given input distribution (i.i.d. Gaussian),  $\mathbb{E}[t_k t_k^T] = \mathbf{I}_{N \times N}$ . Hence,  $\lim_{N \rightarrow \infty} \mathbb{E}[L_{ij}] = \lim_{N \rightarrow \infty} \mathbb{E}[\frac{1}{N} \frac{1}{n} \mathbf{h}_i^T \mathbf{h}_j] = \lim_{N \rightarrow \infty} \frac{1}{N} \mathbf{w}_i^T \mathbf{w}_j = Q_{ij}$ , and we have the lemma.  $\square$

From [13] we know that  $Q$  is a block diagonal matrix where each "block" (or "group" used interchangeably henceforth) refers to the set of student hidden nodes that represent (explain/learn) one particular teacher hidden node.

## C PROOF OF THE THEOREMS

*Proof of Theorem 1 and 2.* Let  $H_R = \{h_{i_1}, \dots, h_{i_{k_n}}\}$  be the set of selected nodes by DPP Node pruning method. Recall from [13] that every student hidden node specializes in learning a teacher node. Denote  $t(h)$  to be the teacher node learnt by  $h$ .  $S_m \subseteq H_R$  be the set of selected hidden nodes of the pruned network which learnt the  $m^{\text{th}}$  teacher node, i.e.,  $S_m = \{h \in H_R | t(h) = t_m\}$  ( $t_m$  is the  $m^{\text{th}}$  teacher node). Hence,  $prn = |\{\mathbb{1}(|S_m| > 0) | 1 \leq m \leq M\}|$  is the number of teacher nodes explained by the pruned network and W.L.O.G. we can assume that  $t_1, \dots, t_{prn}$  are those set of teacher nodes. Let  $l_1, \dots, l_{prn}$  be the number of student nodes in the pruned network which learn the corresponding teacher node. Note that,  $\sum_{i=1}^{prn} l_i = k_n$  and  $l_i \leq Z$  (where  $Z$  is the number of student nodes dedicated to learn a single teacher node in the unpruned network) for all  $i$ . Applying Lemma 2 directly we can see that the GE for the pruned network is

$$\frac{(v^*)^2}{6} \left[ \sum_{i=1}^{prn} \left(1 - \frac{l_i}{Z}\right)^2 \right] + \frac{(M - prn)(v^*)^2}{6} \quad (17)$$

The first part of (17) is the GE for the group whose corresponding teacher node is partially explained and the second part accounts for the GE due to unexplained teacher nodes (number of such teacher nodes are  $M - prn$ ). From Lemma 1 we know that the expected kernel matrix for DPP Node pruning is the order parameter  $Q$  and it becomes a block diagonal matrix after the training converges, where size of each block is  $Z$  (which is also the number of student nodes dedicated to learn a single teacher node in the unpruned network). Because of the block diagonal property of the DPP kernel matrix, at most 1 student hidden node will be chosen from each block, i.e.,  $l_i = 1 \forall i$ . Hence,  $prn = k_n$ . From Lemma 2 we can see that the GE of node pruned network only depends on the number of student node survived in each block after pruning, and, for DPP node pruning, it is always 1 (given  $k_n \leq M$ ). This is why there is no expectation in the GE term. So for DPP node pruning the GE is,

$$\epsilon_{k_n}^{DPPNode}(f) = (v^*)^2 \left[ \frac{k_n}{6} \left(1 - \frac{1}{Z}\right)^2 + \frac{M - k_n}{6} \right].$$

Each of the  $k_n$  student nodes in the pruned network learns a different teacher node. Consider one such teacher node and call it  $t_i$ . In the unpruned network, there are  $Z$  student hidden nodes which learn a single teacher node  $t_i$ , only one of which survives after DPP node pruning. The first part of the error is due to the removal of student nodes ( $Z - 1$  student nodes for each  $t_i$ ). However, these errors can be retrieved by reweighting the survived student node. On the contrary, there are  $M - k_n$  teacher nodes which don't have any representative (some student hidden node from the set of student nodes which specialized in this particular teacher node) in the pruned network. And the error (second part of the GE) due to those nodes can not be retrieved even after reweighting. Hence, the GE after reweighting becomes,

$$(M - k_n) \times \frac{(v^*)^2}{6}$$

Thus, we have the Theorem 1.

Next, we will prove Theorem 2. We will show, for any network pruned by Random Node, the GE is more than the expected GE of DPP Node pruning. Recall the randomly pruned network  $f$  discussed in the beginning of the proof. From Lemma 2 we can see that for node pruning the GE only depends on the number of nodes survived in each block. From (17) we have,

$$\begin{aligned} & \epsilon_{k_n}^{RandNode}(f) \\ &= \frac{(v^*)^2}{6} \left[ \sum_{i=1}^{prn} \left(1 - \frac{l_i}{Z}\right)^2 \right] + \frac{(M - prn)(v^*)^2}{6} \\ &= \frac{(M - k_n)(v^*)^2}{6} + \sum_{i=1}^{prn} \left[ (l_i - 1) \frac{(v^*)^2}{6} + \frac{(v^*)^2}{6} \left(1 - \frac{l_i}{Z}\right)^2 \right] \\ &\geq \frac{(M - k_n)(v^*)^2}{6} + \sum_{i=1}^{prn} l_i \frac{(v^*)^2}{6} \left(1 - \frac{1}{Z}\right)^2 \\ &= \frac{(M - k_n)(v^*)^2}{6} + l \frac{(v^*)^2}{6} \left(1 - \frac{1}{Z}\right)^2 \\ &= \epsilon_{k_n}^{DPPNode}(f) \end{aligned} \quad (18)$$

where (18) follows from the inequality below:

$$(l_i - 1) \frac{(v^*)^2}{6} + \frac{(v^*)^2}{6} \left(1 - \frac{l_i}{Z}\right)^2 = l_i \frac{(v^*)^2}{6} \left[1 + \frac{1}{Z^2} - \frac{2}{Z}\right] \geq l_i \frac{(v^*)^2}{6} \left(1 - \frac{1}{Z}\right)^2$$

which proves the first part of Theorem 2. The proof for the reweighted network is similar.

In case of importance node pruning, the nodes with lowest absolute value of outgoing edges are dropped. Following [13] the outgoing weights of all the hidden teacher nodes are equal (we call it  $v^*$ ). Also, from Lemma 3 we see that the sum of the weights of the outgoing edges of the student nodes which learn the same teacher node add up to the outgoing edge weight of the corresponding teacher hidden node. Moreover, we assume the ansatz  $v_i = v_j$  when  $i, j \in G_n$ , where  $G_n$  denotes the set of student nodes which learn the same teacher node  $t_n$ . Hence, we can see that all the outgoing edges are approximately similar. We also verify this fact experimentally. Therefore, this defines an approximately uniform distribution on the set of hidden nodes. Hence, this is almost same as random node pruning and so the result follows from Theorem 2.  $\square$

**Remark 5.** *The comparison between performance of importance node pruning and DIVNET depends on the fact that all the outgoing edges of the teacher hidden nodes are equal. However, when the outgoing weights are not equal the importance pruning first selects student hidden nodes from a group whose corresponding teacher node has the highest weight. Once all the student nodes are selected from that group then it selects the group whose corresponding teacher node has second highest outgoing edge weight and the process continues. Because of this approach, even without reweighting a complete information about the teacher node is preserved in the pruned network. However, in DPP node pruning one candidate from each group (representing a particular teacher node) is selected first. But if a member is selected from a group then the reweighting method can recover the complete lost information for the corresponding group. Hence, DIVNET is able to preserve information about more number of teacher hidden nodes than importance pruning which results in better performance.*

*Proof of Theorem 3.* In this theorem, we will give the GE of the expected network pruned by the Random Edge method. Pruning is performed on the edges between input layer and the hidden layer. Hence, the order parameter changes. From Lemma 4, we have the order parameters of the expected network (call these  $Q', R', T'$ ). However, the weights of the second layer remain unchanged. Putting these values in (14), (15) and (16) we have,

$$\begin{aligned} f_1(Q') &= \frac{1}{\pi} \sum_{i,k} v_i v_k \arcsin \frac{Q'_{ik}}{\sqrt{1+Q'_{ii}} \sqrt{1+Q'_{kk}}} \\ &= \frac{M(v^*)^2}{\pi} \arcsin \frac{c^2}{1+c} + \frac{M(v^*)^2}{Z\pi} \left[ \arcsin \frac{c}{1+c} - \arcsin \frac{c^2}{1+c} \right] \end{aligned} \quad (19)$$

and,

$$\begin{aligned} f_3(R', Q', T') &= \frac{2}{\pi} \sum_{i,n} v_i v_n^* \arcsin \frac{R'_{in}}{\sqrt{1+Q'_{ii}} \sqrt{1+T'_{nn}}} \\ &= \frac{2M(v^*)^2}{\pi} \arcsin \frac{c}{\sqrt{2(1+c)}} \end{aligned} \quad (20)$$

Therefore, the GE of the expected network after Random Edge pruning is,

$$\frac{M(v^*)^2}{\pi} \left[ \arcsin \frac{c^2}{1+c} + \frac{\pi}{6} - 2 \arcsin \frac{c}{\sqrt{2(1+c)}} \right] + \frac{M(v^*)^2}{Z\pi} \left[ \arcsin \frac{c}{1+c} - \arcsin \frac{c^2}{1+c} \right]$$

This proves the first part of the theorem.  $\square$

*Proof of Theorem 4.* Theorem 1 and 3 provide the closed form of the GE after DPP node pruning and random node pruning respectively. Using this closed form we plot  $\epsilon_{k_n}^{DPPNode}(f) - \epsilon_c^{RandEdge}(f)$  in Figure 4 A. Here  $k_n$  and  $c$  satisfy (4), i.e., parameter count is same after two kinds of pruning. We can see for  $Z \geq 4$  this value is  $\geq 0$  given  $0 \leq c \leq 1.0/Z$ , which proves the theorem.  $\square$

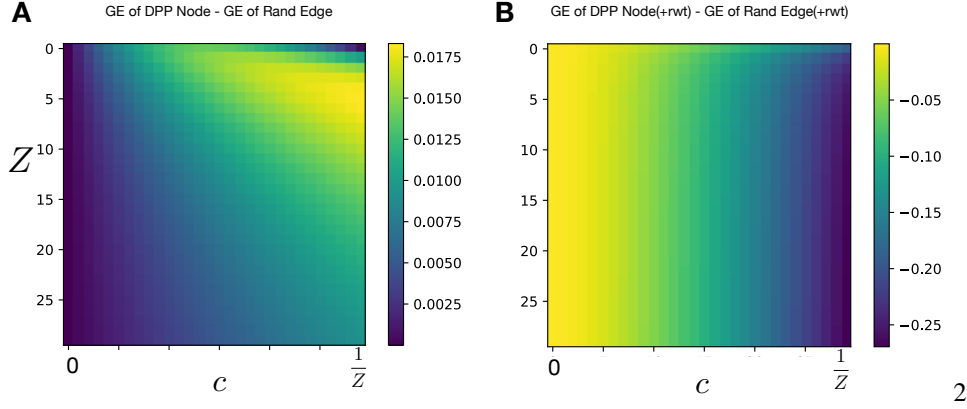


Figure 4: **(A)** Difference between the GE of DPP node pruning and Random edge pruning for  $4 \leq Z \leq 30$ . The matrix consist of only nonzero entries which proves that random edge pruning performs better than DPP node pruning when parameter count is same. **(B)** Difference between the GE of DPP node pruning with reweighting and Random edge pruning with reweighting for  $4 \leq Z \leq 30$ . The matrix consist of only negative entries which proves that random edge pruning can never perform better than DPP node pruning when reweighting is applied in the second layer.

**Remark 6.** Our results hold for  $Z \geq 4$ , where  $Z$  is the number of student nodes which learn the same teacher node. This is because in DPP node pruning at most 1 student node survives per group. As a result for larger  $Z$  the lost information per group is higher (in the scale of  $(1 - \frac{1}{Z})^2$ ).

Next we state the impossibility result as discussed in main text. We will show that, no reweighting scheme in the second layer for random edge pruning which is based on scaling can beat DPP node pruning after reweighting. Formally we have the following:

**Theorem 5.** Assume (A1) – (A7). Let  $k_n$  and  $c$  satisfy (4), and  $0 \leq c \leq \frac{1}{Z}$  and  $Z \geq 4$ . Assume the reweighting scheme for random edge in second layer such that,  $\hat{v}_i = Av_i$ . Then  $\forall A \in \mathbb{R}$  we have,

$$\hat{\epsilon}_{k_n}^{DPPNode}(f) \leq \hat{\epsilon}_c^{RandEdge}(\mathbb{E}[f]) \quad (21)$$

*Proof of Theorem 5.* From Theorem 1 we know that the GE after reweighting the DPP node pruned network is

$$\frac{(v^*)^2}{6} (M - k_n) = \frac{M(v^*)^2}{6} (1 - Zc) \quad (22)$$

where  $c$  satisfies (4). Now for the given reweighting scheme in the hypothesis the GE for random edge pruning will be,

$$\frac{M(v^*)^2}{\pi} \left[ A^2 \left( \frac{1}{Z} \arcsin \frac{c}{1+c} + \left(1 - \frac{1}{Z}\right) \arcsin \frac{c^2}{1+c} \right) + \frac{\pi}{6} - 2A \arcsin \frac{c}{\sqrt{2(1+c)}} \right] \quad (23)$$

(23) can be viewed as a quadratic equation of  $A$  whose minimum correspond to the best reweighting scheme in the scaling family. In Figure 4 B we compare this minimum with (22). Formally we plotted  $\hat{\epsilon}_{k_n}^{DPPNode}(f) - \hat{\epsilon}_c^{RandEdge}(\mathbb{E}[f])$ . It can be seen that this value is  $-ve$  for all  $0 \leq c \leq \frac{1}{Z}$ , which implies GE of reweighted DPP node pruned network is always lower than reweighted random edge pruned network.  $\square$

## D PROOF OF LEMMAS

**Lemma 2.** Assume (A1)-(A7). Let  $t_1, \dots, t_M$  denote the teacher hidden nodes and  $l_1, \dots, l_M$  denote the number of student hidden nodes in a node pruned network which learnt the corresponding teacher node. If  $\sum_{m=1}^M l_m \leq M$ , then the GE of this node pruned network is,

$$\frac{(v^*)^2}{6} \left[ \sum_{m=1}^M \left(1 - \frac{l_m}{Z}\right)^2 \right].$$

*Proof.* Let  $G_1, \dots, G_M$  be the subsets of student nodes such that all student nodes in  $G_m$  learn the  $m^{\text{th}}$  teacher node. From the assumption we have,  $|G_m| = Z$  for all  $m$ . After pruning, a subset  $P_m \subseteq G_m$  is chosen, and  $|P_m| = l_m$ . Denote the order parameters of the pruned network as  $Q', R', T'$ . For node pruning we can see that

$$Q'_{ik} = \begin{cases} Q_{ik} & \text{if } \exists m \text{ s.t. } h_i \in P_m \text{ and } h_k \in P_m \\ 0 & \text{otherwise} \end{cases}$$

Also, for the unpruned network we have

$$Q_{ik} = \begin{cases} 1 & \text{if } \exists m \text{ s.t. } h_i \in G_m \text{ and } h_k \in G_m \\ 0 & \text{otherwise} \end{cases}$$

Now from (13) we can break down the GE into three parts. From (14), (15) and (16) we have,.

$$\begin{aligned} f_1(Q') &= \frac{1}{\pi} \sum_{i,k} v_i v_k \arcsin \frac{Q'_{ik}}{\sqrt{1+Q'_{ii}} \sqrt{1+Q'_{kk}}}, \\ &= \frac{1}{\pi} \sum_{n=1}^M \sum_{i,k \in P_n} v_i v_k \arcsin \frac{1}{2}, \end{aligned} \quad (24)$$

$$\begin{aligned} &= \frac{1}{\pi} \sum_{n=1}^M \sum_{i,k \in P_n} v_i v_k \frac{\pi}{6}, \\ &= \frac{1}{6} \sum_{n=1}^M \left( \sum_{i \in P_n} v_i \right)^2, \\ &= \frac{(v^*)^2}{6} \sum_{n=1}^M \left( \frac{l_i}{Z} \right)^2 \end{aligned} \quad (25)$$

(24) follows from the fact that  $h_i$  and  $h_k$  belong to the same group  $G_n$ . So we have,

$$\frac{Q'_{ik}}{\sqrt{1+Q'_{ii}} \sqrt{1+Q'_{kk}}} = \frac{1}{\sqrt{2}\sqrt{2}} = \frac{1}{2}$$

We can also see that (25) follows from Lemma 3 and the ansatz  $v_i = v_j$  when  $i, j \in G_n$ . The order parameters  $T_{nm}$  doesn't change after pruning, and so we have,

$$\begin{aligned} f_2(T') &= \frac{1}{\pi} \sum_{n,m} v_n^* v_m^* \arcsin \frac{T_{nm}}{\sqrt{1+T_{nn}} \sqrt{1+T_{mm}}}, \\ &= \frac{1}{6} \sum_{n=1}^M (v^*)^2 \end{aligned} \quad (26)$$

And similarly,

$$\begin{aligned} f_3(R', Q', T') &= \frac{2}{\pi} \sum_{i,n} v_i v_n^* \arcsin \frac{R'_{in}}{\sqrt{1+Q'_{ii}} \sqrt{1+T'_{nn}}}, \\ &= \frac{2}{\pi} \sum_{n=1}^M v_n^* \sum_{i \in P_n} v_i \arcsin \frac{1}{2}, \\ &= \frac{2}{6} \sum_{n=1}^M v_n^* \sum_{i \in P_n} v_i. \end{aligned} \quad (27)$$

Then from (25),(26) and (27) the GE of node pruning is,

$$\frac{(v^*)^2}{6} \left[ \sum_{m=1}^M \left( 1 - \frac{l_m}{Z} \right)^2 \right]. \quad (28)$$

Hence we have the lemma.  $\square$

Intuitively, this lemma states that for teacher hidden node  $t_n$  if  $l_n$  student hidden nodes survive after node pruning, then the fraction of information lost due to the deletion of nodes is  $1 - \frac{l_n}{Z}$ , where  $Z$  is the number of student nodes learn a particular teacher node in the unpruned network.

**Lemma 3.** *Let  $v^*$  denotes the weight of the second layer of the teacher network and  $\{v_1, \dots, v_K\}$  be the weights of the student network after convergence. Then in the noiseless case for all  $n$  we have,*

$$v^* = \sum_{i \in G_n} v_i$$

*Proof of Lemma 3.* From (S36) of [13] we have,

$$\begin{aligned} \frac{dv_i}{dt} &= \eta_v \left[ \sum_{n=1}^M v_n^* I_2(i, n) - \sum_{j=1}^K v_j I_2(i, j) \right] \\ &= \eta_v \arcsin \frac{1}{2} \left[ v^* - \sum_{j \in G_n} v_j \right] \end{aligned}$$

Hence, a fixed point (in terms of  $v_i$ 's) of the ODE is,

$$\{(v_1, \dots, v_K) \mid \sum_{i \in G_n} v_i = v^*, \forall 1 \leq n \leq M\}$$

□

Intuitively, this lemma states that the sum of the outgoing edges of the student hidden nodes which learn a particular teacher hidden node is approximately equal to the weight of the outgoing edge of that teacher hidden node.

**Lemma 4.** *Let  $Q, R, T$  are the order parameters of the unpruned network, and  $Q', R', T'$  are the respective order parameters after applying the Random Edge pruning where  $c$  fraction of the edges are kept. Then we have the following:*

•

$$\mathbb{E}[Q'_{ik}] = \begin{cases} cQ_{ik} & \text{if } i = k \\ c^2Q_{ik} & \text{otherwise} \end{cases}$$

- $\mathbb{E}[R'_{st}] = cR_{st}$
- $T'_{mn} = T_{mn}$

*Proof.* In case of Random Edge pruning each edge is kept with probability  $c$ . Then we have,

$$\mathbb{E}[Q'_{st}] = \frac{1}{N} \sum_{i=1}^N c \cdot w_{is} \times c \cdot w_{it} = c^2 \frac{1}{N} \sum_{i=1}^N w_{is} w_{it} = c^2 Q_{st}$$

and

$$\mathbb{E}[Q'_{ss}] = \frac{1}{N} \sum_{i=1}^N c^2 \cdot w_{ss}^2 = c^2 Q_{ss}.$$

Similarly,

$$\mathbb{E}[R'_{st}] = \frac{1}{N} \sum_{i=1}^N c \cdot w_{is} w_{it}^* = cR_{st}$$

The teacher node is not affected by the pruning. So  $T$  is not modified by the pruning process. This proves the lemma. □

Intuitively, this lemma states that the order parameters of the pruned network using random edge pruning is a scaled version of the order parameters of the unpruned networks. However, the scaling of diagonal elements are different from that of off-diagonal elements (for more see Figure 6 A).



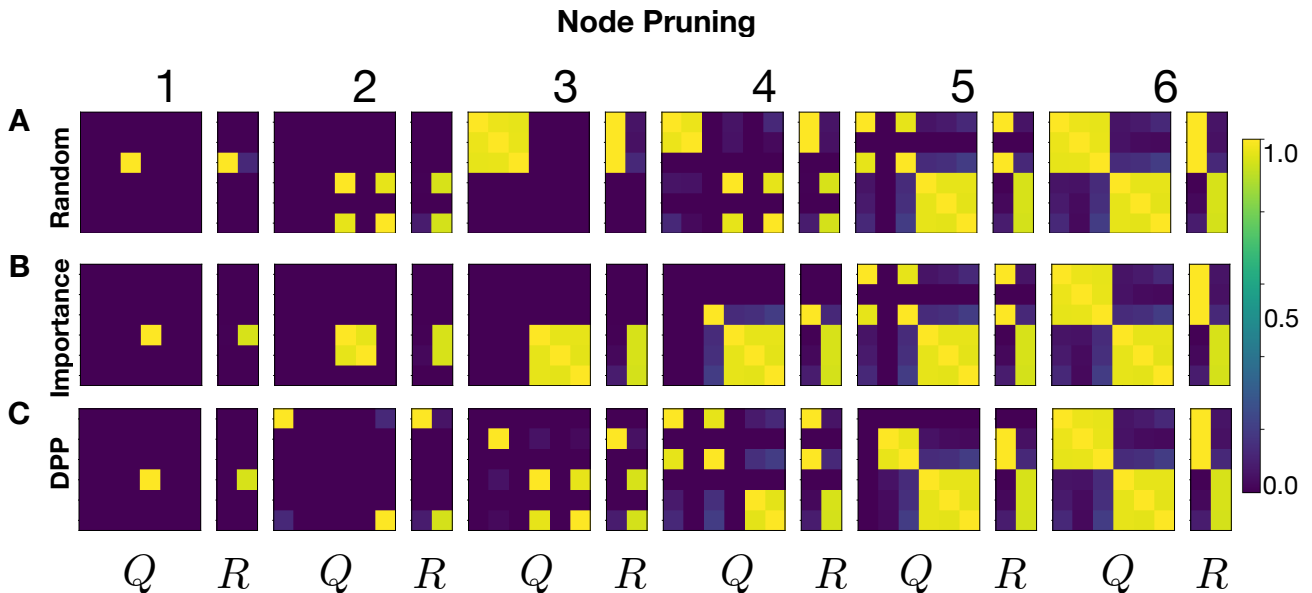


Figure 5: Order parameters after different node pruning methods in the teacher student setup. For this example, number of student hidden nodes  $M = 2$  and number of teacher hidden nodes  $K = 6$ . From [34] we know that the first 3 student nodes (call 1<sup>st</sup> group) learn one teacher node, and the next 3 (call 2<sup>nd</sup> group) learn the second teacher node. Recall that  $k_n$  is the number of student hidden nodes survived after node pruning. In this figure each row represents a particular node pruning method and each column (Q, R) shows results for different choices of  $k_n$  (left to right goes from most pruned to unpruned network). (A) In case of random node pruning when  $k_n = 2$ , two student node survives from the 2<sup>nd</sup> group after pruning. As a result, information about the 1<sup>st</sup> teacher node is completely lost in the pruned network. (B) Importance pruning keeps a student hidden node depending on its outgoing edge weights. The outgoing edge weights of each group is almost equally distributed among themselves, and they sum up to the second layer weight of corresponding teacher node (see Lemma 3). As all the group size is equal (3 for this example), importance node pruning first selects node from the group whose corresponding second layer teacher weight is highest. In our example, it is the second group and hence for  $k_n = 1, 2, 3$ , it selects node from the second group. Once a group is exhausted, it then selects from another group according to the aforementioned policy and so on. (C) For DPP node pruning when  $k_n = 2$ , two student hidden nodes are chosen from different groups which preserve information about both the teacher nodes. It can also be shown that, in case of node pruning, if at least one representative from a group survives after pruning, then the reweighting can recover the complete information about that block. Hence, in teacher student framework DPP node pruning performs the best among the node pruning methods especially after reweighting.

## E SIMULATION DETAILS

In total, 10 rounds of simulations are run for each of the 5 pruning methods, and we report the average and standard deviations (as error bars). The standard deviations are negligible (in the magnitude of  $10^{-3}$ ). A *round* is the entire process of generating a new teacher network with datasets, training the student from scratch, performing pruning and finally testing with the pruned network. For DPP and random methods, we sampled 100 masks per round and reported the average performance in each round. Given  $M = 2$  and  $K = 6$ , we tried pruning with  $[1, 2, 3, 4, 5]$  nodes (and the equivalent number of edges) left in the student, respectively. We keep the total number of weights same to compare different pruning methods. The node-to-edge ratio, given  $N = 500$ ,  $K = 6$ , and  $M = 2$ , is  $[1 : 83, 2 : 166, 3 : 250, 4 : 333, 5 : 417]$ . This is calculated, for the teacher-student setup (single output node) specifically, as  $k_e = \frac{k_n(1+N)-K}{K}$ . We grid-searched  $\eta$  in the range of  $[0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7]$  and found 0.50 to be the optimal. We used  $\beta = 0.3$  for all DPP node kernel calculations in all simulations.

## F HYPERPARAMETERS FOR REAL DATASETS

Besides the hyperparameters and setup we proposed in Section 5.1 on the synthetic dataset, we report the hyperparameters used for the results on the MNIST and CIFAR10. As stated in Section 5.2, we used that exact same experiment setup (network architectures, error thresholds, etc.) as in [34] for fair and consistent comparisons. We used SGD optimizers, a

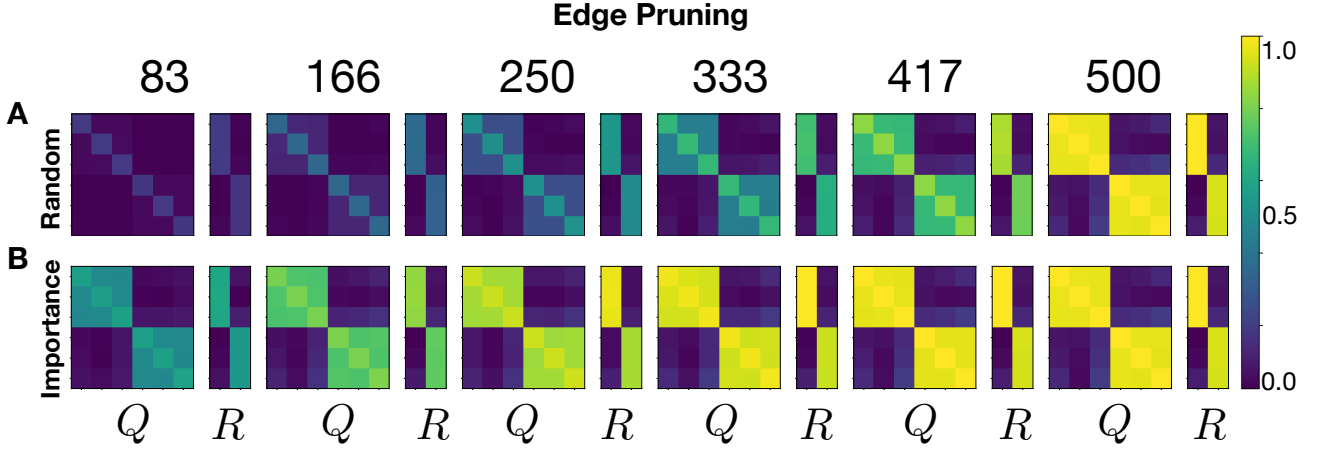


Figure 6: Order parameters after different edge pruning methods in the teacher student setup. For this example, number of student hidden nodes  $M = 2$  and number of teacher hidden nodes  $K = 6$ . From [34] we know that the first 3 student nodes (call  $1^{st}$  group) learn one teacher node, and the next 3 (call  $2^{nd}$  group) learn the second teacher node. Recall that  $k_e$  is the number of incoming edges for each student hidden nodes survived after edge pruning. In this figure each row represents a particular edge pruning method and each column (Q, R) shows results for different choices of  $k_e$  (left to right goes from most pruned to unpruned network). **(A)** In case of random edge pruning, the expected order parameters have the form described in Lemma 4. **(B)** Order parameters for importance edge pruning. For importance edge pruning, the edges with lowest absolute values are removed. As the input dimension goes to infinity, the order parameters of the pruned network are close to that of the unpruned network ( $k_e = 500$ ). In particular, for any fix  $k_e$ , let  $Q_{k_e}^{imp}$  be the order parameter of the pruned network when importance pruning is used.  $Q_{k_e}^{rand}$  is defined similarly. Our simulations show that,  $\|Q_{unpruned} - Q_{k_e}^{imp}\| \leq \|Q_{unpruned} - Q_{k_e}^{rand}\|$ . This is why the blocks in the  $Q$  matrix are the brightest in case of importance pruning. Hence, importance edge pruning performs the best without reweighting.

learning rate of 0.001, and a momentum of 0.9 for training on both datasets. For MNIST, the training batch size was 1000. For CIFAR10, the training batch size was 128. All pruning methods were performed 10 times, and we report the means and standard deviations in Figure 3 (with reweighting).

The node-to-edge ratio for pruning, which keeps the number of parameters in the pruned network the same, is [397 : 614, 472 : 921, 548 : 1228, 623 : 1536, 699 : 1843, 774 : 2150, 849 : 2457, 925 : 2764] for CIFAR10 and [256 : 156, 287 : 235, 317 : 313, 348 : 392, 378 : 470, 409 : 548, 439 : 627, 470 : 705] for MNIST, given the network architecture in Table 1 of [34]. These ratios correspond to 20% to 90% of the edges left for each node, as shown on the x-axis of Figure 3. These node-to-edge ratios are calculated based on the conversion equation in Section 5.2. We used  $\beta = 10/|T|$  where  $|T|$  is the size of the training dataset for all DPP node and edge kernel calculations on real data, following the choice of [34].

## G TABLES AND FIGURES

Table 3 shows the experimental results on the synthetic data with the setup discussed in main text. For all the node-to-edge ratios in (4), given  $K = 6$  and  $M = 2$ , we calculated the mean square GEs for both the noiseless and noisy case ( $\sigma = 0.25$ ). We sampled 100 masks per simulation, and there are in total 10 rounds of simulations. As mentioned earlier, DPP methods are stable, and the standard deviations are in the magnitude of  $10^{-3}$  for all ratios.

Table 3: The mean square GE on synthetic data for all pruning methods. The left-most row indicates the percentage of parameters left in the network. For specific node-to-edge ratio, see 4. The upper table shows the noiseless case, and the lower shows the noisy case ( $\sigma = 0.25$ ). We also observed the implicit regularization effects of pruning proposed by [34]

% OF PARAMETERS	DPP NODE	RAND. EDGE	RAND. NODE	IMP. EDGE	IMP. NODE
17.0%	3.737± 0.009	3.451± 0.011	3.978 ± 0.016	1.911	3.760
33.0%	2.310± 0.012	2.300± 0.015	2.800 ± 0.035	0.814	2.719
50.0%	1.438± 0.015	1.402± 0.006	1.748 ± 0.036	0.311	1.540
67.0%	0.740± 0.017	0.730± 0.006	1.046 ± 0.018	0.110	0.721
83.0%	0.258± 0.008	0.204± 0.005	0.540 ± 0.010	0.040	0.360
ORIGINAL TEST LOSS: 0.051 (NOISELESS)					
17.0%	4.000± 0.005	3.769± 0.012	4.188 ± 0.001	1.963	4.167
33.0%	2.622± 0.015	2.558± 0.011	3.041 ± 0.024	0.905	2.910
50.0%	1.633± 0.002	1.675± 0.010	2.023 ± 0.035	0.450	2.031
67.0%	0.890± 0.018	1.007± 0.007	1.269 ± 0.022	0.271	1.144
83.0%	0.394± 0.001	0.490± 0.003	0.643 ± 0.002	0.253	0.659
ORIGINAL TEST LOSS: 0.241 ( $\sigma = 0.25$ )					