

Data-Efficient Structured Pruning via Submodular Optimization

Marwa El Halabi^{*1} Suraj Srinivas² Simon Lacoste-Julien^{1,3}

¹Samsung, SAIT AI Lab, Montreal

²Idiap Research Institute & EPFL

³Mila, Université de Montreal, Canada CIFAR AI Chair

Abstract

Structured pruning is an effective approach for compressing large pre-trained neural networks without significantly affecting their performance, which involves removing redundant regular regions of weights. However, current structured pruning methods are highly empirical in nature, do not provide any theoretical guarantees, and often require fine-tuning, which makes them inapplicable in the limited-data regime. We propose a principled data-efficient structured pruning method based on submodular optimization. In particular, for a given layer, we select neurons/channels to prune and corresponding new weights for the next layer, that minimize the change in the next layer’s input induced by pruning. We show that this selection problem is a weakly submodular maximization problem, thus it can be provably approximated using an efficient greedy algorithm. Our method is one of the few in the literature that uses only a limited-number of training data and no labels. Our experimental results demonstrate that our method outperforms popular baseline methods in various one-shot pruning settings.

1 Introduction

As modern neural networks (NN) grow increasingly large, with some models reaching billions of parameters (McGuffie and Newhouse, 2020), they require an increasingly large amount of memory, power, hardware, and inference time, which makes it necessary to compress them. This is especially important for models deployed on resource-constrained devices like mobile phones and smart speakers, and for latency-critical applications such as self-driving cars. Fortunately, the literature on network pruning has observed that parameters of deep NNs have a lot of redundancy, and thus they can be compressed without significantly affecting performance.

Several orthogonal approaches exist to compress NNs with the objective of reducing parameter count, or decreasing inference time. One family of methods approximate model weights using quantization or hashing (Gong et al., 2014; Courbariaux et al., 2015), while another set of methods use low-rank approximation or tensor factorization (Denil et al., 2013; Sainath et al., 2013; Jaderberg et al., 2014; Denton et al., 2014; Lebedev et al., 2015; Kim et al., 2015; Zhang et al., 2016; Su et al., 2018). In another class of methods called knowledge distillation, a small network is trained to mimic a much larger network (Bucila et al., 2006; Hinton et al., 2015). Other methods employ sparsity and group-sparsity regularisation during training, to induce sparse weights (Collins and Kohli, 2014; Voita et al., 2019). In this work, we follow the network pruning approach, where the redundant units (weights, neurons or filters/channels) of a pre-trained NN are removed; see (Kuzmin et al.,

^{*}part of this work was done while a postdoc at MIT, CSAIL.

2019; Blalock et al., 2020; Hoefler et al., 2021) for recent surveys. It is also possible to combine two or more of these types of compression methods to compound their benefits, see e.g., (Kuzmin et al., 2019, Section 4.3.4).

Pruning methods date back to at least the 1990s with the works of (LeCun et al., 1990; Hassibi et al., 1993). Existing pruning methods fall into two main categories: unstructured pruning methods which prune individual weights leading to irregular sparsity patterns, and structured pruning methods which prune regular regions of weights, such as neurons, channels, or attention heads. Structured pruning methods are generally preferable as the resulting pruned models can work with off-the-shelf hardware or kernels, as opposed to models pruned with unstructured pruning which require specialized ones.

Related work The goal in structured pruning is to select groups of parameters to prune that will minimize the decrease in performance. As minimizing the true decrease in accuracy is intractable, one has to replace it with a proxy more amenable to efficient optimization. Various proxy choices has been proposed in the literature. Some works prune neurons/channels individually based on their input or output weights norm (He et al., 2014; Li et al., 2017), or activations entropy (He et al., 2014), or product of activations and gradients of the loss w.r.t the activations (Molchanov et al., 2017), or the similarity of their weights to other neurons in the same layer (Srinivas and Babu, 2015). Such methods are efficient and easy to implement, but fail to capture higher-order interactions between the pruned parameters.

Closer to our approach are methods that aim to prune neurons/channels that minimize the change induced by pruning in the output of the layer being pruned, or its input to the next layer. Although these criteria are easier to optimize than the accuracy drop, they still yield an intractable combinatorial problem. Zhuang et al. (2018) propose to prune channels that minimize the change in the output, before the activation function, of the layer being pruned, and an additional discrimination-aware loss. They use a heuristic greedy algorithm to solve the resulting problem. The pruning method in (Ye et al., 2020) is also based on minimizing the change in the output of the layer being pruned, but after the activation function. They employ a greedy forward selection algorithm that provably solves the ℓ_1 -relaxation of the problem. Luo et al. (2017) and He et al. (2017) both propose to prune neurons/channels that minimize the change in the input to the next layer. Luo et al. (2017) solve the resulting problem using a heuristic backward greedy algorithm, while He et al. (2017) solve the ℓ_1 -relaxation of the problem using alternating minimization. Most of these methods simultaneously optimize the weights of the layer being pruned (Zhuang et al., 2018) or the next one (He et al., 2017), or do so in a separate step at the end (Luo et al., 2017). Mariet and Sra (2015) depart from the usual strategy of pruning parameters whose removal influences the network the least. They instead select a subset of diverse neurons to keep in each layer, then optimize the weights of the next layer to minimize the change in the input to the next layer.

The majority of existing structured pruning methods are heuristics that do not offer any theoretical guarantees. One exception is the work of (Ye et al., 2020), but their guarantee is w.r.t to the ℓ_1 -relaxation of the problem. Moreover, most pruning methods rely on fine-tuning at least for a few epochs to boost the performance of the pruned NN. Fine-tuning is not possible in the limited-data regime, where only few training data is available, or data labels are unavailable. Results in (Mariet and Sra, 2015) suggest that optimizing the weights to fuse information from the pruned neurons into the remaining ones can provide a similar boost to performance as fine-tuning, without the need for data labels.

Algorithm 1 GREEDY

```
1: Input: Ground set  $V$ , set function  $F : 2^V \rightarrow \mathbb{R}_+$ , budget  $k \in \mathbb{N}_+$ 
2:  $S \leftarrow \emptyset$ 
3: while  $|S| < k$  do
4:    $i^* \leftarrow \arg \max_{i \in V \setminus S} F(i \mid S)$ 
5:    $S \leftarrow S \cup \{i^*\}$ 
6: end while
7: Output:  $S$ 
```

Our contributions In this paper, we propose a structured pruning method, which simultaneously selects neurons to prune and new weights for the next layer, that minimize the change in the next layer’s input induced by pruning. We refer to the replacement of the next layer’s weights with the new optimal weights as “reweighting”. The resulting subset selection problem is intractable, but we show that it can be formulated as a weakly submodular maximization problem (see definition in Section 2). This formulation allows us to use the standard greedy selection algorithm to obtain $(1 - e^{-\gamma})$ -approximation to the optimal solution, where γ is non-zero if we use sufficient training data.

Our method uses only limited training data and no labels. Similar to (Mariet and Sra, 2015), we observe that reweighting provide a significant boost in performance not only to our method, but also to other baselines we consider. However unlike (Mariet and Sra, 2015), we only use a small fraction of the training data, around $\sim 1\%$ in our experiments. We further adapt our method to prune (1) any regular regions of weights; we focus in particular on pruning channels in convolution neural networks (CNNs), and (2) multiple layers in the network, by pruning each layer independently or sequentially. Our experimental results demonstrate that our method often outperforms popular baselines, and matches their performance even after fine-tuning.

2 Preliminaries

We begin by introducing our notation and some relevant background from submodular optimization.

Notation: Given a ground set $V = \{1, 2, \dots, d\}$ and a set function $F : 2^V \rightarrow \mathbb{R}_+$, we denote the *marginal gain* of adding a set $I \subseteq V$ to another set $S \subseteq V$ by $F(I \mid S) = F(S \cup I) - F(S)$, which quantifies the change in value when adding I to S . The cardinality of a set S is written as $|S|$. Given a vector $x \in \mathbb{R}^d$, we denote its support set by $\text{supp}(x) = \{i \in V \mid x_i \neq 0\}$, and its ℓ_2 -norm by $\|x\|_2$. Given a matrix $X \in \mathbb{R}^{d' \times d}$, we denote its i -th column by X_i , and its Frobenius norm by $\|X\|_F$. Given a set $S \subseteq V$, X_S is the matrix with columns X_i for all $i \in S$, and 0 otherwise, and $\mathbf{1}_S$ is the indicator vector of S , with $[\mathbf{1}_S]_i = 1$ for all $i \in S$, and 0 otherwise.

Weakly submodular maximization: A set function F is *submodular* if it has diminishing marginal gains: $F(i \mid S) \geq F(i \mid T)$ for all $S \subseteq T$, $i \in V \setminus T$. We say that F is *normalized* if $F(\emptyset) = 0$, and non-decreasing if $F(S) \leq F(T)$ for all $S \subseteq T$.

Given a non-decreasing submodular function F , selecting a set $S \subseteq V$ with cardinality $|S| \leq k$ that maximize $F(S)$ can be done efficiently using the GREEDY algorithm (Alg. 1). The returned solution is guaranteed to satisfy $F(\hat{S}) \geq (1 - 1/e) \max_{|S| \leq k} F(S)$ (Nemhauser et al., 1978). In general though

maximizing a non-submodular function over a cardinality constraint is NP-Hard (Natarajan, 1995). However, Das and Kempe (2011) introduced a notion of *weak submodularity* which is sufficient to obtain provable bounds of the GREEDY algorithm’s performance.

Definition 2.1. Given a set function $F : 2^V \rightarrow \mathbb{R}$, $U \subseteq V, k \in \mathbb{N}_+$, we say that F is $\gamma_{U,k}$ -weakly submodular, with $\gamma_{U,k} > 0$ if

$$\gamma_{U,k} F(S|L) \leq \sum_{i \in S} F(i|L),$$

for every two disjoint sets $L, S \subseteq V$, such that $L \subseteq U, |S| \leq k$.

The parameter $\gamma_{U,k}$ is called the *submodularity ratio* of F . It characterizes how close a set function is to being submodular. If F is non-decreasing then $\gamma_{U,k} \in [0, 1]$, and F is submodular if and only if $\gamma_{U,k} = 1$ for all $U \subseteq V, k \in \mathbb{N}_+$. Given a non-decreasing $\gamma_{\hat{S},k}$ -weakly submodular function F , the Greedy algorithm is guaranteed to return a solution $F(\hat{S}) \geq (1 - e^{-\gamma_{\hat{S},k}}) \max_{|S| \leq k} F(S)$ (Das and Kempe, 2011). Hence, the closer F is to being submodular, the better is the approximation guarantee.

3 Reweighted input change pruning

In this section, we introduce our approach for pruning neurons in a single layer.

Given a large pre-trained NN, n training data samples, and a layer ℓ with n_ℓ neurons, our goal is to select a small number k out of the n_ℓ neurons to keep, and prune the rest, in a way that influences the network the least. One way to achieve this is by minimizing the change in input to the next layer, $\ell + 1$, induced by pruning. However, simply throwing away the activations from the dropped neurons is wasteful. Instead, we optimize the weights of the next layer to reconstruct the inputs from the remaining neurons.

Formally, let $A^\ell \in \mathbb{R}^{n \times n_\ell}$ be the activation matrix of layer ℓ with columns $a_1^\ell, \dots, a_{n_\ell}^\ell$, where $a_i^\ell \in \mathbb{R}^n$ is the vector of activations of the i th neuron in layer ℓ for each training input, and let $W^{\ell+1} \in \mathbb{R}^{n_\ell \times n_{\ell+1}}$ be the weight matrix of layer $\ell + 1$ with columns $w_1^{\ell+1}, \dots, w_{n_{\ell+1}}^{\ell+1}$, where $w_i^{\ell+1} \in \mathbb{R}^{n_\ell}$ is the vector of weights connecting the i th neuron in layer $\ell + 1$ to the neurons in layer ℓ .

When a neuron is pruned in layer ℓ , the corresponding column of weights in W^ℓ and row in $W^{\ell+1}$ are removed. Pruning $n_\ell - k$ neurons in layer ℓ reduces the number of parameters and computation cost by $(n_\ell - k)/n_\ell$ for both layer ℓ and $\ell + 1$.

Let $V_\ell = \{1, \dots, n_\ell\}$. Given a set $S \subseteq V_\ell$, we denote by A_S^ℓ the matrix with columns a_i^ℓ for all $i \in S$, and 0 otherwise. That is, A_S^ℓ is the activation matrix of layer ℓ after pruning. We choose a set of neurons $S \subseteq V_\ell$ to keep and new weights $\tilde{W}^{\ell+1} \in \mathbb{R}^{n_\ell \times n_{\ell+1}}$ that minimize:

$$\min_{|S| \leq k, \tilde{W}^{\ell+1} \in \mathbb{R}^{n_\ell \times n_{\ell+1}}} \|A^\ell W^{\ell+1} - A_S^\ell \tilde{W}^{\ell+1}\|_F^2 \quad (1)$$

Note that $A^\ell W^{\ell+1}$ are the original inputs of layer $\ell + 1$, and $A_S^\ell \tilde{W}^{\ell+1}$ are the inputs after pruning and reweighting, i.e., replacing the weights $W^{\ell+1}$ of layer $\ell + 1$ with the new weights $\tilde{W}^{\ell+1}$.

3.1 Greedy selection

Solving Problem (1) exactly is NP-Hard (Natarajan, 1995). However, we show below that it can be formulated as a weakly submodular maximization problem, hence it can be efficiently approximated. Let

$$F(S) = \|A^\ell W^{\ell+1}\|_F^2 - \min_{\tilde{W}^{\ell+1}} \|A^\ell W^{\ell+1} - A_S^\ell \tilde{W}^{\ell+1}\|_F^2, \quad (2)$$

then Problem (1) is equivalent to $\max_{|S| \leq k} F(S)$.

Proposition 3.1. *Given $U \subseteq V, k \in \mathbb{N}_+$, F is a normalized non-decreasing $\gamma_{U,k}$ -weakly submodular function, with*

$$\gamma_{U,k} \geq \frac{\min_{\|z\|_2=1, \|z\|_0 \leq |U|+k} \|A^\ell z\|_2^2}{\max_{\|z\|_2=1, \|z\|_0 \leq |U|+1} \|A^\ell z\|_2^2}.$$

The proof of Proposition 3.1 follows by writing F as the sum of $n_{\ell+1}$ sparse linear regression problems $F(S) = \sum_{m=1}^{n_{\ell+1}} \|A^\ell w_m^{\ell+1}\|_2^2 - \min_{\text{supp}(\tilde{w}_m) \subseteq S} \|A^\ell w_m^{\ell+1} - A_S^\ell \tilde{w}_m\|_2^2$, and from the relation established in (Elenberg et al., 2016; Das and Kempe, 2011) between weak submodularity and sparse eigenvalues of the covariance matrix (see Appendix A.1).

We use the GREEDY algorithm to select a set $\hat{S} \subseteq V_\ell$ of k neurons to keep in layer ℓ . As discussed in Section 2, the returned solution is guaranteed to satisfy

$$F(\hat{S}) \geq (1 - e^{-\gamma_{\hat{S},k}}) \max_{|S| \leq k} F(S) \quad (3)$$

The submodularity ratio $\gamma_{\hat{S},k}$ is non-zero if any $2k$ columns of A^ℓ are linearly independent. If the number of training data is larger than the number of neurons, i.e., $n > n_\ell$, this is likely to be satisfied for any $k \leq n_\ell/2$; it is unlikely for the activations of a neuron to match exactly a linear combination of other neurons. We verify that this is indeed the case in our experiments.

We show in Appendix B that F satisfies an even stronger notion of approximate submodularity than weak submodularity, which implies a better approximation guarantee for GREEDY than the one provided in Eq. (3). Though, this requires a somewhat stronger assumption: any $k+1$ columns of A^ℓ should be linearly independent and all rows of $W^{\ell+1}$ should be linearly independent. In particular, we would need that $n_\ell \leq n_{\ell+1}$, which is not always satisfied.

3.2 Reweighting

For a fixed $S \subseteq V_\ell$, the reweighted input change $\|A^\ell W^{\ell+1} - A_S^\ell \tilde{W}^{\ell+1}\|_F^2$ is minimized by setting

$$\tilde{W}^{\ell+1} = x^S(A^\ell) W^{\ell+1}, \quad (4)$$

where $x^S(A^\ell) \in \mathbb{R}^{n_\ell \times n_\ell}$ is the matrix with columns $x^S(a_j^\ell)$ such that

$$x^S(a_j^\ell) \in \arg \min_{\text{supp}(x) \subseteq S} \|a_j^\ell - Ax\|_2^2 \text{ for all } j \in V_\ell. \quad (5)$$

Note that $x^S(a_j) = \mathbf{1}_j$ for all $j \in S$, and that $x^S(A^\ell)$ has zero rows for all $i \notin S$. Hence, the new weights are given by $\tilde{w}_{im}^{\ell+1} = w_{im}^{\ell+1} + \sum_{j \notin S} [x^S(A^\ell)]_{ij} w_{jm}^{\ell+1}$ for all $i \in S$, and $\tilde{w}_{im}^{\ell+1} = 0$ for all $i \notin S, m \in V_{\ell+1}$. Namely, the new weights merge the weights from the dropped neurons into the selected ones.

This is the same reweighting procedure introduced in (Mariet and Sra, 2015). But instead of applying it only at the end to the selected neurons \hat{S} , it is implicitly done at each iteration of our pruning method, as it is required to evaluate F . We discuss next how this can be done efficiently.

3.3 Cost

Each iteration of GREEDY requires $O(n_\ell)$ function evaluations of F . Computing $F(S)$ from scratch needs $O(k \cdot (n_\ell \cdot n_{\ell+1} + n \cdot (n_\ell + n_{\ell+1})))$ time, so a naive implementation of GREEDY is too expensive. The following Proposition outlines how we can efficiently evaluate $F(S + i)$ given that $F(S)$ was computed in the previous iteration.

Proposition 3.2. *Given $S \subseteq V_\ell$ such that $|S| \leq k$, $i \notin S$, let $\text{proj}_S(a_j^\ell) = A_S^\ell x^S(a_j^\ell)$ be the projection of a_j^ℓ onto the column space of A_S^ℓ , $R_S(a_i^\ell) = a_i^\ell - \text{proj}_S(a_i^\ell)$ and $\text{proj}_{R_S(a_i^\ell)}(a_j^\ell) \in \arg \min_{z=R_S(a_i^\ell), \gamma, \gamma \in \mathbb{R}} \|a_j^\ell - z\|_2^2$ the corresponding residual and the projection of a_j^ℓ onto it. We can write*

$$F(i|S) = \sum_{m=1}^{n_{\ell+1}} \|\text{proj}_{R_S(a_i^\ell)}(A_{V \setminus S}^\ell w_m^{\ell+1})\|_2^2,$$

where $\text{proj}_{R_S(a_i)}(A_{V \setminus S}^\ell)$ is the matrix with columns $\text{proj}_S(a_j^\ell)$ for all $j \notin S$, 0 otherwise. Assuming $F(S)$, $\text{proj}_S(a_j^\ell)$ and $x^S(a_j^\ell)$ for all $j \notin S$ were computed in the previous iteration, we can compute $F(S + i)$, $\text{proj}_{S+i}(a_j^\ell)$ and $x^{S+i}(a_j^\ell)$ for all $j \notin (S + i)$ in

$$O(n_\ell \cdot (n_{\ell+1} + n + k)) \text{ time.}$$

Computing the optimal weights in Eq. (4) at the end of GREEDY can then be done in $O(k \cdot n_\ell \cdot n_{\ell+1})$ time.

The proof is given in Appendix A.2, and relies on using optimality conditions to construct the least squares solution $x^{S+i}(a_j^\ell)$ from $x^S(a_j^\ell)$.

In total GREEDY's runtime is then $O(k \cdot (n_\ell)^2 \cdot (n_{\ell+1} + n + k))$. Using a faster variant of GREEDY called STOCHASTIC-GREEDY (Mirzasoleiman et al., 2015) further reduces the cost to $O(\log(1/\epsilon) \cdot (n_\ell)^2 \cdot (n_{\ell+1} + n + k))$ while maintaining almost the same approximation guarantee $(1 - e^{-\gamma_{S,k}} - \epsilon)$ in expectation.¹ Note also that computing the solutions for different budgets $k' \leq k$ can be done at the cost of one by running GREEDY with budget k .

4 Pruning regular regions of neurons

In this section, we discuss how to adapt our approach to pruning regular regions of neurons. This is easily achieved by mapping any set of regular regions to the corresponding set of neurons, then applying the same method in Section 3. In particular, we focus on pruning channels in CNNs.

Given a layer ℓ with n_ℓ output channels, let $X^\ell \in \mathbb{R}^{n \cdot p_\ell \times n_\ell \times r_h \times r_w}$ be its activations for each output channel and training input, where p_ℓ is number of patches obtained by applying a filter of size $r_h \times r_w$, and let $F^{\ell+1} \in \mathbb{R}^{n_{\ell+1} \times n_\ell \times r_h \times r_w}$ be the weights of layer $\ell + 1$, corresponding to n_ℓ filters of size $r_h \times r_w$ for each of its output channels.

¹The results in (Mirzasoleiman et al., 2015) only apply to submodular functions, but it is straightforward to extend them to weakly submodular functions.

When an output channel is pruned in layer ℓ , the corresponding weights in F^ℓ and $F^{\ell+1}$ are removed. Pruning $n^\ell - k$ output channels in layer ℓ reduces the number of parameters and computation cost by $(n^\ell - k)/n^\ell$ for both layer ℓ and $\ell + 1$. If layer ℓ is followed by a batch norm layer, the weights therein corresponding to the pruned channels are also removed.

We arrange the activations $X_c^\ell \in \mathbb{R}^{n \cdot p_\ell \times r_h \cdot r_w}$ of each channel c into $r_h r_w$ columns of $A^\ell \in \mathbb{R}^{n \cdot p_\ell \times n_\ell \cdot r_h \cdot r_w}$, i.e., $A^\ell = [X_1^\ell, \dots, X_{n_\ell}^\ell]$. Similarly, we arrange the weights $F_c^{\ell+1} \in \mathbb{R}^{n_{\ell+1} \times r_h \times r_w}$ of each channel c into $r_h \cdot r_w$ rows of $W^{\ell+1} \in \mathbb{R}^{n_\ell \cdot r_h \cdot r_w \times n_{\ell+1}}$, i.e., $(W^{\ell+1})^\top = [(F_1^\ell)^\top, \dots, (F_{n_\ell}^\ell)^\top]$. Recall that $V_\ell = \{1, \dots, n_\ell\}$, and let $V'_\ell = \{1, \dots, r_h r_w n_\ell\}$. We define a function $M : 2^{V_\ell} \rightarrow 2^{V'_\ell}$ which maps every channel c to its corresponding $r_h r_w$ columns A^ℓ . Let $G(S) = F(M(S))$, with F defined in Eq. (2), then minimizing the reweighted input change $\|A^\ell W^{\ell+1} - A_{M(S)}^\ell \tilde{W}^{\ell+1}\|_F^2$ with a budget k is equivalent to $\max_{|S| \leq k} G(S)$. The following proposition shows that this remains a weakly submodular maximization problem.

Proposition 4.1. *Given $U \subseteq V_\ell, k \in \mathbb{N}_+$, G is a normalized non-decreasing $\gamma_{U,k}$ -weakly submodular function, with*

$$\gamma_{U,k} \geq \frac{\min_{\|z\|_2=1, \|z\|_0 \leq r_h r_w (|U|+k)} \|A^\ell z\|_2^2}{\max_{\|z\|_2=1, \|z\|_0 \leq r_h r_w (|U|+1)} \|A^\ell z\|_2^2}.$$

Proof sketch. G is $\gamma_{U,k}$ -weakly submodular iff F satisfies

$$\gamma_{U,k} F(M(S)|M(L)) \leq \sum_{i \in S} F(M(i)|M(L)),$$

for every two disjoint sets $L, S \subseteq V_\ell$, such that $L \subseteq U, |S| \leq k$. We extend the relation established in (Elenberg et al., 2016; Das and Kempe, 2011) between weak submodularity and sparse eigenvalues of the covariance matrix to this case. The proof then follows in a similar way to Proposition 3.1. \square

As before, we use the GREEDY algorithm, with function G , to select a set $\hat{S} \subseteq V_\ell$ of k channels to keep in layer ℓ . We get the same approximation guarantee $G(\hat{S}) \geq (1 - e^{-\gamma_{\hat{S},k}}) \max_{|S| \leq k} G(S)$. The submodularity ratio $\gamma_{\hat{S},k}$ is non-zero if any $2kr_h r_w$ columns of A^ℓ are linearly independent. In our experiments, we observe that linear independence only holds for small $k \leq 0.2n_\ell$ in certain layers. This is due to the correlation between patches which overlap. To remedy that, we experimented with using only $r_h r_w$ random patches from each image, instead of using all patches. This indeed raises the rank of A^ℓ , such that the range where linear independence is satisfied is around $k \leq 0.4n_\ell$. However, the obtained results are very similar to the ones with all patches, we thus omit them. Note that having linear dependence does not necessarily imply that $\gamma_{\hat{S},k} = 0$; our method still performs well in these cases.

For a fixed $S \subseteq V_\ell$, the optimal weights are again given by $\tilde{W}^{\ell+1} = x^{M(S)}(A^\ell)W^{\ell+1}$. The cost of running GREEDY and reweighting is the same as before (see Appendix A.2).

5 Pruning multiple layers

In this section, we explain how to apply our pruning method to prune multiple layers of a NN.

5.1 Reweighted input change pruning variants

We consider three variants of our method: LAYERINCHANGE, SEQINCHANGE, and ASYMINCHANGE. In LAYERINCHANGE, we prune each layer independently, i.e., we apply exactly the method in

Section 3 or 4, according to the layer’s type. This is the fastest variant; it has the same cost as pruning a single layer, as each layer can be pruned in parallel, and it only requires one forward pass to get the activations of all layers. However, it does not take into account the effect of pruning one layer on subsequent layers.

In SEQINCHANGE, we prune each layer sequentially, starting from the earliest layer to latest one. For each layer ℓ , we apply our method with A^ℓ replaced by the *updated* activations B^ℓ after having pruned previous layers, i.e., we solve $\min_{|S| \leq k, \tilde{W}^{\ell+1} \in \mathbb{R}^{n_\ell \times n_{\ell+1}}} \|B^\ell W^{\ell+1} - B_S^\ell \tilde{W}^{\ell+1}\|_F^2$.

In ASYMINCHANGE, we also prune each layer sequentially, but to avoid the accumulation of error, we use an asymmetric formulation of the reweighted input change, where instead of approximating the *updated* input $B^\ell W^{\ell+1}$, we approximate the *original* input $A^\ell W^{\ell+1}$, i.e., we solve $\min_{|S| \leq k, \tilde{W}^{\ell+1} \in \mathbb{R}^{n_\ell \times n_{\ell+1}}} \|A^\ell W^{\ell+1} - B_S^\ell \tilde{W}^{\ell+1}\|_F^2$. This problem is still a weakly submodular maximization problem, with the same submodularity ratio given in Propositions 3.1 and 4.1, with A^ℓ replaced by B^ℓ therein (see Appendix A.1). Hence, the same approximation guarantee as in the symmetric formulation holds here. Moreover, a better approximation guarantee can again be obtained with a stronger notion of approximate submodularity, under stronger assumptions (see Appendix B). The cost of running GREEDY with the asymmetric formulation and reweighting is also the same as before (see Appendix A.2).

We evaluate all three variants in our experiments. As expected, ASYMINCHANGE usually performs the best, and LAYERINCHANGE the worst. But in some settings, LAYERINCHANGE yields the best results (see Fig. 2).

5.2 Per-layer budget selection

Another important design choice is how much to prune in each layer, given a desired global compression ratio (see Appendix F for the effect of this choice on performance). In our experiments, we use the compression ratio selection method introduced in (Kuzmin et al., 2019, Section 3.4.1), which can be applied to any layerwise pruning method, thus enabling us to have a fair comparison.

Given a network with L layers to prune, let $k = \alpha \sum_{\ell=1}^L n_\ell$, for some $\alpha \in [0, 1]$, be the global number of neurons/channels to keep. We want to select for each layer ℓ , the number of neurons/channels $k_\ell = \alpha_\ell n_\ell$ to keep, with α_ℓ chosen from a set of possible values, e.g., $\alpha_\ell \in \{0.05, 0.1, \dots, 1\}$. We define a layerwise accuracy metric $P_\ell(k_\ell)$ as the accuracy obtained after pruning layer ℓ , with a budget k_ℓ , while other layers are kept intact, evaluated on a verification set. We set aside a subset of the training set to use as a verification set. Let P_{orig} be the original model accuracy. We select the per-layer budgets that minimize the per-layer accuracy drop:

$$\min_{k_1, \dots, k_L} \{ \tau : \forall \ell \in [L], P_\ell(k_\ell) \geq P_{\text{orig}} - \tau, \sum_{\ell=1}^L k_\ell \leq k \}$$

Alternatively, another simple strategy is to prune each layer until the perlayer error (the reweighted input change in our case) reaches some threshold ϵ , and vary ϵ to obtain the desired global compression ratio, as done in (Zhuang et al., 2018; Ye et al., 2020).

6 Empirical Evaluation

In this section, we examine the performance of our proposed pruning method in the limited-data regime. To that end, we focus on one-shot pruning, in which a pre-trained model has to be compressed

in a single step, without any fine-tuning. We study the effect of fine-tuning in Appendix E.

We compare the three variants of our method, LAYERINCHANGE, SEQINCHANGE, and ASYMINCHANGE, with the following baselines:

- RANDOM: prunes randomly selected neurons/channels globally across layers in the network.
- LAYERRANDOM: prunes randomly selected neurons/channels in each layer.
- WEIGHTNORM: prunes neurons/channels with the lowest output weights ℓ_1 -norm, normalized by the number of neurons/channels in their layer. This is the “onorm” function proposed in (He et al., 2014) for neuron pruning, which provides the best results among the three neuron’s importance criteria considered therein.
- LAYERWEIGHTNORM: prunes neurons/channels with the lowest output weights ℓ_1 -norm, in each layer. This is the layerwise version of WEIGHTNORM. It was proposed in (Li et al., 2017) for channel pruning.
- ACTGRAD: prunes neurons/channels with the lowest (activations \times gradients), averaged over the training data, with layerwise ℓ_2 -normalization. This criteria is derived in (Molchanov et al., 2017) from the first order Taylor approximation of the change in loss induced by pruning. It is the best performing criteria among the ones evaluated in (Molchanov et al., 2017).
- LAYERACTGRAD: prunes neurons/channels with the lowest (activations \times gradients), averaged over the training data, in each layer. This is the layerwise variant of ACTGRAD.

We evaluate the performance of these methods on neuron pruning of the LeNet model (LeCun et al., 1989) on the MNIST dataset (Lecun et al., 1998), and on both neuron and channel pruning of the VGG11 model (Simonyan and Zisserman, 2015) on the CIFAR-10 dataset (Krizhevsky et al., 2009). We also consider a transfer learning setting, where we pretrain VGG11 on CIFAR-10, then fine-tune, prune and evaluate it on MNIST.

We implemented all the compared pruning methods using Pytorch (Paszke et al., 2017). Our code builds on the open source ShrinkBench library introduced in (Blalock et al., 2020), and uses the code from (Buschjäger et al., 2020) for GREEDY. We use the implementation of LeNet included in ShrinkBench (Blalock et al., 2020), and a modified version of the implementation of VGG11 provided in (Phan, 2021), where we changed the number of neurons in the first two layers to 128. The code for reproducing all experiments is available at <https://github.com/marwash25/subpruning>.

To compute the gradients and activations used for pruning in ACTGRAD, LAYERACTGRAD, and our method’s variants, we use four batches of size 128 of training images, i.e., $n = 512$, which corresponds to $\sim 1\%$ of the training data in MNIST and CIFAR10. We report top-1 accuracy results evaluated on the validation set, as we vary the fraction α of prunable (i.e., belonging to layers being pruned) neurons/channels kept. Accordingly, when pruning L layers, the global budget is set to $k = \alpha \sum_{\ell=1}^L n_{\ell}$, where n_{ℓ} is the number of neuron/channels in layer ℓ . Unless otherwise specified, we use the per-layer budget selection method described in Section 5.2 for all the pruning methods. We use a subset of the training set, of the same size as the validation set, as verification set for the budget selection method. For fair comparison and to disentangle the benefit of using our pruning method from the benefit of reweighting (Section 3.2), we report results with reweighting applied to all pruning methods, or none of them. Though, we will focus our analysis on the more interesting results with reweighting, with the plots without reweighting mostly serving as a demonstration of the benefit of reweighting.

Results are averaged over five random runs, with standard deviations plotted as error bars. Top-5 accuracy results are reported in Appendix D.2. For additional details on the experimental set-up, see Appendix C.

6.1 Neuron Pruning

We focus in this section on neuron pruning. We pre-train LeNet model on MNIST for 200 epochs, with a batch size of 128, using SGD with Nestrov momentum 0.9 and a fixed learning rate of 1×10^{-3} . The resulting model achieves 97.75% top-1 accuracy. We prune two of the three linear layers in LeNet, where $n_\ell = 120, n_{\ell+1} = 84$ in the first layer, and $n_\ell = 84, n_{\ell+1} = 10$ in the second one. We do not prune the last classifier layer.

Figure 1 shows the top-1 accuracy for different fractions of prunable neurons kept. All three variants of our method consistently outperform other baselines, with ASYMINCHANGE doing the best, and LAYERINCHANGE the worst. We also observe that reweighting significantly improves performance for all methods. It also reduces variance in results; the right subplot has larger error bars than the left one.

We also pre-train VGG11 model on CIFAR-10 for 200 epochs, with a batch size of 128, using Adam with $\beta_1 = 0.9, \beta_2 = 0.99$, a weight decay of 5×10^{-4} , and an initial learning rate of 1×10^{-3} , dropped by 0.1 at epochs 100 and 150. The resulting model achieves 90.11% top-1 accuracy. We prune two of the three linear layers in VGG11, where $n_\ell = n_{\ell+1} = 128$ in the first layer, and $n_\ell = 128, n_{\ell+1} = 10$ in the second one. We again do not prune the last classifier layer.

Figure 2 shows the top-1 accuracy for different fractions of prunable neurons kept. The three variants of our method again consistently outperform other baselines. But surprisingly LAYERINCHANGE is the best variant here. The random methods perform surprisingly well here, even better than other baselines. Though in all other settings we consider, they perform badly. As before, reweighting benefits all methods and reduces variance.

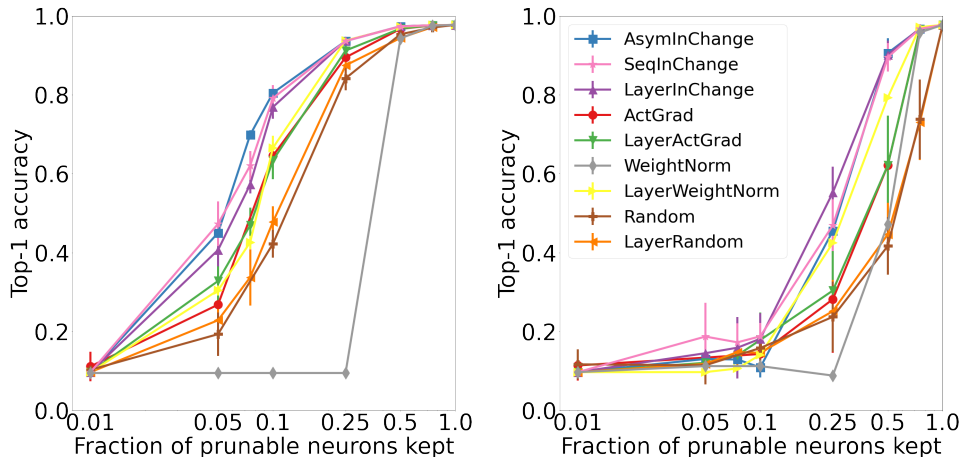


Figure 1: Top-1 Accuracy of different neuron pruning methods on MNIST, after pruning two linear layers in LeNet model, with different fractions of remaining neurons (in log-scale), with (left) and without (right) reweighting.

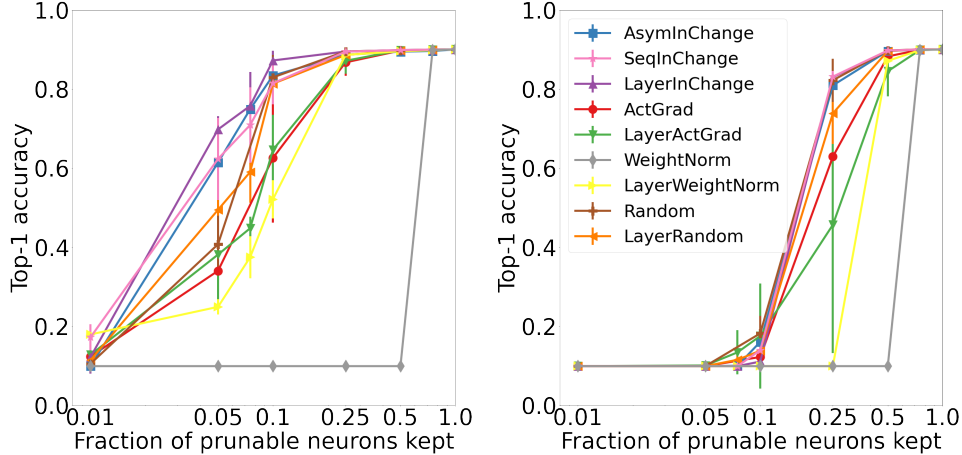


Figure 2: Top-1 Accuracy of different neuron pruning methods on CIFAR10, after pruning two linear layers in VGG11 model, with different fractions of remaining neurons (in log-scale), with (left) and without (right) reweighting.

6.2 Channel Pruning

We focus in this section on channel pruning. We construct the activations and weights matrices A^ℓ and $W^{\ell+1}$ as described in Section 4. We use the same pretrained VGG11 model from Section 6.1. We prune the first two of the eight convolution layers of VGG11, where $n_\ell = 64, n_{\ell+1} = 128$ in the first layer, $n_\ell = 128, n_{\ell+1} = 256$ in the second one, and the kernel size is 3×3 , i.e., $r_h = r_w = 3$. We present results with four convolution layers pruned in Appendix D.1, which are similar to the ones presented here.

Figure 3 shows the top-1 accuracy for different fractions of prunable channels kept. The three variants of our method are among the best performing methods, though their performance is matched by LAYERACTGRAD for most fractions, except very small ones. However, LAYERACTGRAD performs worst in all other settings we consider. SEQINCHANGE is the best performing variant here, with the two others yielding similar results. As in the neuron pruning case, reweighting benefits all methods and reduces variance. Pruning convolution layers seems to affect performance more than pruning linear layers, as we observe a significant drop in accuracy at fraction 0.1 in Figure 3, which doesn't happen until fraction 0.05 in Figures 1 and 2. This can be explained by the fact that pruning channels results in removing more weights than pruning neurons.

6.3 Transfer learning

In this section, we consider a transfer learning setting, where given a large pre-trained network, our goal is to obtain a small model on a different but related task to the one it was trained on. We upsample MNIST images to match the shape of CIFAR-10 images, and use the same VGG11 model pretrained on CIFAR-10 as in Sections 6.1 and 6.2. We fine-tune it for 10 epochs on MNIST with the same setup used during pre-training, and obtain 98.09% top-1 accuracy. We refer to the resulting model as VGG11-MNIST. For the per-layer budget selection, we use the layerwise accuracy metrics $P_\ell(k_\ell)$ computed on CIFAR-10. We prune the same two convolution layers as in Section 6.2. Results with four convolution layers pruned are presented in Appendix D.1.

Figure 4 shows the top-1 accuracy for different fractions of prunable channels kept. The three

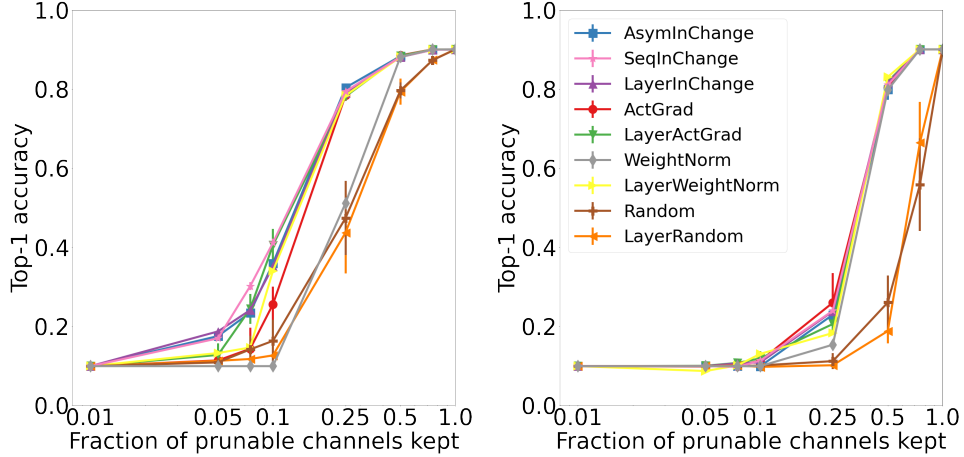


Figure 3: Top-1 Accuracy of different channel pruning methods on CIFAR10, after pruning two convolution layers in VGG11 model, with different fractions of remaining channels (in log-scale), with (left) and without (right) reweighting.

variants of our method outperform other baselines significantly at fraction 0.05 and slightly at fraction 0.075. ASYMINCHANGE is the best performing variant here, with the two others yielding matching results. Here also reweighting improves the performance of all methods and reduces variance. Note that pruning is easier in this setting, as it is not necessary to have a network as large as VGG11 to achieve good performance on MNIST. This allows us to maintain an accuracy close to the original model up to fraction 0.05.

6.4 Discussion

We summarize our observations from the empirical results:

- Our proposed pruning method outperforms common baselines like WEIGHTNORM, ACTGRAD, RANDOM, and their layerwise variants, in various one shot pruning settings. Our results also illustrate the robustness of our method, as it reliably yields the best results in various settings, while other baselines perform well in some settings but not in others.
- Among the three variants of our method, ASYMINCHANGE usually performs the best, and LAYERINCHANGE the worst. But in some settings, LAYERINCHANGE yields the best results (see Fig. 2).
- In all one shot pruning settings we considered, reweighting significantly improves performance for all methods, and reduces variance in results.
- The choice of how much to prune in each layer given a global budget can have a drastic effect on performance, as illustrated in Appendix F.
- Fine-tuning boosts performance even more than reweighting, as illustrated in Appendix E. Though reweighting often still helps even when fine-tuning is used. Hence, both fine-tuning and reweighting should be used when possible, i.e., in settings where data is not limited. In fact, for the relatively small models we consider, the choice of the pruning method is less important when both fine-tuning and reweighting are used, since all methods perform fairly similarly, except WEIGHTNORM and the random selection methods.

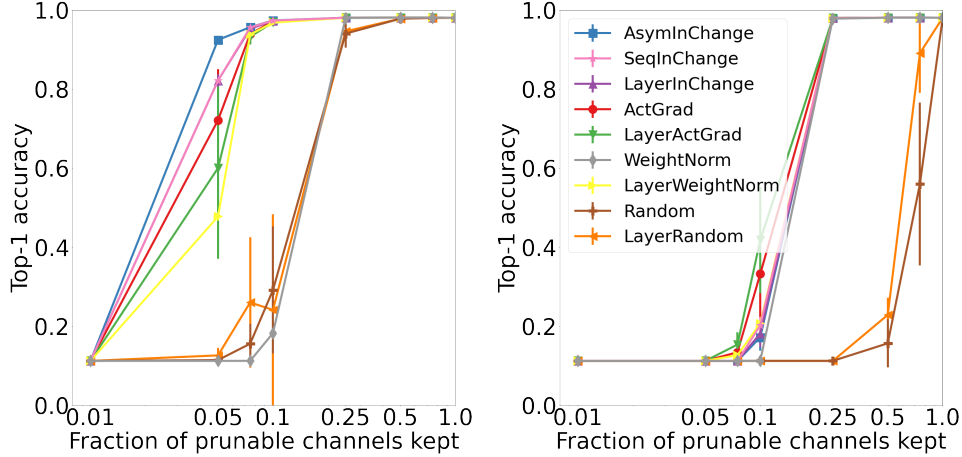


Figure 4: Top-1 Accuracy of different channel pruning methods on MNIST, after pruning two convolution layers in VGG11-MNIST model, with different fractions of remaining channels (in log-scale), with (left) and without (right) reweighting.

7 Conclusion

We proposed a data-efficient structured pruning method, based on submodular optimization. By casting the layerwise subset selection problem as a weakly submodular optimization problem, we are able to use the GREEDY algorithm to provably approximate it. Empirically, our method outperforms common baselines on various neuron and channel one-shot pruning tasks.

Acknowledgements

We thank Stefanie Jegelka, Debadeepta Dey, Jose Gallego-Posada for their helpful discussions, and Boris Knyazev for his help with running experiments. We also acknowledge the MIT SuperCloud and Lincoln Laboratory Supercomputing Center (supercloud.mit.edu), Compute Canada (www.computeCanada.ca), WestGrid (www.westgrid.ca), and the Mila IDT team for providing HPC resources that have contributed to the research results reported within this paper. This research was partially supported by the Canada CIFAR AI Chair Program. Simon Lacoste-Julien is a CIFAR Associate Fellow in the Learning Machines & Brains program.

References

- Andrew An Bian, Joachim M Buhmann, Andreas Krause, and Sebastian Tschiatschek. Guarantees for greedy maximization of non-submodular functions with applications. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pages 498–507. JMLR. org, 2017.
- 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.
- Cristian Bucila, Rich Caruana, and Alexandru Niculescu-Mizil. Model compression. In *Proceedings of the 12th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, KDD ’06, page 535–541, New York, NY, USA, 2006. Association for Computing Machinery. ISBN 1595933395. doi: 10.1145/1150402.1150464. URL <https://doi.org/10.1145/1150402.1150464>.

- Sebastian Buschjäger, Philipp-Jan Honysz, and Katharina Morik. Very fast streaming submodular function maximization, 2020.
- Maxwell D Collins and Pushmeet Kohli. Memory bounded deep convolutional networks. *arXiv preprint arXiv:1412.1442*, 2014.
- Matthieu Courbariaux, Yoshua Bengio, and Jean-Pierre David. Binaryconnect: Training deep neural networks with binary weights during propagations. In *Advances in neural information processing systems*, pages 3123–3131, 2015.
- A. Das and D. Kempe. Submodular meets spectral: Greedy algorithms for subset selection, sparse approximation and dictionary selection. *arXiv preprint arXiv:1102.3975*, 2011.
- Misha Denil, Babak Shakibi, Laurent Dinh, Marc Aurelio Ranzato, and Nando de Freitas. Predicting parameters in deep learning. In C. J. C. Burges, L. Bottou, M. Welling, Z. Ghahramani, and K. Q. Weinberger, editors, *Advances in Neural Information Processing Systems*, volume 26. Curran Associates, Inc., 2013. URL <https://proceedings.neurips.cc/paper/2013/file/7fec306d1e665bc9c748b5d2b99a6e97-Paper.pdf>.
- Emily L Denton, Wojciech Zaremba, Joan Bruna, Yann LeCun, and Rob Fergus. Exploiting linear structure within convolutional networks for efficient evaluation. In *Advances in neural information processing systems*, pages 1269–1277, 2014.
- M. El Halabi, F. Bach, and V Cevher. Combinatorial penalties: Structure preserved by convex relaxations. *Proceedings of the 21st International Conference on Artificial Intelligence and Statistics*, 2018.
- Ethan R Elenberg, Rajiv Khanna, Alexandros G Dimakis, and Sahand Negahban. Restricted strong convexity implies weak submodularity. *arXiv preprint arXiv:1612.00804*, 2016.
- Yunchao Gong, Liu Liu, Ming Yang, and Lubomir Bourdev. Compressing deep convolutional networks using vector quantization. *arXiv preprint arXiv:1412.6115*, 2014.
- B. Hassibi, D.G. Stork, and G.J. Wolff. Optimal brain surgeon and general network pruning. In *IEEE International Conference on Neural Networks*, pages 293–299 vol.1, 1993. doi: 10.1109/ICNN.1993.298572.
- 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.
- 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.
- Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network. *Neural Information Processing Systems (NeurIPS) Workshops*, 2015.
- Torsten Hoefer, Dan Alistarh, Tal Ben-Nun, Nikoli Dryden, and Alexandra Peste. Sparsity in deep learning: Pruning and growth for efficient inference and training in neural networks. *arXiv preprint arXiv:2102.00554*, 2021.
- Max Jaderberg, Andrea Vedaldi, and Andrew Zisserman. Speeding up convolutional neural networks with low rank expansions. In *Proceedings of the British Machine Vision Conference*. BMVA Press, 2014. doi: <http://dx.doi.org/10.5244/C.28.88>.

- Yong-Deok Kim, Eunhyeok Park, Sungjoo Yoo, Taelim Choi, Lu Yang, and Dongjun Shin. Compression of deep convolutional neural networks for fast and low power mobile applications. *arXiv preprint arXiv:1511.06530*, 2015.
- Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009.
- Andrey Kuzmin, Markus Nagel, Saurabh Pitre, Sandeep Pendyam, Tijmen Blankevoort, and Max Welling. Taxonomy and evaluation of structured compression of convolutional neural networks. *arXiv preprint arXiv:1912.09802*, 2019.
- Vadim Lebedev, Yaroslav Ganin, Maksim Rakhuba, Ivan V Oseledets, and Victor S Lempitsky. Speeding-up convolutional neural networks using fine-tuned cp-decomposition. In *International Conference on Learning Representations*, 2015.
- Y Lecun, C Cortes, and C Burges. The mnist database of handwritten digits, 1998.
- Yann LeCun, Bernhard Boser, John S Denker, Donnie Henderson, Richard E Howard, Wayne Hubbard, and Lawrence D Jackel. Backpropagation applied to handwritten zip code recognition. *Neural computation*, 1(4):541–551, 1989.
- Yann LeCun, John Denker, and Sara Solla. Optimal brain damage. In D. Touretzky, editor, *Advances in Neural Information Processing Systems*, volume 2. Morgan-Kaufmann, 1990. URL <https://proceedings.neurips.cc/paper/1989/file/6c9882bbac1c7093bd25041881277658-Paper.pdf>.
- Benny Lehmann, Daniel Lehmann, and Noam Nisan. Combinatorial auctions with decreasing marginal utilities. *Games and Economic Behavior*, 55(2):270–296, 2006.
- Hao Li, Asim Kadav, Igor Durdanovic, Hanan Samet, and Hans Peter Graf. Pruning filters for efficient convnets. In *5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings*. OpenReview.net, 2017. URL <https://openreview.net/forum?id=rJqFGTslg>.
- Jian-Hao Luo, Jianxin Wu, and Weiyao Lin. Thinet: A filter level pruning method for deep neural network compression. In *Proceedings of the IEEE international conference on computer vision*, pages 5058–5066, 2017.
- Zelda Mariet and Suvrit Sra. Diversity networks: Neural network compression using determinantal point processes. *arXiv preprint arXiv:1511.05077*, 2015.
- Kris McGuffie and Alex Newhouse. The radicalization risks of gpt-3 and advanced neural language models. *arXiv preprint arXiv:2009.06807*, 2020.
- Baharan Mirzasoleiman, Ashwinkumar Badanidiyuru, Amin Karbasi, Jan Vondrák, and Andreas Krause. Lazier than lazy greedy. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 29, 2015.
- Pavlo Molchanov, Stephen Tyree, Tero Karras, Timo Aila, and Jan Kautz. Pruning convolutional neural networks for resource efficient inference. *ICLR*, 2017.
- Balas Kausik Natarajan. Sparse approximate solutions to linear systems. *SIAM journal on computing*, 24(2):227–234, 1995.
- G.L. Nemhauser, L.A. Wolsey, and M.L. Fisher. An analysis of approximations for maximizing submodular set functions — I. *Mathematical Programming*, 14(1):265–294, 1978.

- Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito, Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. Automatic differentiation in pytorch. 2017.
- Huy Phan. huyvnphan/pytorch_cifar10, January 2021. URL <https://doi.org/10.5281/zenodo.4431043>.
- 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, 2013. doi: 10.1109/ICASSP.2013.6638949.
- Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. In Yoshua Bengio and Yann LeCun, editors, *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*, 2015. URL <http://arxiv.org/abs/1409.1556>.
- Suraj Srinivas and R. Venkatesh Babu. Data-free parameter pruning for deep neural networks. In *Proceedings of the British Machine Vision Conference (BMVC)*, pages 31.1–31.12. BMVA Press, September 2015.
- Jiahao Su, Jingling Li, Bobby Bhattacharjee, and Furong Huang. Tensorial neural networks: Generalization of neural networks and application to model compression. *arXiv preprint arXiv:1805.10352*, 2018.
- Maxim Sviridenko, Jan Vondrák, and Justin Ward. Optimal approximation for submodular and supermodular optimization with bounded curvature. *Mathematics of Operations Research*, 42(4): 1197–1218, 2017.
- Elena Voita, David Talbot, Fedor Moiseev, Rico Sennrich, and Ivan Titov. Analyzing multi-head self-attention: Specialized heads do the heavy lifting, the rest can be pruned. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 5797–5808, 2019.
- Mao Ye, Chengyue Gong, Lizhen Nie, Denny Zhou, Adam Klivans, and Qiang Liu. Good subnetworks provably exist: Pruning via greedy forward selection. *ICML*, 2020.
- Xiangyu Zhang, Jianhua Zou, Kaiming He, and Jian Sun. Accelerating very deep convolutional networks for classification and detection. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 38(10):1943–1955, 2016. doi: 10.1109/TPAMI.2015.2502579.
- Zhuangwei Zhuang, Minghui Tan, Bohan Zhuang, Jing Liu, Yong Guo, Qingyao Wu, Junzhou Huang, and Jinhui Zhu. Discrimination-aware channel pruning for deep neural networks. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 31. Curran Associates, Inc., 2018. URL <https://proceedings.neurips.cc/paper/2018/file/55a7cf9c71f1c9c495413f934dd1a158-Paper.pdf>.

A Missing proofs

Recall that $F(S) = \|A^\ell W^{\ell+1}\|_F^2 - \min_{\tilde{W}^{\ell+1} \in \mathbb{R}^{n_\ell \times n_{\ell+1}}} \|A^\ell W^{\ell+1} - A_S^\ell \tilde{W}^{\ell+1}\|_F^2$, and $G(S) = F(M(S))$, where M maps each channel to its corresponding columns in A^ℓ . We denote by $\tilde{F}(S)$ the objective corresponding to the asymmetric formulation introduced in Section 5.1, i.e., $\tilde{F}(S) = \|A^\ell W^{\ell+1}\|_F^2 - \min_{\tilde{W}^{\ell+1} \in \mathbb{R}^{n_\ell \times n_{\ell+1}}} \|A^\ell W^{\ell+1} - B_S^\ell \tilde{W}^{\ell+1}\|_F^2$, and similarly $\tilde{G}(S) = \tilde{F}(M(S))$, where M maps each channel to its corresponding columns in A^ℓ .

We introduce some notation that will be used throughout the Appendix. Given any matrix D and vector y , we denote by $x^S(y) \in \arg \min_{\text{supp}(x) \subseteq S} \frac{1}{2} \|y - Dx\|_2^2$ the vector of optimal regression coefficients, and by $\text{proj}_S(y) = Dx^S(y)$, $R^S(y) = y - \text{proj}_S(y)$ the corresponding projection and residual.

A.1 Proof of Proposition 3.1 and Proposition 4.1 and their extension to the asymmetric formulation

In this section, we prove that F, G , and their asymmetric variants \tilde{F}, \tilde{G} are all non-decreasing weakly submodular functions. We start by reviewing the definition of restricted smoothness (RSM) and restricted strong convexity (RSC).

Definition A.1 (RSM/RSC). Given a differentiable function $\ell : \mathbb{R}^d \rightarrow \mathbb{R}$ and $\Omega \subset \mathbb{R}^d \times \mathbb{R}^d$, ℓ is μ_Ω -RSC and ν_Ω -RSM if $\frac{\mu_\Omega}{2} \|x - y\|_2^2 \leq \ell(y) - \ell(x) - \langle \nabla \ell(x), y - x \rangle \leq \frac{\nu_\Omega}{2} \|x - y\|_2^2$, $\forall (x, y) \in \Omega$.

If ℓ is RSC/RSM on $\Omega = \{(x, y) : \|x\|_0 \leq k, \|y\|_0 \leq k, \|x - y\|_0 \leq k\}$, we denote by μ_k, ν_k the corresponding RSC and RSM parameters.

Proposition 3.1. *Given $U \subseteq V, k \in \mathbb{N}_+$, F is a normalized non-decreasing $\gamma_{U,k}$ -weakly submodular function, with*

$$\gamma_{U,k} \geq \frac{\min_{\|z\|_2=1, \|z\|_0 \leq |U|+k} \|A^\ell z\|_2^2}{\max_{\|z\|_2=1, \|z\|_0 \leq |U|+1} \|A^\ell z\|_2^2}.$$

Proof. We can write $F(S) = \sum_{m=1}^{n_{\ell+1}} F_m(S) := \ell_m(0) - \min_{\text{supp}(\tilde{w}_m) \subseteq S} \ell_m(\tilde{w}_m)$, where $\ell_m(\tilde{w}_m) = \|A^\ell w_m^{\ell+1} - A^\ell \tilde{w}_m\|_2^2$. The function $F_m(S)$ is then $\gamma_{U,k}$ -weakly submodular with $\gamma_{U,k} \geq \frac{\mu_{|U|+k}}{\nu_{|U|+1}}$ (Elenberg et al., 2016), where $\mu_{|U|+k}$ and $\nu_{|U|+1}$ are the RSC and RSM parameters of ℓ_m , given by $\mu_{|U|+k} = \min_{\|z\|_2=1, \|z\|_0 \leq |U|+k} \|A^\ell z\|_2^2$, and $\nu_{|U|+1} = \max_{\|z\|_2=1, \|z\|_0 \leq |U|+1} \|A^\ell z\|_2^2$. It follows then that F is also $\gamma_{U,k}$ -weakly submodular. It is easy to check that F is also normalized and non-decreasing. \square

Proposition 4.1. *Given $U \subseteq V_\ell, k \in \mathbb{N}_+$, G is a normalized non-decreasing $\gamma_{U,k}$ -weakly submodular function, with*

$$\gamma_{U,k} \geq \frac{\min_{\|z\|_2=1, \|z\|_0 \leq r_h r_w (|U|+k)} \|A^\ell z\|_2^2}{\max_{\|z\|_2=1, \|z\|_0 \leq r_h r_w (|U|+1)} \|A^\ell z\|_2^2}.$$

Proof. By definition, G is $\gamma_{U,k}$ -weakly submodular iff F satisfies

$$\gamma_{U,k} F(M(S)|M(L)) \leq \sum_{i \in S} F(M(i)|M(L)),$$

for every two disjoint sets $L, S \subseteq V_\ell$, such that $L \subseteq U, |S| \leq k$. We extend the relation established in (Elenberg et al., 2016) between weak submodularity and RSC/RSM parameters to this case.

Let $S' = M(S)$, $L' = M(L)$, $I' = M(i)$, and $k' = r_h r_w k$. As before, we can write $G(S) = F(S') = \sum_{m=1}^{n_{\ell}+1} F_m(S') := \ell_m(0) - \min_{\text{supp}(\tilde{w}_m) \subseteq S'} \ell_m(\tilde{w}_m)$, where $\ell_m(\tilde{w}_m) = \|A^\ell w_m^{\ell+1} - A^\ell \tilde{w}_m\|_2^2$. We denote by μ_k and ν_k the RSC and RSM parameters of ℓ_m , given by $\mu_k = \min_{\|z\|_2=1, \|z\|_0 \leq k} \|A^\ell z\|_2^2$, and $\nu_k = \max_{\|z\|_2=1, \|z\|_0 \leq k} \|A^\ell z\|_2^2$. To simplify notation, we use $x^S := x^S(A^\ell w_m^{\ell+1})$.

For every two disjoint sets $L, S \subseteq V_\ell$, such that $L \subseteq U$, $|S| \leq k$, we have:

$$\begin{aligned} 0 \leq F_m(S'|L') &= \ell_m(x^{L'}) - \ell_m(x^{S' \cup L'}) \\ &\leq -\langle \nabla \ell_m(x^{L'}), x^{S' \cup L'} - x^{L'} \rangle - \frac{\mu_{|L'|+k'}}{2} \|x^{S' \cup L'} - x^{L'}\|_2^2 \\ &\leq \max_{\text{supp}(x) \subseteq S' \cup L'} -\langle \nabla \ell_m(x^{L'}), x - x^{L'} \rangle - \frac{\mu_{|L'|+k'}}{2} \|x - x^{L'}\|_2^2 \end{aligned}$$

By setting $x = x^{L'} - \frac{[\nabla \ell_m(x^{L'})]_{S'}}{\mu_{|L'|+k'}}$, we get $G(S'|L') \leq \frac{\|[\nabla \ell_m(x^{L'})]_{S'}\|_2^2}{2\mu_{|L'|+k'}}$.

Given any $i \in S$, $I' = M(i)$, we have

$$\begin{aligned} F_m(I'|L') &= \ell_m(x^{L'}) - \ell_m(x^{I' \cup L'}) \\ &\geq \ell_m(x^{L'}) - \ell_m(x^{L'} - \frac{[\nabla \ell_m(x^{L'})]_{I'}}{\nu_{|L'|+|I'|}}) \\ &\geq \langle \nabla \ell_m(x^{L'}), \frac{[\nabla \ell_m(x^{L'})]_{I'}}{\nu_{|L'|+|I'|}} \rangle - \frac{\nu_{|L'|+k'}}{2} \left\| \frac{[\nabla \ell_m(x^{L'})]_{I'}}{\nu_{|L'|+|I'|}} \right\|_2^2 \\ &= \frac{\|[\nabla \ell_m(x^{L'})]_{I'}\|_2^2}{2\nu_{|L'|+|I'|}} \end{aligned}$$

Hence,

$$\begin{aligned} G(S|L) &\leq \sum_{m=1}^{n_{\ell}+1} \frac{\|[\nabla \ell_m(x^{L'})]_{S'}\|_2^2}{2\mu_{|L'|+k'}} = \sum_{i \in S, I' = M(i)} \frac{\|[\nabla \ell_m(x^{L'})]_{I'}\|_2^2}{2\mu_{|L'|+k'}} \\ &= \sum_{i \in S, I' = M(i)} \frac{\nu_{|L'|+|I'|}}{\mu_{|L'|+k'}} \frac{\|[\nabla \ell_m(x^{L'})]_{I'}\|_2^2}{2\nu_{|L'|+|I'|}} = \frac{\nu_{|L'|+k'}}{\mu_{|L'|+k'}} \sum_{i \in S} G(i|L). \end{aligned}$$

We thus have $\gamma_{U,k} \geq \frac{\mu_{|U'|+k'}}{\nu_{|U'|+|I'|}}$. □

Both Proposition 3.1 and Proposition 4.1 apply also to the asymmetric variants, using exactly the same proofs.

Proposition A.2. *Given $U \subseteq V$, $k \in \mathbb{N}_+$, \tilde{F} is a normalized non-decreasing $\gamma_{U,k}$ -weakly submodular function, with*

$$\gamma_{U,k} \geq \frac{\min_{\|z\|_2=1, \|z\|_0 \leq |U|+k} \|B^\ell z\|_2^2}{\max_{\|z\|_2=1, \|z\|_0 \leq |U|+1} \|B^\ell z\|_2^2}.$$

Proposition A.3. *Given $U \subseteq V_\ell$, $k \in \mathbb{N}_+$, G is a normalized non-decreasing $\gamma_{U,k}$ -weakly submodular function, with*

$$\gamma_{U,k} \geq \frac{\min_{\|z\|_2=1, \|z\|_0 \leq r_h r_w (|U|+k)} \|A^\ell z\|_2^2}{\max_{\|z\|_2=1, \|z\|_0 \leq r_h r_w (|U|+1)} \|A^\ell z\|_2^2}.$$

A.2 Proof of Proposition 3.2 and its extension to other variants

In this section, we investigate the cost of applying GREEDY with F, G and their asymmetric variants \tilde{F}, \tilde{G} . To that end, we need the following key lemmas showing how to update the least squares solutions and the function values after adding one or more elements.

Lemma A.4. *Given a matrix D , vector y , and a vector of optimal regression coefficients $x^S(y) \in \arg \min_{\text{supp}(x) \subseteq S} \frac{1}{2} \|y - Dx\|_2^2$, we have for all $S \subseteq V, i \notin S$:*

$$x^{S \cup i}(y) = (x^S(y) - x^S(d_i) \gamma^{S,i}(y)) + \gamma^{S,i}(y) \mathbf{1}_i \in \arg \min_{\text{supp}(x) \subseteq S \cup i} \frac{1}{2} \|y - Dx\|_2^2,$$

where $\gamma^{S,i}(y) \in \arg \min_{\gamma \in \mathbb{R}} \frac{1}{2} \|y - R^S(d_i) \gamma\|_2^2$. Hence, $\text{proj}_{S \cup i}(y) = \text{proj}_S(y) + \text{proj}_{R^S(d_i)}(y)$, where $\text{proj}_{R^S(d_i)}(y) = R^S(d_i) \gamma^{S,i}(y)$.

Similarly, for $I \subseteq V \setminus S$, let $R^S(D_I)$ be the matrix with columns $R^S(d_i)$, $x^S(D_I)$ the matrix with columns $x^S(d_i)$, and $\gamma^{S,I}(y) \in \arg \min_{\gamma \in \mathbb{R}^{|I|}} \frac{1}{2} \|y - R^S(D_I) \gamma\|_2^2$, then

$$x^{S \cup I}(y) = (x^S(y) - x^S(D_I) \gamma^{S,I}(y)) + e_I \gamma^{S,I}(y) \in \arg \min_{\text{supp}(x) \subseteq S \cup I} \frac{1}{2} \|y - Dx\|_2^2,$$

where $e_I \in \mathbb{R}^{|V| \times |I|}$ is the matrix with $[e_I]_{i,i} = 1$ for all $i \in I$, and 0 elsewhere. Hence, $\text{proj}_{S \cup I}(y) = \text{proj}_S(y) + \text{proj}_{R^S(D_I)}(y)$, where $\text{proj}_{R^S(D_I)}(y) = R^S(D_I) \gamma^{S,I}(y)$.

Proof. By optimality conditions, we have:

$$D_S^\top (D_S x^S(y) - y) = 0 \tag{6}$$

$$D_S^\top (D_S x^S(d_i) - d_i) = 0 \Rightarrow -D_S^\top R^S(d_i) = 0 \tag{7}$$

$$R^S(d_i)^\top (R^S(d_i) \gamma^{S,i}(y) - y) = 0 \tag{8}$$

We prove that $\hat{x}^{S \cup i}(y) = (x^S(y) - x^S(d_i) \gamma^{S,i}(y)) + \gamma^{S,i}(y) \mathbf{1}_i$ satisfies the optimality conditions on $x^{S \cup i}(y)$, hence $\hat{x}^{S \cup i}(y) = x^{S \cup i}(y)$.

We have $D_{S \cup i} \hat{x}^{S \cup i}(y) = D_S x^S(y) + R^S(d_i) \gamma^{S,i}(y)$, then

$$D_S^\top (D_{S \cup i} \hat{x}^{S \cup i}(y) - y) = D_S^\top (D_S x^S(y) - y) + D_S^\top R^S(d_i) \gamma^{S,i}(y) = 0$$

and

$$\begin{aligned} d_i^\top (D_{S \cup i} \hat{x}^{S \cup i}(y) - y) &= (R^S(d_i) + D_S x^S(d_i))^\top (D_S x^S(y) + R^S(d_i) \gamma^{S,i}(y) - y) \\ &= R^S(d_i)^\top D_S x^S(y) + R^S(d_i)^\top (R^S(d_i) \gamma^{S,i}(y) - y) \\ &\quad + (D_S x^S(d_i))^\top (D_S x^S(y) - y) + (D_S x^S(d_i))^\top R^S(d_i) \gamma^{S,i}(y) \\ &= 0 \end{aligned}$$

The proof for the case where we add multiple indices at once follows similarly. \square

Lemma A.5. *For all $S \subseteq V_\ell, i \notin S$, let $R_S(b_i^\ell) = b_i^\ell - \text{proj}_S(b_i^\ell)$, $\text{proj}_{R_S(b_i^\ell)}(y) = R_S(b_i^\ell) \gamma^{S,i}(y)$ with $\gamma^{S,i}(y) \in \arg \min_{\gamma \in \mathbb{R}} \|y - R_S(b_i^\ell) \gamma\|_2^2$. We can write the marginal gain of adding i to S w.r.t \tilde{F} as:*

$$\tilde{F}(i|S) = \sum_{m=1}^{n^{\ell+1}} \|\text{proj}_{R_S(b_i^\ell)}(A^\ell) w_m^{\ell+1}\|_2^2,$$

where $\text{proj}_{R_S(b_i^\ell)}(A^\ell)$ is the matrix with columns $\text{proj}_{R_S(b_i^\ell)}(a_j^\ell)$ for all $j \in V_\ell$. Similarly, for all $S \subseteq V_\ell, I \subseteq V_\ell \setminus S$, let $R_S(B_I^\ell) = B_I^\ell - \text{proj}_S(B_I^\ell)$, $\text{proj}_{R_S(B_I^\ell)}(y) = R_S(B_I^\ell)\gamma^{S,I}(y)$ with $\gamma^{S,I}(y) \in \arg \min_{\gamma \in \mathbb{R}^{|I|}} \|y - R_S(B_I^\ell)\gamma\|_2^2$. We can write the marginal gain of adding I to S w.r.t \tilde{F} as:

$$\tilde{F}(I|S) = \sum_{m=1}^{n^{\ell+1}} \|\text{proj}_{R_S(B_I^\ell)}(A^\ell)w_m^{\ell+1}\|_2^2,$$

where $\text{proj}_{R_S(B_I^\ell)}(A^\ell)$ is the matrix with columns $\text{proj}_{R_S(B_I^\ell)}(a_j^\ell)$ for all $j \in V_\ell$.

Proof. We prove the claim for the case where we add several elements. The case where we add a single element then follows as a special case. For a fixed $S \subseteq V_\ell$, the reweighted asymmetric input change $\|A^\ell W^{\ell+1} - B_S^\ell \tilde{W}^{\ell+1}\|_F^2$ is minimized by setting $\tilde{W}^{\ell+1} = x^S(A^\ell)W^{\ell+1}$, where $x^S(A^\ell) \in \mathbb{R}^{n_\ell \times n_\ell}$ is the matrix with columns $x^S(a_j^\ell)$ such that

$$x^S(a_j^\ell) \in \arg \min_{\text{supp}(x) \subseteq S} \|a_j^\ell - B^\ell x\|_2^2 \text{ for all } j \in V_\ell. \quad (9)$$

Plugging $\tilde{W}^{\ell+1}$ into the expression of $\tilde{F}(S)$ yields $\tilde{F}(S) = \|A^\ell W^{\ell+1}\|_F^2 - \|(A^\ell - \text{proj}_S(A^\ell))W^{\ell+1}\|_F^2$, where $\text{proj}_S(A^\ell) = B_S^\ell x^S(A^\ell)$. For every $m \in \{1, \dots, n_{\ell+1}\}$, we have:

$$\begin{aligned} & \|(\text{proj}_{S \cup I}(A^\ell) - A^\ell)w_m^{\ell+1}\|_2^2 - \|(\text{proj}_S(A^\ell) - A^\ell)w_m^{\ell+1}\|_2^2 \\ &= \|(\text{proj}_S(A^\ell) + \text{proj}_{R_S(B_I)}(A^\ell) - A^\ell)w_m^{\ell+1}\|_2^2 - \|(\text{proj}_S(A^\ell) - A^\ell)w_m^{\ell+1}\|_2^2 \quad (\text{by Lemma A.4}) \\ &= \|\text{proj}_{R_S(B_I)}(A^\ell)w_m^{\ell+1}\|_2^2 - 2\langle (A^\ell - \text{proj}_S(A^\ell))w_m^{\ell+1}, \text{proj}_{R_S(B_I)}(A^\ell)w_m^{\ell+1} \rangle \\ &= \|\text{proj}_{R_S(B_I)}(A^\ell)w_m^{\ell+1}\|_2^2 - 2\langle \text{proj}_{R_S(B_I)}(A^\ell)w_m^{\ell+1}, \text{proj}_{R_S(B_I)}(A^\ell)w_m^{\ell+1} \rangle \\ &= -\|\text{proj}_{R_S(B_I)}(A^\ell)w_m^{\ell+1}\|_2^2 \end{aligned}$$

where the second to last equality holds because $y - \text{proj}_S(y) - \text{proj}_{R_S(B_I)}(y)$ and $\text{proj}_{R_S(B_I)}(y')$ are orthogonal by optimality conditions (see proof of Lemma A.4):

$$\langle y - \text{proj}_S(y) - \text{proj}_{R_S(B_I)}(y), \text{proj}_{R_S(B_I)}(y') \rangle = \langle y - B_S x^S(y) - R_S(B_I)\gamma^{S,I}(y), R_S(B_I)\gamma^{S,I}(y') \rangle = 0.$$

Hence, $\tilde{F}(I|S) = \sum_{m=1}^{n^{\ell+1}} \|\text{proj}_{R_S(B_I)}(A^\ell)w_m^{\ell+1}\|_2^2$. In particular, if $I = \{i\}$, $\tilde{F}(i|S) = \sum_{m=1}^{n^{\ell+1}} \|\text{proj}_{R_S(b_i^\ell)}(A^\ell)w_m^{\ell+1}\|_2^2$. \square

Lemma A.6. For all $S \subseteq V_\ell, i \notin S$, let $R_S(a_i^\ell) = a_i^\ell - \text{proj}_S(a_i^\ell)$, $\text{proj}_{R_S(a_i^\ell)}(y) = R_S(a_i^\ell)\gamma^{S,i}(y)$ with $\gamma^{S,i}(y) \in \arg \min_{\gamma \in \mathbb{R}} \|y - R_S(a_i^\ell)\gamma\|_2^2$. We can write the marginal gain of adding i to S w.r.t F as:

$$F(i|S) = \sum_{m=1}^{n^{\ell+1}} \|\text{proj}_{R_S(a_i^\ell)}(A_{V \setminus S}^\ell)w_m^{\ell+1}\|_2^2,$$

where $\text{proj}_{R_S(a_i^\ell)}(A_{V \setminus S}^\ell)$ is the matrix with columns $\text{proj}_{R_S(a_i^\ell)}(a_j^\ell)$ for all $j \in V \setminus S$, 0 otherwise. Similarly, for all $S \subseteq V_\ell, I \subseteq V_\ell \setminus S$, let $R_S(A_I^\ell) = A_I^\ell - \text{proj}_S(A_I^\ell)$, $\text{proj}_{R_S(A_I^\ell)}(y) = R_S(A_I^\ell)\gamma^{S,I}(y)$ with $\gamma^{S,I}(y) \in \arg \min_{\gamma \in \mathbb{R}^{|I|}} \|y - R_S(A_I^\ell)\gamma\|_2^2$. We can write the marginal gain of adding I to S w.r.t F as:

$$F(I|S) = \sum_{m=1}^{n^{\ell+1}} \|\text{proj}_{R_S(A_I^\ell)}(A_{V \setminus S}^\ell)w_m^{\ell+1}\|_2^2,$$

where $\text{proj}_{R_S(A_I^\ell)}(A_{V \setminus S}^\ell)$ is the matrix with columns $\text{proj}_{R_S(A_I^\ell)}(a_j^\ell)$ for all $j \in V \setminus S$, 0 otherwise.

Proof. Setting $B^\ell = A^\ell$ in Lemma A.5, we get $F(I|S) = \sum_{m=1}^{n_\ell+1} \|\text{proj}_{R_S(A_i^\ell)}(A^\ell)w_m^{\ell+1}\|_2^2$. Note that for all $i \in I, j \in S$, a_j^ℓ and $R_S(a_i^\ell)$ are orthogonal, and hence $\text{proj}_{R_S(A_i^\ell)}(a_j^\ell) = 0$, by optimality conditions (see proof of Lemma A.4). It follows then that $F(I|S) = \sum_{m=1}^{n_\ell+1} \|\text{proj}_{R_S(A_i^\ell)}(A_{V \setminus S}^\ell)w_m^{\ell+1}\|_2^2$. In particular, if $I = \{i\}$, $F(i|S) = \sum_{m=1}^{n_\ell+1} \|\text{proj}_{R_S(a_i^\ell)}(A_{V \setminus S}^\ell)w_m^{\ell+1}\|_2^2$. \square

Proposition A.7. *Given $S \subseteq V_\ell$ such that $|S| \leq k$, $i \notin S$, let $\text{proj}_S(a_j^\ell) = A_S^\ell x^S(a_j^\ell)$. Assuming $\tilde{F}(S), \text{proj}_S(b_j^\ell), x^S(b_j^\ell)$ for all $j \notin S$, and $x^S(a_j^\ell)$ for all $j \in V_\ell$, were computed in the previous iteration, we can compute $\tilde{F}(S+i), \text{proj}_{S+i}(b_j^\ell), x^{S+i}(b_j^\ell)$ for all $j \notin (S+i)$, and $x^{S+i}(a_j^\ell)$ for all $j \in V_\ell$, in*

$$O(n_\ell \cdot (n_{\ell+1} + n + k)) \text{ time.}$$

Computing the optimal weights in Eq. (4) at the end of GREEDY can then be done in $O(k \cdot n_\ell \cdot n_{\ell+1})$ time.

Proof. By Lemma A.5, we can update the function value using

$$\tilde{F}(i|S) = \sum_{m=1}^{n_{\ell+1}} \|\text{proj}_{R_S(b_i^\ell)}(A^\ell)w_m^{\ell+1}\|_2^2 = \sum_{m=1}^{n_{\ell+1}} \left(\sum_{j \in V_\ell} \gamma^{S,i}(a_j^\ell) w_{jm}^{\ell+1} \right)^2 \|R_S(b_i^\ell)\|_2^2.$$

This requires $O(n)$ to compute $R_S(b_i^\ell)$ and its norm, $O(n_\ell \cdot n)$ to compute $\gamma^{S,i}(a_j^\ell) = \frac{R_S(b_i^\ell)^\top a_j^\ell}{\|R_S(b_i^\ell)\|_2^2}$ for all $j \in V_\ell$, and an additional $O(n_\ell \cdot n_{\ell+1})$ to finally evaluate $\tilde{F}(S+i)$. We also need $O(|S^c| \cdot n)$ to update $\text{proj}_{S \cup i}(b_j) = \text{proj}_S(b_j) + \text{proj}_{R_S(b_i^\ell)}(b_j)$ (by Lemma A.4), using $\text{proj}_{R_S(b_i^\ell)}(b_j) = R_S(b_i^\ell) \gamma^{S,i}(b_j^\ell)$ for all $j \in V_\ell \setminus S \cup i$, $O(|S^c| \cdot |S|)$ to update $x^{S \cup i}(b_j^\ell) = x^S(b_j^\ell) + (\mathbf{1}_i - x^S(b_i)) \gamma^{S,i}(b_j^\ell)$ (by Lemma A.4) for all $j \in V_\ell \setminus S \cup i$, and $O(n_\ell \cdot |S|)$ to update $x^{S \cup i}(a_j^\ell) = x^S(a_j^\ell) + (\mathbf{1}_i - x^S(b_i)) \gamma^{S,i}(a_j^\ell)$ for all $j \in V_\ell$. So in total, we need $O(n_\ell \cdot (n_{\ell+1} + n + k))$. Computing the new weights $\tilde{W}^{\ell+1} = x^S(A^\ell)W^{\ell+1}$ at the end can be done in $O(n_\ell \cdot n_{\ell+1} \cdot |S|) = O(n_\ell \cdot n_{\ell+1} \cdot k)$. \square

Proposition 3.2. *Given $S \subseteq V_\ell$ such that $|S| \leq k$, $i \notin S$, let $\text{proj}_S(a_j^\ell) = A_S^\ell x^S(a_j^\ell)$ be the projection of a_j^ℓ onto the column space of A_S^ℓ , $R_S(a_i^\ell) = a_i^\ell - \text{proj}_S(a_i^\ell)$ and $\text{proj}_{R_S(a_i^\ell)}(a_j^\ell) \in \arg \min_{z=R_S(a_i^\ell)\gamma, \gamma \in \mathbb{R}} \|a_j^\ell - z\|_2^2$ the corresponding residual and the projection of a_j^ℓ onto it. We can write*

$$F(i|S) = \sum_{m=1}^{n_{\ell+1}} \|\text{proj}_{R_S(a_i^\ell)}(A_{V \setminus S}^\ell)w_m^{\ell+1}\|_2^2,$$

where $\text{proj}_{R_S(a_i^\ell)}(A_{V \setminus S}^\ell)$ is the matrix with columns $\text{proj}_S(a_j^\ell)$ for all $j \notin S$, 0 otherwise. Assuming $F(S), \text{proj}_S(a_j^\ell)$ and $x^S(a_j^\ell)$ for all $j \notin S$ were computed in the previous iteration, we can compute $F(S+i), \text{proj}_{S+i}(a_j^\ell)$ and $x^{S+i}(a_j^\ell)$ for all $j \notin (S+i)$ in

$$O(n_\ell \cdot (n_{\ell+1} + n + k)) \text{ time.}$$

Computing the optimal weights in Eq. (4) at the end of GREEDY can then be done in $O(k \cdot n_\ell \cdot n_{\ell+1})$ time.

Proof. The proof follows from Lemma A.6 and A.4 in the same way as in Proposition A.7. \square

Proposition 3.2 and Proposition A.7 apply also to G and \tilde{G} respectively, since $|M(i)| = O(1)$.

B Stronger notion of approximate submodularity

In this section, we show that F and \tilde{F} satisfy stronger properties than the weak submodularity discussed in Section 3.1, which lead to a stronger approximation guarantee for GREEDY. These properties do not necessarily hold for G and \tilde{G} .

B.1 Additional preliminaries

We start by reviewing some preliminaries. Recall that a set function F is *submodular* if it has diminishing marginal gains: $F(i | S) \geq F(i | T)$ for all $S \subseteq T$, $i \in V \setminus T$. If $-F$ is submodular, then F is said to be *supermodular*, i.e., F satisfies $F(i | S) \leq F(i | T)$, for all $S \subseteq T$, $i \in V \setminus T$. When F is both submodular and supermodular, it is said to be *modular*.

Relaxed notions of submodularity/supermodularity, called *weak DR-submodularity/supermodularity*, were introduced in (Lehmann et al., 2006) and (Bian et al., 2017), respectively.

Definition B.1 (Weak DR-sub/supermodularity). A set function F is α_k -weakly DR-submodular, with $k \in \mathbb{N}_+$, $\alpha_k > 0$, if

$$F(i|S) \geq \alpha_k F(i|T), \text{ for all } S \subseteq T, i \in V \setminus T, |T| \leq k.$$

Similarly, F is β_k -weakly DR-supermodular, with $k \in \mathbb{N}_+$, $\beta > 0$, if

$$F(i|T) \geq \beta_k F(i|S), \text{ for all } S \subseteq T, i \in V \setminus T, |T| \leq k.$$

We say that F is (α_k, β_k) -weakly DR-modular if it satisfies both properties.

The parameters α_k, β_k characterize how close a set function is to being submodular and supermodular, respectively. If F is non-decreasing, then $\alpha_k, \beta_k \in [0, 1]$, F is submodular (supermodular) if and only if $\alpha_k = 1$ ($\beta_k = 1$) for all $k \in \mathbb{N}_+$, and modular if and only if both $\alpha_k = \beta_k = 1$ for all $k \in \mathbb{N}_+$. The notion of weak DR-submodularity is a stronger notion of approximate submodularity than weak submodularity, as $\gamma_{S,k} \geq \alpha_{|S|+k-1}$ for all $S \subseteq V, k \in \mathbb{N}_+$ (El Halabi et al., 2018, Prop. 8). This implies that GREEDY achieves a $(1 - e^{-\alpha_{2k-1}})$ -approximation when F is α_{2k-1} -weakly DR-submodular.

A stronger approximation guarantee can be obtained with the notion of *total curvature* introduced in (Sviridenko et al., 2017), which is a stronger notion of approximate submodularity than weak DR-modularity.

Definition B.2 (Total curvature). Given a set function F , we define its total curvature c_k where $k \in \mathbb{N}_+$, as

$$c_k = 1 - \min_{|S| \leq k, |T| \leq k, i \in V \setminus T} \frac{F(i|S)}{F(i|T)}.$$

Note that if F has total curvature c_k , then F is $(1 - c_k, 1 - c_k)$ -weakly DR-modular. Given a non-decreasing function F with total curvature c_k , the GREEDY algorithm is guaranteed to return a solution $F(\hat{S}) \geq (1 - c_k) \max_{|S| \leq k} F(S)$ (Sviridenko et al., 2017, Theorem 6).

B.2 Approximate modularity of reweighted input change

The reweighted input change objective F is closely related to the column subset selection objective (the latter is a special case of F where $W^{\ell+1}$ is the identity matrix), whose total curvature was shown to be related to the condition number of A^ℓ in (Sviridenko et al., 2017).

We show in Propositions B.3 and B.4 that the total curvatures of \tilde{F} and F are related to the condition number of $A^\ell W^{\ell+1}$.

Proposition B.3. *Given $k \in \mathbb{N}_+$, \tilde{F} is a normalized non-decreasing α_k -weakly DR-submodular function, with*

$$\alpha_k \geq \frac{\min_{\|z\|_2=1} \|(A^\ell W^{\ell+1})^\top z\|_2^2}{\max_{\|z\|_2=1} \|(A^\ell W^{\ell+1})^\top z\|_2^2}.$$

Moreover, if any collection of $k+1$ columns of B^ℓ are linearly independent, \tilde{F} is also α_k -weakly DR-supermodular and has total curvature $1 - \alpha_k$.

Proof. We adapt the proof from (Sviridenko et al., 2017, Lemma 6). For all $S \subseteq V, i \in V \setminus S$, we have $F(i|S) = \sum_{m=1}^{n_{\ell+1}} \|\text{proj}_{R_S(b_i^\ell)}(A^\ell w_m^{\ell+1})\|_2^2$ by Lemma A.5. For all $j \notin S$, we have $\text{proj}_{R_S(b_i^\ell)}(a_j^\ell) = R_S(b_i^\ell) \frac{R_S(b_i^\ell)^\top a_j^\ell}{\|R_S(b_i^\ell)\|^2}$ if $\|R_S(b_i^\ell)\| > 0$, and 0 otherwise, by optimality conditions. Hence, we can write for all i such that $\|R_S(b_i^\ell)\| > 0$,

$$\tilde{F}(i|S) = \sum_{m=1}^{n_{\ell+1}} \|\text{proj}_{R_S(b_i)}(A^\ell) w_m^{\ell+1}\|_2^2 = \sum_{m=1}^{n_{\ell+1}} \left\| \sum_{j \in V} w_{jm}^{\ell+1} R_S(b_i^\ell) \frac{R_S(b_i^\ell)^\top a_j^\ell}{\|R_S(b_i^\ell)\|^2} \right\|_2^2 = \|(A^\ell W^{\ell+1})^\top \frac{R_S(b_i^\ell)}{\|R_S(b_i^\ell)\|}\|_2^2$$

Hence,

$$\min_{\|z\|_2=1} \|(A^\ell W^{\ell+1})^\top z\|_2^2 \leq \tilde{F}(i|S) \leq \max_{\|z\|_2=1} \|(A^\ell W^{\ell+1})^\top z\|_2^2$$

Let $v_j = x^S(B^\ell)_{ij}$ for $j \in S$, $v_i = -1$, and $z = v/\|v\|_2$, then $\|v\|_2 \geq 1$, and

$$\|R_S(b_i^\ell)\|^2 = \|B^\ell v\|_2^2 \geq \|B^\ell z\|_2^2 \geq \min_{\|z\|_2 \leq 1, \|z\|_0 \leq |S|+1} \|B^\ell z\|_2^2$$

The bound on α_k then follows by noting that $\|R_S(b_i^\ell)\| \geq \|R_T(b_i^\ell)\|$ for all $S \subseteq T$. The rest of the proposition follows by noting that if any collection of $k+1$ columns of B^ℓ are linearly independent, then $\|R_S(b_i^\ell)\| \geq \min_{\|z\|_2 \leq 1, \|z\|_0 \leq k+1} \|B^\ell z\|_2^2 > 0$ for any S such that $|S| \leq k$. \square

As discussed in Section B.1, Proposition B.3 implies that GREEDY achieves a $(1 - e^{-\alpha_{2k-1}})$ -approximation with $\tilde{F}(S)$, where α_k is non-zero if all rows of $A^\ell W^{\ell+1}$ are linearly independent. If in addition any $k+1$ columns of B^ℓ are linearly independent, then GREEDY achieves an α_k -approximation.

Proposition B.4. *Given $k \in \mathbb{N}_+$, F is a normalized non-decreasing α_k -weakly DR-submodular function, with*

$$\alpha_k \geq \frac{\max\{\min_{\|z\|_2=1} \|(A^\ell W^{\ell+1})^\top z\|_2^2, \min_{\|z\|_2=1} \|(W^{\ell+1})^\top z\|_2^2 \min_{\|z\|_2 \leq 1, \|z\|_0 \leq k+1} \|A^\ell z\|_2^2\}}{\max_{\|z\|_2=1} \|(A^\ell W^{\ell+1})^\top z\|_2^2}.$$

Moreover, if any collection of $k+1$ columns of A^ℓ are linearly independent, F is also α_k -weakly DR-supermodular and has total curvature $1 - \alpha_k$.

Proof. Setting $B^\ell = A^\ell$ in Lemma A.5, we get

$$\alpha_k \geq \frac{\min_{\|z\|_2=1} \|(A^\ell W^{\ell+1})^\top z\|_2^2}{\max_{\|z\|_2=1} \|(A^\ell W^{\ell+1})^\top z\|_2^2}$$

To obtain the second lower bound, we note that by Lemma A.6 we can write for all $S \subseteq V, i \in V \setminus S$ such that $\|R_S(a_i^\ell)\| > 0$,

$$F(i|S) = \sum_{m=1}^{n_{\ell+1}} \left(\sum_{j \notin S} w_{jm}^{\ell+1} \frac{R_S(a_i^\ell)^\top a_j^\ell}{\|R_S(a_i^\ell)\|^2} \right)^2 \|R_S(a_i^\ell)\|_2^2 = \|(W^{\ell+1})^\top (A_{V \setminus S}^\ell)^\top \frac{R_S(a_i^\ell)}{\|R_S(a_i^\ell)\|}\|_2^2 \|R_S(a_i^\ell)\|_2^2.$$

Note that $\gamma^{S,i}(a_i) = \frac{R_S(a_i^\ell)^\top a_i^\ell}{\|R_S(a_i^\ell)\|^2} = 1$ by optimality conditions ($R_S(a_i^\ell)^\top a_i^\ell = R_S(a_i^\ell)^\top (R_S(a_i^\ell) + A_S x^S(a_i)) = \|R_S(a_i^\ell)\|_2^2$; see proof of Lemma A.4), hence $\|(A_{V \setminus S}^\ell)^\top \frac{R_S(a_i^\ell)}{\|R_S(a_i^\ell)\|}\|_2^2 \geq 1$ and $\|(W^{\ell+1})^\top (A_{V \setminus S}^\ell)^\top \frac{R_S(a_i^\ell)}{\|R_S(a_i^\ell)\|}\|_2^2 \geq \min_{\|z\|_2=1, \|z\|_0 \leq |V \setminus S|} \|(W^{\ell+1})^\top z\|_2^2$. Let $v_j = x^S(A^\ell)_{ij}$ for $j \in S$, $v_i = -1$, and $z = v/\|v\|_2$, then $\|v\|_2 \geq 1$, and

$$\|R_S(a_i^\ell)\|^2 = \|A^\ell v\|_2^2 \geq \|A^\ell z\|_2^2 \geq \min_{\|z\|_2 \leq 1, \|z\|_0 \leq |S|+1} \|A^\ell z\|_2^2$$

We thus have $F(i|S) \geq \min_{\|z\|_2=1, \|z\|_0 \leq |V \setminus S|} \|(W^{\ell+1})^\top z\|_2^2 \min_{\|z\|_2 \leq 1, \|z\|_0 \leq |S|+1} \|A^\ell z\|_2^2$.

The bound on α_k then follows by noting that $\|R_S(a_i^\ell)\| \geq \|R_T(a_i^\ell)\|$ for all $S \subseteq T$. The rest of the proposition follows by noting that if any collection of $k+1$ columns of A^ℓ are linearly independent, then $\|R_S(a_i^\ell)\| \geq \min_{\|z\|_2 \leq 1, \|z\|_0 \leq k+1} \|A^\ell z\|_2^2 > 0$ for any S such that $|S| \leq k$. \square

As discussed in Section B.1, Proposition B.3 implies that GREEDY achieves a $(1 - e^{-\alpha_{2k-1}})$ -approximation with $F(S)$, where α_k is non-zero if all rows of $A^\ell W^{\ell+1}$ are linearly independent. Moreover, if any $k+1$ columns of A^ℓ are linearly independent and all rows of $W^{\ell+1}$ are linearly independent, then GREEDY achieves an α_k -approximation.

It is worth noting that if $W^{\ell+1}$ is identity matrix, i.e., F is the column subset selection function, then Proposition B.4 implies a stronger result than (Sviridenko et al., 2017, Lemma 6) for some cases. In particular, we have

$$\alpha_k \geq \frac{\max\{\min_{\|z\|_2 \leq 1, \|z\|_0 \leq k+1} \|A^\ell z\|_2^2, \min_{\|z\|_2=1} \|(A^\ell)^\top z\|_2^2\}}{\max_{\|z\|_2 \leq 1} \|(A^\ell)^\top z\|_2^2},$$

which implies that F is weakly DR-submodular if any k columns of A^ℓ are linearly independent, or if all rows of A^ℓ are linearly independent.

C Experimental setup

Random seeds: 42, 43, 44, 45, 46

Pruning setup:

- Number of Batches: 4
- Batch size: 128
- Values used for fraction of prunable neurons/channels kept: $\alpha \in \{0.01, 0.05, 0.075, 0.1, 0.25, 0.5, 0.75, 1.0\}$
- Values used for per-layer fraction selection: $\alpha_\ell \in \{0.01, 0.05, 0.075, 0.1, 0.15, 0.2, 0.25, 0.3, 0.35, 0.4, 0.45, 0.5, 0.55, 0.6, 0.65, 0.7, 0.75, 0.8, 0.85, 0.9, 0.95, 1.0\}$

Training setup for LeNet on MNIST:

- Batch size: 128
- Epochs for pre-training: 200
- Epochs for fine-tuning: 10
- Optimizer for pre-training: SGD with Nesterov momentum 0.9
- Optimizer for fine-tuning: Adam with $\beta_1 = 0.9$ and $\beta_2 = 0.99$
- Initial learning rate: 1×10^{-3}
- Learning rate schedule: Fixed

Training setup for VGG11 on CIFAR10:

- Batch size: 128
- Epochs for pre-training: 200
- Epochs for fine-tuning: 20
- Optimizer for pre-training: Adam with $\beta_1 = 0.9$ and $\beta_2 = 0.99$
- Optimizer for fine-tuning: Adam with $\beta_1 = 0.9$ and $\beta_2 = 0.99$
- Initial learning rate: 1×10^{-3}
- Weight decay: 5×10^{-4}
- Learning rate schedule for pre-training: learning rate dropped by 0.1 at epochs 100 and 150
- Learning rate schedule for fine-tuning: learning rate dropped by 0.1 at epochs 10 and 15

The setup for fine-tuning VGG11 on MNIST, both before and after pruning, is the same one used for pre-training it on CIFAR10, as outline above, but with 10 epochs.

D Additional results

D.1 Other channel pruning results

In this section, we present additional channel pruning results. We first consider the same setup as in Section 6.2, but we prune four of the convolutional layers of VGG11 (features.0, 8, 11, 18), instead of two. In these layers, $n_\ell = 64, 256, 512, 512$, and $n_{\ell+1} = 128, 256, 512, 512$, respectively. Figure 5 shows the top-1 and top-5 accuracy for different fractions of prunable channels kept. The sequential variants of our method are again among the best performing methods. Their performance is matched only by LAYERWEIGHTNORM, which performs worst in all other settings we consider. LAYERINCHANGE performs slightly worst here.

Next we consider the same transfer learning setup as in Section 6.2, but we again prune four of the convolutional layers of VGG11-MNIST (features.0, 8, 11, 18), instead of two. Figure 6 shows the top-1 and top-5 accuracy for different fractions of prunable channels kept. All variants of our method perform similarly here, and are again among the best performing methods. Their performance is matched only by ACTGRAD, which also performs worst in all other settings we consider.

Note that the results in Figures 5 and 6 are better than their counterparts in Figures 3 and 4, where we prune less layers. This is because the current four layers being pruned include a late

layer (features.18), which can be pruned almost entirely without affecting the model’s accuracy. Hence, given a global budget, the per-layer budget selection assigns most of the budget to the first three layers, which yield better results than when we are forced to prune a large fraction of early layers. We further highlight the effect of per-layer budget selection on performance in Appendix F. Reweighting again helps all methods in both pruning tasks considered here.

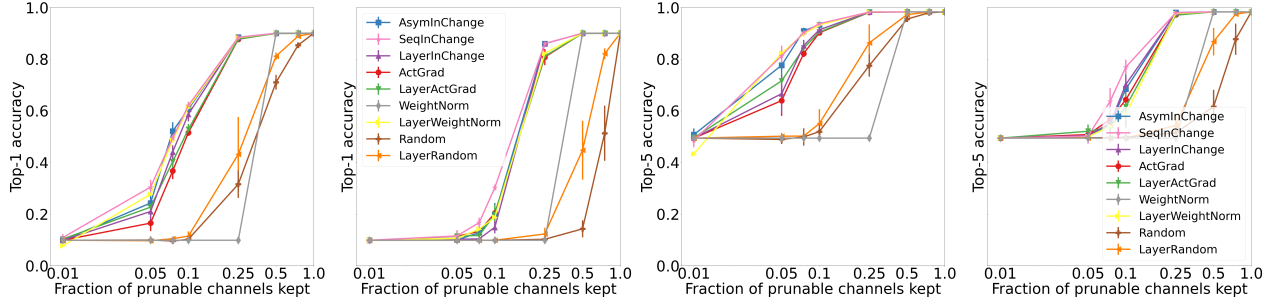


Figure 5: Top-1 and Top-5 Accuracy of different channel pruning methods on CIFAR10, after pruning four convolution layers in VGG11 model, with different fractions of remaining channels (in log-scale), with (left) and without (right) reweighting.

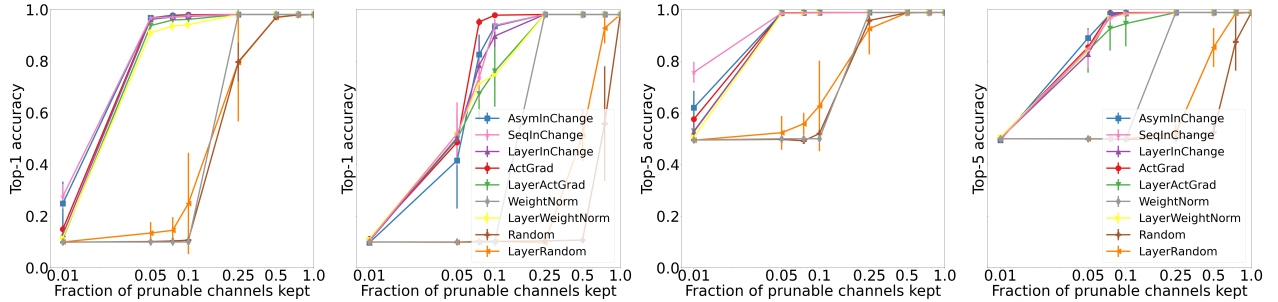


Figure 6: Top-1 and Top-5 Accuracy of different channel pruning methods on MNIST, after pruning four convolution layers in VGG11-MNIST model, with different fractions of remaining channels (in log-scale), with (left) and without (right) reweighting.

D.2 Top-5 Accuracy results of Section 6

For completeness, we report in Figures 7, 8, 9, and 10 the top-5 accuracy results of the experiments presented in Section 6.

E Effect of fine-tuning

In this section, we study the effect of fine-tuning. To that end, we report in Figures 11, 12, 13, 14, 15, and 16, the top-1 and top-5 accuracy results of all the pruning tasks considered in Section 6 and Appendix D.1 after fine-tuning.

We fine-tune for 10 epochs in all MNIST experiments, and for 20 epochs for all CIFAR-10 experiments. We do not fine-tune at fraction 1 (i.e., when nothing is pruned). As expected, fine-tuning provides a significant boost in performance to all methods, even more than reweighting. Though reweighting still improves performance even when fine-tuning is used. However, in some cases, it can lead to

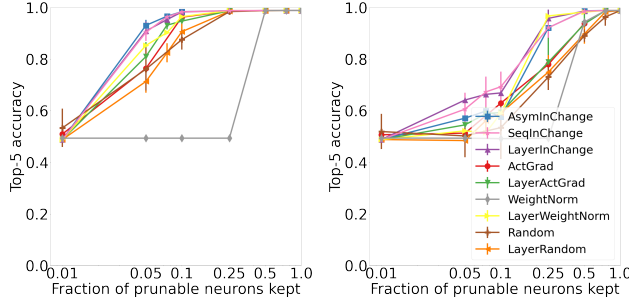


Figure 7: Top-5 Accuracy of different neuron pruning methods on MNIST, after pruning two linear layers in LeNet model, with different fractions of remaining neurons, with (left) and without (right) reweighting.

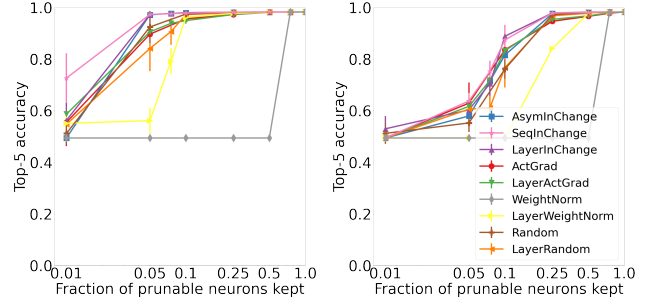


Figure 8: Top-5 Accuracy of different neuron pruning methods on CIFAR10, after pruning two linear layers in VGG11 model, with different fractions of remaining neurons (in log-scale), with (left) and without (right) reweighting.

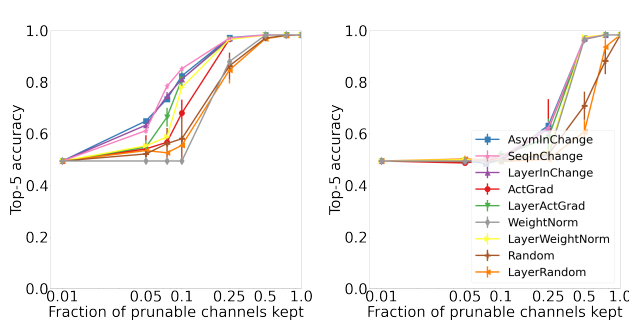


Figure 9: Top-5 Accuracy of different channel pruning methods on CIFAR10, after pruning two convolution layers in VGG11 model, with different fractions of remaining channels (in log-scale), with (left) and without (right) reweighting.

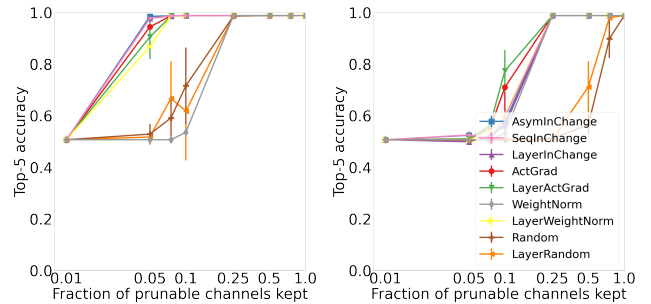


Figure 10: Top-5 Accuracy of different channel pruning methods on MNIST, after pruning two convolution layers in VGG11-MNIST model, with different fractions of remaining channels (in log-scale), with (left) and without (right) reweighting.

slightly worst results, as in Figure 13; our methods have slightly lower accuracy when reweighting is used, for fractions larger than 0.05. When both fine-tuning and reweighting are used, all methods perform fairly similarly, except WEIGHTNORM which consistently performs worst, and random pruning methods which sometimes perform worst (see Figures 13 and 14).

F Importance of per-layer budget selection

In this section, we study the effect of per-layer budget selection on accuracy. To that end, we consider the same setup as in Section 6.2, but instead of pruning the first two convolutional layers of VGG11, we prune the first and second to last layer. Since the second to last layer in VGG11 (features.22) has little effect on accuracy when pruned, we expect the choice of how the global budget is distributed on the two layers to have a significant impact on performance. Figure 17 shows the top-1 accuracy for different fractions of prunable channels kept, when the per-layer budget selection from Section 5.2 is used, while Figure 18 shows the results when equal fractions of channels kept are used in each layer. As expected, all layerwise methods perform much more poorly with equal

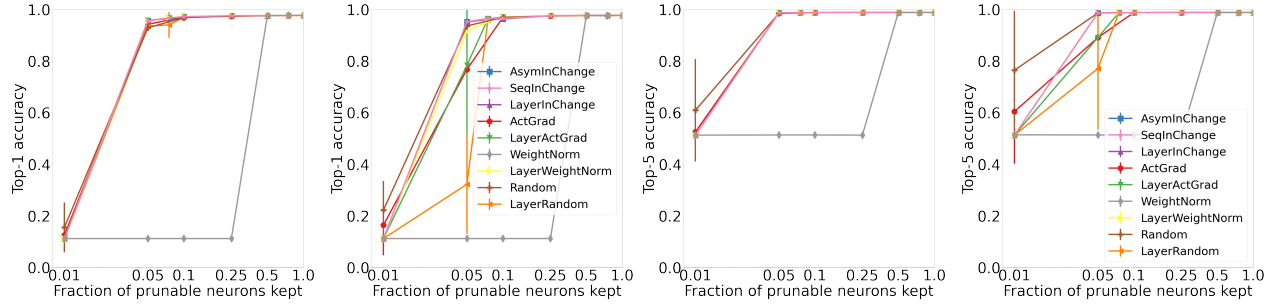


Figure 11: Top-1 and Top-5 Accuracy of different neuron pruning methods on MNIST, after pruning two linear layers in LeNet model, and fine-tuning for 10 epochs, with different fractions of remaining neurons (in log-scale), with (left) and without (right) reweighting.

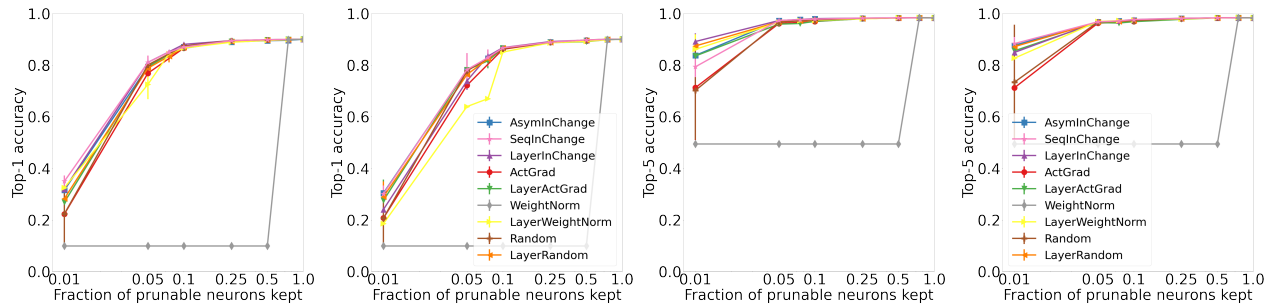


Figure 12: Top-1 and Top-5 Accuracy of different neuron pruning methods on CIFAR10, after pruning two linear layers in VGG11 model, and fine-tuning for 20 epochs, with different fractions of remaining neurons (in log-scale), with (left) and without (right) reweighting.

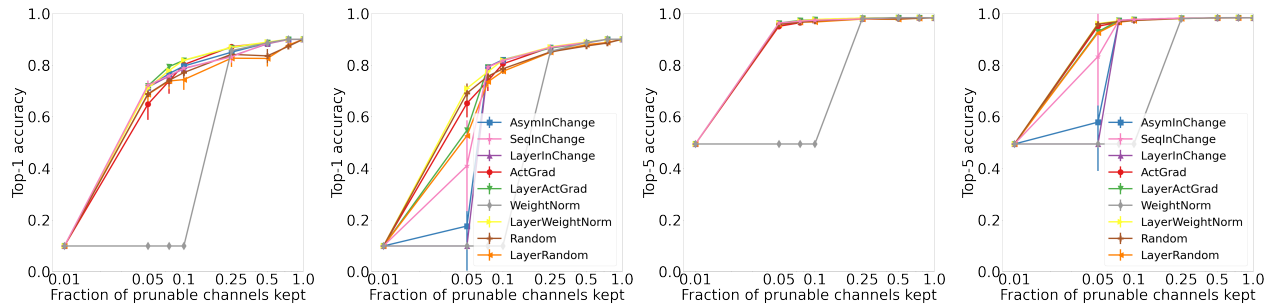


Figure 13: Top-1 and Top-5 Accuracy of different channel pruning methods on CIFAR10, after pruning two convolution layers in VGG11 model, and fine-tuning for 20 epochs, with different fractions of remaining channels (in log-scale), with (left) and without (right) reweighting.

per-layer fractions, both with and without fine-tuning. Though, the difference is less drastic when fine-tuning is used.

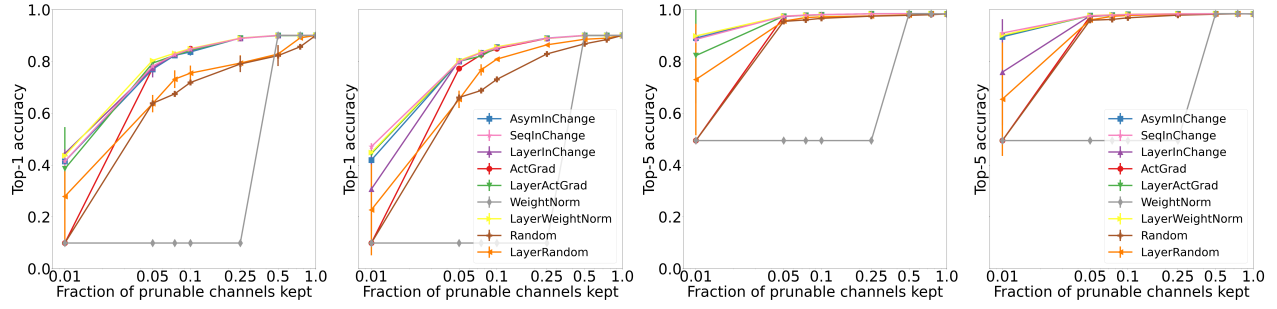


Figure 14: Top-1 and Top-5 Accuracy of different channel pruning methods on CIFAR10, after pruning four convolution layers in VGG11 model, and fine-tuning for 20 epochs, with different fractions of remaining channels (in log-scale), with (left) and without (right) reweighting.

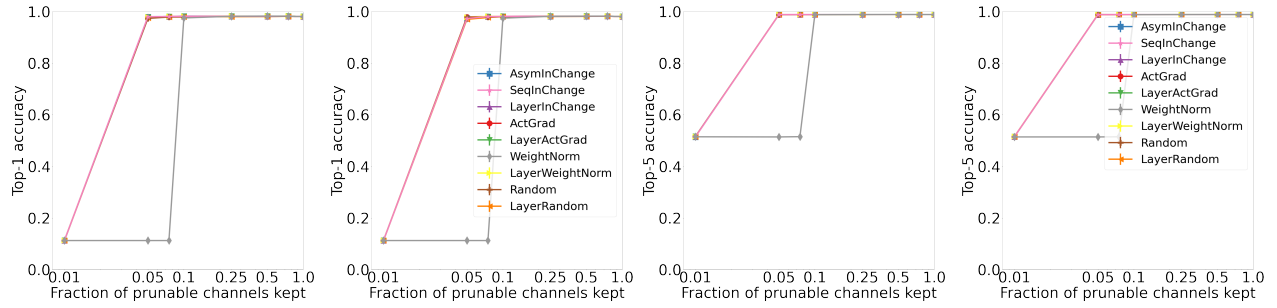


Figure 15: Top-1 and Top-5 Accuracy of different channel pruning methods on MNIST, after pruning two convolution layers in VGG11-MNIST model, and fine-tuning for 10 epochs, with different fractions of remaining channels (in log-scale), with (left) and without (right) reweighting.

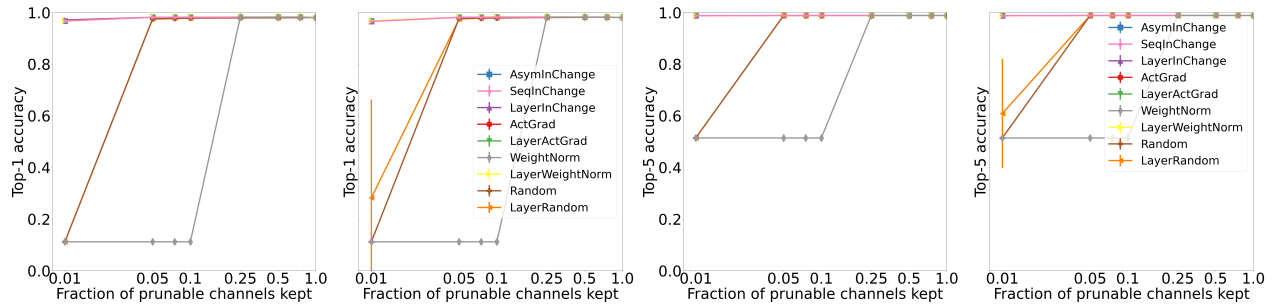


Figure 16: Top-1 and Top-5 Accuracy of different channel pruning methods on MNIST, after pruning four convolution layers in VGG11-MNIST model, and fine-tuning for 10 epochs, with different fractions of remaining channels (in log-scale), with (left) and without (right) reweighting.

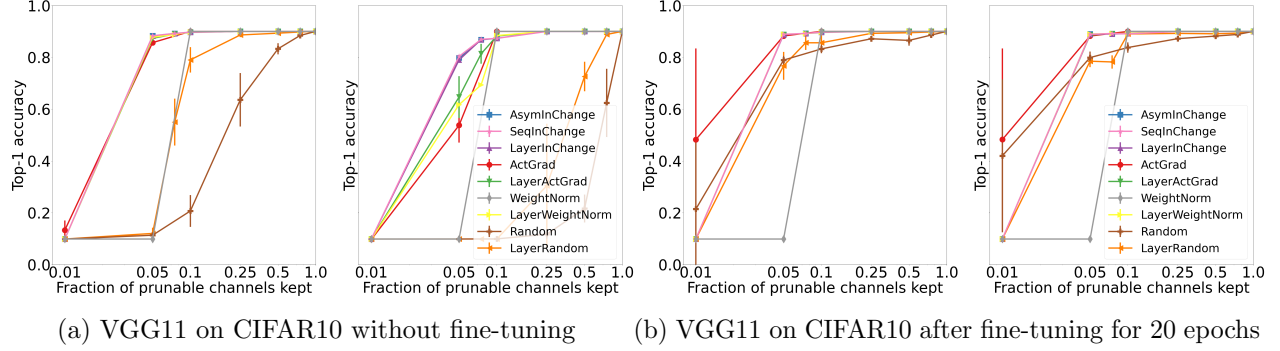


Figure 17: Top-1 Accuracy of different channel pruning methods on CIFAR10, after pruning the first and second to last convolution layers in VGG11 model, with different fractions of remaining channels (in log-scale), with (a, b - left) and without (a, b - right) reweighting, with per-layer fractions selected using the selection method discussed in Section 5.2.

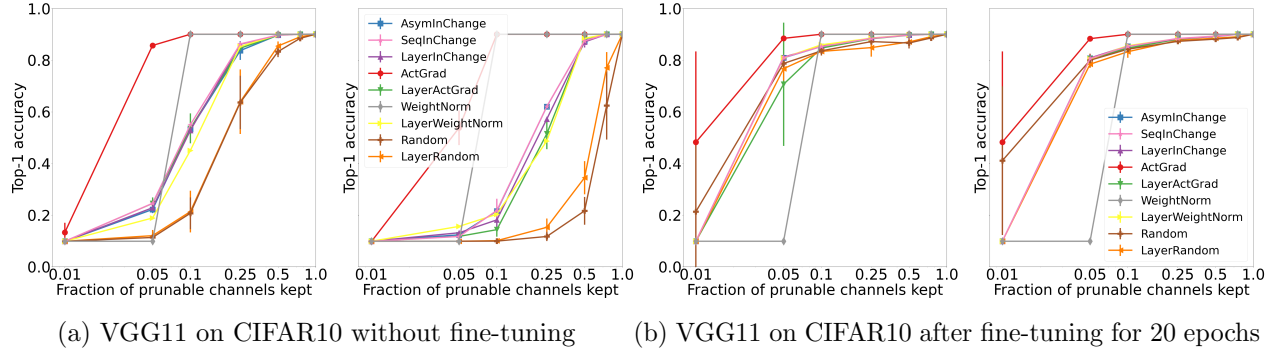


Figure 18: Top-1 Accuracy of different channel pruning methods on CIFAR10, after pruning the first and second to last convolution layers in VGG11 model, with different fractions of remaining channels (in log-scale), with (a, b - left) and without (a, b - right) reweighting, with equal per-layer fractions.