

On Causal Inference for Data-free Structured Pruning

Martin Ferianc¹, Anush Sankaran², Olivier Mastropietro³, Ehsan Saboori³, and Quentin Cappart⁴

¹Department of Electronic and Electrical Engineering, University College London, London, UK WC1E 7JE, martin.ferianc.19@ucl.ac.uk

²Independent Researcher, anush.sankaran@gmail.com

³Deeplite Inc., Montreal, QC, Canada H3C 2G9, {olivier, ehsan}@deeplite.ai

⁴Department of Computer Engineering and Software Engineering, Polytechnique Montréal, Montreal, QC, Canada H3T 1J4, quentin.cappart@polymtl.ca

Abstract

Neural networks (NNs) are making a large impact both on research and industry. Nevertheless, as NNs' accuracy increases, it is followed by an expansion in their size, required number of compute operations and energy consumption. Increase in resource consumption results in NNs' reduced adoption rate and real-world deployment impracticality. Therefore, NNs need to be compressed to make them available to a wider audience and at the same time decrease their runtime costs. In this work, we approach this challenge from a causal inference perspective, and we propose a scoring mechanism to facilitate structured pruning of NNs. The approach is based on measuring mutual information under a maximum entropy perturbation, sequentially propagated through the NN. We demonstrate the method's performance on two datasets and various NNs' sizes, and we show that our approach achieves competitive performance under challenging conditions.

Introduction

Neural networks (NNs) have been successfully deployed in several applications, such as computer vision (Zhai et al. 2021) or natural language processing (Brown et al. 2020). The accuracy in these tasks increases with ongoing development, however, so does the model size and power consumption (Guo et al. 2019). Novel NNs' increasing size, complexity, and energy demands limit their deployment to low-cost compute platforms.

On one hand, hardware optimizations have been proposed to ease the deployment of demanding NNs, but these are usually targeting on a specific pair of NN architecture and a hardware platform (Chen et al. 2020). On the other hand, software optimizations, such as structured pruning (Hoeffler et al. 2021) have been proposed, which can be applied to a variety of NNs to make them smaller. In structured pruning, a certain neuron in the NN is removed completely, saving computation time, reducing their memory, and energy consumption across multiple hardware platforms. For example, by compression and subsequent fine-tuning, ResNet-18's compute operations' count can be reduced by $7\times$ and its memory foot-print by $4.5\times$ (Sankaran et al. 2021). However, most of the methods for structured pruning are based on some heuristics which are not conclusive in the context of structurally examining NNs. Additionally, pruning might require access to the original data on which the NN was trained, for further fine-tuning.

To address these challenges, in this research work we propose first steps towards a data-free approach to structured pruning which is facilitated through causal inference. In this approach, we evaluate the importance of neurons by measuring mutual information (MI) under a maximum entropy perturbation (MEP) propagated through the NN. We demonstrate our performance and generalizability on various fully-connected NN architectures on two datasets. In our evaluation, the proposed method produces marginal improvements over the existing work and it hopes to pave new directions for research into causal inference applied to optimize NNs.

Related Work

Causal Inference and Information Bottleneck

Our method is primarily inspired by the work of Mattsson, Michaud, and Hoel (2020) who have proposed a suite of metrics based on information theory; to quantify and track changes in the causal structure of NNs. They introduced, the notion of *effective information* which is the MI between layer input and output following a local layer-wise MEP. We build on this notion and introduce several changes. First, we sample a random intervention only at the input of the NN and we measure the MI with respect to the output of the previous layer obtained by propagating the intervention throughout the net. Second, we pick a different maximum entropy distribution, a Gaussian instead of uniform, that more closely reflects real-world data. Third, we combine the different measurements per neural connection, and we use them to score the likeliness of that neuron for structured pruning. Additional concept related to this work is *information bottleneck* (Tishby, Pereira, and Bialek 2000), which measures MI with respect to the information plane and propagating data through the network. They have shown that at a certain point in the NN, the NN minimizes MI between input and output. In this work, we contemplate that if Gaussian noise is propagated through the net, the neurons which maximize the MI between input and output are preferred with respect to generalization on the test data.

Structured Data-free Pruning

Hoeffler et al. (2021) completed a comprehensive survey of NN pruning methods and in this work we will focus on those that: do not require data to prune and focus on struc-

tured pruning. Srinivas and Babu (2015) proposed a *data-free pruning* (DFP) method that examines the importance of different neurons based on their similarity through the magnitude of their weights. Their method iteratively examines, prunes and updates this similarity along with the weights of the NN. Mussay et al. (2019) proposed a data-independent way for pruning neurons in an NN through coreset approximation in its preceding layers. (Wang et al. 2019) developed *correlation-based pruning* (COP), which can detect the redundant neurons efficiently through removing the ones which are the most correlated with the others. Moreover, Narendra et al. (2018) developed a method to reason over NN as a structured causal model, nevertheless, this method is data-bound. Lastly, Ganesh, Corso, and Sekeh (2021) introduced *MINT* which is based on measuring MI with respect to data, however, without considering the notion of causal inference or MEP. With respect to the related work, our method also wants to appeal to users who seek data-free pruning methods, potentially due to privacy-related constraints. Nevertheless, our method differs by avoiding the usage of heuristics, such as: the weight magnitude or correlation. It relies on examining the causal structure in the NN, rather than deterministic heuristics.

Method

Without using the train data, NNs and their internal connectivity have been often described through heuristics, such as correlations and magnitude of the connecting weights for the individual neurons (Hoeffler et al. 2021). As the depth and width of the NNs increase, these metrics become less transparent and less interpretable in feature space. Additionally, there is no clear link between these heuristics and the causal structure by which the NN makes decisions and generalizes beyond train data. As it has already been argued, generalizability must be a function of a NN’s causal structure since it reflects how the trained NN responds to unseen or even not-yet-defined future inputs (Mattsson, Michaud, and Hoel 2020). Therefore, from a causal perspective, the neurons which are identified to be more important in the architecture should be preserved and the ones that are identified less important could be removed. This paradigm paves the way for observing the causal structure, identifying important connections and subsequent structured pruning in NNs, replacing heuristics, to achieve better generalization.

In this work, we propose a perturbation-based approach to study the causal structure of the NN which enables us to quantify the significance of each neuron in the NN. The method performs an intervention $do(\mathbf{x}_0)$ at the input level of the NN. However, instead of choosing a single type of intervention we opt for a maximum entropy distribution - a Gaussian distribution, which covers the space of all potential interventions with a fixed variance, and it is used to sample the input $\mathbf{x}_0 \sim \mathcal{N}(0, 1)$. Subsequently, the output of the input layer is propagated to deeper layers $i; i \in 1, \dots, L$ up to the entire depth L of the net. Under this perturbation, we propose to measure MI between the input and output pairs of layers, to measure the strength of their causal interactions (Mattsson, Michaud, and Hoel 2020). Unlike the standard MI, which is a measure of correlation, all mutual bits

Algorithm 1: Inference of MI scores for structured pruning.

```

1: Phase 1: Record and normalize inputs and outputs
2: Sample  $S$  samples for NN input  $\mathbf{x}_0 \sim \mathcal{N}(0, 1)$ 
3: Pass  $\mathbf{x}_0$  and cache layers’ outputs  $\mathbf{x}_i; i \in [1, L]$ 
4: Infer range of activations for each layer  $i \in [0, L]$ 
5: Clamp and normalize each  $\mathbf{x}_i; i \in [0, L]$ 
6: Phase 2: Infer input-output mutual information
7: Initialize MI as an empty list
8: for  $(\mathbf{x}_{i-1} \in \mathbb{R}^{S \times N}, \mathbf{x}_i \in \mathbb{R}^{S \times M}); i \in [1, L]$  do
9:   Initialize  $\mathbf{m}i_i \in \mathbb{R}^{N \times M}$  as zeros
10:  for  $n$  in range  $N$  do
11:    for  $m$  in range  $M$  do
12:       $\mathbf{h}_i^{n,m}$  = joint histogram for  $\mathbf{x}_{i-1}^n, \mathbf{x}_i^m$ 
13:       $\mathbf{m}i_i[n, m]$  = Mutual information for  $\mathbf{h}_i^{n,m}$ 
14:    end for
15:  end for
16:  Add  $\mathbf{m}i_i$  to MI to be used during structured pruning
17: end for
18: Phase 3: Pruning
19: for layer  $i; i \in [0, L - 1]$  do
20:    $\mathbf{m}i = \text{MI}[i]$ 
21:   Set  $\mathbf{m}i$  indices of already pruned neurons to zero
22:    $\text{Scores} = \text{Sum } \mathbf{m}i$  row-wise
23:   Sort  $\text{Scores}$ 
24:   Remove neuron connections for lowest  $\text{Scores}$ 
25: end for

```

with a noise injection will be caused by that noise (Mattsson, Michaud, and Hoel 2020). We hypothesize that the connections that preserve the most information on average under MEP are the strongest, with the most impact on the generalization performance, and they should be preserved in case of pruning. The MI is measured individually per input-output connection, and later summed for the given output neuron, for computational simplicity. However, the individual computation and summation imply independence with respect to the input connections for a particular neuron. This is a limitation that manifests itself with increasing depth of the net.

The Algorithm 1 summarizes the method. In general, it consists of three phases: 1. recording and normalization of NN’s layers’ inputs and outputs, 2. inference of MI from the records and 3. using the MI to determine which neurons to prune. The first phase starts with propagating the random noise through the trained NN, while caching, clamping and normalizing the outputs of all layers between $[0, 1]$, for the inferred maximums and minimums of their paired activations. In practice, the normalization was used to standardize the relative bin size of the input-output histograms that are being used for the MI estimation.

Using the cached inputs and outputs for all layers $i \in 1, \dots, L$, the algorithm then proceeds to the second phase to estimate per-connection MI for each layer i . For example, if considering an input from the previous layer $\mathbf{x}_{i-1} \in \mathbb{R}^{S \times N}$ with N nodes and S total samples and the output of the current layer, after an activation was applied, $\mathbf{x}_i \in \mathbb{R}^{S \times M}$ with M output nodes, we calculate the MI for a connection between n^{th} and m^{th} neuron and the i^{th}

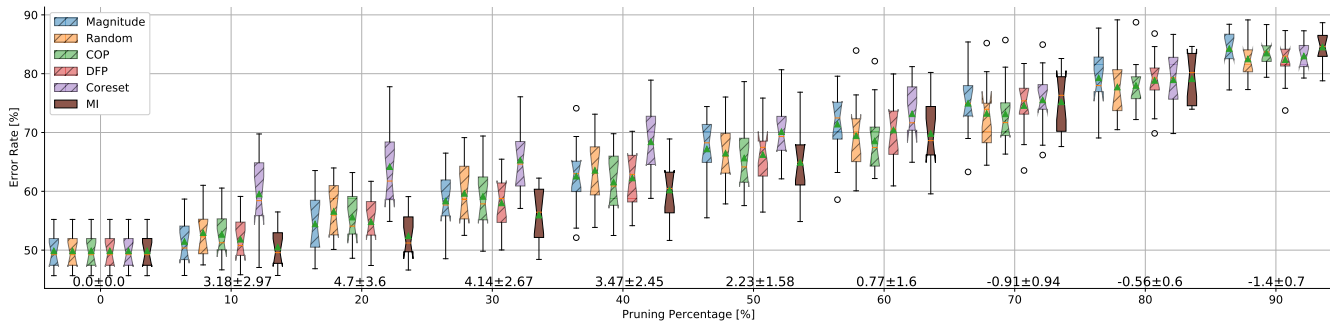


Figure 1: CIFAR-10 test dataset results, each box is aggregation of all 12 networks pruned with respect to the set percentage.

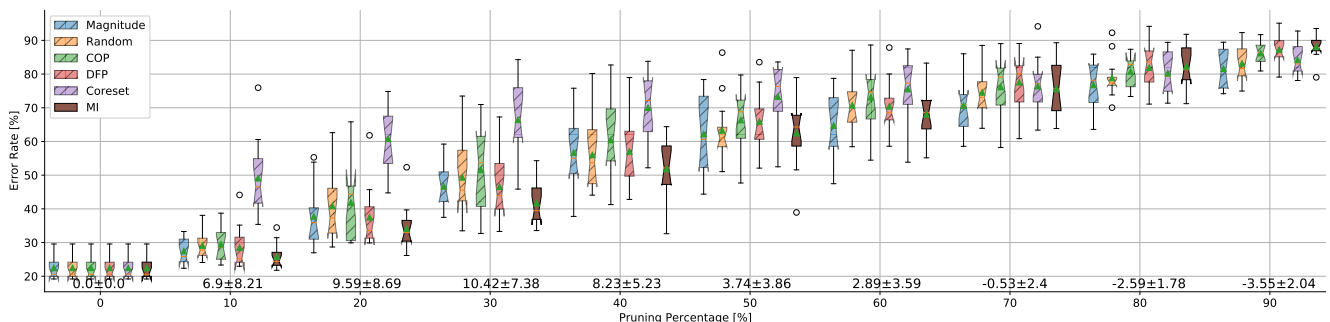


Figure 2: SVHN test dataset results, each box is aggregation of all 12 networks pruned with respect to the set percentage.

layer as $mi_i^{n,m}(\mathbf{x}_{i-1}^n, \mathbf{x}_i^m | do(\mathbf{x}_0))$; $\mathbf{x}_0 \sim \mathcal{N}(0, 1)$, where mi stands for MI estimation based on a histogram (Fraser and Swinney 1986). In detail, the MI for node pair n, m and layer i as $mi_i^{n,m}$ is computed as shown in Equation 1:

$$mi_i^{n,m}(\mathbf{x}_{i-1}^n, \mathbf{x}_i^m | do(\mathbf{x}_0)) = \sum_{u \in \mathbf{x}_{i-1}^n} \sum_{v \in \mathbf{x}_i^m} p_{x_{i-1}^n, x_i^m}(u, v) \log \left(\frac{p_{x_{i-1}^n, x_i^m}(u, v)}{p_{x_{i-1}^n}(u) p_{x_i^m}(v)} \right) \quad (1)$$

The different probability densities $p(\cdot)$ are captured by histograms with $B \times B$ bins with respect to their 2D space, defined by S normalized and cached samples of $\mathbf{x}_{i-1}^n, \mathbf{x}_i^m$ for the given input intervention $do(\mathbf{x}_0)$. This process is repeated with respect to all $n \in N$ for a particular m^{th} node out of M nodes and for every layer $i \in 1, \dots, L$ with the same \mathbf{x}_0 . These scores are pre-computed and cached until the pruning phase.

In the pruning phase, the final significance score for layer i and some node m is determined by summation of the pre-recorded $mi_i^{n,m}$ with respect to all existing connections $n \in N$, where the n^{th} link was not cancelled by the pruning in the previous layer. The scores for all $m \in M$ for a given layer i are then sorted and the neurons with the smallest score are pruned, moving to the next layer. Hence, the overall algorithm focuses on preserving the neurons that were observed with the highest input-output MI under the MEP and we hypothesize that they should have the most impact on the generalization performance of the NN.

Experiments

To validate the proposed method, we performed comprehensive experiments involving two datasets and various network depths and widths. All networks were trained with respect to 200 epochs, learning rate set initially to $1e^{-3}$ and exponentially decayed with an Adam optimizer and weight decay set to $1e^{-4}$. We conducted experiments with respect to CIFAR-10 and SVHN, to vary the complexity of the datasets, without any data augmentations except normalization. For both datasets we trained in total 12 networks, paired with ReLU activations, with $\{1, 2, 3\}$ hidden layers and $\{64, 128, 192, 256\}$ channels for CIFAR-10 or $\{16, 32, 48, 64\}$ channels for SVHN, giving 12 model combinations for each dataset. The models were arguably small, where it can be assumed that each neuron has certain importance and there are no or few inactive neurons. Therefore, the pruning methods need to be careful about scoring the neurons, since removing even a single neuron will affect the algorithmic performance. In terms of pruning, we ask each compared method: magnitude-based (He, Zhang, and Sun 2017), Random, COP, DFP, coreset or ours (MI) to provide a relative importance score for all hidden neurons in an NN. We used publicly available implementations of the respective methods, except DFP which we reimplemented. We adopted a linearly increasing pruning schedule with respect to depth of a layer with some maximum percentage, omitting the input or output layers. For example, if we set the pruning rate to 30% and the network has 2 hidden layers, we would prune 15% of neurons in the 1st hidden layer and

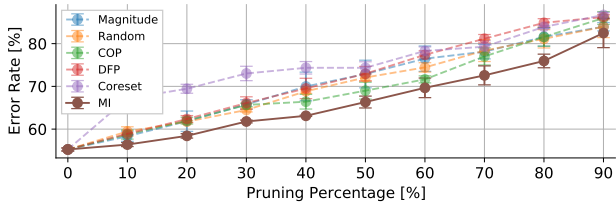


Figure 3: Network with one hidden layer with 64 channels for CIFAR-10.

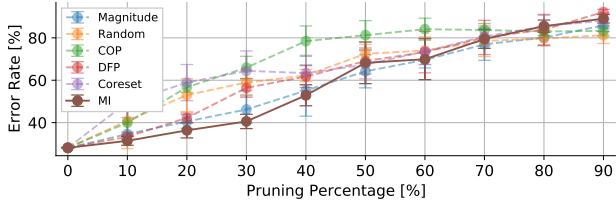


Figure 4: Network with one hidden layer with 16 channels for SVHN.

30% in the 2nd hidden layer depending on the lowest scores given by each method. We used $S = 5000$ samples for MI estimation with $B = 32$.

Aggregated Results

In Figures 1 and 2 we demonstrate the results showing varying error rate across different limiting pruning percentages. Each box represents aggregated results from 12 benchmarked models pruned with respect to the limiting percentage. We chose this form of presentation to demonstrate the versatility of our method and related work across different network depths and widths. As it can be seen with respect to both datasets, our method’s error rates increase less in comparison to the related work across a range of different architectures and pruning percentages, especially in lower pruning percentages. Quantitatively, the benefits of our method are shown at the bottom of the plots. The numbers symbolize the mean and the standard deviation of reduction in average error rate across the 12 architectures and the related work compared to the proposed method. Note that, the evaluated networks were small and every connection can be hypothetically significant to generalization and the error increases irrespective of selected method after setting the bounding percentage to approximately 70%.

Detailed Results

In Figures 3 and 4 we present the results with respect to the smallest and most challenging architectures in our experiments with only one hidden layer. All experiments were repeated 3 times with different random seeds to observe mean and standard deviation for robustness. As it can be seen, MI was able to more concretely identify the significant neurons, resulting in lower average error rates, mainly for CIFAR-10.

Additionally, in Tables 1 and 2 we demonstrate the Spearman correlation and Kendall τ ranking correlation with respect to magnitude-based pruning, which is a well-

	CIFAR-10		SVHN	
	Correlation	Kendall τ	Correlation	Kendall τ
Layer 1	0.86 ± 0.006	0.78 ± 0.01	0.5 ± 0.06	0.3 ± 0.05
Layer 2	0.5 ± 0.1	-0.15 ± 0.04	-0.29 ± 0.05	-0.2 ± 0.03
Layer 3	0.6 ± 0.06	0.02 ± 0.04	-0.08 ± 0.68	-0.3 ± 0.16
Layer 4	0.9 ± 0.0	0.47 ± 0.02	0.74 ± 0.1	0.05 ± 0.1

Table 1: Ranking similarity to magnitude-based score for the deepest and widest network variants.

	CIFAR-10		SVHN	
	Correlation	Kendall τ	Correlation	Kendall τ
Layer 1	0.48 ± 0.07	0.32 ± 0.05	0.11 ± 0.08	0.06 ± 0.12
Layer 2	-0.43 ± 0.02	-0.23 ± 0.05	-0.01 ± 0.43	0.05 ± 0.42

Table 2: Ranking similarity to magnitude-based score for the shallowest and thinnest network variants.

established baseline, to provide deeper insight into the proposed method. As it can be seen, the method is partially correlated to the magnitude of the weights connecting that neuron to the rest of the NN. However, looking simultaneously at the Kendall τ comparing weight magnitude and our score, it can be seen that the overall ranking is completely different. These results demonstrate that causal inference and MI in general are vital for deeper understanding of the structure of the NN and there is only a relatively weak connection to the weights’ magnitude.

Challenging Settings

During the experimentation we made several observations and we also encountered challenging experimental and deployment settings. In initial experiments with respect to MNIST and FashionMNIST, we noticed that the proposed method underperformed. We associated this with respect to atypicality of those datasets, containing predominately zeros resulting in skewed input data distributions, far from the Gaussian assumption. Moreover, we noticed that the performance of the method degrades with increasing the depth of the network. We associated this with respect to the MI estimation, where the independence assumption does not hold and thus MI is extremely challenging to estimate.

Conclusion

In this work, we presented empirical first steps towards a causal inference-based approach for data-free structured NN pruning. We evaluated the proposed methodology with respect to different NN structures on two real-world datasets. Additionally, we detailed positive cases for pruning as well as challenging conditions. In the future work, we aim to extend the current framework with respect to complex networks, specifically convolutional NNs, to promote its practicality in real-world applications. Moreover, we want to investigate the methodology with respect to additional tasks or larger datasets.

Acknowledgements

This work was completed, while Martin Ferianc was an intern and Anush Sankaran was a research scientist at Deeplite Inc. This research was supported by MITACS/IT26487. This funding source had no role in the design of this study and did not have any role during its execution, analyses, interpretation of the data, or decision to submit results. Lastly, we thank ITCI'22 reviewers for feedback that helped us to improve the paper.

References

- Brown, T. B.; Mann, B.; Ryder, N.; Subbiah, M.; Kaplan, J.; Dhariwal, P.; Neelakantan, A.; Shyam, P.; Sastry, G.; Askell, A.; et al. 2020. Language models are few-shot learners. *arXiv preprint arXiv:2005.14165*.
- Chen, Y.; Xie, Y.; Song, L.; Chen, F.; and Tang, T. 2020. A survey of accelerator architectures for deep neural networks. *Engineering*, 6(3): 264–274.
- Fraser, A. M.; and Swinney, H. L. 1986. Independent coordinates for strange attractors from mutual information. *Physical review A*, 33(2): 1134.
- Ganesh, M. R.; Corso, J. J.; and Sekeh, S. Y. 2021. MINT: Deep Network Compression via Mutual Information-based Neuron Trimming. In *Proceedings of the International Conference on Pattern Recognition (ICPR)*, 8251–8258. IEEE.
- Guo, Q.; Chen, S.; Xie, X.; Ma, L.; Hu, Q.; Liu, H.; Liu, Y.; Zhao, J.; and Li, X. 2019. An empirical study towards characterizing deep learning development and deployment across different frameworks and platforms. In *Proceedings of the International Conference on Automated Software Engineering (ASE)*, 810–822. IEEE.
- He, Y.; Zhang, X.; and Sun, J. 2017. Channel pruning for accelerating very deep neural networks. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 1389–1397. IEEE.
- Hoeffler, T.; Alistarh, D.; Ben-Nun, T.; Dryden, N.; and Peste, A. 2021. Sparsity in Deep Learning: Pruning and growth for efficient inference and training in neural networks. *arXiv preprint arXiv:2102.00554*.
- Mattsson, S.; Michaud, E. J.; and Hoel, E. 2020. Examining the causal structures of deep neural networks using information theory. *arXiv preprint arXiv:2010.13871*.
- Mussay, B.; Osadchy, M.; Braverman, V.; Zhou, S.; and Feldman, D. 2019. Data-independent neural pruning via coresets. *arXiv preprint arXiv:1907.04018*.
- Narendra, T.; Sankaran, A.; Vijaykeerthy, D.; and Mani, S. 2018. Explaining deep learning models using causal inference. *arXiv preprint arXiv:1811.04376*.
- Sankaran, A.; Mastropietro, O.; Saboori, E.; Idris, Y.; Sawyer, D.; Askari Hemmat, M. H.; and Hacene, G. B. 2021. Deeplite Neutrino: An End-to-End Framework for Constrained Deep Learning Model Optimization. *arXiv preprint arXiv:2101.04073*.
- Srinivas, S.; and Babu, R. V. 2015. Data-free parameter pruning for deep neural networks. *arXiv preprint arXiv:1507.06149*.
- Tishby, N.; Pereira, F. C.; and Bialek, W. 2000. The information bottleneck method. *arXiv preprint physics/0004057*.
- Wang, W.; Fu, C.; Guo, J.; Cai, D.; and He, X. 2019. Cop: Customized deep model compression via regularized correlation-based filter-level pruning. *arXiv preprint arXiv:1906.10337*.
- Zhai, X.; Kolesnikov, A.; Houlsby, N.; and Beyer, L. 2021. Scaling vision transformers. *arXiv preprint arXiv:2106.04560*.