

Magnitude and Uncertainty Pruning Criterion for Neural Networks

Vinnie Ko
University of Oslo
vinniebk@math.uio.no

Stefan Oehmcke
University of Copenhagen
stefan.oehmcke@di.ku.dk

Fabian Gieseke
University of Copenhagen
fabian.gieseke@di.ku.dk

Abstract—Neural networks have achieved dramatic improvements in recent years and depict the state-of-the-art methods for many real-world tasks nowadays. One drawback is, however, that many of these models are overparameterized, which makes them both computationally and memory intensive. Furthermore, overparameterization can also lead to undesired overfitting side-effects. Inspired by recently proposed magnitude-based pruning schemes and the Wald test from the field of statistics, we introduce a novel magnitude and uncertainty (M&U) pruning criterion that helps to lessen such shortcomings. One important advantage of our M&U pruning criterion is that it is scale-invariant, a phenomenon that the magnitude-based pruning criterion suffers from. In addition, we present a “pseudo bootstrap” scheme, which can efficiently estimate the uncertainty of the weights by using their update information during training. Our experimental evaluation, which is based on various neural network architectures and datasets, shows that our new criterion leads to more compressed models compared to models that are solely based on magnitude-based pruning criteria, with, at the same time, less loss in predictive power.

Index Terms—Neural network compression, pruning, overparameterization, Wald test

I. INTRODUCTION

In recent years, deep neural networks have achieved state-of-the-art performance for a broad range of tasks and have, hence, gained big popularity in many fields like computer vision and natural language processing [1]. Typically, these powerful models have millions or even billions of parameters, which require significant computational and memory resources when deployed. For the training phase, this computational bottleneck can usually be overcome by using powerful hardware such as graphics processing units, which renders it possible to train large models based on a lot of data in a reasonable amount of time. However, during the inference phase, once the trained model has been deployed, one is often faced with very limited computational resources. For instance, mobile or edge computing devices have little to no hardware acceleration available. In addition, many applications of neural networks—such as face recognition in video and autonomous driving—require the inference to be done with minimal latency.

A prominent approach to deal with these scenarios is *pruning*. Here, one usually generates a powerful model without any restrictions on its size and on its computational requirements during the training phase. Afterwards, one “compresses” the model such that it fits into the (main) memory of the

device that should be used for the inference phase. Pruning neural networks this way can be a critical task for many real-world applications. In addition, [2] show that there is significant redundancy in the parameterization of many deep learning models. This overparameterization not only requires unnecessary computational resources, but can also lead to an overfitting of the models. Various approaches have been suggested in the literature that aim at decreasing the size and the computational footprint of neural networks [3, 4, 5, 6, 7]; see [8] for an overview. One of the most successful works for pruning neural networks is provided by [3], who prune network parameters based on their magnitude, measured as absolute value.

Contribution: Inspired by the *Wald test* from the field of statistics [9], we introduce a novel magnitude and uncertainty (M&U) pruning criterion, which depicts a generalization of both, the magnitude based criterion used by [3] and the Wald test statistic. One of the biggest advantages above the magnitude based pruning criteria [3, 10, 11, 12, 13, 14, 15] is that our M&U pruning criterion is scale-invariant. Furthermore, it can be easily applied to many papers on pruning strategies that exploit a magnitude based criterion. Furthermore, we also propose a new “pseudo bootstrap” scheme that efficiently estimates the weight uncertainties based on update information obtained during the training process. Our experiments show that using our M&U pruning criterion outperforms the previous pruning criterion and results in compressed models with less loss in predictive power.

II. RELATED WORK

The attempt to prune neural networks started with [16] and [17], who used a second-order Taylor expansion to determine the saliency of each parameter. Their methods require the computation of (the inverse of) the Hessian matrix regarding the loss function. [18] consider pruning as a statistical model selection problem and apply diverse model selection tools from the field of maximum likelihood theory.

Since the recent success of deep learning, there has been a number of new approaches for reducing the computation and memory footprint of neural networks. One approach, quantization, aims to reduce the numerical precision required to represent each parameter (or neuron). For example, [19], [20], and [21] compress networks by using 8-bit precision, 3-bit

precision, and vector quantizations, respectively. There are more quantization based methods that expanded on these approaches [22, 23, 24, 25]. Other strategies to reduce the computational burden resort to low-rank decompositions of tensors [26, 27, 28]. In another category, [29] uses distilling approach to obtain a small ‘student’ network that can mimic what the bigger ‘teacher’ network does.

The strategy that has received the most attention in recent years is probably the pruning of parameters based on certain criteria. Particularly, [3] achieve state-of-the-art results by pruning weights based on their magnitudes (i.e., based on the absolute values of the parameters) and by retraining the network afterwards. The same authors further improve on their results by combining their absolute value based pruning scheme with quantizing and Huffman encoding techniques [10]. Extensions of this pruning strategy have also been proposed for neural machine translation (NMT) [30] and for recurrent neural networks (RNN) [31], respectively. Yet, those two follow-up works keep the same criterion as the one originally proposed by [3] for ranking the weights to be pruned, i.e., they use the absolute values of the weights. Finally, there are also attempts to prune specific classes of weights [15, 32, 33, 34]. These approaches have the advantage that the resulting pruned networks have less irregular network connections and are, hence, better suited for parallel implementations.

III. BACKGROUND

In this section, we briefly recap the theory behind statistical hypothesis testing and look into the similarities between the Wald test statistic and the magnitude based criterion [3].

A. Asymptotic Normality

Consider a regression model g with independent and dependent variable pairs $(\mathbf{x}_1, \mathbf{y}_1), \dots, (\mathbf{x}_n, \mathbf{y}_n)$. Further, let $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n]^T$ denote a tensor containing all n observations of independent variables. The arbitrary regression model g can then be written as

$$E[\mathbf{Y}] = g(\mathbf{X}, \mathbf{w}), \quad (1)$$

where $\mathbf{w} = [w_1, \dots, w_p]^T$ are the parameters of the model, which are commonly called *weights* in case of neural networks. The parameter p indicates the total number of parameters of the model.

For now, we assume that there exists a unique set of optimal parameters \mathbf{w}_* , which will give the best approximation of \mathbf{y} . We can then obtain $\hat{\mathbf{w}}$, our estimate of \mathbf{w}_* , by minimizing a loss function with the optimization technique of choice (e.g., stochastic gradient descent, mini-batch gradient descent, momentum based methods, etc.).

When the loss function is the negative log-likelihood, or an equivalent quantity such as cross entropy, gradient descent can be seen as a numerical realization of maximum likelihood (ML) estimation. In this case, we can apply the results from classic literature on ML theory [35, 36, 37] and we

have for $n \rightarrow \infty$, the following convergence in distribution property:

$$\sqrt{n}(\hat{\mathbf{w}} - \mathbf{w}_*) \xrightarrow{d} \mathcal{N}(\mathbf{0}, \mathbf{V}_{\text{ML}}). \quad (2)$$

Here \mathbf{V}_{ML} is the covariance matrix as defined in ML theory. In the field of statistics, this type of convergence to a normal distribution is called *asymptotic normality* [38]. In the model robust case, \mathbf{V}_{ML} is equal to the so-called sandwich estimator

$$\mathbf{V}_{\text{ML}} = \mathcal{I}_{\mathbf{w}}^{-1} \mathbf{K}_{\mathbf{w}} \mathcal{I}_{\mathbf{w}}^{-1}, \quad (3)$$

where $\mathcal{I}_{\mathbf{w}}$ is the Fisher information and where $\mathbf{K}_{\mathbf{w}}$ is the covariance matrix of the score function [37]. In statistics, the score function is defined as the first derivative of the log-likelihood function with respect to its parameters. It indicates how sensitive a likelihood function is w.r.t. its parameters. For exact definitions of the quantities in Equation (3), we refer to the literature in maximum likelihood theory, for instance by [37]. When one assumes that the model is true, in other words, that the current candidate model is the true model that generated the data, we have $\mathbf{K}_{\mathbf{w}} = \mathcal{I}_{\mathbf{w}}$ and Equation (3) simplifies to

$$\mathbf{V}_{\text{ML}} = \mathcal{I}_{\mathbf{w}}^{-1}. \quad (4)$$

For more details on the impact of this ‘true model’ assumption and the comparison between Equations (3) and (4), we refer to the corresponding literature [39, 40, 41].

Even though stochastic gradient descent (SGD) and its modified versions can be considered as numerical realization of ML estimation, the details of SGD lead to properties different from the ones obtained via ML. The recent work by [42] obtains, for the first time, a full analytical characterization of the asymptotic behavior of SGD procedures, including an asymptotic normality and analytical expression for the covariance matrix, which is different from that of ML estimators. They also show the loss of asymptotic statistical efficiency for SGD estimators. Although the work of [42] is a great step forward towards understanding the asymptotic properties of SGD estimators, its result cannot be directly applied to non-vanilla versions of SGD without necessary extension work. For example, combining the SGD estimator with training strategies like Dropout [43] and Dropconnect [44] or using momentum based modifications such as RMSprop [45] or Adam [46] would invalidate the asymptotic covariance matrix from [42].

Based on [37] and [42], it is not unreasonable to assume that the modified versions of the SGD estimator are still consistent and have asymptotic normality with a finite, but unknown covariance matrix \mathbf{V} , i.e., that we have

$$\sqrt{n}(\hat{\mathbf{w}} - \mathbf{w}_*) \xrightarrow{d} \mathcal{N}(\mathbf{0}, \mathbf{V}). \quad (5)$$

B. Statistical Hypothesis Testing

With the asymptotic normality property as defined in Section III-A, one can perform statistical hypothesis testing on the parameters \mathbf{w} . The natural null hypothesis would be

$$H_0 : w_j = 0, \quad (6)$$

i.e., is the j -th parameter of the model significantly different from 0? The three commonly used statistical hypothesis tests are the *Wald test* [9], the *likelihood ratio test* [47], and the *score test* [48]. Note that all three tests are asymptotically equivalent [49]. Yet, the Wald test has a practical advantage over the other two since it does not require nested models for testing the null hypothesis on individual weights. Here, “nested” means that one can obtain one model by constraining some of the parameters of another model.

Given Equations (5) and (6), the Wald test statistic for the parameter w_j is given by

$$z = \frac{\hat{w}_j}{\text{SE}(\hat{w}_j)} = \frac{\hat{w}_j}{\sqrt{\hat{\mathbf{V}}_{j,j}}}, \quad (7)$$

where $\text{SE}(\hat{w}_j)$ indicates the standard error of \hat{w}_j and $\hat{\mathbf{V}}$ the estimate of the covariance matrix from Equation (5).¹ The test statistic z follows the standard normal distribution [49]. The corresponding p -value can be obtained via $p = 2(1 - \Phi(|z|))$, where Φ is the cumulative distribution function of the standard normal distribution.

C. Wald Test and Absolute Value Based Pruning

In statistics, one of the most canonical ways of removing redundant parameters is by performing null hypothesis testing as described above and by dropping the parameters that have a p -value above a certain threshold (e.g., 0.05 or 0.01). Since $p = 2(1 - \Phi(|z|))$ is a monotonic function of $|z|$, one can skip the step of converting $|z|$ into a p -value. Instead, one can directly use

$$|z| = \frac{|\hat{w}_j|}{\text{SE}(\hat{w}_j)} = \frac{|\hat{w}_j|}{\sqrt{\hat{\mathbf{V}}_{j,j}}}, \quad (8)$$

and can compare it with the threshold values in $|z|$ -scale, which can be obtained by transforming the threshold p -values via the function $\Phi^{-1}(1-p/2)$. Here, we can realize two things. Firstly, Equation (8) is highly similar to the magnitude based pruning criterion suggested by [3] and also used by [10], [30], and [31]. In those works, only the absolute value of each parameter is considered and the parameters with absolute values below a certain threshold are removed. We denote this criterion as the ABS pruning criterion from now on. The only difference between the Wald test statistic (Equation (8)) and the ABS pruning criterion is the standard error term in the denominator.

In this sense, the Wald test statistic can actually be interpreted as evaluating parameters based on their absolute value, while compensating for their uncertainty. For example, if \hat{w}_j has high uncertainty, there is higher chance that we observe a large value of \hat{w}_j “by chance”. Dividing by the standard error compensates for this uncertainty.

¹Since we defined \mathbf{V} as a finite, but unknown covariance matrix, its estimate $\hat{\mathbf{V}}$ is an unknown and arbitrary estimate.

IV. MAGNITUDE AND UNCERTAINTY PRUNING CRITERION

Inspired by Wald test, we evaluate the importance of parameters by taking both their magnitudes and their uncertainties into account. By using this combined criterion, we can prune to the desired proportion of parameters in a model, specific layer, feature map, or any selected part of a model. In addition, we also introduce a “pseudo bootstrap” approach to efficiently estimate the uncertainties of the weights by keeping track of their changes during the training process.

A. Definition

Based on our observation of the similarity between the Wald test statistic and the ABS pruning criterion in Section III-C as well as based on the interpretation of the Wald test statistic, we define our M&U pruning criterion as

$$\tau_j = \frac{|\hat{w}_j|}{\lambda + \tilde{\sigma}_{\hat{w}_j}}. \quad (9)$$

Here, $\tilde{\sigma}_{\hat{w}_j}$ is an uncertainty estimate of \hat{w}_j , a quantity that should reflect how uncertain the estimated parameter \hat{w}_j is. We will cover possible choices for $\tilde{\sigma}_{\hat{w}_j}$ in Sections IV-C and IV-D. The hyperparameter $0 \leq \lambda$ is a user-defined parameter that determines the balance between magnitude and uncertainty. When $\lambda \rightarrow \infty$, no uncertainty is taken into account and the M&U pruning criterion becomes equal to the ABS pruning criterion used by [3]. As λ decreases, more and more uncertainty is taken into account. In case of $\lambda = 0$, the criterion τ_j has the same functional form as the Wald test statistic. If one chooses $\tilde{\sigma}_{\hat{w}_j}$ such that it is an analytically justified estimate of the standard error, τ_j is simply the Wald test statistic. This implies that M&U pruning criterion can be seen as a generalization of the Wald test statistic and the ABS pruning criterion.

B. Scale Invariant Property

Pruning strategies that are based on the magnitude based pruning criterion [3, 10, 11, 12, 13, 14, 15, 30, 31] suffer from the so-called scale issue by their nature.

Consider a very simple model (e.g. linear regression) with 2 input variables X_1 and X_2 . Assume that the weights are $w_1 = 10$ and $w_2 = 0.1$. One can be easily tricked to think that X_1 is more important than X_2 . However, if we divide X_2 by 1000 (e.g. m to km), w_2 will get much bigger and this will change which weight is pruned by the magnitude based pruning criterion. In contrast, akin to the Wald test statistic, our M&U pruning criterion cancels out this change of scale effect by standardizing with $\tilde{\sigma}_{\hat{w}_j}$. Moreover, we introduce in Section IV-E a reparametrization trick that makes the hyperparameter λ robust to the change of scale.

C. Choosing Uncertainty Estimates

The first thing to note is that $\tilde{\sigma}_{\hat{w}_j}$ is not the standard error, but concerns a much more general concept of uncertainty estimate. Ideally, the standard error is a very suitable candidate as an uncertainty estimate $\tilde{\sigma}_{\hat{w}_j}$, like in the original Wald test statistic formula. However, there is at the moment no

asymptotic theories developed for non-convex neural networks with commonly used estimation methods such as Adam [46] or RMSprop [45].

In rare cases, when the optimization task induced by a neural network is convex (or piecewise convex) with the negative log-likelihood, or an equivalent quantity such as cross entropy as loss function, one can take analytically derived asymptotic variance formulas such as Equation (3), Equation (4), or other ones [42] and estimate these uncertainty measures via their sample equivalent (e.g., by replacing the expectation with a sum over the samples). However, estimating those quantities has a number of numerical bottlenecks. For example, to obtain the sample equivalent of \mathcal{I}_w , one needs to calculate the Hessian matrix, which requires the second derivative of the log-likelihood function with respect to all the parameters in the model. Although modern deep learning software offers highly optimized automatic differentiation for the first derivative, many of the automatic differentiation implementations do not support the second derivative. In case they do support it, the corresponding implementations are not as highly optimized as those for the first derivative.

Another computational bottleneck is the computation of \mathbf{K}_w . Computation of this matrix requires the first derivative of the log-likelihood with respect to the model parameters, for every sample separately. This is highly inefficient for most deep learning software packages since optimized for mini-batch operations. Yet, the biggest computational obstacle is that one needs to take the inverse of \mathcal{I}_w . For a network with, say, a million parameters, inverting this matrix is infeasible in practice even using significant computational resources.

D. Pseudo Bootstrap

We introduce two possible alternatives to compute the uncertainty estimate $\tilde{\sigma}_{\hat{w}_j}$. The first alternative is to make use of bootstrapping [50]. There have been extensive studies done on both theoretical and numerical properties of this approach [51, 52, 53, 54] and it is a canonical method in statistics for estimating the standard error in case an analytical estimate is not feasible. The disadvantage is that we need to train the same model multiple times, which is a big increase in the computational costs.²

The second alternative is a novel estimation scheme, called *pseudo bootstrap*, which uses uncertainty information gathered during the training process. This idea stems from the fact that the training process with mini-batch SGD has similarities with the bootstrap resampling strategy: Each mini-batch we draw during training can be seen as a bootstrap sample with replacement and of smaller size. In addition, the iterative weight update process of mini-batch SGD can be seen as training with bootstrap samples, where each bootstrap training instance uses the parameter values from the previous instance as parameter initialization.

²Since neural networks are usually trained with large number of observations, the Jackknife resampling method [55], which requires repeating the training process equal to the number of total data points, is not suitable.

Algorithm 1 *Pseudo bootstrap*, our proposed novel algorithm for a fast estimation of the uncertainty of each weight in neural network. The algorithm described below computes a single weight w_j and has to be repeated for all the weights in the neural network.

Require: $B \in \mathbb{N}$: Weight collection size

Require: $\mathbf{w} = [w_1, \dots, w_p]$: Parameter vector of the given neural network.

Require: n_{iter} : The total number of weight update iterations that is going to be made during the training phase. Note that this is the total number of iterations and not the total number of epochs.

```

1:  $\mathbf{w}_{j,\{B\}}$ : An empty vector of length  $B$ .
2:  $i_{\text{iter}} \leftarrow 0$  (Initialize training iteration count)
3:  $i_B \leftarrow 0$  (Initialize weight collecting count)
4: while training the target model with mini-batch SGD do
5:    $i_{\text{iter}} \leftarrow i_{\text{iter}} + 1$ 
6:   if  $n_{\text{iter}} - B < i_{\text{iter}} \leq n_{\text{iter}}$  then
7:      $i_B \leftarrow i_B + 1$ 
8:      $\mathbf{w}_{j,\{B\}}[i_B] \leftarrow w_j$  (Record current value of  $w_j$ .)
9:   end if
10: end while
11:  $\tilde{\sigma}_{\hat{w}_j} = \text{std}(\mathbf{w}_{j,\{B\}})$  (Compute uncertainty estimate of  $\hat{w}_j$  by using its ‘fluctuations’ during the training.)
12: return  $\tilde{\sigma}_{\hat{w}_j}$ 

```

The pseudo bootstrap approach is sketched in Algorithm 1: To obtain the uncertainty estimate $\tilde{\sigma}_{\hat{w}_j}$, the algorithm monitors how \hat{w}_j has changed during the last B iterations of mini-batch SGD training and by storing that information in the ‘weight collection’ $\mathbf{w}_{j,\{B\}} = [w_{j,1}, \dots, w_{j,B}]$ (see Lines 4–10). After the training, we simply obtain the uncertainty estimate via $\tilde{\sigma}_{\hat{w}_j} = \text{std}(\mathbf{w}_{j,\{B\}})$ (see Line 11). Note that other ‘collecting strategies’ such as recording weight values through the entire training process with a certain interval between them have been tested. However, the best results were obtained by recording the weight values during the last iterations.

E. Reparameterization of λ

The hyperparameter λ determines the tradeoff between magnitude and uncertainty. It is therefore crucial to consider reasonable values for λ . Note that the actual impact of λ on the balance between magnitude and uncertainty depends largely on which scale w_j is. For instance, if we decide to use batch normalization and the absolute value of w_j suddenly becomes much smaller, then a previously ‘optimal’ value for λ (e.g., $\lambda = 1$) might not be optimal anymore. To deal with this problem, we suggest to reparametrize λ via

$$\lambda = \lambda^* \cdot \text{std}(\mathbf{w}_{\{\text{layer}_k\}}) \quad (10)$$

when one prunes within layer k . Here, λ^* is a fixed global assignment for λ and $\text{std}(\mathbf{w}_{\{\text{layer}_k\}})$ depicts the standard deviation of the weights within the k -th layer.³

V. EXPERIMENTS

We implemented our M&U pruning criterion in PyTorch [56] and evaluated our approach via four experiments that were based on the Fashion MNIST, CIFAR-10, MNIST, and ImageNet datasets, respectively. Experiment 1 was repeated 100 times, Experiment 2 and 3 were repeated 20 times,

³Naturally, instead of considering $\mathbf{w}_{\{\text{layer}_k\}}$, one can also consider the standard deviation based on all the weights via $\mathbf{w}_{\{\text{model}\}}$ such that λ is scaled by the whole model (or by a specific part of a certain layer).

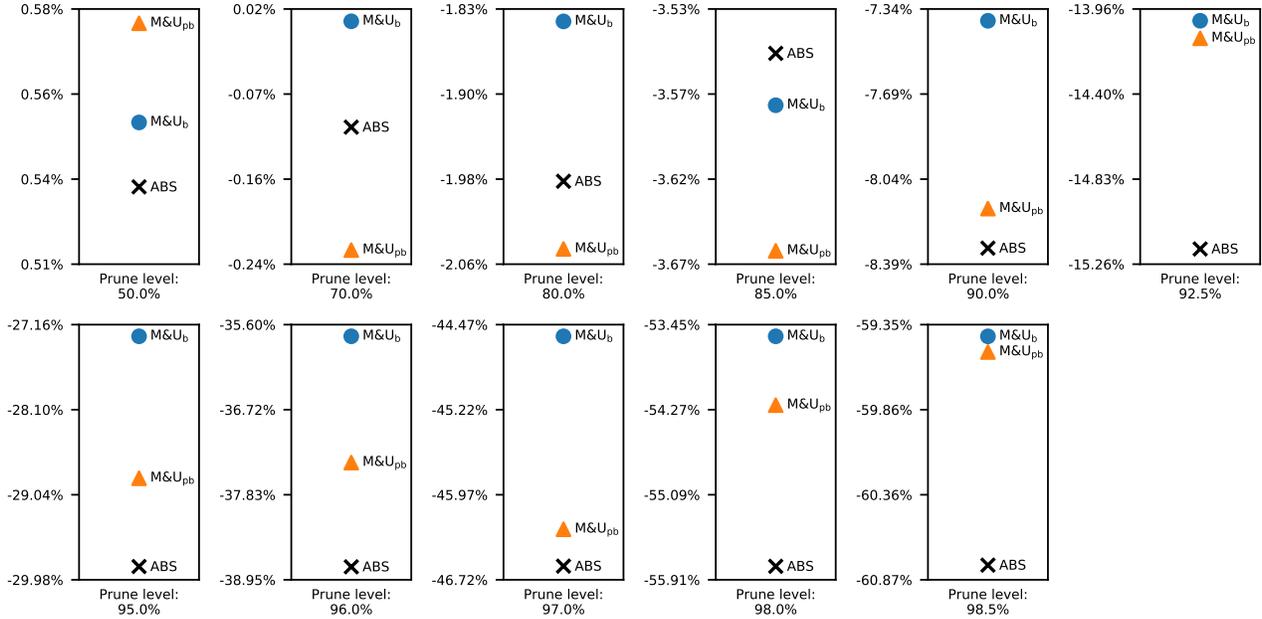


Fig. 1. Result of Experiment 1. The y -axis of each plot shows the loss in test accuracy due to pruning (with retraining afterwards). It was computed as the test accuracy of the pruned model subtracted by the test accuracy of the unpruned original model. Blue circles with ‘M&U_b’ indicate that our M&U pruning method with $\tilde{\sigma}_{\hat{w}_j}^{\text{boot}}$ was used. Orange triangles with ‘M&U_{pb}’ mean that our M&U pruning method with $\tilde{\sigma}_{\hat{w}_j}^{\text{pboot}}$ was used. Black crosses with ‘ABS’ mean that the ABS pruning criterion from [3] was used. The test accuracy of the unpruned model was 81.21%.

and Experiment 4 was repeated 10 times. The reported results are averaged results over those repetitions.

Since the majority of pruning strategies utilize the magnitude based pruning criterion, it is most informative to compare our method directly to the ABS criterion from [3]. Our M&U pruning criterion can be easily extended to replace the pruning criterion part of, for example, automatic gradual pruning [12], structured pruning [13, 15], dynamic network surgery [11] and more. Note further that our M&U pruning criterion should be seen as a companion rather than a competitor to other model compression approaches such as distilling [29] and quantization [19], since it can be combined with those techniques.

A. Experiment 1: Comparison of Uncertainty Estimates

The Fashion MNIST dataset consists of 28x28 grayscale images, each with a label from one of 10 classes [57]. To make the bootstrap estimation computationally feasible, we kept the size of this experiment as small as possible. Thus, we only used a balanced subset of 6,000 examples from the training data. From this subset, 600 examples were set apart for validation during training. For testing, the entire test set was used. We employed a very small feedforward convolutional neural network (CNN) with two convolution layers and three fully connected layers, combined with max pooling, ReLU activation function and dropout. The first fully connected layer got a feature map of size 6 x 6 x 6 as input. The two hidden layers had 32 and 48 neurons, respectively, and the model had

a total amount of 9,660 parameters. We chose cross entropy as loss function and the model was trained for 500 epochs via RMSprop.

Firstly, since the total number of parameters was small, we tried to estimate the covariance matrix of w with analytical expressions (3) and (4). However, \mathcal{I}_w and K_w were not invertible most of the time. As alternative, we obtained the pseudo inverse of those matrices. Yet, the pseudo inverse contained almost always negative diagonal elements. As a sanity check, we estimated the expectation of the score vector and many of its elements had values that were far away from 0. This suggests that the maximum likelihood theory did not hold in this case. This was expected, since many of the assumptions from [36] and [42], including convexity, were broken. Secondly, we estimated $\tilde{\sigma}_{\hat{w}_j}$ with bootstrap. We drew 100 bootstrap datasets and trained 100 independent instances of the model. With the weights from those 100 trained models, we obtained the bootstrap based uncertainty estimate $\tilde{\sigma}_{\hat{w}_j}^{\text{boot}}$. Lastly, we obtained the pseudo bootstrap based uncertainty estimate $\tilde{\sigma}_{\hat{w}_j}^{\text{pboot}}$, with the method described in Section IV-D and $B = 200$. We tried different values of λ^* and empirically choose $\lambda^* = 10^{-1}$ for $\tilde{\sigma}_{\hat{w}_j}^{\text{boot}}$ and $\lambda^* = 10^{-4}$ for $\tilde{\sigma}_{\hat{w}_j}^{\text{pboot}}$.

With the obtained $\tilde{\sigma}_{\hat{w}_j}$ and λ , we pruned $a\%$ of parameters by using our M&U pruning criterion from each fully connected layer. The values of a that were tested can be found in Figure 1. For comparison, we also pruned the model by using the ABS pruning criterion [3]. After pruning, each model was



Fig. 2. Result of Experiment 2. The y -axis is the same as in Figure 1. Orange triangles with ‘M&U_{pb}’ mean that our M&U pruning method with $\tilde{\sigma}_{\hat{w}_j}^{\text{pboot}}$ was used. Black crosses with ‘ABS’ mean that the ABS pruning criterion from [3] was used. The test accuracy of the unpruned model was 79.98%.

retrained for 200 epochs.

Figure 1 shows the results of this experiment. One can observe that our M&U pruning criterion either with $\tilde{\sigma}_{\hat{w}_j}^{\text{boot}}$ or $\tilde{\sigma}_{\hat{w}_j}^{\text{pboot}}$ yielded a better (i.e., smaller) test accuracy drop compared to the ABS pruning criterion. The difference was rather small when the pruning level was relatively low, but the gap between the two increased as the pruning level increased. Using $\tilde{\sigma}_{\hat{w}_j}^{\text{boot}}$ gave the most optimal result (defeated ABS 10 out of 11 times), but using $\tilde{\sigma}_{\hat{w}_j}^{\text{pboot}}$ also gave better results than the ABS pruning criterion (defeated ABS 8 out of 11 times).

B. Experiment 2: CIFAR-10

We further examined the performance of our M&U pruning criterion on the CIFAR-10 dataset [58]. Unlike Experiment 1, we utilized the whole training set, where 5,000 examples of the training set were used as validation set during training. We employed a smaller version of AlexNet [59]; the model has three convolution layers with ReLU activation function and max pooling, which yield feature maps of size 128*4*4 as output. The network also has four hidden layers with 512, 512, 512, and 256 neurons. As the rapid drop in test accuracy suggests, this model has relatively low degree of overparameterization. We chose cross entropy as loss function and the model was trained for 1,000 epochs with RMSprop. Due to large amount of computation required for $\tilde{\sigma}_{\hat{w}_j}^{\text{boot}}$, we only used $\tilde{\sigma}_{\hat{w}_j}^{\text{pboot}}$ with $B = 200$. We empirically chose $\lambda^* = 1$.

We pruned parameters of each fully connected layer by using our M&U pruning criterion or the ABS pruning criterion; the pruning levels are shown in Figure 2. After pruning, the

model was retrained for 200 epochs. Figure 2 shows the results of this experiment. Our M&U pruning criterion with $\tilde{\sigma}_{\hat{w}_j}^{\text{pboot}}$ gave a better test accuracy than the ABS criterion in 9 out of 11 times.

C. Experiment 3: Overparameterization

To explore the case of high overparameterization, we used the MNIST dataset [60] with a MLP model. The MLP model had three hidden layers with 512, 1024, and 512 neurons, respectively. The model was trained for 1,000 epochs with RMSprop.

Due to the large amount of computation required for $\tilde{\sigma}_{\hat{w}_j}^{\text{boot}}$, we only used $\tilde{\sigma}_{\hat{w}_j}^{\text{pboot}}$ with $B = 200$. We empirically chose $\lambda^* = 1$. For this experiment, we used the iterative pruning strategy described in [3]. We pruned the parameters of each fully connected layer by using our M&U pruning criterion or the ABS pruning criterion with the pruning levels shown in Figure 3. After pruning, the model was retrained for 200 epochs. Figure 3 shows the results of this experiment. Since the model is highly overparameterized, the performance drop was minimal, even when we pruned 99% of the weights. In fact, pruning up to 95% of the weights increased the test accuracy of the model, possibly because pruning decreased overfitting. Our M&U pruning criterion with $\tilde{\sigma}_{\hat{w}_j}^{\text{pboot}}$ gave a better test accuracy than using the ABS pruning criterion in 10 out of 11 times.

Note that since the model was highly overparameterized, one could prune a large number of the weights without sacrificing too much predictive power. For example, Figure 3

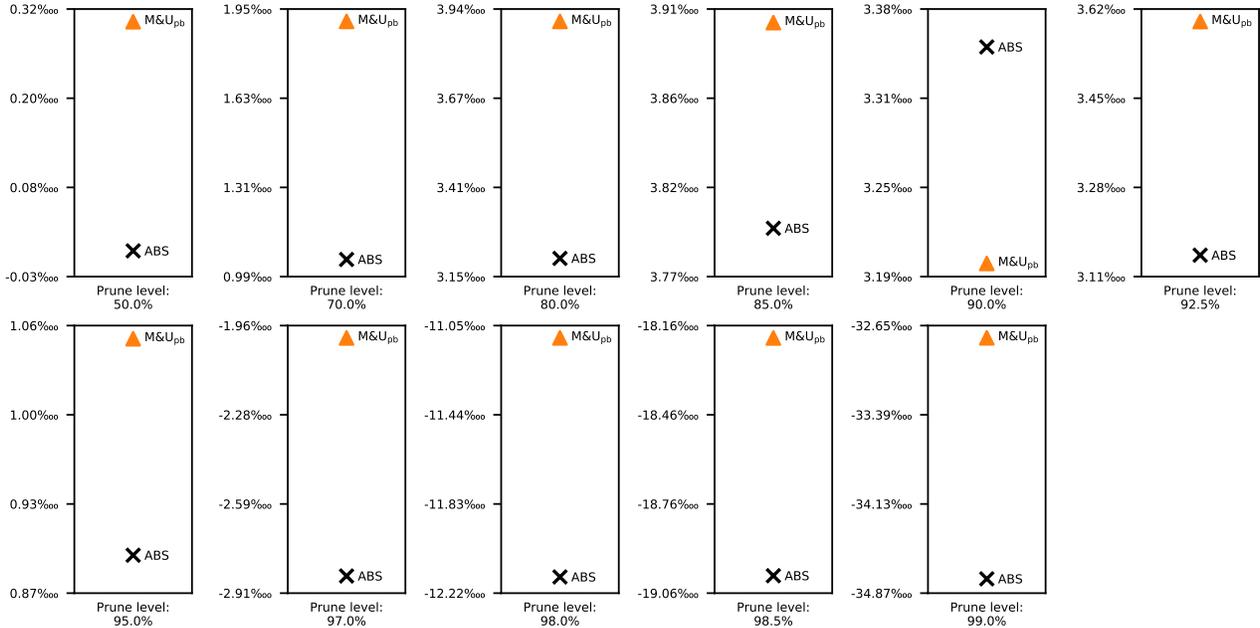


Fig. 3. Result of Experiment 3. The y -axis is the same as in Figure 1. Note that the unit of the y -axis is basis point, and not percent. Orange triangles with ‘M&U_{pb}’ mean that our M&U pruning method with $\tilde{\sigma}_{\tilde{w}_j}^{\text{pboot}}$ was used. Black crosses with ‘ABS’ mean that the ABS pruning criterion from [3] was used. The test accuracy of the unpruned model was 98.5160%.

shows that could prune 99% of the weights while loosing only 0.33% of the test accuracy (note that the unit of the y -axis in Figure 3 is basis point, and not percent). The performance difference between our M&U pruning criterion and the ABS pruning criterion might seem small due to this minimal test accuracy drop throughout all pruning levels. But when one takes this into account, one can see that our M&U pruning criterion still outperforms the ABS criterion consistently.

D. Experiment 4: ImageNet

ImageNet ILSVRC-12 dataset [61] has 1.28M training images, spread over 1,000 classes. To make 10 repetition of the experiment feasible, we used a subset that consists of 100 classes. We used VGG-11 [62] architecture, shrank with the reduced number of classes. Like in Experiment 3, we used $\tilde{\sigma}_{\tilde{w}_j}^{\text{pboot}}$ with $B = 200$. We empirically chose $\lambda^* = 10^{-3}$. We pruned the parameters of each fully connected layer by using our M&U pruning criterion or the ABS pruning criterion with the pruning levels shown in Figure 4. After pruning, the model was retrained for 40 epochs. Figure 4 shows the results of this experiment. Our M&U pruning criterion with $\tilde{\sigma}_{\tilde{w}_j}^{\text{pboot}}$ gave a better test accuracy than using the ABS pruning criterion in 6 out of 6 times.

VI. CONCLUSION AND DISCUSSION

Motivated by the Wald test statistic, we derived the M&U pruning criterion for neural networks. Our M&U pruning criterion takes both magnitude and uncertainty into account and can balance between them through the hyperparameter

λ . It can be considered as a generalization of both the Wald test statistic and the magnitude based criterion [3]. Our M&U pruning criterion is free from the scale variant issue, which the magnitude based pruning criterion suffer from. We also suggested a reparametrization of λ that makes finding reasonable values for λ easier. Further, we introduced ‘pseudo bootstrap’, a very efficient method of measuring uncertainty of weights by tracking their updates during training. The experiments suggest that our M&U pruning criterion outperforms the criterion suggested by [3] most of the time, in terms of test accuracy of the pruned model, with the best result achieved by using $\tilde{\sigma}_{\tilde{w}_j}^{\text{pboot}}$. Thus, using our M&U pruning criterion can lead to models that require less memory and at smaller computational cost, while minimizing the loss of predictive power. In addition, Our MnU pruning criterion can easily be applied to the pruning strategies that utilize the magnitude based pruning criterion.

In the future, we wish to investigate numerical and theoretical properties of pseudo bootstrap and compare it to that of the ‘usual’ bootstrap. Due to limited amount of GPU resources, we had to keep our models relatively small. We would like to investigate the performance of our method on bigger and deeper neural network models. Further, it will be fruitful to combine our pruning criterion with other types of compression methods such as quantization, Huffman coding, and regularization.

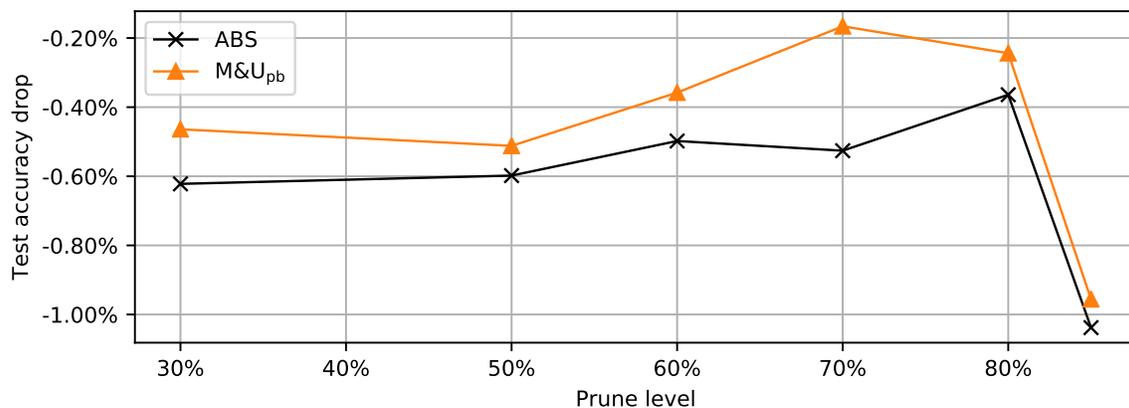


Fig. 4. Result of Experiment 4. Since the drop in test accuracy was at similar level across all the pruning levels we tested, a line plot is used. Orange triangles with ‘M&U_{pb}’ mean that our M&U pruning method with $\tilde{\sigma}_{\hat{w}_j}^{\text{pboot}}$ was used. Black crosses with ‘ABS’ mean that the ABS pruning criterion from [3] was used. The test accuracy of the unpruned model was 72.6820%.

ACKNOWLEDGMENT

This work is supported by Google Cloud Platform through their research credits program.

We acknowledge support from the Danish Industry Foundation through the Industrial Data Analysis Service (IDAS).

REFERENCES

- [1] I. Goodfellow, Y. Bengio, and A. Courville, *Deep Learning*. MIT press, 2016.
- [2] M. Denil, B. Shakibi, L. Dinh, N. De Freitas *et al.*, “Predicting parameters in deep learning,” in *Advances in neural information processing systems (NeurIPS)*, 2013, pp. 2148–2156.
- [3] S. Han, J. Pool, J. Tran, and W. Dally, “Learning both weights and connections for efficient neural network,” in *Advances in neural information processing systems (NeurIPS)*, 2015, pp. 1135–1143.
- [4] C. Louizos, M. Welling, and D. P. Kingma, “Learning sparse neural networks through l_0 regularization,” in *International Conference on Learning Representations (ICLR)*, 2018.
- [5] W. Chen, J. Wilson, S. Tyree, K. Weinberger, and Y. Chen, “Compressing neural networks with the hashing trick,” in *International Conference on Machine Learning (ICML)*, 2015, pp. 2285–2294.
- [6] S. Srinivas and R. V. Babu, “Data-free parameter pruning for deep neural networks,” in *British Machine Vision Conference (BMVC)*, 2015.
- [7] K. Ullrich, E. Meeds, and M. Welling, “Soft weight-sharing for neural network compression,” in *International Conference on Learning Representations (ICLR)*, 2017.
- [8] Y. Cheng, D. Wang, P. Zhou, and T. Zhang, “A survey of model compression and acceleration for deep neural networks,” *arXiv preprint arXiv:1710.09282*, 2017.
- [9] S. D. Silvey, “The lagrangian multiplier test,” *The Annals of Mathematical Statistics*, vol. 30, no. 2, pp. 389–407, 1959.
- [10] S. Han, H. Mao, and W. J. Dally, “Deep compression: Compressing deep neural networks with pruning, trained quantization and huffman coding,” in *International Conference on Learning Representations (ICLR)*, 2016.
- [11] Y. Guo, A. Yao, and Y. Chen, “Dynamic network surgery for efficient dnns,” in *Advances in neural information processing systems (NeurIPS)*, 2016, pp. 1379–1387.
- [12] M. Zhu and S. Gupta, “To prune, or not to prune: exploring the efficacy of pruning for model compression,” *arXiv preprint arXiv:1710.01878*, 2017.
- [13] J.-H. Luo, J. Wu, and W. Lin, “Thinnet: A filter level pruning method for deep neural network compression,” in *International Conference on Computer Vision (ICCV)*, 2017, pp. 5058–5066.
- [14] W. Wen, C. Wu, Y. Wang, Y. Chen, and H. Li, “Learning structured sparsity in deep neural networks,” in *Advances in neural information processing systems (NeurIPS)*, 2016, pp. 2074–2082.
- [15] H. Li, A. Kadav, I. Durdanovic, H. Samet, and H. P. Graf, “Pruning filters for efficient convnets,” in *International Conference on Learning Representations (ICLR)*, 2017.
- [16] Y. LeCun, J. S. Denker, and S. A. Solla, “Optimal brain damage,” in *Advances in neural information processing systems (NeurIPS)*, 1990, pp. 598–605.
- [17] B. Hassibi, D. G. Stork, and G. J. Wolff, “Optimal brain surgeon and general network pruning,” in *IEEE international conference on neural networks*, 1993, pp. 293–299.
- [18] U. Anders and O. Korn, “Model selection in neural networks,” *Neural networks*, vol. 12, no. 2, pp. 309–323, 1999.
- [19] V. Vanhoucke, A. Senior, and M. Z. Mao, “Improving the speed of neural networks on CPUs,” 2011.

- [20] K. Hwang and W. Sung, "Fixed-point feedforward deep neural network design using weights+ 1, 0, and- 1," in *2014 IEEE Workshop on Signal Processing Systems (SiPS)*. IEEE, 2014, pp. 1–6.
- [21] Y. Gong, L. Liu, M. Yang, and L. Bourdev, "Compressing deep convolutional networks using vector quantization," *arXiv preprint arXiv:1412.6115*, 2014.
- [22] M. Courbariaux, Y. Bengio, and J.-P. David, "Binaryconnect: Training deep neural networks with binary weights during propagations," in *Advances in neural information processing systems (NeurIPS)*, 2015, pp. 3123–3131.
- [23] I. Hubara, M. Courbariaux, D. Soudry, R. El-Yaniv, and Y. Bengio, "Quantized neural networks: Training neural networks with low precision weights and activations," *The Journal of Machine Learning Research (JMLR)*, vol. 18, no. 1, pp. 6869–6898, 2017.
- [24] Z. Lin, M. Courbariaux, R. Memisevic, and Y. Bengio, "Neural networks with few multiplications," in *International Conference on Learning Representations (ICLR)*, 2016.
- [25] M. Rastegari, V. Ordonez, J. Redmon, and A. Farhadi, "Xnor-net: Imagenet classification using binary convolutional neural networks," in *European Conference on Computer Vision (ECCV)*. Springer, 2016, pp. 525–542.
- [26] E. L. Denton, W. Zaremba, J. Bruna, Y. LeCun, and R. Fergus, "Exploiting linear structure within convolutional networks for efficient evaluation," in *Advances in neural information processing systems (NeurIPS)*, 2014, pp. 1269–1277.
- [27] V. Lebedev, Y. Ganin, M. Rakhuba, I. Oseledets, and V. Lempitsky, "Speeding-up convolutional neural networks using fine-tuned cp-decomposition," *arXiv preprint arXiv:1412.6553*, 2014.
- [28] M. Jaderberg, A. Vedaldi, and A. Zisserman, "Speeding up convolutional neural networks with low rank expansions," *arXiv preprint arXiv:1405.3866*, 2014.
- [29] G. Hinton, O. Vinyals, and J. Dean, "Distilling the knowledge in a neural network," in *NIPS Deep Learning and Representation Learning Workshop*, 2015. [Online]. Available: <http://arxiv.org/abs/1503.02531>
- [30] A. See, M.-T. Luong, and C. D. Manning, "Compression of neural machine translation models via pruning," in *The SIGNLL Conference on Computational Natural Language Learning (CoNLL)*, 2016.
- [31] S. Narang, E. Elsen, G. Diamos, and S. Sengupta, "Exploring sparsity in recurrent neural networks," in *International Conference on Learning Representations (ICLR)*, 2017.
- [32] S. Anwar, K. Hwang, and W. Sung, "Structured pruning of deep convolutional neural networks," *ACM Journal on Emerging Technologies in Computing Systems (JETC)*, vol. 13, no. 3, p. 32, 2017.
- [33] S. Changpinyo, M. Sandler, and A. Zhmoginov, "The power of sparsity in convolutional neural networks," *arXiv preprint arXiv:1702.06257*, 2017.
- [34] V. Lebedev and V. Lempitsky, "Fast convnets using group-wise brain damage," in *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016, pp. 2554–2564.
- [35] G. Casella and R. L. Berger, *Statistical Inference*. Duxbury Pacific Grove, CA, 2002.
- [36] L. Le Cam, "Maximum likelihood: an introduction," *International Statistical Review*, vol. 58, no. 2, pp. 153–171, 1990.
- [37] H. White, "Maximum likelihood estimation of misspecified models," *Econometrica: Journal of the Econometric Society*, pp. 1–25, 1982.
- [38] A. W. Van der Vaart, *Asymptotic Statistics*. Cambridge University Press, 2000.
- [39] J. W. Hardin, "The sandwich estimate of variance," in *Maximum likelihood estimation of misspecified models: Twenty years later*. Emerald Group Publishing Limited, 2003, pp. 45–73.
- [40] G. Kauermann and R. J. Carroll, "The sandwich variance estimator: Efficiency properties and coverage probability of confidence intervals," 2000.
- [41] —, "A note on the efficiency of sandwich covariance matrix estimation," *Journal of the American Statistical Association*, vol. 96, no. 456, pp. 1387–1396, 2001.
- [42] P. Toulis, E. M. Airoldi *et al.*, "Asymptotic and finite-sample properties of estimators based on stochastic gradients," *The Annals of Statistics*, vol. 45, no. 4, pp. 1694–1727, 2017.
- [43] N. Srivastava, G. Hinton, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov, "Dropout: a simple way to prevent neural networks from overfitting," *The Journal of Machine Learning Research (JMLR)*, vol. 15, no. 1, pp. 1929–1958, 2014.
- [44] L. Wan, M. Zeiler, S. Zhang, Y. Le Cun, and R. Fergus, "Regularization of neural networks using dropconnect," in *International Conference on Machine Learning (ICML)*, 2013, pp. 1058–1066.
- [45] T. Tieleman and G. Hinton, "Lecture 6.5-rmsprop: Divide the gradient by a running average of its recent magnitude," *COURSERA: Neural networks for machine learning*, vol. 4, no. 2, pp. 26–31, 2012.
- [46] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," in *International Conference on Learning Representations (ICLR)*, 2015.
- [47] J. Neyman and E. S. Pearson, "On the problem of the most efficient tests of statistical hypotheses," *Philosophical Transactions of the Royal Society of London. Series A, Containing Papers of a Mathematical or Physical Character*, vol. 231, no. 694-706, pp. 289–337, 1933.
- [48] A. K. Bera and Y. Biliias, "Rao's score, Neyman's $c(\alpha)$ and Silvey's LM tests: an essay on historical developments and some new results," *Journal of Statistical Planning and Inference*, vol. 97, no. 1, pp. 9–44, 2001.
- [49] R. F. Engle, "Wald, likelihood ratio, and lagrange multiplier tests in econometrics," *Handbook of econometrics*, vol. 2, pp. 775–826, 1984.
- [50] B. Efron and R. J. Tibshirani, *An Introduction to the*

Bootstrap. CRC press, 1994.

- [51] P. J. Bickel, D. A. Freedman *et al.*, “Some asymptotic theory for the bootstrap,” *The annals of statistics*, vol. 9, no. 6, pp. 1196–1217, 1981.
- [52] B. Efron, “Bootstrap methods: another look at the jackknife,” in *Breakthroughs in statistics*. Springer, 1992, pp. 569–593.
- [53] A. C. Davison and D. V. Hinkley, *Bootstrap methods and their application*. Cambridge university press, 1997, vol. 1.
- [54] J. Shao and D. Tu, *The jackknife and bootstrap*. Springer Science & Business Media, 2012.
- [55] B. Efron, *The Jackknife, the Bootstrap, and Other Resampling Plans*. Siam, 1982, vol. 38.
- [56] A. Paszke, S. Gross, S. Chintala, G. Chanan, E. Yang, Z. DeVito, Z. Lin, A. Desmaison, L. Antiga, and A. Lerer, “Automatic differentiation in pytorch,” 2017.
- [57] H. Xiao, K. Rasul, and R. Vollgraf, “Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms,” *arXiv preprint arXiv:1708.07747*, 2017.
- [58] A. Krizhevsky and G. Hinton, “Learning multiple layers of features from tiny images,” Citeseer, Tech. Rep., 2009.
- [59] A. Krizhevsky, I. Sutskever, and G. E. Hinton, “Imagenet classification with deep convolutional neural networks,” in *Advances in neural information processing systems (NeurIPS)*, 2012, pp. 1097–1105.
- [60] Y. LeCun, “The mnist database of handwritten digits,” <http://yann.lecun.com/exdb/mnist/>, 1998.
- [61] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, A. C. Berg, and L. Fei-Fei, “ImageNet Large Scale Visual Recognition Challenge,” *International Journal of Computer Vision (IJCV)*, vol. 115, no. 3, pp. 211–252, 2015.
- [62] K. Simonyan and A. Zisserman, “Very deep convolutional networks for large-scale image recognition,” in *International Conference on Learning Representations (ICLR)*, 2015.