# Pruning by Explaining: A Novel Criterion for Deep Neural Network Pruning

Seul-Ki Yeom[a], Philipp Seegerer[a], Sebastian Lapuschkin[b], Simon Wiedemann[b],
Klaus-Robert Müller[a,c,d], Wojciech Samek[b,*]

[a]*Machine Learning Group, Technische Universität Berlin, 10587 Berlin, Germany*
[b]*Machine Learning Group, Fraunhofer Heinrich Hertz Institute, 10587 Berlin, Germany*
[c]*Department of Brain and Cognitive Engineering, Korea University, Seoul 136-713, Korea*
[d]*Max Planck Institut für Informatik, 66123 Saarbrücken, Germany*

## Abstract

The success of convolutional neural networks (CNNs) in various applications is accompanied by a significant increase in computation and parameter storage costs. Recent efforts to reduce these overheads involve pruning and compressing the weights of various layers while at the same time aiming to not sacrifice performance. In this paper, we propose a novel criterion for CNN pruning inspired by neural network interpretability: The most relevant elements, i.e. weights or filters, are automatically found using their relevance score in the sense of explainable AI (XAI). By that we for the first time link the two disconnected lines of interpretability and model compression research. We show in particular that our proposed method can efficiently prune transfer-learned CNN models where networks pre-trained on large corpora are adapted to specialized tasks. To this end, the method is evaluated on a broad range of computer vision datasets. Notably, our novel criterion is not only competitive or better compared to state-of-the-art pruning criteria when successive retraining is performed, but clearly outperforms these previous criteria in the common application setting where the data of the task to be transferred to are very scarce and no retraining is possible. Our method can iteratively compress the model while maintaining or even improving accuracy. At the same time, it has a computational cost in the order of gradient computation and is comparatively simple to apply without the need for tuning hyperparameters for pruning.

*Keywords:* Pruning, Layer-wise Relevance Propagation (LRP), Convolutional Neural Network (CNN), Interpretation of Models, Explainable AI (XAI)

## 1. Introduction

Deep CNNs have become an indispensable tool for a wide range of applications [1], such as image classification, speech recognition [2], natural language processing [3], chemistry

---

[4], neuroscience [5], medicine [6] and even are applied for playing games such as Go, poker or super smash bros [7]. They have achieved high predictive performance, at times even outperforming humans. Furthermore, in specialized domains where limited training data is available, e.g., due to the cost and difficulty of data generation (medical imaging from fMRI, EEG, PET etc.), transfer learning can improve the CNN performance by extracting the knowledge from the source tasks and applying it to a target task which has limited training data [8].

However, the high predictive performance of CNNs often comes at the expense of high storage and computational costs, which are directly related to energy expenditure demanded by the elaborate architecture of the fine-tuned network. These deep architectures are composed of millions of parameters to be trained, leading to overparameterization (i.e. having more parameters than training samples) of the model [9]. The run-times are typically dominated by the evaluation of convolutional layers, while dense layers are cheap but memory-heavy [10]. For instance, the VGG-16 model [11] has approximately 138 million parameters, taking up more than 500MB in storage space, and needs 15.5 billion floating-point operations (FLOPs) to classify a single image. Note that overparametrization is helpful for an efficient and successful training of neural networks, however, once the trained and well generalizing network structure is established, pruning can help to reduce redundancy while still maintaining good performance [12].

Reducing a model's storage requirements and computational cost becomes critical for a broader applicability, e.g., in embedded systems, autonomous agents, mobile devices, or edge devices [13, 14, 15, 16]. Neural network pruning has a decades long history with interest in both academia and industry [12, 14] aiming to detect and eliminate the subset of the network elements (i.e. weights or filters) that include connections and neurons less important w.r.t. the network's intended task. For network pruning, it is very crucial to decide how to quantify the "relevant" subset of the parameters in the current state for deletion. To address this issue, previous researches have proposed specific criteria based on for instance Taylor expansion, weight, gradient, etc. to reduce complexity and computations costs in the network and related works are introduced in Section 2.

From the practical point of view, the full capacity (in terms of weights and filters) of an overparameterized model may not be required, e.g., when

1. parts of the model lie dormant after training (i.e., are permanently "switched off"),

2. a user is not interested in the model's full array of possible outputs, which is a common scenario in transfer learning (e.g. the user only has use for 2 out of 10 available network outputs), or

3. a user lacks data and resources for fine-tuning and running the overparameterized model.

In these scenarios the redundant parts of the model will still occupy space on disk and in system memory, and simultaneously information will be propagated throughout the redundant parts of the model, consuming energy and increasing runtime. Thus, criteria able to stably and significantly reduce the computational complexity of deep neural networks broadly across applications would be highly versatile for all practitioners.

In this paper, we propose a novel pruning framework based on layer-wise relevance propagation (LRP) [17]. LRP was originally developed as an explanation method to assign importance scores, so called *relevance*, to the different input dimensions of a neural network that reflect the contribution of an input dimension to the models decision and has been applied to different fields of computer vision (e.g., [18, 19]). The relevance is backpropagated from the output to the input and hereby assigned to each element of the deep model. Since relevance scores are computed for every layer and neuron from the model output to the input, these relevance scores essentially reflect the importance of every single element of a model and its contribution to the information flow through the network – a natural candidate to be used as pruning criterion. The LRP criterion can be motivated theoretically through the concept of deep Taylor decomposition (c.f. [20]). Moreover, LRP is scalable and easy to apply, and has been implemented in software frameworks such as [21]. Furthermore, it has linear computational cost like backpropagation.

We systematically evaluate the compression efficacy of the LRP criterion compared to common pruning criteria on two different scenarios.

**Scenario 1**: We focus on pruning of pre-trained CNNs including subsequent fine-tuning. This is the usual setting in CNN pruning and requires a sufficiently large amount of data and computational power.

**Scenario 2**: We focus on another scenario, in which a model was pre-trained and needs to be transferred to a related problem but the data available for the new task is too scarce for a proper fine-tuning and/or the time consumption, computational power or energy consumption is constrained. Such transfer learning with restrictions is common in mobile or embedded applications.

Our experimental results on various benchmark datasets and two different popular CNN architectures show that the LRP criterion for pruning is more scalable and efficient and leads to better performance than existing criteria regardless of data types and model architectures if retraining is performed (scenario 1). Especially, if retraining is prohibited due to external constraints after pruning, the LRP criterion clearly outperforms previous criteria on all datasets (scenario 2). Finally we would like to note that our proposed pruning framework is not limited to LRP, but can be also used with other explanation techniques.

The rest of this paper is organized as follows: Section 2 summarizes related works for network compression and introduces the typical criteria for network pruning. Section 3 describes the framework and details of our approach. The experimental results are illustrated and discussed in Section 4. Section 5 gives conclusions and an outlook to future work.

## 2. Related Work

Research efforts have been made in the field of network compression and acceleration. For instance, network quantization methods have been proposed for storage space compression by decreasing the number of possible and unique values for the parameters [13, 22, 23]. Tensor decomposition approaches decompose network matrices into several smaller ones to estimate the informative parameters of the deep CNNs with low-rank approximation/factorization [24, 25, 26, 27]. More recently, [28] also propose a framework of architecture distillation based on layer-wise replacement, called LightweightNet for memory and time saving. Algorithms for

designing efficient models focus more on acceleration instead of compression by optimizing convolution operations or architectures directly [29, 30].

Network pruning approaches remove redundant or irrelevant elements — i.e. nodes, filters, or layers — from the model which are not critical for performance [12, 14, 31]. Network pruning is robust to various settings and easily and cheaply gives reasonable compression rates while not (or minimally) hurting the model accuracy. Also it can support both training from scratch and transfer learning from pre-trained models. Early works have shown that network pruning is effective in reducing network complexity and simultaneously addressing over-fitting problems. Current network pruning techniques make weights or channels sparse by removing non-informative connections and require an appropriate criterion for identifying which elements of the model are not relevant for solving a problem. Thus, it is very crucial to decide how to quantify the relevance of the parameters (i.e. weights or channels) from the current sample in stochastic learning for deletion without sacrificing predictive performance. In previous studies, pruning criteria have been proposed based on the magnitude of their 1) weights, 2) gradients, 3) Taylor expansion/derivative, and 4) other criteria, as described in the following section.

**Taylor expansion:** Early approaches towards neural network pruning — optimal brain damage [12] and optimal brain surgeon [32] — leveraged a second-order Taylor expansion based on the Hessian matrix of the loss function to select parameters for deletion. However, computing the inverse of Hessian is accompanied by high computational effort. The work of [33, 34] used a first-order Taylor expansion as a criterion to approximate the change of loss in the objective function as an effect of pruning away network elements.

**Gradient:** [35] proposes a sparsified back propagation approach for neural network training using the magnitude of the gradient to find essential and non-essential features in MLP and LSTM models, which can be used for pruning as a criterion. [36] proposed a hierarchical global pruning strategy by calculating the mean gradient of feature maps in each layer. It is adopted between the layers with similar sensitivity.

**Weight:** A recent trend is to prune redundant, non-informative weights in pre-trained CNN models. [37] and [38] proposed pruning the weights whose magnitude is below a certain threshold and to subsequently fine-tune with a $l_1$-norm regularization. This pruning strategy has been used on fully-connected layers and introduced sparse connections with BLAS libraries, supporting specialized hardware to achieve its acceleration. In the same context, Structured Sparsity Learning (SSL) added group sparsity regularization to penalize unimportant parameters by removing some weights [39]. [40] proposed a one-shot channel pruning method using the $l_1$ norm of weights for filter selection, provided that those channels with smaller weights always produce weaker activations. Recently, channel pruning alternatively used LASSO regression based channel selection and feature map reconstruction to prune filters [41]. [42] performed structured pruning in convolutional layers by considering strided sparsity of feature maps and kernels to avoid the need for custom hardware and uses particle filters to decide the importance of connections and paths. In contrast to previous pruning studies for deep deterministic models, [43] proposed a pruning approach for deep probabilistic models by using the mask of the weights. [44] proposed a fusion approach to combine with weight-based channel pruning and network quantization. [45] proposed a weight-based network pruning approach using the modified $l_{1/2}$ penalty to increase the sparsity of the pre-trained models. More recently, [46] proposed evolutionary paradigm for weight-based pruning and

gradient-based growing to reduce the network heuristically.

**Other criteria:** [47] proposed the Neuron Importance Score Propagation (NISP) algorithm to propagate the importance scores of final responses before the softmax, classification layer in the network. The method is based on — in contrast to our proposed metric — a layer-independent pruning process which does not consider global importance in the network. And [48, 49] proposed thinet that is a data-driven statistical channel pruning technique based on the statistics information computed from the next layer.

## 3. LRP-Based Network Pruning

A feedforward DNN consists of neurons established in a sequence of multiple layers of neurons, where each neuron receives the input data from the previous layer and propagates its output to every neuron in the next layer, using a non-linear mapping. Network pruning aims to sparsify the elements by eliminating weights or filters that are non-informative (according to a certain criterion). We specifically focus our experiments on transfer learning, where the parameters of a network pre-trained on a *source* domain is subsequently fine-tuned on a *target* domain, i.e., the final data or prediction task [8]. Here, the pruning procedure can be described as follows:

---

**Network Pruning**:

  *1.* Given a pre-trained model in the target domain

  *2.* Define a pruning criterion

  *3.* Repeatedly prune the network as follows:

    i. For each layer,

        a. For each element (weight/filter), evaluate the importance according to the pruning criterion (compute magnitudes)
        *b. optional:* Globally scale the magnitudes with regularization (e.g. $l_p$-norm)

    ii. Sort the magnitudes for all the layers throughout the network

    iii. Prune the least important elements and their inputs and outputs

    *iv. optional:* Further fine-tune to compensate performance degradation

  *4.* Stop pruning if the model is reduced to a desired amount of model size or performance

---

Even though most approaches use an identical process, choosing a suitable pruning criterion to quantify the importance of model parameters for deletion while minimizing performance drop (Step 3) is of critical importance, governing the success of the approach.

*3.1. Layer-wise Relevance Propagation*

In this paper, we propose a novel criterion for pruning neural network elements: the *relevance* quantity computed with LRP [17]. LRP decomposes a classification decision into

contributions called "relevances" of each network element to the overall classification score. When computed for the input dimensions of a CNN and visualized as a heatmap, these relevances highlight parts of the input that are important for the classification decision. LRP thus originally served as a tool to interpret non-linear learning machines and has been applied as such in various fields, amongst others for general image recognition [18], medical imaging [19] and natural language processing [50]. However, the direct linkage of the relevances to the classifier output makes LRP not only attractive for model explaining, but can also naturally serve as pruning criterion (see section 4.1).

The main characteristic of LRP is a backward pass through the network during which the network output is redistributed to all elements of the network in a layer-by-layer fashion. This backward pass is structurally similar to gradient backpropagation and has therefore a similar runtime. The redistribution is based on a *conservation principle* such that the relevances can immediately be interpreted as the contribution that an element makes to the network output, hence establishing a direct connection to the network output and thus its predictive performance. Therefore, as a pruning criterion, the method is efficient (similar runtime as gradient backpropagation) and easily scalable to generic network structures. Independent of the type of neural network layer — that is pooling, fully-connected, convolutional layers — LRP allows to quantify the importance of elements throughout the network [40].

### 3.2. LRP-based Pruning

The procedure of LRP-based pruning is summarized in Figure 1. In the first phase, a standard forward pass is performed by the network and the activations at each layer are collected. In the second phase, the score obtained at the output of the network $f(\mathbf{x})$ is propagated backwards through the network according to LRP propagation rules [17]. In the third phase, the current model is pruned by eliminating the irrelevant neurons/filters and is further fine-tuned.

LRP is based on a layer-wise conservation principle that allows the propagated quantity (e.g. relevance for a predicted class) to be preserved between neurons of two adjacent layers. Let $R_i^{(l)}$ be the relevance of neuron $i$ at layer $l$ and $R_j^{(l+1)}$ be the relevance of neuron $j$ at the next layer $l+1$. Stricter definitions of conservation that involve only subsets of neurons can further impose that relevance is locally redistributed in the lower layers and we define $R_{i \leftarrow j}^{(l)}$ as the share of $R_j^{(l+1)}$ that is redistributed to neuron $i$ in the lower layer. The conservation property always satisfies

$$\sum_i R_{i \leftarrow j}^{(l)} = R_j^{(l+1)} \ , \tag{1}$$

where the sum runs over all neurons $i$ of the (during inference) preceding layer $l$. When using relevance as a pruning criterion, this property helps to preserve its quantity layer-by-layer, regardless of hidden layer size and the number of iteratively pruned neurons for each layer. At each layer $l$, we can extract node $i$'s global importance as its attributed relevance $R_i^{(l)}$.

In this paper, we specifically adopt relevance quantities computed with the LRP-$\alpha_1\beta_0$-rule as pruning criterion. The LRP-$\alpha\beta$-rule was developed with feedforward-DNNs with ReLU activations in mind and assumes positive (pre-softmax) logit activations $f_{\text{logit}}(\mathbf{x}) > 0$ for decomposition. The rule has been shown to work well in practice in such a setting [51].
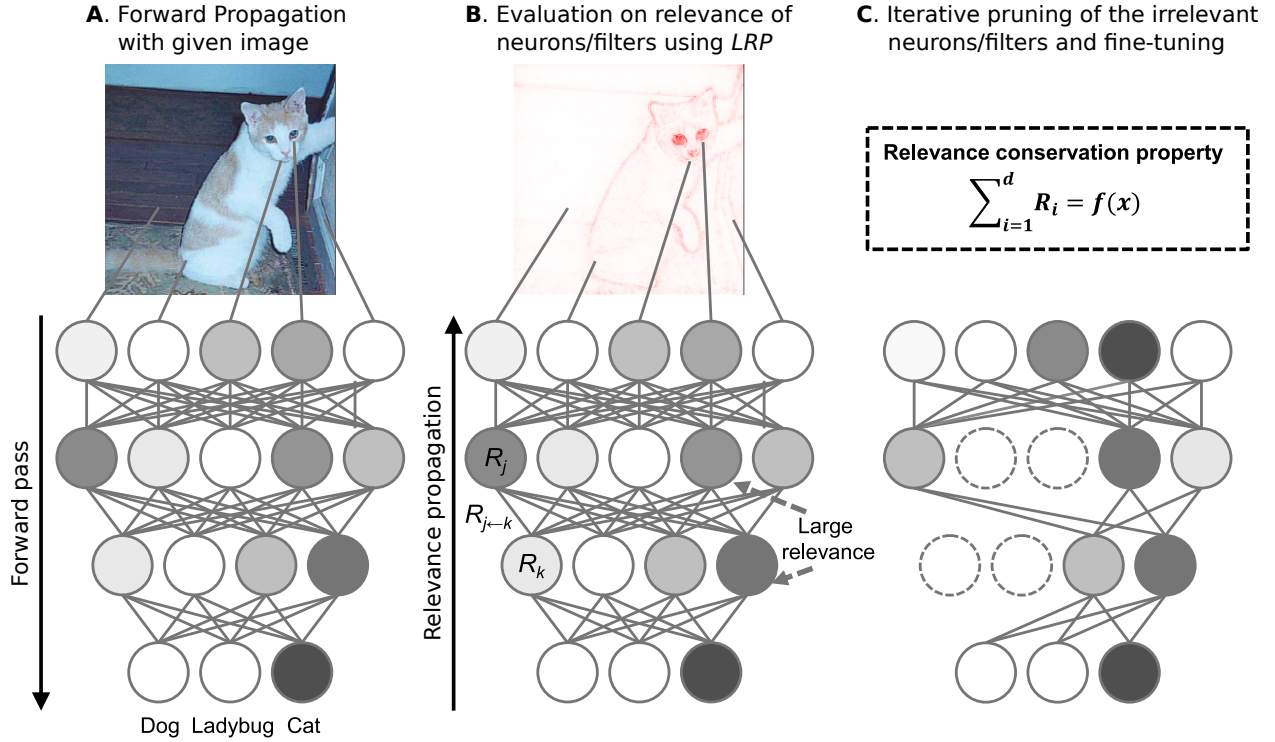
6

Figure 1: Illustration of LRP-based sequential process for pruning. **A.** Forward propagation of a given image (i.e. cat) through a pre-trained model. **B.** Evaluation on relevance for weights/filters using LRP, **C.** Iterative pruning by eliminating the least relevant elements (depicted by circles) and fine-tuning if necessary. The elements can be individual neurons, filters, or other arbitrary grouping of parameters, depending on the model architecture.

The propagation rule performs two separate relevance propagation steps per layer: one exclusively considering activatory parts of the forward propagated quantities (i.e. all $a_i^{(l)} w_{ij} > 0$) and another only processing the inhibitory parts ($a_i^{(l)} w_{ij} < 0$) which are subsequently merged in a sum with components weighted by $\alpha$ and $\beta$ (s.t. $\alpha + \beta = 1$) respectively.

By selecting $\alpha = 1$, the propagation rule expresses as

$$R_i^{(l)} = \sum_j \frac{a_i^{(l)} w_{ij}^+}{\sum_i a_i^{(l)} w_{ij}^+} R_j^{(l+1)} \quad , \tag{2}$$

where $R_i^{(l)}$ denotes relevance attributed to $i^{th}$ neuron at layer $l$, as an aggregation of downward-propagated relevance messages $R_{i \leftarrow j}^{(l,l+1)}$. The variables $w_{ij}^+$ and $a_i^{(l)}$ indicate the positive part of the weight filter and the (strictly positive) activation of the $i^{th}$ neuron at layer $l$, respectively. Note that a choice of $\alpha = 1$ only decomposes w.r.t. the parts of the inference signal *supporting* the model decision *for* the class of interest.

Note that Equation (2) is *locally conservative*, i.e. no quantity of relevance gets lost or injected during the distribution of $R_j$ where each term of the sum corresponds to a relevance message $R_{j \leftarrow k}$. For this reason, LRP has the following technical advantages over other pruning techniques such as gradient-based or activation-based methods:

1. Localized relevance conservation implicitly ensures layer-wisely regularized global redistribution of importances from each network element.

2. By summing relevance over within each (convolutional) filter channel, the LRP-based criterion is directly applicable as a measure of total relevance per node/filter, without requiring a post-hoc layer-wise renormalization, e.g., via $l_p$ norm.

3. The use of relevance scores is not restricted to a global application of pruning but can be easily applied to locally and (neuron- or filter-)group-wise constrained pruning without regularization. Different strategies for selecting (sub-)parts of the model might still be considered, e.g., applying different weightings/priorities for pruning different parts of the model: Should the aim of pruning be the reduction of FLOPs required during inference, one would prefer to focus on primarily pruning elements of the convolutional layers. In case the aim is a reduction of the memory requirement, pruning should focus on the fully-connected layers instead.

Algorithm 1 provides an overview of the LRP-based pruning process.

## 4. Experiments

In the following, we will first attempt to intuitively illuminate the properties of different pruning criteria — that is: weight magnitude, Taylor, gradient and LRP — at hand of a series of toy datasets.

---

**Algorithm 1** LRP-based pruning

---

1: **Input:** pre-trained model $f$,
2:       data $\mathbf{x}$, pruning rate $pr$,
3:       total number of weights/filters $m$
4: **while** pruning **do**
5:    Initialization of storage, $S$
6:    $R = f(\mathbf{x})$           // Forward processing
7:    **for** $layer$ **in** $f[:: -1]$ **do**
8:      $R = \mathrm{LRP}(layer, R)$    // LRP processing
9:      **if** $layer$ **is** '*conv*' **then**
10:        Compute filter-wise sum of $R$
11:      **end if**
12:      Stack $R$ in $S$
13:    **end for**
14:    **for all** $j$ **where** $S_j < m \times pr$ **do**
15:      Remove the $j^{th}$ weight/filter and its consecutive connections
16:    **end for**
17:    Fine-tuning (optional)
18: **end while**
19: **Output:** pruned model $f'$

---

We then show the effectiveness of the LRP criterion for pruning on current and widely-used Image recognition benchmark datasets — i.e. the Scene 15 [52], Event 8 [53], Cats and Dogs [54], Oxford Flower 102 [55], Cifar 10 [56], and ILSVRC2012 [57] datasets — and two pre-trained feed-forward deep neural network architectures, the AlexNet and the VGG-16.

The first scenario focuses specifically on pruning of pre-trained CNNs with subsequent fine-tuning, as it is common in pruning research [33]. We compare our method with several state-of-the-art criteria to demonstrate the effectiveness of LRP as a pruning criterion in CNNs. In the second scenario, we tested whether the proposed pruning criterion also works well if only a very limited number of samples is available for pruning the model. This is relevant in case of devices with limited computational power, energy and storage such as mobile devices or embedded applications.

*4.1. Pruning Toy Models*

Firstly, we systematically compare the properties and effectiveness of the different pruning criteria on toy datasets. In order to evaluate the criteria under the different data distributions, We tested four pruning criteria on three toy datasets ("moon", "circle" and "multi-class" ($k = 4$)), as generated using the respective functions from the Scikit-Learn machine learning toolkit for python [58] and shown in Figure 2. We generated a 2D dataset which consists of 1000 training samples and 5 test samples per class. We constructed the model as follows,

1. Linear ($2 \times 1000$), ReLU

2. Dropout with 0.5

3. Linear ($1000 \times 1000$), ReLU

4. Linear ($1000 \times 1000$), ReLU

5. Linear ($1000 \times k$)

During pruning, we pruned 1000 of 3000 neurons that have the least relevance for prediction according to each criterion. After pruning, we observed the change in the decision boundaries and re-evaluated classification accuracy using the original training samples across criteria.

On the first toy example ("moon" dataset, first row in Figure 2), the original model accuracy is 99.9%. After pruning, the accuracies are 99.6% (weight), 74.5% (Taylor), 76.9% (gradient), and 99.85% (LRP), respectively. On the second example ("circle" dataset, second row), the original accuracy is 100.0%. After pruning, the accuracies are 97.1% (weight), 97.25% (Taylor), 75.05% (gradient), and 100.0% (LRP), respectively. On the third example ("multi-class" dataset, third row), the original model accuracy is 94.95%. After pruning, the accuracies are 91.0% (weight), 85.15% (Taylor), 84.93% (gradient), and 91.3% (LRP), respectively. (See the Table 1)

These results show that among the considered pruning criteria, pruning with the relevance-criterion based on LRP best preserves the *effective* core of the example models, and thus prediction performance in every toy setting.

In contrast to the other heuristics-based pruning criteria, LRP directly relates relevance to the classification output, thus allows to safely remove the unimportant (w.r.t. classification)
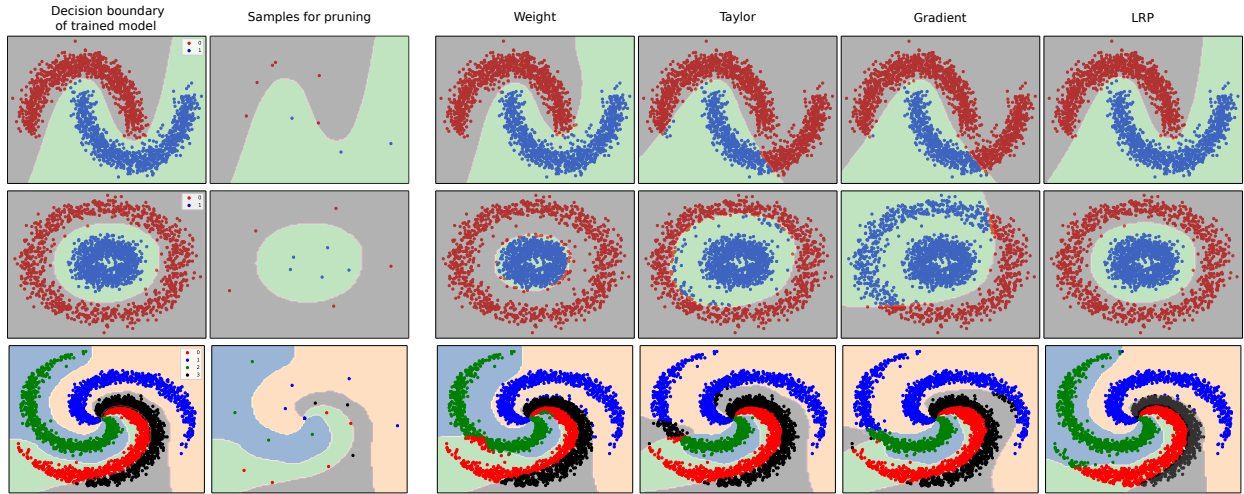
Figure 2: Performance comparison of the criteria on toy datasets (moon, circle, and multi-class dataset) *1st column*: scatter plot and decision boundary of trained model, *2nd column*: data samples for pruning, *3rd to 6th columns*: changed decision boundaries of different criteria.

Table 1: Comparison of accuracy in each criterion with pruned models on toy datasets (moon, circle, and multi-class dataset)

| Dataset | Original | Weight | Taylor expansion | Gradient | LRP |
|---|---|---|---|---|---|
| Moon | 0.9990 | 0.9960 | 0.7450 | 0.7690 | **0.9985** |
| Circle | 1.0000 | 0.9710 | 0.9725 | 0.7505 | **1.0000** |
| Multi-class | 0.9495 | 0.9100 | 0.8515 | 0.8493 | **0.9130** |

elements. Furthermore, Figure 2 shows that when pruning 1000 out of 3000 neurons, LRP-based pruning results in only minimal change in the decision boundary, compared to the other criteria. One can also clearly see that compared to weight- and LRP-based criteria, models pruned by gradient-based criteria misclassify a large part of samples.

## 4.2. Pruning Deep Image Classifiers for Large-scale Benchmark Data

The performance of the pruning criteria is evaluated on the CNNs VGG-16 and AlexNet pre-trained on ILSVRC2012 that are popular in model compression research [59]. VGG-16 consists of 13 convolutional layers with 4224 filters and 3 fully-connected layers and AlexNet contains 5 convolutional layers with 1552 filters and 3 fully-connected layers. In dense layers, there exist 8192 + N (# of classes) of neurons, respectively. In terms of complexity of the model, the pre-trained VGG-16 and AlexNet on ImageNet originally consist of 138.36/60.97 million of parameters and 154.7/7.27 GMAC for FLOPs, respectively [60]. Experiments are performed within the *PyTorch* and *torchvision* frameworks [61] under *Intel(R) Xeon(R) CPU E5-2660 2.20GHz* and *NVIDIA Tesla P100 with 12GB* for GPU processing. We evaluated the criteria on six public datasets (Scene 15 [52], Event 8 [53], Cats and Dogs [54], Oxford Flower 102 [55], Cifar 10 [56], and ILSVRC2012 [57]). For more detail on the datasets and the preprocessing, see Appendix Appendix A.
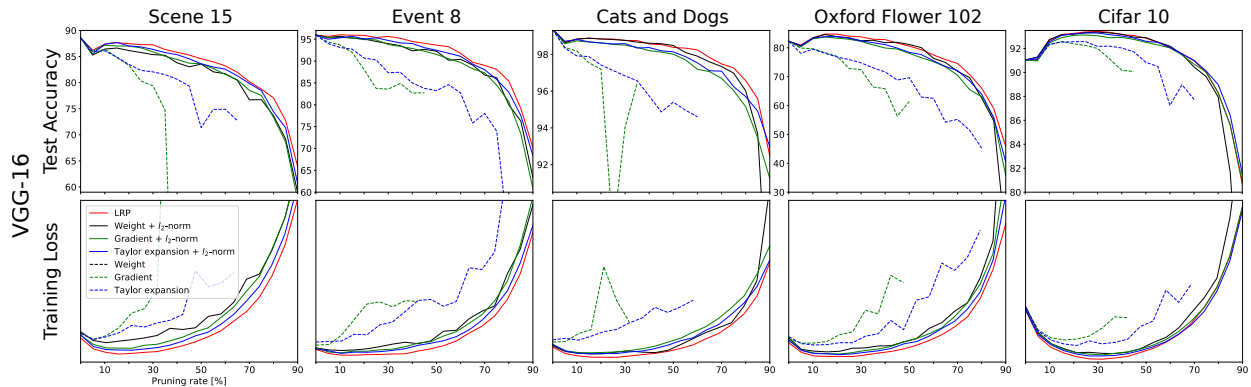
Figure 3: Performance comparison of training loss and test accuracy in different criteria as pruning rate increases on VGG-16 with five datasets.

We fine-tuned the model for 200 epochs with constant learning rate 0.001 and batch size of 20. We used the Stochastic Gradient Descent (SGD) optimizer with momentum of 0.9. In addition, we also apply dropout to the fully-connected layers with probability of 0.5. Fine-tuning and pruning are performed on the training set, while results are evaluated on each test dataset. Throughout the experiments, we iteratively prune 5% of all the filters in the network by eliminating neurons including their input and output connections. In Scenario 1, we subsequently fine-tune and re-evaluate the model to account for dependency across parameters and regain performance, as it is common.
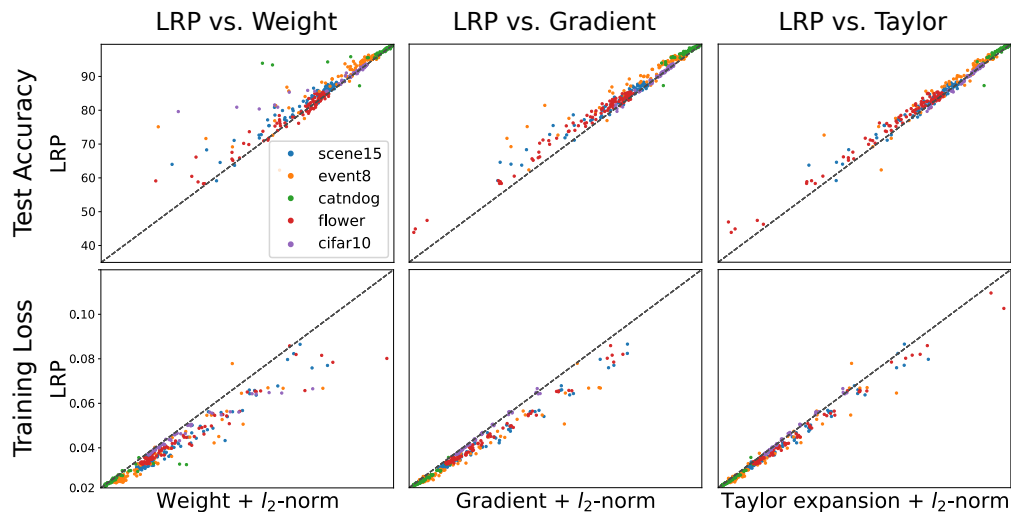


Figure 4: Performance comparison of the proposed method (i.e. LRP) and other criteria on VGG-16 with five datasets. Please note that weight cannot prune the model without $l_p$ regularization.

11

Table 2: Performance comparison of criteria (Weight, Taylor, Gradient with $l_2$-norm, and LRP) on **VGG-16** with five datasets.

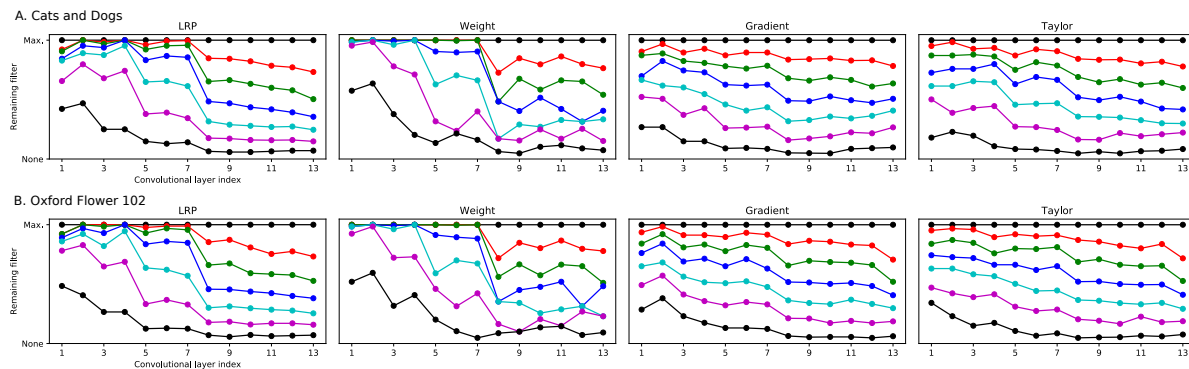| Dataset | Criterion | Training loss | Accuracy (%) | Pruning rate | Parameters ($\times 10^7$) | FLOPs (GMAC) |
|---|---|---|---|---|---|---|
| | VGG-16 | 0.0209 | 88.59 | | 119.61 | 15.50 |
| | Weight | 0.0227 | 82.07 | | 56.17 | 8.03 |
| Scene 15 | Taylor expansion | 0.0176 | 83.00 | 55 % | 53.10 | **4.66** |
| | Gradient | 0.0190 | 82.72 | | 53.01 | 4.81 |
| | LRP | **0.0162** | **83.99** | | **49.67** | 6.94 |
| | VGG-16 | 0.0085 | 95.95 | | 119.58 | 15.50 |
| | Weight | 0.0135 | 90.19 | | 56.78 | 8.10 |
| Event 8 | Taylor expansion | 0.0101 | 91.79 | 55 % | 48.48 | 5.21 |
| | Gradient | 0.0118 | 90.55 | | 50.25 | **5.05** |
| | LRP | **0.0083** | **93.29** | | **47.35** | 7.57 |
| | VGG-16 | 0.0019 | 99.36 | | 119.55 | 15.50 |
| | Weight | 0.0050 | 97.90 | | 47.47 | 7.02 |
| Cats and Dogs | Taylor expansion | 0.0051 | 97.54 | 60 % | 51.19 | 3.86 |
| | Gradient | 0.0057 | 97.19 | | 57.27 | **3.68** |
| | LRP | **0.0044** | **98.24** | | **43.75** | 6.49 |
| | VGG-16 | 0.0369 | 82.26 | | 119.96 | 15.50 |
| | Weigh | 0.0383 | 71.84 | | 39.34 | 5.48 |
| Oxford Flower 102 | Taylor expansion | 0.0327 | 72.11 | 70 % | 41.37 | **2.38** |
| | Gradient | 0.0354 | 70.53 | | 42.68 | 2.45 |
| | LRP | **0.0296** | **74.59** | | **37.54** | 4.50 |
| | VGG-16 | 0.0157 | 91.04 | | 119.59 | 15.50 |
| | Weight | 0.0183 | 93.36 | | **74.55** | 11.70 |
| Cifar 10 | Taylor expansion | 0.0176 | 93.29 | 30 % | 97.30 | **8.14** |
| | Gradient | 0.0180 | 93.05 | | 97.33 | 8.24 |
| | LRP | **0.0171** | **93.42** | | 89.20 | 9.93 |



Figure 5: Transition of the number of the remaining convolutional filters in different pruning rates (0 (black), 15 (red), 30 (green), 45 (blue), 60 (cyan), 75 (magenta), and 90% (black) from top to bottom) across criteria.

### 4.2.1. Scenario 1: Pruning with Fine-tuning

On the first scenario, we retrain the model after each iteration of pruning in order to regain lost performance. We then evaluate the performance of the different pruning criteria after each pruning-retraining-step.

That is, we quantify the importance of each filter by the magnitude of the respective criterion and iteratively prune the 5% of filters (w.r.t. the original number of filters in the model) rated least important in each pruning step. Then, we compute and record the training loss, test accuracy, number of remaining parameters and total estimated FLOPs. We assume

that the least important filters should have only little influence on the prediction and thus incur the lowest performance drop if they are removed from the network.

Figures 3 and A.8 depict the training loss (bottom) and test accuracy (top) as increasing the pruning rate in VGG-16 and AlexNet after fine-tuning for each dataset and each criterion. At each pruning iteration, we remove 5% of the weights/filters in the entire network based on the magnitude of the different criteria. It is observed that LRP achieves higher test accuracies as well as the lower training losses compared to other criteria for the entire range of pruning iterations on every dataset (see Figures 4 and A.7). These results demonstrate that the performance of LRP-based pruning is stable and independent of the chosen dataset. Apart from performance, regularization by layer must be a critical constraint which obstructs the expansion of the criterion toward several pruning strategies such as local pruning, global pruning, etc. Except for the LRP criterion, all criteria perform substantially worse without $l_p$ regularization compared to those with $l_p$ regularization and result in unexpected interruptions during the pruning process due to the biased redistribution of importance in the network.

Table 2 shows the predictive performance of the different criteria in terms of training loss, test accuracy, number of remaining parameters and FLOPs. Except for Cifar10, the highest compression rate (i.e. lowest number of parameters) could be achieved by the proposed LRP-based criterion (column "#Parameters" in Table 2). However, in terms of FLOPs, the proposed criterion only outperformed the weight criterion, but not the Taylor and Gradient criteria (see column "FLOPs" in Table 2). This is due to the fact that a reduction in number of FLOPs depends on the location where pruning is applied within the network: Figure 5 shows that the LRP and weight criteria focus the pruning on upper layers closer to the model output, whereas the Taylor and Gradient criteria focus more on the lower layers.

Throughout the pruning process usually a gradual decrease in performance can be observed. However, with the Event 8, Oxford Flower 102, Cifar 10 datasets, pruning leads to an initial performance increase, until a pruning rate of approx. 30% is reached. This behavior has been reported before in the literature and might stem from improvements of the model structure through elimination of filters related to classes in the source dataset (i.e., ILSVRC2012) that are not present in the target dataset anymore [62].

Table A.4 and Figure A.8 similarly show that LRP achieves the highest test accuracy and lowest training loss for nearly all pruning ratios with almost every dataset. Furthermore, as discussed in Section 3, due to the widespread use of LRP, we do not take into account dataset and model type, which could be a powerful advantage going forward.

Figure 5 shows the number of the remaining convolutional filters for each iteration. We observe that, on the one hand, as pruning rate increases, the convolutional filters in earlier layers that are associated with very generic features, such as edge and blob detectors, tend to generally be preserved as opposed to those in latter layers which are associated with abstract, task-specific features. On the other hand, the LRP- and weight-criterion first keep the filters in early layers in the beginning, but later aggressively prune filters near the input which now have lost functionality as input to later layers, compared to the gradient-based criteria such as gradient and Taylor-based approaches. Although gradient-based criteria also adopt the greedy layer-by-layer approach, we can see that gradient-based criteria pruned the less important filters almost uniformly across all the layers due to re-normalization of the criterion in each iteration. However, this result contrasts with previous gradient-based works [33, 35] that have shown that neurons deemed unimportant in earlier layers, contribute significantly compared
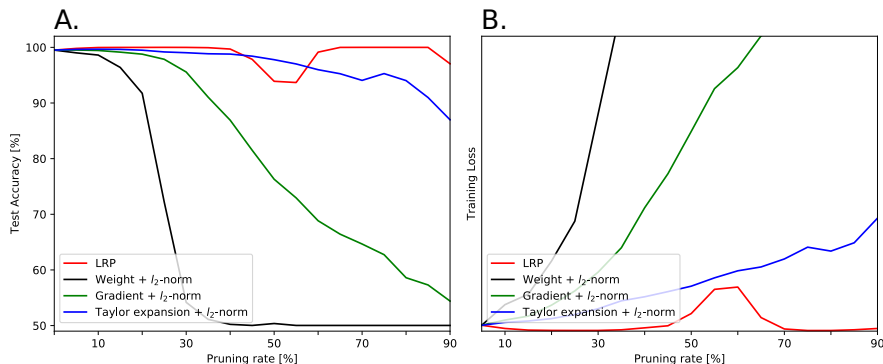
Figure 6: Performance comparison of pruning without fine-tuning for VGG-16 with small samples from Cats&Dogs dataset, as a means for domain adaption.

to neurons deemed important in latter layers. In contrast to this, LRP can efficiently preserve neurons in the early layers — as long as they serve a purpose — despite of iterative global pruning.

### 4.2.2. Scenario 2: Pruning without Fine-tuning

In this section, we evaluate whether pruning works well if only a (very) limited number of samples is available for quantifying the pruning criteria. To the best of our knowledge, there are no previous studies that show the performance of pruning approaches when acting w.r.t. very small amounts of data. With large amounts of data available (and even though we can expect reasonable performance after pruning), an iterative pruning and fine-tuning procedure of the network can amount to a very time consuming and computationally heavy process. From a practical point of view, this issue becomes a significant problem, e.g. with limited computational resources (mobile devices or in general; consumer-level hardware) and reference data (private photo collections), where capable and effective one-shot pruning approaches are desired and only little leeway (or none at all) for post-pruning fine-tuning strategies is available.

To investigate whether pruning is possible also in this scenario, we performed experiments with a relatively small number of data on the 1) Cats and Dogs and 2) ILVSRC 2012 datasets. On the Cats and Dogs dataset, we only used 10 samples each from the "cat" and "dog" classes to prune the (on ImageNet) pre-trained VGG-16 network with the goal of domain/dataset adaption. The binary classification (i.e. "cat" vs. "dog") is a subtask of ImageNet and corresponding output neurons can be identified by its WordNet [1] associations. On the ILSVRC 2012 dataset, we randomly chose $k = 3$ classes for model specialization, selected $n = 10$ images per class from the training set and used them to compare the different pruning criteria. For each criterion, we used the same selection of classes and samples. In these experiments, *we do not fine-tune* the models after each pruning iteration, in contrast to Section 4.2.1. Performance is averaged over 20 random selections of classes and samples to account for randomness. Please note that before pruning, we firstly reconstructed output

---

[1]http://www.image-net.org/archive/wordnet.is_a.txt

layers which has the size of N × k by eliminating the redundant network outputs of 1000 - $k$.

Furthermore, as our target datasets are relatively small and only have an extremely reduced set of target classes, the pruned models would still be very heavy w.r.t. memory requirements if the pruning process would be limited to the convolutional layers, as in Section 4.2.1.

More specifically, while convolutional layers dominantly constitute the source of computation cost, fully connected layers are proven to be more redundant [40]. In this respect, we applied pruning procedures in both fully connected layers and convolutional layers.

For pruning, we iterate a sequence of first pruning filters from the convolutional layers, followed by a step of pruning filters/neurons from the model's fully connected layers.

Table 3 indicates the performances of each criterion for classifying a small number of classes ($k = 3$) from the ILSVRC2012 dataset. During pruning at fully-connected layers, no significant difference across different pruning ratios can be observed. Without further fine-tuning, pruning weights/filters at the fully connected layers can retain performance efficiently.

However, there is a subtle difference between LRP and other criteria with increasing pruning ratio of convolutional layers (LRP vs. Taylor with $l_2$-norm: up to of 9.6 %, LRP vs. gradient with $l_2$-norm: up to 28.0 %, LRP vs. weight with $l_2$-norm: up to 27.1 %).

Moreover, pruning convolutional layers needs to be carefully managed compared to pruning fully connected layers. We can observe that LRP is applicable for pruning any layer type (i.e. fully connected, convolutional, pooling, etc.) efficiently. Additionally, as mentioned in 3.1, our method can be applied to general network architectures because it it can automatically measure the importance of weights or filters in a global (network-wise) context without further normalization.

Figure 6 shows the test accuracy and training loss as function of the pruning ratio, in context a domain adaption task towards the Cats&Dogs dataset. As the pruning ratio increases, we can see that even without fine-tuning, using LRP as pruning criterion can keep the test accuracy and training loss relatively not only stable, but close to 100% and 0 respectively, given the extreme scarcity of data in this experiment. In contrast, the performance decreases significantly when using the $l_2$-norm based criteria (i.e. weight, gradient and Taylor). Initially, the performance is even slightly increasing when pruning with LRP. During iterative pruning, unexpected changes in accuracy with LRP (2 of 20 iterations) have been shown around 50 - 55% pruning ratio, but accuracy is regained quickly again. By pruning 90% of convolutional filters in the network using our proposed method, we can have 1) greatly reduced computational cost, 2) faster forward and backward processing, and 3) a lighter model even in the small sample case, all while adapting an off-the-shelf pre-trained ImageNet model towards a dog-vs.-cat classification task.

## 5. Conclusion

Modern CNNs typically have a high capacity with millions of parameters as this allows to obtain good optimization results in the training process. After training, however, high inference costs remain, despite the fact that the number of effective parameters in the deep model is actually significantly lower (see e.g. [63]). To alleviate this, pruning aims at compressing and accelerating the given models without sacrificing much predictive performance. In this paper,

Table 3: Performance of pruning of convolutional and fully connected layers on VGG-16 without fine-tuning for a random subset of the classes from ILSVRC 2012 ($k = 3$) based on LRP(left top), Taylor(right top), gradient(left bottom), and weight(right bottom).

**LRP (left top)**

| % of Remaining Dense Layers | % of Remaining Convolutional Layers | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | 100 | 95 | 90 | 85 | 80 | 75 | 70 | 65 | 60 | 55 |
| 100 | 97.33 | 96.56 | 94.94 | 91.94 | 87.83 | 83.61 | 78.39 | 69.28 | 60.72 | 54.56 |
| 95 | 97.06 | 96.44 | 95.06 | 91.94 | 87.83 | 83.39 | 78.33 | 69.28 | 60.78 | 54.67 |
| 90 | 96.94 | 96.39 | 94.89 | 92.00 | 88.00 | 83.72 | 78.17 | 69.11 | 60.94 | 54.83 |
| 85 | 97.06 | 96.44 | 94.83 | 91.67 | 87.78 | 83.56 | 78.33 | 69.50 | 61.28 | 55.06 |
| 80 | 96.89 | 96.44 | 94.72 | 91.78 | 87.78 | 83.28 | 78.06 | 69.11 | 61.33 | 55.39 |
| 75 | 96.83 | 96.22 | 94.72 | 91.72 | 87.67 | 83.06 | 78.00 | 69.00 | 61.28 | 55.28 |
| 70 | 96.89 | 96.44 | 94.72 | 91.78 | 87.89 | 83.17 | 78.00 | 69.28 | 61.39 | 55.06 |
| 65 | 96.56 | 96.17 | 94.78 | 91.61 | 87.89 | 83.39 | 78.17 | 69.33 | 61.72 | 55.22 |
| 60 | 96.50 | 96.11 | 94.50 | 91.50 | 87.94 | 84.06 | 78.22 | 70.11 | 61.72 | 55.44 |
| 55 | 96.44 | 96.06 | 94.22 | 91.44 | 87.72 | 83.83 | 78.50 | 70.22 | 62.50 | 55.67 |
| 50 | 96.50 | 95.83 | 93.94 | 91.11 | 87.83 | 83.94 | 78.39 | 70.11 | 62.83 | 56.56 |
| 45 | 96.39 | 95.78 | 94.00 | 90.89 | 87.61 | 84.11 | 78.11 | 70.28 | 63.11 | 56.72 |
| 40 | 96.33 | 95.50 | 93.94 | 90.83 | 87.33 | 83.17 | 78.11 | 69.67 | 63.50 | 57.78 |
| 35 | 95.78 | 95.44 | 93.50 | 90.72 | 87.06 | 83.39 | 77.83 | 69.39 | 63.28 | 57.72 |
| 30 | 95.94 | 95.50 | 93.44 | 90.33 | 86.44 | 82.56 | 77.17 | 69.44 | 62.33 | 56.56 |
| 25 | 95.72 | 95.22 | 93.11 | 89.33 | 84.50 | 80.94 | 76.50 | 68.22 | 61.56 | 55.94 |
| 20 | 95.83 | 94.83 | 92.56 | 89.06 | 84.89 | 79.67 | 75.67 | 66.78 | 59.72 | 54.28 |
| 15 | 95.50 | 94.39 | 91.72 | 87.78 | 83.83 | 78.89 | 73.94 | 64.28 | 56.67 | 52.06 |

**Taylor (right top)**

| % of Remaining Dense Layers | % of Remaining Convolutional Layer | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | 100 | 95 | 90 | 85 | 80 | 75 | 70 | 65 | 60 | 55 |
| 100 | 97.33 | 95.74 | 92.62 | 89.59 | 85.18 | 80.10 | 70.62 | 62.62 | 53.90 | 50.41 |
| 95 | 97.33 | 95.74 | 92.62 | 89.59 | 85.13 | 80.10 | 70.62 | 62.56 | 54.00 | 50.31 |
| 90 | 97.33 | 95.74 | 92.56 | 89.59 | 85.13 | 80.21 | 70.72 | 62.56 | 53.95 | 50.31 |
| 85 | 97.13 | 95.79 | 92.56 | 89.59 | 85.13 | 80.10 | 70.62 | 62.62 | 53.79 | 50.36 |
| 80 | 97.13 | 95.49 | 92.62 | 89.69 | 85.13 | 79.95 | 70.62 | 62.62 | 53.64 | 50.31 |
| 75 | 96.97 | 95.49 | 92.67 | 89.64 | 85.03 | 79.95 | 70.62 | 62.67 | 53.69 | 50.21 |
| 70 | 97.08 | 95.44 | 92.62 | 89.79 | 85.28 | 80.15 | 70.51 | 62.62 | 53.79 | 50.21 |
| 65 | 97.03 | 95.49 | 92.97 | 89.23 | 85.38 | 80.05 | 69.95 | 62.46 | 53.44 | 50.10 |
| 60 | 96.97 | 95.69 | 92.92 | 89.18 | 84.97 | 79.85 | 70.05 | 62.46 | 53.64 | 50.15 |
| 55 | 96.82 | 95.59 | 93.33 | 89.54 | 85.44 | 80.10 | 70.36 | 62.10 | 53.38 | 49.79 |
| 50 | 96.92 | 95.38 | 93.38 | 89.74 | 85.95 | 80.15 | 70.67 | 62.05 | 53.49 | 49.64 |
| 45 | 96.72 | 95.64 | 93.74 | 89.79 | 84.72 | 80.21 | 70.21 | 61.38 | 53.59 | 49.54 |
| 40 | 96.56 | 95.44 | 93.90 | 89.85 | 85.38 | 80.77 | 69.85 | 62.15 | 53.85 | 49.08 |
| 35 | 96.72 | 95.64 | 94.10 | 89.64 | 85.90 | 80.92 | 70.00 | 62.21 | 53.38 | 49.33 |
| 30 | 96.46 | 95.33 | 93.85 | 89.59 | 86.41 | 81.44 | 70.10 | 61.33 | 53.23 | 48.05 |
| 25 | 96.26 | 95.74 | 93.59 | 90.00 | 86.15 | 81.49 | 71.44 | 59.90 | 53.23 | 48.77 |
| 20 | 96.36 | 94.36 | 92.10 | 89.44 | 83.69 | 81.69 | 70.92 | 59.59 | 51.44 | 46.77 |
| 15 | 94.92 | 93.69 | 90.72 | 88.00 | 86.05 | 79.23 | 68.05 | 57.28 | 49.95 | 45.95 |

**gradient (left bottom)**

| % of Remaining Dense Layer | % of Remaining Convolutional Layer | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | 100 | 95 | 90 | 85 | 80 | 75 | 70 | 65 | 60 | 55 |
| 100 | 97.33 | 95.74 | 91.28 | 83.28 | 74.56 | 65.69 | 56.77 | 49.64 | 44.41 | 40.31 |
| 95 | 97.28 | 95.74 | 91.28 | 83.23 | 74.51 | 65.69 | 56.72 | 49.54 | 44.46 | 40.31 |
| 90 | 97.23 | 95.74 | 91.23 | 83.13 | 74.46 | 65.74 | 56.67 | 49.64 | 44.46 | 40.31 |
| 85 | 97.28 | 95.74 | 91.03 | 83.08 | 74.46 | 65.64 | 56.62 | 49.64 | 44.46 | 40.31 |
| 80 | 96.92 | 95.64 | 91.08 | 83.28 | 74.41 | 65.59 | 56.67 | 49.74 | 44.46 | 40.31 |
| 75 | 96.97 | 95.44 | 91.13 | 83.64 | 74.10 | 65.33 | 56.67 | 49.69 | 44.36 | 40.31 |
| 70 | 96.82 | 95.44 | 91.33 | 83.44 | 74.26 | 65.49 | 56.41 | 49.85 | 44.51 | 40.31 |
| 65 | 96.92 | 95.28 | 91.18 | 83.08 | 73.69 | 65.18 | 55.69 | 49.49 | 44.46 | 40.26 |
| 60 | 97.03 | 95.18 | 91.03 | 83.13 | 73.23 | 64.92 | 55.44 | 48.82 | 44.41 | 40.10 |
| 55 | 96.82 | 95.03 | 91.03 | 82.92 | 73.44 | 64.97 | 55.44 | 48.41 | 44.46 | 39.95 |
| 50 | 97.03 | 95.59 | 91.18 | 83.08 | 72.87 | 64.82 | 55.03 | 48.21 | 43.85 | 40.41 |
| 45 | 96.62 | 94.97 | 91.13 | 82.77 | 71.90 | 64.46 | 55.03 | 47.69 | 43.49 | 39.59 |
| 40 | 96.56 | 95.13 | 90.97 | 82.31 | 71.85 | 63.03 | 53.49 | 46.41 | 42.41 | 39.03 |
| 35 | 96.46 | 95.08 | 90.46 | 81.49 | 71.64 | 62.62 | 52.87 | 46.21 | 42.10 | 38.56 |
| 30 | 96.56 | 94.21 | 90.92 | 80.72 | 70.41 | 61.33 | 51.79 | 45.23 | 40.92 | 37.18 |
| 25 | 95.38 | 94.97 | 90.26 | 79.74 | 69.44 | 60.10 | 51.13 | 44.82 | 40.41 | 36.56 |
| 20 | 95.64 | 94.31 | 89.08 | 77.33 | 68.31 | 57.90 | 47.64 | 43.59 | 38.92 | 36.15 |
| 15 | 93.69 | 93.74 | 88.26 | 74.77 | 67.03 | 55.90 | 46.10 | 43.38 | 38.26 | 35.38 |

**weight (right bottom)**

| % of Remaining Dense Layer | % of Remaining Convolutional Layer | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | 100 | 95 | 90 | 85 | 80 | 75 | 70 | 65 | 60 | 55 |
| 100 | 97.33 | 96.97 | 96.10 | 93.64 | 88.31 | 77.54 | 61.74 | 53.08 | 48.31 | 40.15 |
| 95 | 97.33 | 96.82 | 95.85 | 93.38 | 87.90 | 77.79 | 62.56 | 52.82 | 47.08 | 42.21 |
| 90 | 97.44 | 96.82 | 95.79 | 93.33 | 88.15 | 78.62 | 63.90 | 52.21 | 45.49 | 42.46 |
| 85 | 97.28 | 96.92 | 95.79 | 92.72 | 87.90 | 78.97 | 65.79 | 51.38 | 46.00 | 41.23 |
| 80 | 97.23 | 96.97 | 95.64 | 92.82 | 88.51 | 78.87 | 65.85 | 50.00 | 45.64 | 40.15 |
| 75 | 96.82 | 97.03 | 95.85 | 92.87 | 87.33 | 77.69 | 65.23 | 49.13 | 43.64 | 39.59 |
| 70 | 97.18 | 96.62 | 95.54 | 92.56 | 87.69 | 77.74 | 65.49 | 46.82 | 42.62 | 40.31 |
| 65 | 96.82 | 96.56 | 95.18 | 92.41 | 86.56 | 78.05 | 64.56 | 46.46 | 42.05 | 38.77 |
| 60 | 96.31 | 96.31 | 95.38 | 91.95 | 85.79 | 77.03 | 62.56 | 45.18 | 42.51 | 37.44 |
| 55 | 95.44 | 95.85 | 94.46 | 91.28 | 85.28 | 74.51 | 62.10 | 44.51 | 41.59 | 38.31 |
| 50 | 94.97 | 95.08 | 93.69 | 89.90 | 83.08 | 71.69 | 60.67 | 43.28 | 39.90 | 37.33 |
| 45 | 94.10 | 94.46 | 93.03 | 89.49 | 82.56 | 71.23 | 59.44 | 43.23 | 43.13 | 37.03 |
| 40 | 93.64 | 94.26 | 92.97 | 90.00 | 81.08 | 70.31 | 57.23 | 42.62 | 39.33 | 38.62 |
| 35 | 92.82 | 93.69 | 92.87 | 88.92 | 78.62 | 68.82 | 58.26 | 46.15 | 40.62 | 37.90 |
| 30 | 90.62 | 91.64 | 90.72 | 87.08 | 76.10 | 66.05 | 55.18 | 43.64 | 42.05 | 36.72 |
| 25 | 90.05 | 91.90 | 90.00 | 85.49 | 74.56 | 64.72 | 52.05 | 42.15 | 40.56 | 33.69 |
| 20 | 87.69 | 88.62 | 86.31 | 79.79 | 70.41 | 63.13 | 51.69 | 41.08 | 39.28 | 35.54 |
| 15 | 85.69 | 85.74 | 82.82 | 76.72 | 62.21 | 57.38 | 48.92 | 40.36 | 37.18 | 37.85 |

we have proposed a novel criterion for iterative pruning of CNNs based on the explanation method LRP, linking for the first time two so far disconnected lines of research. LRP has a clearly defined meaning, namely the contribution of an individual network element, i.e. weight or filter, to the network output. Removing elements according to low LRP scores thus means discarding all aspects in the model that do not contribute relevance to its decision making. Hence, as a criterion, the computed relevance scores can easily and cheaply give efficient compression rates without further postprocessing, such as per-layer normalization. Besides, technically LRP is scalable to general network structures and its computational cost

is similar to the one of a gradient backward pass.

In our experiments, the LRP criterion has shown favorable compression performance on a variety of datasets both with and without retraining after pruning. Especially when pruning without retraining, our results for small datasets suggest that the LRP criterion outperforms the state of the art and therefore, its application is especially recommended in transfer learning settings where only a small target dataset is available.

In addition to pruning, the same method can be used to visually interpret the model and explain individual decisions as intuitive relevance heatmaps. Therefore, in future work, we propose to use these heatmaps to elucidate and explain which image features are most strongly affected by pruning to additionally avoid that the pruning process leads to undesired Clever Hans phenomena [18]. Also, we would like to further investigate LRP-based pruning with other modern neural network architectures such as GoogLeNet, ResNet, DenseNet, etc. in order to see how to effectively prune the weights/filters in the residual dense blocks of the desired networks.

## Acknowledgements

## References

[1] J. Gu, Z. Wang, J. Kuen, L. Ma, A. Shahroudy, B. Shuai, T. Liu, X. Wang, G. Wang, J. Cai, T. Chen, Recent advances in convolutional neural networks, Pattern Recognition 77 (2018) 354–377.

[2] O. Abdel-Hamid, A. Mohamed, H. Jiang, L. Deng, G. Penn, D. Yu, Convolutional neural networks for speech recognition, IEEE/ACM Transactions on Audio, Speech, and Language Processing 22 (2014) 1533–1545.

[3] T. Young, D. Hazarika, S. Poria, E. Cambria, Recent trends in deep learning based natural language processing, IEEE Computational Intelligence Magazine 13 (2018) 55–75.

[4] K. T. Schütt, F. Arbabzadah, S. Chmiela, K. R. Müller, A. Tkatchenko, Quantum-chemical insights from deep tensor neural networks, Nature Communications 8 (2017) 13890.

[5] R. T. Schirrmeister, J. T. Springenberg, L. D. J. Fiederer, M. Glasstetter, K. Eggensperger, M. Tangermann, F. Hutter, W. Burgard, T. Ball, Deep learning with convolutional neural networks for eeg decoding and visualization, Human Brain Mapping 38 (2017) 5391–5420.

[6] P. Jurmeister, M. Bockmayr, P. Seegerer, T. Bockmayr, D. Treue, G. Montavon, C. Vollbrecht, A. Arnold, D. Teichmann, K. Bressem, U. Schüller, M. von Laffert, K. Müller, D. Capper, F. Klauschen, Machine learning analysis of dna methylation profiles distinguishes primary lung squamous cell carcinomas from head and neck metastases, Science Translational Medicine 11 (2019).

[7] D. Silver, A. Huang, C. J. Maddison, A. Guez, L. Sifre, G. van den Driessche, J. Schrittwieser, I. Antonoglou, V. Panneershelvam, M. Lanctot, S. Dieleman, D. Grewe, J. Nham, N. Kalchbrenner, I. Sutskever, T. P. Lillicrap, M. Leach, K. Kavukcuoglu, T. Graepel, D. Hassabis, Mastering the game of go with deep neural networks and tree search, Nature 529 (2016) 484–489.

[8] S. J. Pan, Q. Yang, A survey on transfer learning, IEEE Transactions on Knowledge and Data Engineering 22 (2010) 1345–1359.

[9] M. Denil, B. Shakibi, L. Dinh, M. Ranzato, N. de Freitas, Predicting parameters in deep learning, in: Advances in Neural Information Processing Systems (NIPS), 2013, pp. 2148–2156.

[10] V. Sze, Y. Chen, T. Yang, J. S. Emer, Efficient processing of deep neural networks: A tutorial and survey, Proceedings of the IEEE 105 (2017) 2295–2329.

[11] K. Simonyan, A. Zisserman, Very deep convolutional networks for large-scale image recognition, in: International Conference on Learning Representations, (ICLR), 2015.

[12] Y. LeCun, J. S. Denker, S. A. Solla, Optimal brain damage, in: Advances in Neural Information Processing Systems (NIPS), 1989, pp. 598–605.

[13] S. Wiedemann, K.-R. Müller, W. Samek, Compact and computationally efficient representation of deep neural networks, IEEE Transactions on Neural Networks and Learning Systems (2019).

[14] Y. Cheng, D. Wang, P. Zhou, T. Zhang, Model compression and acceleration for deep neural networks: The principles, progress, and challenges, IEEE Signal Processing Magazine 35 (2018) 126–136.

[15] J. Pei, L. Deng, S. Song, M. Zhao, Y. Zhang, S. Wu, G. Wang, Z. Zou, Z. Wu, W. He, et al., Towards artificial general intelligence with hybrid tianjic chip architecture, Nature 572 (2019) 106–111.

[16] Y. Tu, Y. Lin, Deep neural network compression technique towards efficient digital signal modulation recognition in edge device, IEEE Access 7 (2019) 58113–58119.

[17] S. Bach, A. Binder, G. Montavon, F. Klauschen, K.-R. Müller, W. Samek, On pixel-wise explanations for non-linear classifier decisions by layer-wise relevance propagation, PLoS ONE 10 (2015) e0130140.

[18] S. Lapuschkin, S. Wäldchen, A. Binder, G. Montavon, W. Samek, K.-R. Müller, Unmasking clever hans predictors and assessing what machines really learn, Nature Communications 10 (2019) 1096.

[19] M. Hägele, P. Seegerer, S. Lapuschkin, M. Bockmayr, W. Samek, F. Klauschen, K.-R. Müller, A. Binder, Resolving challenges in deep learning-based analyses of histopathological images using explanation methods, arXiv:1908.06943 (2019).

[20] G. Montavon, W. Samek, K.-R. Müller, Methods for interpreting and understanding deep neural networks, Digital Signal Processing 73 (2018) 1–15.

[21] M. Alber, S. Lapuschkin, P. Seegerer, M. Hägele, K. T. Schütt, G. Montavon, W. Samek, K.-R. Müller, S. Dähne, P.-J. Kindermans, iNNvestigate neural networks!, Journal of Machine Learning Research 20 (2019) 1–8.

[22] F. Tung, G. Mori, Deep neural network compression by in-parallel pruning-quantization, IEEE Transactions on Pattern Analysis and Machine Intelligence (2018) 1–1.

[23] M. Courbariaux, Y. Bengio, J. David, Binaryconnect: Training deep neural networks with binary weights during propagations, in: Advances in Neural Information Processing Systems (NIPS), 2015, pp. 3123–3131.

[24] E. L. Denton, W. Zaremba, J. Bruna, Y. LeCun, R. Fergus, Exploiting linear structure within convolutional networks for efficient evaluation, in: Advances in Neural Information Processing Systems (NIPS), 2014, pp. 1269–1277.

[25] H. Ding, K. Chen, Q. Huo, Compressing CNN-DBLSTM models for OCR with teacher-student learning and tucker decomposition, Pattern Recognition 96 (2019).

[26] X. Xiao, L. Jin, Y. Yang, W. Yang, J. Sun, T. Chang, Building fast and compact convolutional neural networks for offline handwritten chinese character recognition, Pattern Recognition 72 (2017) 72–81.

[27] K. Guo, X. Xie, X. Xu, X. Xing, Compressing by learning in a low-rank and sparse decomposition form, IEEE Access 7 (2019) 150823–150832.

[28] T. Xu, P. Yang, X. Zhang, C. Liu, Lightweightnet: Toward fast and lightweight convolutional neural networks via architecture distillation, Pattern Recognition 88 (2019) 272–284.

[29] X. Zhang, X. Zhou, M. Lin, J. Sun, Shufflenet: An extremely efficient convolutional neural network for mobile devices, in: IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2018, pp. 6848–6856.

[30] A. G. Howard, M. Zhu, B. Chen, D. Kalenichenko, W. Wang, T. Weyand, M. Andreetto, H. Adam, Mobilenets: Efficient convolutional neural networks for mobile vision applications, arXiv:1704.04861 (2017).

[31] P. Molchanov, A. Mallya, S. Tyree, I. Frosio, J. Kautz, Importance estimation for neural network pruning, in: IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2019, pp. 11264–11272.

[32] B. Hassibi, D. G. Stork, Second order derivatives for network pruning: Optimal brain surgeon, in: Advances in Neural Information Processing Systems (NIPS), 1992, pp. 164–171.

[33] P. Molchanov, S. Tyree, T. Karras, T. Aila, J. Kautz, Pruning convolutional neural networks for resource efficient transfer learning, in: Proceedings of the International Conference on Learning Representations (ICLR), 2017.

[34] C. Yu, J. Wang, Y. Chen, X. Qin, Transfer channel pruning for compressing deep domain adaptation models, International Journal of Machine Learning and Cybernetics 10 (2019) 3129–3144.

[35] X. Sun, X. Ren, S. Ma, H. Wang, meprop: Sparsified back propagation for accelerated deep learning with reduced overfitting, in: International Conference on Machine Learning (ICML), 2017, pp. 3299–3308.

[36] C. Liu, H. Wu, Channel pruning based on mean gradient for accelerating convolutional neural networks, Signal Processing 156 (2019) 84–91.

[37] S. Han, J. Pool, J. Tran, W. J. Dally, Learning both weights and connections for efficient neural network, in: Advances in Neural Information Processing Systems (NIPS), 2015, pp. 1135–1143.

[38] S. Han, X. Liu, H. Mao, J. Pu, A. Pedram, M. A. Horowitz, W. J. Dally, EIE: efficient inference engine on compressed deep neural network, in: International Symposium on Computer Architecture (ISCA), 2016, pp. 243–254.

[39] W. Wen, C. Wu, Y. Wang, Y. Chen, H. Li, Learning structured sparsity in deep neural networks, in: Advances in Neural Information Processing Systems (NIPS), 2016, pp. 2074–2082.

[40] H. Li, A. Kadav, I. Durdanovic, H. Samet, H. P. Graf, Pruning filters for efficient convnets, in: International Conference on Learning Representations, (ICLR), 2017.

[41] Y. He, X. Zhang, J. Sun, Channel pruning for accelerating very deep neural networks, in: IEEE International Conference on Computer Vision (ICCV), 2017, pp. 1398–1406.

[42] S. Anwar, K. Hwang, W. Sung, Structured pruning of deep convolutional neural networks, ACM Journal on Emerging Technologies in Computing Systems (JETC) 13 (2017) 32:1–32:18.

[43] C. Zhang, Q. Zhao, C. L. P. Chen, W. Liu, Deep compression of probabilistic graphical networks, Pattern Recognition 96 (2019).

[44] J. Gan, W. Wang, K. Lu, Compressing the cnn architecture for in-air handwritten chinese character recognition, Pattern Recognition Letters 129 (2020) 190 – 197.

[45] J. Chang, J. Sha, Prune deep neural networks with the modified $L_{1/2}$ penalty, IEEE Access 7 (2019) 2273–2280.

[46] X. Dai, H. Yin, N. K. Jha, Nest: A neural network synthesis tool based on a grow-and-prune paradigm, IEEE Transactions on Computers 68 (2019) 1487–1497.

[47] R. Yu, A. Li, C. Chen, J. Lai, V. I. Morariu, X. Han, M. Gao, C. Lin, L. S. Davis, NISP: pruning networks using neuron importance score propagation, in: IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2018, pp. 9194–9203.

[48] J. Luo, J. Wu, W. Lin, Thinet: A filter level pruning method for deep neural network compression, in: IEEE International Conference on Computer Vision (ICCV), 2017, pp. 5068–5076.

[49] J. Luo, H. Zhang, H. Zhou, C. Xie, J. Wu, W. Lin, Thinet: Pruning CNN filters for a thinner net, IEEE Transactions on Pattern Analysis and Machine Intelligence 41 (2019) 2525–2538.

[50] L. Arras, F. Horn, G. Montavon, K.-R. Müller, W. Samek, "what is relevant in a text document?": An interpretable machine learning approach, PLoS ONE 12 (2017) e0181142.

[51] W. Samek, A. Binder, G. Montavon, S. Lapuschkin, K.-R. Müller, Evaluating the visualization of what a deep neural network has learned, IEEE Transactions on Neural Networks and Learning Systems 28 (2017) 2660–2673.

[52] S. Lazebnik, C. Schmid, J. Ponce, Beyond bags of features: Spatial pyramid matching for recognizing natural scene categories, in: IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2006, pp. 2169–2178.

[53] L. Li, F. Li, What, where and who? classifying events by scene and object recognition, in: IEEE International Conference on Computer Vision (ICCV), 2007, pp. 1–8.

[54] J. Elson, J. R. Douceur, J. Howell, J. Saul, Asirra: a CAPTCHA that exploits interest-aligned manual image categorization, in: Proceedings of the 2007 ACM Conference on Computer and Communications Security (CCS), 2007, pp. 366–374.

[55] M. Nilsback, A. Zisserman, Automated flower classification over a large number of classes, in: Sixth Indian Conference on Computer Vision, Graphics & Image Processing (ICVGIP), 2008, pp. 722–729.

[56] A. Krizhevsky, Learning multiple layers of features from tiny images, Technical Report, Citeseer, 2009.

[57] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. S. Bernstein, A. C. Berg, F. Li, Imagenet large scale visual recognition challenge, International Journal of Computer Vision 115 (2015) 211–252.

[58] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, E. Duchesnay, Scikit-learn: Machine learning in Python, Journal of Machine Learning Research 12 (2011) 2825–2830.

[59] H. Wang, Q. Zhang, Y. Wang, H. Hu, Structured probabilistic pruning for convolutional neural network acceleration, in: British Machine Vision Conference (BMVC), 2018, p. 149.

[60] S. H. HasanPour, M. Rouhani, M. Fayyaz, M. Sabokrou, Lets keep it simple, using simple architectures to outperform deeper and more complex architectures, arXiv:1608.06037 (2016).

[61] A. Paszke, S. Gross, S. Chintala, G. Chanan, E. Yang, Z. DeVito, Z. Lin, A. Desmaison, L. Antiga, A. Lerer, Automatic differentiation in pytorch, in: NIPS-W, 2017.

[62] J. Liu, Y. Wang, Y. Qiao, Sparse deep transfer learning for convolutional neural network, in: AAAI Conference on Artificial Intelligence, 2017, pp. 2245–2251.

[63] N. Murata, S. Yoshizawa, S.-i. Amari, Network information criterion-determining the number of hidden units for an artificial neural network model, IEEE Transactions on Neural Networks 5 (1994) 865–872.

## Appendix A. Data pre-processing

During fine-tuning, images are resized to 256×256 and randomly cropped to 224×224 pixels, and then horizontally flipped with a random chance of 50% for data augmentation. For testing, images are resized to 224×224 pixels.

**Scene 15:** The Scene 15 dataset contains about 4500 images and consists of 15 natural scene categories obtained from COREL collection, Google image search and personal photographs [52]. We fine-tuned two different models on 20% of the images from each class and achieved initial Top-1 accuracy of 88.59% for VGG-16 and 85.48% for AlexNet, respectively.

**Event 8:** Event-8 consists of 8 sports event categories by integrating scene and object recognition [53]. We use 40% of the dataset's images for fine-tuning and the remaining 60% for testing. We adopted the common data augmentation method as in [37].

**Cats and Dogs:** This is the Asirra dataset provided by Microsoft Research (from Kaggle) The given dataset for the competition (KMLC-Challenge-1) [54]. Training dataset contains 12,500 colored images of dogs and 12,500 colored images of cats, while containing 12,500 test images. We reached initial accuracies of 99.36% for VGG-16 and 96.84% for AlexNet based on transfer learning approach.

**Oxford Flower 102:** Oxford 102 flowers contains 102 species of flower categories found in the UK, which is a collection with over 2000 training and 6100 test images [55]. We fine-tuned models with pre-trained VGG-16 on ImageNet for transfer learning and the initial performance is 82.26%.

**Cifar 10:** This contains 50,000 training images and 10,000 test images spanning 10 categories of objects [56]. The resolution of each image is 32×32 pixels and therefore we resize the images as 224×224 pixels.

**ILSVRC 2012:** In order to show the effectiveness of the pruning criteria on small sample case, we pruned and tested on randomly selected data among 150,000 photographs, collected from flickr and other search engines, hand-labeled with the presence or absence of 1000 object categories [57].
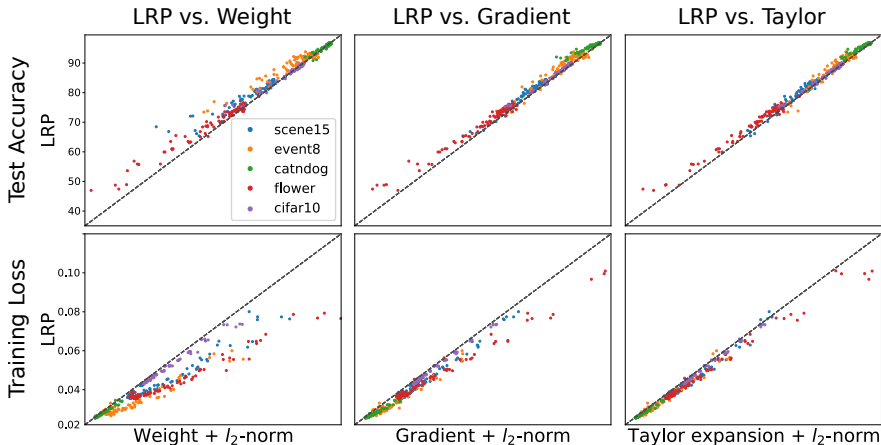


Figure A.7: Performance comparison of the proposed method (i.e. LRP) and other criteria on AlexNet with five datasets.
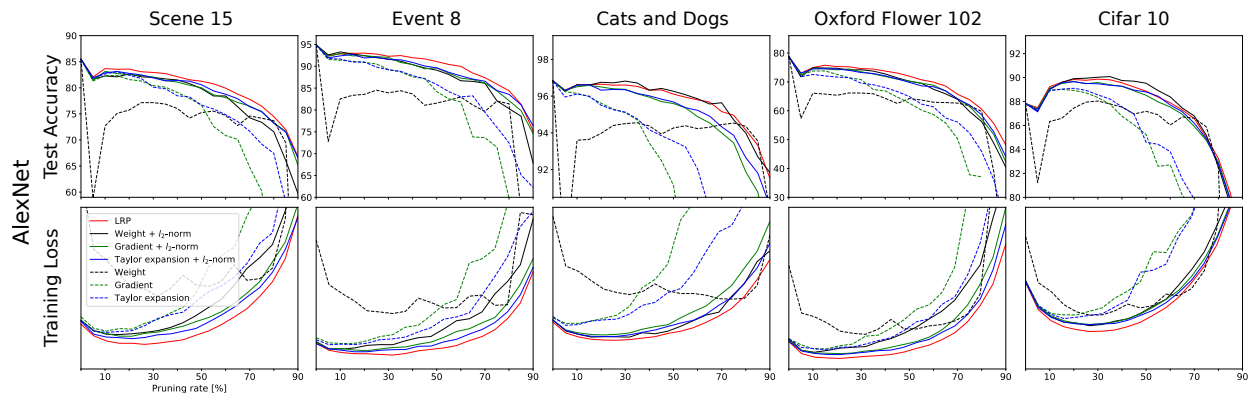
Figure A.8: Performance comparison of training loss and test accuracy in different criteria as pruning rate increases on AlexNet with five datasets.

Table A.4: Performance comparison of criteria (Weight, Taylor, Gradient with $l_2$-norm, and LRP) on **AlexNet** with five datasets. Please note that initial performances of AlexNet models are lower than those of the deeper VGG-16 models.

| Dataset | Criterion | Training loss | Accuracy (%) | Pruning rate | Parameters ($\times 10^7$) | FLOPs (MMAC) |
|---|---|---|---|---|---|---|
| Scene 15 | AlexNet | 0.0249 | 85.48 | | 54.60 | 711.51 |
| | Weight | 0.0278 | 78.43 | | 35.19 | 304.88 |
| | Taylor expansion | 0.0231 | 79.40 | 55 % | 33.79 | 229.95 |
| | Gradient | 0.0246 | 78.63 | | **33.29** | **225.19** |
| | LRP | **0.0202** | **80.76** | | 33.93 | 277.98 |
| Event 8 | AlexNet | 0.0116 | 94.89 | | 54.57 | 711.48 |
| | Weight | 0.0167 | 88.10 | | 34.25 | 301.19 |
| | Taylor expansion | 0.0121 | 88.62 | 55 % | 33.00 | 241.18 |
| | Gradient | 0.0137 | 88.42 | | 33.29 | **238.36** |
| | LRP | **0.0100** | **90.41** | | **32.79** | 291.95 |
| Cats and Dogs | AlexNet | 0.0050 | 96.84 | | 54.54 | 711.46 |
| | Weight | 0.0077 | **95.86** | | 32.73 | 264.04 |
| | Taylor expansion | 0.0078 | 95.23 | 60 % | 33.50 | 199.88 |
| | Gradient | 0.0087 | 94.89 | | 33.97 | **190.84** |
| | LRP | **0.0070** | 95.81 | | **32.66** | 240.19 |
| Oxford Flower 102 | AlexNet | 0.0515 | 78.74 | | 54.95 | 711.87 |
| | Weight | 0.0483 | 63.93 | | **28.13** | 192.69 |
| | Taylor expansion | 0.0339 | 64.10 | 70 % | 29.19 | **132.34** |
| | Gradient | 0.0377 | 64.11 | | 28.72 | 141.82 |
| | LRP | **0.0301** | **65.69** | | 28.91 | 161.35 |
| Cifar 10 | AlexNet | 0.0195 | 87.83 | | 54.58 | 711.49 |
| | Weight | 0.0244 | **90.03** | | 48.31 | 477.16 |
| | Taylor expansion | 0.0246 | 89.48 | 30 % | 48.23 | **371.48** |
| | Gradient | 0.0246 | 89.58 | | **46.31** | 395.43 |
| | LRP | **0.0233** | 89.87 | | 48.22 | 402.93 |