

MLCV2017: Multi-Layer Knowledge Transfer for Neural Networks

Lennard Kiehl

`lennard.kiehl@gmail.com`

Roman Remme

`roman.remme@gmx.de`

Abstract

This paper investigates an extension to the method of knowledge distillation between models originally proposed by Hinton et al. As distillation only establishes a link between the last (softmax) layers, we propose additional intermediate links to accelerate training and improve test accuracies. We study this on the example of VGG-16 and an effective “VGG-7” with each stack of layers reduced to one. We show that intermediate links improve test accuracies on CIFAR-10 by 13% compared to training from scratch, although we did not have enough time to fully optimize hyperparameters. This paper was done as a project for the lecture MLCV2017 by Fred Hamprecht.

1. Introduction

The idea of transferring knowledge between different architectures of neural networks, specifically from bigger models to smaller models, has been introduced in [1]. Part of the motivation for this process called distilling is to create a smaller model which is faster at runtime, with the same knowledge as the bigger model. We want to extend on this idea and not only compare the last layers of both networks while training the smaller one but also add links between intermediate layers. This extension is a really canonical one as especially bigger models used for image classification have a lot of their knowledge saved in their convolutional layers which may not completely translate into the final prediction. The popular VGG-16 model introduced in [4] serves as the bigger model and the goal is to distill each group of convolutional layers into only one convolutional layer, for an overview of the architectures see Table 1.3. We investigate how training hyper-parameters influence the process of successfully distilling knowledge.

1.1. Distillation

As it is a prerequisite for distilling to have an already trained model, to make this process worthwhile the trained model should have some kind of disadvantage at inference because distilling allows to transfer knowledge to a smaller model better suited for inference. And as shown in [1]

distillation also works with a fraction of the original training set as well as unlabeled data because the trained model should produce reliable predictions even for unseen data. To transfer knowledge an additional term to the loss function is introduced that links the softmax layers of both networks by calculating their cross entropy. This way the smaller model will not only have to produce the correct label but also the relationship between classes of lower probability. The softmax also has an added temperature dependency for training

$$q_i = \frac{\exp(z_i/T)}{\sum_j \exp(z_j/T)} \quad (1)$$

with logits z_i , predicted class probabilities q_i and temperature T , which would normally be set to 1. Using a higher value for T produces a softer probability distribution and should force the smaller model to optimize better for intermediate relationships between classes. It was shown in [1] that this alone serves as a very good knowledge transfer tool. We will investigate how adding a selection of other links while training will influence the distillation. For an overview of the proposed links see Table 1.3.

1.2. Loss

The loss function for transferring knowledge into the small model is a weighted sum of three terms:

- “hard” loss: cross-entropy between output and correct label at temperature $T = 1$

$$\text{hard} = - \sum_i \text{label}_i(x) \log \text{small}_i(x) \quad (2)$$

- “soft” loss: cross-entropy between output and prediction of big model at temperature T

$$\text{soft} = - \sum_i \text{big}_i(x) \log \text{small}_i(x) \quad (3)$$

- “intermediate” loss: sum of MSE between linked intermediate layers

$$\text{intermediate} = \frac{1}{N} \sum_i^N \text{MSE}(\text{link}_i) \quad (4)$$

The third term is the new part of our approach. The “intermediate” loss averages the MSE over all links that were selected. With this we make sure that the magnitude of the loss stays consistent between experiments with different numbers of links. An anticipated advantage is that training times should be reduced, as gradients do not have to be propagated through the whole network to reach the first layers. Typical factors in the weighted sum are

$$\text{total loss} = \text{hard} + 10 \cdot \text{soft} + 10 \cdot \text{intermediate} \quad (5)$$

The main contributions come equally from “soft” and “intermediate” loss. While still significantly improving distillation the “hard” loss contribution is much less. This is also consistent with the ideas in [1]. We would place the importance of the “intermediate” loss somewhere between these two which is reflected in Equation ??.

1.3. Models

For our experiments we use VGG-16 [4] as the big model. For the small model all stacks of convolutional layers have been replaced by one single convolutional layer (see Table 1.3) and the number of fully connected layers was reduced by one. The similarity between the models is by design and makes it possible to have a maximum of 6 separate links while doing the knowledge transfer. To get the knowledge we want to transfer in the first place, the big model is trained on CIFAR10 [2] with the hyperparameters shown in Table 2. The convolutional layers of the big model are initialized with pre-trained weights on ImageNet [3] while the fully connected layers are initialized randomly. The small model had to be trained from scratch as it is an uncommon architecture. The accuracies in Table 2 serve as our baseline and we expect the accuracy of the small model after distilling to be somewhere between these two test accuracies.

1.4. Dataset

We use CIFAR-10 [2] as the dataset to train both models for all experiments. It consists of 50000 training and 10000 test RGB images of size 32×32 pixels. Each image belongs to one of ten classes. To use the standard VGG architecture with these low-resolution images, they are scaled up to 224×224 pixels. Each image is preprocessed by subtracting the mean RGB value, computed on the training set, from each pixel.

2. Experiments

First we perform an experiment to find out what temperatures is best suited for distillation with our choice of models (results in Table 3). Next we compare multiple combinations of links between intermediate layers to find out if our

Network configurations with linkable layers		
	VGG-16	VGG-7
	input (224×224 RGB image)	
link 1	conv3-64 conv3-64	conv3-64
	maxpool 2×2	
link 2	conv3-128 conv3-128	conv3-128
	maxpool 2×2	
link 3	conv3-256 conv3-256 conv3-256	conv3-256
	maxpool 2×2	
link 4	conv3-512 conv3-512 conv3-512	conv3-512
	maxpool 2×2	
link 5	conv3-512 conv3-512 conv3-512	conv3-512
	maxpool 2×2	
link 6	FC-4096 FC-4096 FC-10	FC-4096 FC-10
	softmax	

Table 1. **Network configurations.** The convolutional layers are denoted as conv(*kernel size*)-(number of channels) and the fully connected layers as FC-(number of output channels). ReLu units are omitted for brevity. The leftmost column gives the links between both networks that are added to the loss function.

	Small model	Big model
Batchsize	40	40
Momentum	0.9	.9
Weight decay	0.01	0.0002
Init learning rate (LR)	0.004	0.004
Epochs between LR decay	10	25
Epochs	25	100
Train accuracy	99.3%	100.0%
Test accuracy	79.1%	91.4%

Table 2. **Baseline Training.** Both network architectures were trained with the given parameters to have a baseline to compare our transfer training to. The conv. layers of the big model had pre-trained weights while the small model was trained from scratch.

approach can improve the knowledge transfer over normal distilling (results in Table 4). For all experiments stochastic gradient descent with momentum as a regularizer is used as an optimizer. Furthermore after every ten epochs the learning rate is decaying by a factor of 10. This way the accuracy should stop changing significantly prior to a drop in learn-

Temperature	Test accuracy
0.6	76.3 %
1	76.5 %
1.5	77.0 %
2	77.4 %
2.5	76.7 %
3	77.2 %
5	73.1 %
10	64.4 %
40	67.3 %

Table 3. **Distillation using only “hard” and “soft” loss.** Test accuracies after distillation using different temperatures for the softmax.

ing rate.

2.1. Temperature of “soft” loss

This experiment is done exactly like [1] describes the process of distillation. That means that the loss function consists of the “hard” and “soft” terms only. This is used to determine the best temperature to test our new approach. The relative weight of the “soft” loss was chosen to be ten times that of the “hard” loss. We found the best temperature T to be 2, see Table 3. The corresponding test accuracy is only 77.4%, which is less than our baseline of 79.1% for the small model. This is due to a lack of further experiments. Optimizing these numbers takes a lot of time and especially compute time. We think that a missing regularization is the main cause and that normally this accuracy should be a little higher than our baseline. This would also be consistent with the findings in [1]. It should be noted that the goal of determining the best temperature has still been achieved and we can proceed with the main experiment.

2.2. Linking intermediate layers

This is the main experiment of this paper. We want to investigate if adding links between intermediate layers while distilling knowledge between models can improve the transfer and have an impact on test accuracies. If this is the case we also want to compare different sets of intermediate links. The results can be seen in Table 4. The relative weight of the “intermediate” loss was chosen to be identical to that of the “soft” loss as we think it should be just as important for the transfer. However to make results with different numbers of links comparable the “intermediate” loss is always an average over all losses that result from links between intermediate layers. If we would not have implemented some kind of mechanism to guarantee a more or less consistent loss, it could have been the case that more links would lead to a difference in magnitude of the “intermediate” loss. First, we use one link at a time. Link 5, the last link in the convolutional part of the network, has the best test accu-

Linked layers	Test accuracy
1	78.3%
2	80.3%
3	83.0%
4	85.6%
5	87.9%
6	86.1%
3, 4	84.8%
2, 3, 4, 5	87.0%
2, 3, 4, 5	87.9%
2, 3, 4, 5, 6	78.5%
1, 2, 3, 4, 5, 6	81.2%

Table 4. **Distillation with added “intermediate” loss.** Test accuracies for different sets of links between intermediate layers ($T = 2$).

racy of 87.9%. This is a 13% improvement over our baseline of 79.1%. Using multiple links between intermediate layers yielded at best equally good results. But since far from all possibilities were explored, it is possible that further improvements could be achieved with the right choice of layers to link.

3. Discussion

This paper investigated an extension of the process called distillation originally proposed by Hinton et al.[1]. Distillation is a method to transfer knowledge from an already trained model to a new one that requires less time and data than training from scratch. Ideally the new model is similar in architecture but with fewer layers to make it faster at runtime than the already available trained model. Distilling knowledge works by linking the softmax layers of both networks by calculating their cross-entropy, making the “knowledge” of relationships in between classes an additional target to the normal onehot encoded label. We extended this idea by adding links between different intermediate layers but using the MSE instead, because this compares activations and not probability distributions. For our experiments we used the popular VGG-16 and a slimmed-down version with each stack of layers being reduced to one, effectively being “VGG-7” (see Table 1.3). From our comparison of different temperatures T for calculating the softmax (Equation 1) in Table 3, we decided to use $T = 2$ for comparing the effectiveness between different sets of intermediate links. The resulting test accuracies can be seen in Table 4. From this we conclude that adding links between intermediate layers while doing knowledge distillation improves the transfer significantly. Not only is the training time reduced by the fact that it doesn’t take multiple epochs for the gradient to reach layers in the front of the model, but the test accuracies on CIFAR-10 [2] im-

proved by 13% to 87.9% over our baseline of 79.1% when training from scratch. While doing these experiments we observed a partial invariance to the choice of hyperparameters for training, leading us to believe that additional links might also make the distillation more robust.

Given the scope and our investment of time in this project, we think that there is still room to improve and investigate. Although our results look promising we propose further experiments with other architectures and datasets. It might also be possible that our tuning of the hyperparameters is off and the observed invariance is not real. With this in mind, we think that knowledge distillation is an idea worth investigating further, especially considering the discrepancy between the ever growing size of top-performing models and the spread of machine learning into small handheld devices. If at least a few future AI services would aim to be independent of a fast internet connection, knowledge distillation into feasible models might be a necessity.

References

- [1] G. Hinton, O. Vinyals, and J. Dean. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531*, 2015.
- [2] A. Krizhevsky and G. Hinton. Learning multiple layers of features from tiny images. *Technical report, University of Toronto*, 2009.
- [3] 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)*, 115(3):211–252, 2015.
- [4] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. *CoRR*, abs/1409.1556, 2014.