

Generalization Ideas in Deep Learning

Marcel Bruckner

Seminar: Optimization and Generalization in Deep Learning

Abstract

Trying to understand the generalization abilities of deep neural networks several capacity measures are experimentally explored. Following the thoughts of Neyshabur et al. [3] different norm-based capacity measures over the network weights and the sharpness as the robustness to perturbations on the parameter space are investigated. The measures are used in different experiments and the results are plotted against each other.

1 Introduction

Huge empirical success of deep neural networks
 Complex, non-convex optimization problem
 Simple stochastic gradient descent (SGD) methods achieve good generalization
 Even in this highly over parameterized setting
 Many local minima with different generalization behavior

2 What is a matrix norm and which do we use

2.1 What is a matrix norm

What does it do to the function

2.2 Why do we use matrix norms as measures for capacity bounds

how controls the norm the model capacity

2.3 Which norms do we use

- l_2 norm with capacity proportional to $\frac{1}{\gamma_{margin}^2} \prod_{i=1}^d 4 \|W_i\|_F^2$ (1)

- l_1 -path norm with capacity proportional to $\frac{1}{\gamma_{margin}^2} \left| \sum_{j \in \prod_{k=0}^d [h_k]} \left| \prod_{i=1}^d 2W_i[j_i, j_{i-1}] \right| \right|^2$ (2)

- l_2 -path norm with capacity proportional to $\frac{1}{\gamma_{margin}^2} \sum_{j \in \prod_{k=0}^d [h_k]} \prod_{i=1}^d 4h_i W_i^2[j_i, j_{i-1}]$ (3)

- spectral norm with capacity proportional to $\frac{1}{\gamma_{margin}^2} \prod_{i=1}^d h_i \|W_i\|_2^2$ (4)

3 What is sharpness

$$\zeta_\alpha(w) = \frac{\max_{|\nu_i| \leq \alpha(|w_i|+1)} \hat{L}(f_{w+\nu}) - \hat{L}(f_w)}{1 + \hat{L}(f_w)} \cong \max_{|\nu_i| \leq \alpha(|w_i|+1)} \hat{L}(f_{w+\nu}) - \hat{L}(f_w) \quad (5)$$

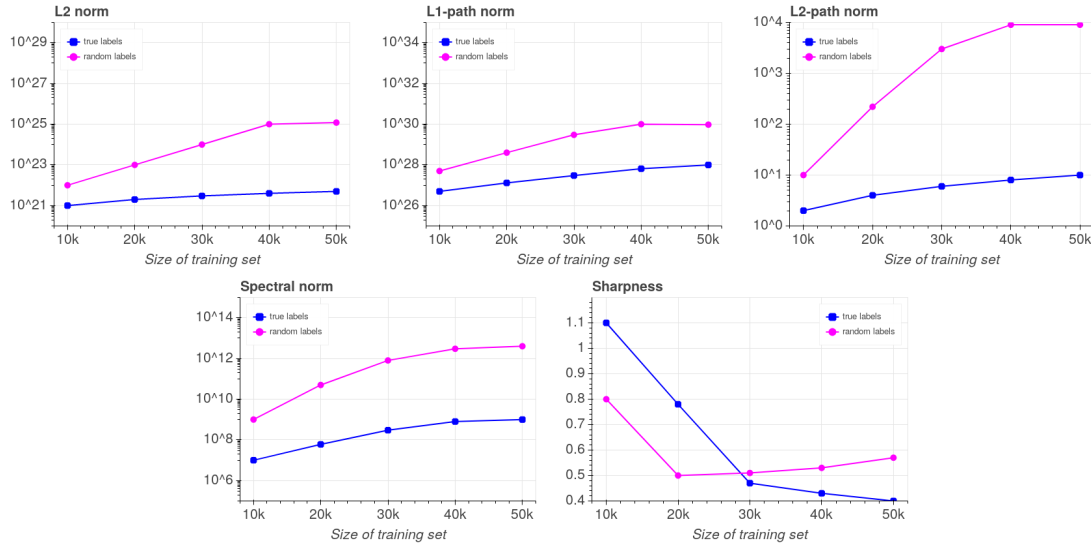


Figure 1: The proposed capacity measures calculated on the VGG model [6] after training on random (magenta) and true (blue) labeled subsets of CIFAR-10. We see only a small increase of the measures on true labels but a huge increase using random labels and the measures on random labels to be bigger for every subset. Reproduced from [3].

4 Empirical study

In this section we discuss the experimental results using the proposed measures and their implications on different phenomena seen in practice regarding the generalization behavior.

4.1 Difference between true and random labeled data

It is possible to obtain zero training error on random labels using the same architecture for which training with real labels leads to good generalization. We would expect the networks learned using real labels (and which generalizes well) to have much lower complexity, under the suggested measure, than those learned using random labels (and which obviously do not generalize well). [3]

Figure 1 shows the results of [3]. We see that training the network on true labels (blue) only results in small increases of the norm-based capacity measures in comparison to the increases when training on random labels (magenta). As the network learns the true functional dependence between the input and output using true labels it only requires small increases in complexity as the subset size increases, whereas it requires more capacity for every newly seen data point when using random labels as it needs to memorize the data point.

We also see that the norm-based capacity measures are all bigger for random labeled data as the network has to memorize the data whereas the dependence between the input and output can be learned with lower capacity on true labeled data.

Looking at the sharpness as a capacity measure in Figure 1 we do not see a direct correlation between its value and the capacity of the network. We expect the model to become less sharp when increasing the network size which can be seen looking at the true labeled data points (blue) but clearly gets violated by the random labeled data points (magenta). Furthermore, we see a lower sharpness for random labels for the two smallest subset sizes which contradicts the intuition that the sharpness should be always lower for true labels.

4.1.1 Problems reproducing results

Following the vague instructions of [3] the VGG-16 network with batch normalization (Fig: 7) is trained on different subsets of the CIFAR-10 (Fig: 7) data set. For every subset a copy with the labels replaced by random labels has been created and the model is trained on both the true and random labeled subset.

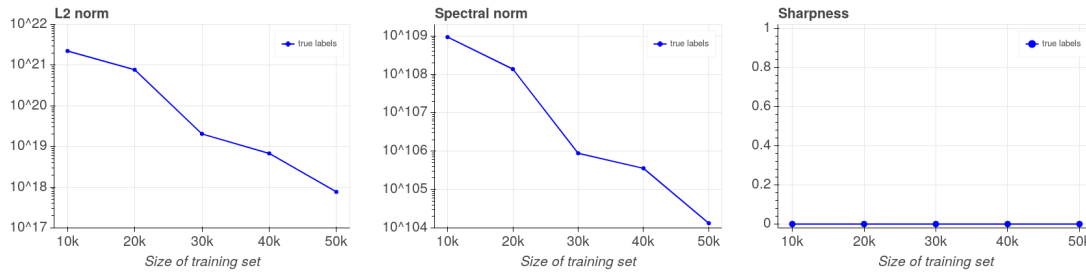


Figure 2: Our own results for the l_2 and spectral norm as well as the sharpness calculated on the VGG-16 model after training on subsets of CIFAR-10. The results show no reproducibility of the results presented by [3].

The huge scope for interpretation of the approach used in [3] caused the non reproducibility of the results which are shown in Figure 2. The problems we face are:

Reproduction of network training Neyshabur et al. [3] only stated that their results were produced using VGG. As VGG is a class of CNN architectures with varying depths and optional batch normalization we performed a grid search over the implementations. Considering the capacity vs. training time trade-off we achieved the best results using VGG-16 with batch normalization. Furthermore, [3] didn't provide information on the used solving strategy. Therefore, we tried simple stochastic gradient descent and Adam but stuck with the later as of the faster convergence.

Calculation of norms The calculation of path norms is trivial for the fully connected case but not for convolutional layers. As we need to enumerate every path through the network we need a way to associate the layer neurons of succeeding layers when applying a convolution. In the short time of this paper we couldn't come up with a feasible implementation which made it impossible to reproduce the path norms.

Training on random labeled data Following the missing details of the network implementation and solving strategy a training on random labeled data was not possible. The network didn't overcome a 1/10 test accuracy on random labels which represents a random choice of the result by the network out of the 10 classes in the CIFAR-10 data set.

Training time The short scope of the seminar enclosing this paper combined with training times ranging from roughly five hours on the smallest subset ranging up to over one day for the biggest one resulted in only small search intervals for the grid search used on the hyper parameters as well for the algorithmic and architectural choices.

Sharpness Sharpness as the maximum over the adversarial perturbations is a maximization problem where the difference between the loss of the network with the added perturbation and the original network gets maximized. This maximization problem is constrained on the values of the perturbations and thus requires a constrained optimization strategy. We couldn't come up with such a strategy in the seminars time and thus tried a random search strategy. Due to the high duration of forward passes over the whole training set this random search didn't provide results in a feasible amount of time.

4.2 Difference between local minima

When training the same architecture, with the same training set, using two different optimization methods (or different algorithmic or parameter choices), one method results in better generalization even though both lead to zero training error. We would expect to see a correlation between the complexity measure and generalization ability among zero-training error models. [3]

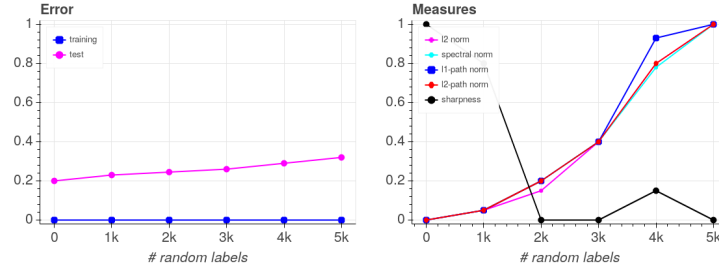


Figure 3: The test and training error on the subset of CIFAR-10 and varying confusion set sizes (left) and the corresponding proposed capacity measures (right). Reproduced from [3].

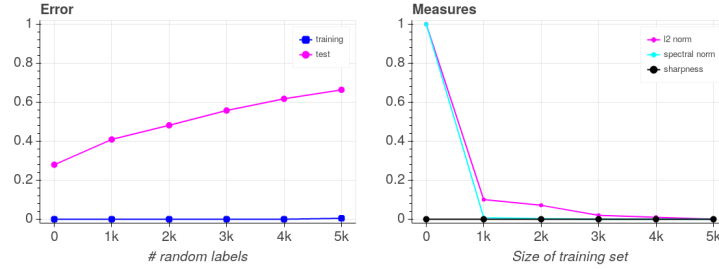


Figure 4: Our own results for the test and training error on the subset of CIFAR-10 and varying confusion set sizes (left) and the corresponding proposed capacity measures (right).

To drive the network towards different parameter choices [3] proposed to train the VGG network [6](Fig: 7) on a subset of 10000 true labeled data points of the CIFAR-10 (Fig: 7) and to add confusion sets with varying sizes. A confusion set is a subset of the same data set where the labels are replaced with random labels. When optimizing the network over the union of the true and random labeled sets the resulting optimum has to be optimal for both of the sets.

We expect to see an increase of the test error with increasing confusion set size as the network not only has to learn the functional dependence between the input and output data but also has to memorize an increasing number of random labeled data points. This memorization prevents the network from generalizing well and can be seen in Figure 3.

Following the increase in test error and the resulting worse generalization behavior we expect the capacity measure to increase with the test error. In Figure 3 we see this behavior for all the norm-based capacity measures which hints towards them being a valid measure for generalization behavior. But again we see the sharpness to have no direct correlation to the test error which hints that it cannot predict generalization.

4.2.1 Problems reproducing results

We faced most of the problems already described in 4.1.1. Especially the missing model and solving strategy details 4.1.1 as well as the calculation of the norms 4.1.1 prevented us from reproducing meaningful results. Nonetheless, the results are displayed in Figure 4 where we can see an increase of the test error as expected when increasing the subset size. The calculated measures do not follow the expected behavior to increase with the worse generalization ability and sharpness wasn't calculable at all.

Training to zero training error not possible For the biggest confusion set size training to zero training error was not possible. We tried different optimization strategies and hyper parameters but couldn't drive the model to zero training error.

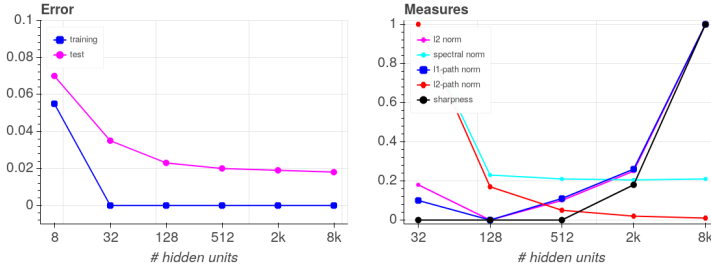


Figure 5: The test and training error on the MNIST handwritten digit data set and varying hidden unit sizes (left) and the corresponding proposed capacity measures (right). Reproduced from [3].

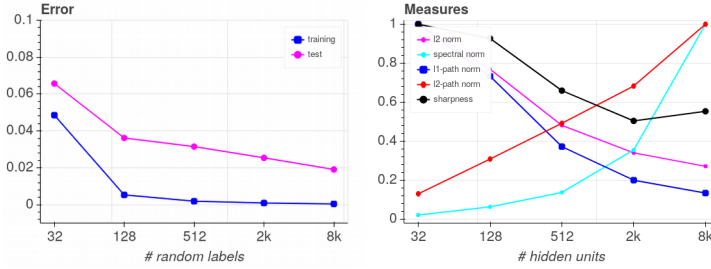


Figure 6: Our own results for the test and training error on the MNIST handwritten digit data set and varying hidden unit sizes (left) and the corresponding proposed capacity measures (right).

4.3 Implications of different hidden unit sizes

Increasing the number of hidden units, thereby increasing the number of parameters, can lead to a decrease in generalization error even when the training error does not decrease. We would expect to see the complexity measure decrease as we increase the number of hidden units. [3]

To investigate into this phenomenon Neyshabur et al. [3] trained a two layer perceptron (Fig: 8) on the MNIST handwritten digit data set (Fig: 8) using different numbers of hidden units. We see in 5 that the networks achieve zero training error for 32 and more hidden units and even though there is no decrease in the training error afterwards we see a decrease in the test error with every increase of the hidden units.

Looking at the capacity measures in Figure 5 we see that the l_2 -path and spectral norm follow the expected behavior and decrease with the increase of hidden units whereas the l_1 -path norm, the l_2 norm and sharpness do not behave as proposed. Reproducing the results of [3] put out the values displayed in Figure 6 where we can see that we could achieve a decrease of the l_2 and l_1 -path norm with the increase of hidden units but see the opposing behavior for the other norms and again no direct correlation between sharpness and generalization behavior.

5 Conclusion

6 Appendix

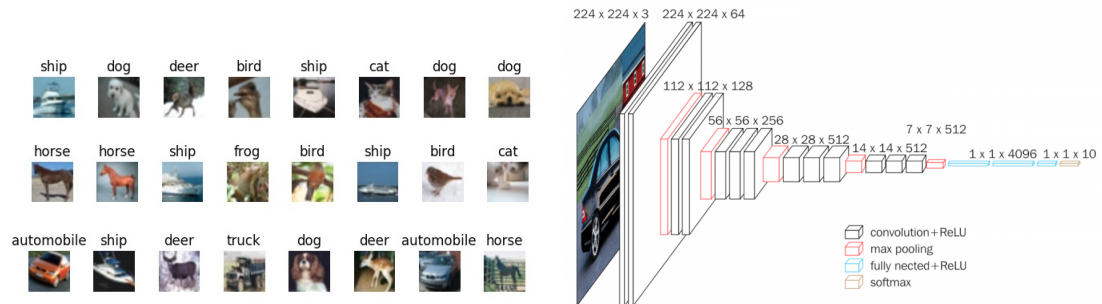


Figure 7: Left: A subset of the CIFAR-10 small image data set with ten classes. Right: The used implementation of the VGG-16 network with batch normalization.

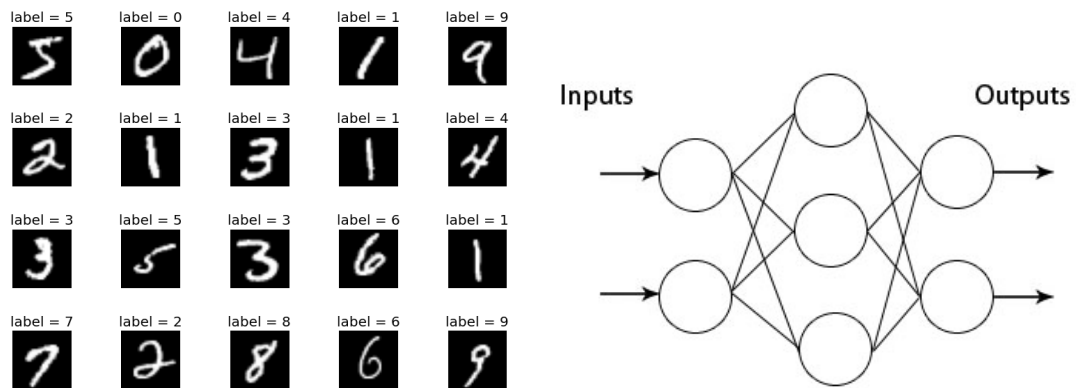


Figure 8: Left: A subset of the MNIST handwritten digit data set. Right: An exemplary two layer perceptron with one input layer, one hidden layer of hidden unit size three and one output layer.

References

- [1] Peter L. Bartlett and Shahar Mendelson. Rademacher and gaussian complexities: Risk bounds and structural results. In David Helmbold and Bob Williamson, editors, *Computational Learning Theory*, pages 224–240, Berlin, Heidelberg, 2001. Springer Berlin Heidelberg.
- [2] Kenji Kawaguchi, Leslie Pack Kaelbling, and Yoshua Bengio. Generalization in deep learning, 2017.
- [3] Behnam Neyshabur, Srinadh Bhojanapalli, David McAllester, and Nati Srebro. Exploring generalization in deep learning. In *Advances in Neural Information Processing Systems*, pages 5947–5956, 2017.
- [4] Behnam Neyshabur, Ruslan Salakhutdinov, and Nathan Srebro. Path-sgd: Path-normalized optimization in deep neural networks. *CoRR*, abs/1506.02617, 2015.
- [5] Behnam Neyshabur, Ryota Tomioka, and Nathan Srebro. Norm-based capacity control in neural networks. *CoRR*, abs/1503.00036, 2015.
- [6] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition, 2014.
- [7] Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. Understanding deep learning requires rethinking generalization. *CoRR*, abs/1611.03530, 2016.