

Exploring the Double Descent Curve in Deep Neural Networks

James Enouen
The Ohio State University
enouen.8@osu.edu

Abstract

This project intended to explore how the double descent curve from [1] is able to generalize to more complex models like deeper neural networks. Honestly, this work failed to get a perfectly comprehensive picture of exactly how the second descent interacts with multiple layers in a network; however, it manages to produce results manifesting the second descent and to give some future directions of further investigation.

I. INTRODUCTION

The goal of this project was to explore the double descent curve in its application to deep neural networks. The original goal was to thoroughly study how the double descent curve emerges throughout different datasets and different training regimens. This paper will walk through the journey of different approaches taken to attempt to find the double descent curve in this setting. Disappointingly, this paper spends more focus on 'how to train deep networks' than the originally desired 'investigating the second descent in deep networks.' In the end, there is no 'fairy tale ending' because none of the results are strongly conclusive; however, the final results indicate success in finding the second descent in the deep setting.

Let us begin with a review of what the double descent curve claims because it is the central focus of this paper. The paper [1] offers a new perspective on the classically believed bias-variance tradeoff. The new perspective suggests that we indeed get better results by adding more and more parameters/ hidden units to our model, despite the local maximum we approach as we reach the interpolation threshold of perfect fitting.

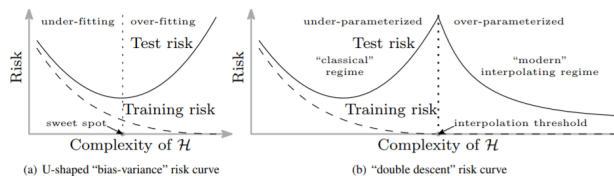


Fig. 1. Double descent curve from [1]

This goes against the traditionally believed idea that a model can overfit the training set and the only way to fix this problem is through regularization techniques. The paper suggests, moreover, that the minimum at the end of this curve generally tends to be lower than the local minimum attained in the bias-variance tradeoff curve. This advocates an overall different look on how to handle training models which is very interesting and consequently deserves study in applications like deep neural networks.

II. INSPIRATION

This project began as a failure to easily replicate the results of the double descent curve in a deep neural network scenario. These original results were trained on a dataset similar to the one below. As can be seen in the tables below, the results are very far away from what one might expect for the double descent curve to react.

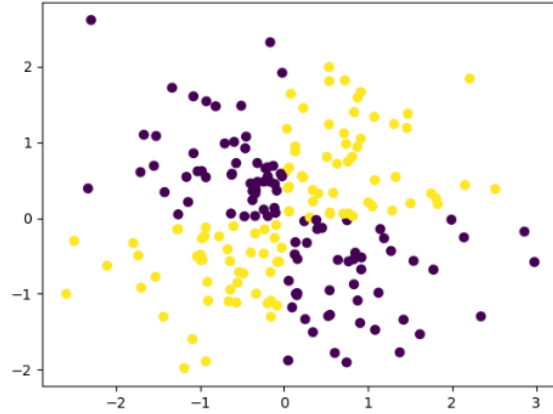


Fig. 2. Depiction of the dataset used

1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20
1	0.433333	0.366667	0.366667	0.433333	0.366667	0.4	0.466667	0.3	0.366667	0.333333	0.433333	0.266667	0.3	0.366667	0.3	0.3	0.3	0.3	0.3
2	0.433333	0.4	0.5	0.266667	0.5	0.333333	0.5	0.466667	0.3	0.333333	0.266667	0.466667	0.366667	0.3	0.333333	0.333333	0.3	0.366667	0.333333
3	0.433333	0.533333	0.366667	0.533333	0.5	0.333333	0.366667	0.3	0.366667	0.433333	0.466667	0.466667	0.366667	0.366667	0.333333	0.333333	0.466667	0.333333	0.333333
4	0.433333	0.466667	0.466667	0.4	0.4	0.566667	0.466667	0.5	0.366667	0.5	0.5	0.333333	0.266667	0.333333	0.4	0.3	0.433333	0.433333	0.366667
5	0.433333	0.466667	0.533333	0.533333	0.533333	0.533333	0.5	0.5	0.533333	0.366667	0.433333	0.5	0.333333	0.33	0.266667	0.366667	0.366667	0.366667	0.333333
6	0.433333	0.433333	0.466667	0.533333	0.566667	0.533333	0.466667	0.533333	0.5	0.366667	0.466667	0.366667	0.3	0.433333	0.466667	0.333333	0.466667	0.366667	0.3
7	0.433333	0.533333	0.533333	0.533333	0.466667	0.466667	0.466667	0.3	0.433333	0.333333	0.333333	0.4	0.533333	0.366667	0.266667	0.333333	0.233333	0.4	0.5
8	0.433333	0.533333	0.533333	0.466667	0.466667	0.4	0.4	0.333333	0.466667	0.433333	0.466667	0.433333	0.366667	0.533333	0.466667	0.3	0.433333	0.3	0.433333
9	0.433333	0.533333	0.466667	0.466667	0.466667	0.633333	0.466667	0.466667	0.366667	0.466667	0.233333	0.4	0.366667	0.366667	0.233333	0.333333	0.4	0.466667	0.366667
10	0.433333	0.533333	0.533333	0.533333	0.466667	0.5	0.366667	0.533333	0.533333	0.533333	0.433333	0.266667	0.4	0.366667	0.3	0.433333	0.4	0.333333	0.5

Fig. 3. Training and testing results for original dataset

This clearly does not depict the pattern we were looking for, Because increasing the length of the network also increases the number of parameters, it should increase the complexity of the model. Hence, the results need to be 'flipped' as can be seen in the following two graphics. This is clearly an issue of the results which were gathered.

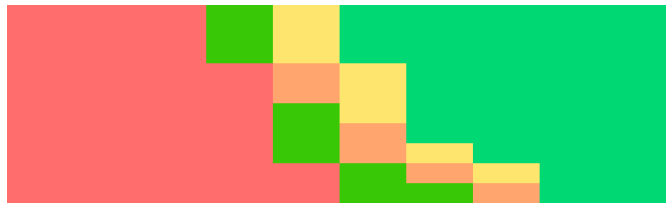


Fig. 4. Observed pattern of the validation errors

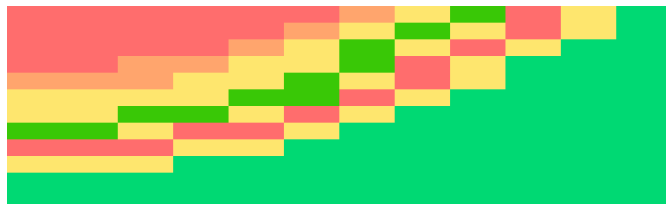


Fig. 5. Expected pattern of the validation errors from double descent perspective

A possibly even more fundamental issue is the fact that the more complex models had trouble scoring even better on the training loss function, which is fairly unreasonable. More advanced models have strictly greater functional complexity. This points to there being an issue in how these models are being trained. The rest of this paper are the steps taken to remedy this issue in light of the ultimate goal of finding the double descent curve in this setting.

III. EXPLORATIVE JOURNEY

A. Beginning Failures

The first steps taken to try to fix this issue was to try the same procedure on different datasets. Trying these on a sphere in two and four dimensions, the following results were generated:

Fig. 6. 2D sphere training loss and testing error

Fig. 7. 4D sphere training loss and testing error

Both of these trials are not a huge step above the work started with. While the axis for these charts are now flipped from the previous images, they still have the issue of being in the wrong direction of the 'complexity' of the model. There are quite a few more failures of these first steps taken towards exemplifying the double descent curve. Notably, these charts have NaN as one of their values. This is due to numerical issues to do with logarithms and exponents. This happens to be because the loss function these training procedures are using is negative log likelihood. In retrospect, this was a little bit of an unreasonable loss function, considering the original paper was using mean squared error and zero-one loss as their functions. It is still interesting to see if this loss function has the potential to see the double descent curve itself, but this will be a future direction which led to these results.

One of the first things that unaligned the training procedure from the double descent paper was the size of the dataset. The goal of the second descent is to be able to perfectly fit the training set and the set being used was probably too large to feasibly do with the number of parameters being used, hence the size of the training dataset was reduced and resulted in the following miniscule changes:

Fig. 8. Smaller dataset - 4D cylinder training loss and testing error

B. Minor Improvements (pretrain)

One of the characteristics that distinguishes the random Fourier features setting (where the double descent curve was discovered) from the neural network setting is the variance between the different layers of the model. The next step taken to fix this training procedure is an attempt to embrace this property that random Fourier features are fixing all previous layers and only training a final layer. For this, the models which got longer and longer by adding hidden layers would start off by using the original layers of the trained previous model. In this way, the goal was to reduce the variability by random initialization greatly and to smooth the results in the direction of longer and longer machines. This definitely made a notable improvement in the training losses of the models; however, it barely made a dent in the large scheme of our goal. This had nowhere near enough impact to 'flip' the results and barely added anything new to the boundary of our diagonal curve:

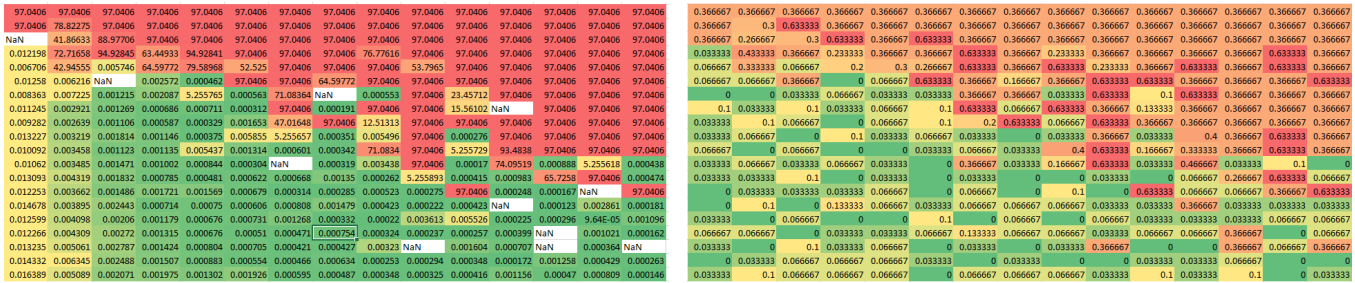


Fig. 9. First pretrained training losses and testing errors

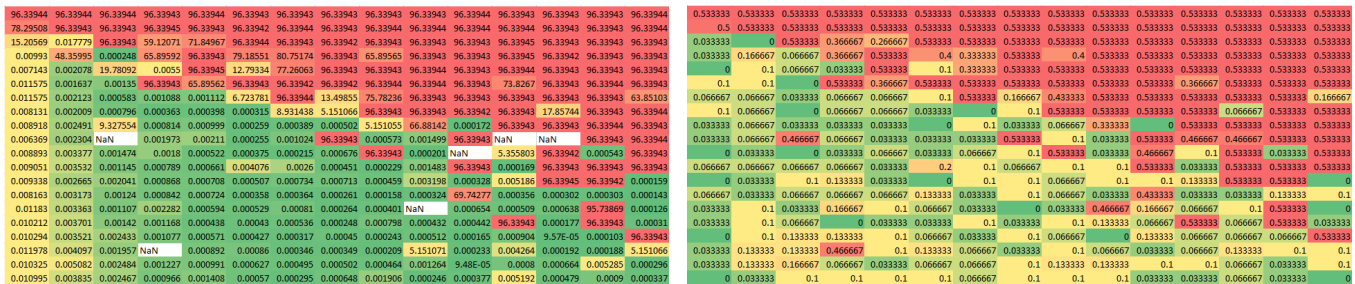


Fig. 10. Second pretrained training losses and testing errors

C. Fair Improvements (MSE)

The next alteration to the training procedure was to use MSE instead of NLL. This had a very noticeable impact on the results. This was the first dataset where the 'flip' was achieved. This is a vast improvement over the previous sets of models for this reason. The training loss curve made sense given the complexities of the models. The results are provided below:

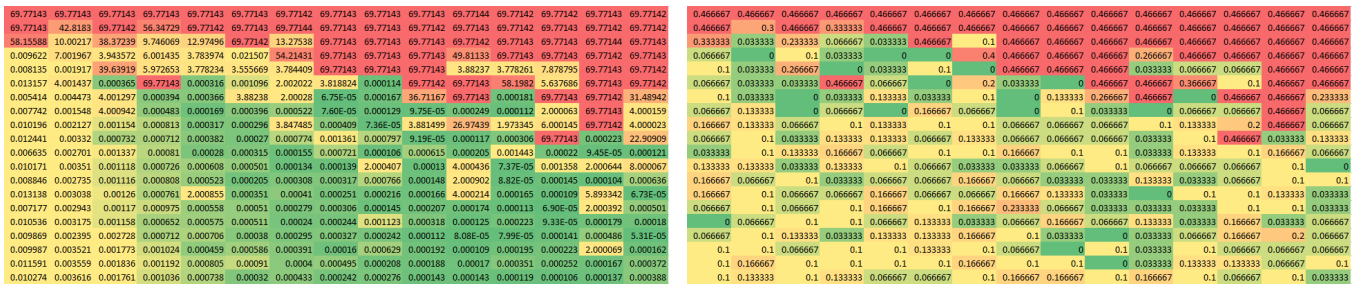


Fig. 11. First MSE training loss and testing error data

The shortcoming of this data is twofold. The first aspect is that neither are actually able to obtain the zero training loss which is required to hit the interpolation threshold. The second aspect which is in part a consequence of the first is how great the variance is in the testing accuracies of these models. There seems to be no clear cut pattern throughout these models. One could argue some small patterns like very good generalizability right after dropping off of the total inability to score well on the training set. This corresponds to the green ridge near the 'cliffs' of red and orange near the top and right of these

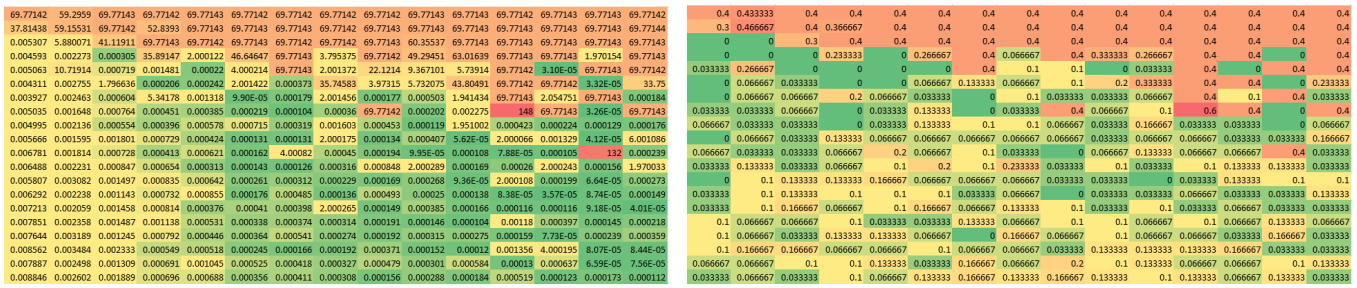


Fig. 12. Second MSE training loss and testing error data

depictions. It is likely that these correspond to the first descent in the double descent curve, but the goal of this project is to see the second.

Additionally, this training procedure was unable to perform well on more complex datasets. The following results were for a six dimensional version of the checkerboard pattern that was used in the first dataset used. Clearly the models were not able to properly learn the training set:

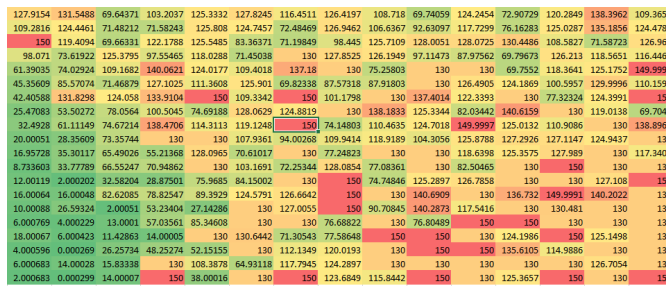


Fig. 13. Failure of MSE on six-dimensional dataset (training loss)

D. Learning Rate Schemes

The work of many others including those in [2] was sufficiently inspiring to alter the training of these models with a new, adaptive learning rate. This paper discusses more advanced learning rate procedures, but the one implemented for the next results is simply a decreasing function over time. It was, however, able to drastically outperform the same procedure with a fixed learning rate. On the same dataset as the last figure from the previous section, the new procedure was able to get the following results:

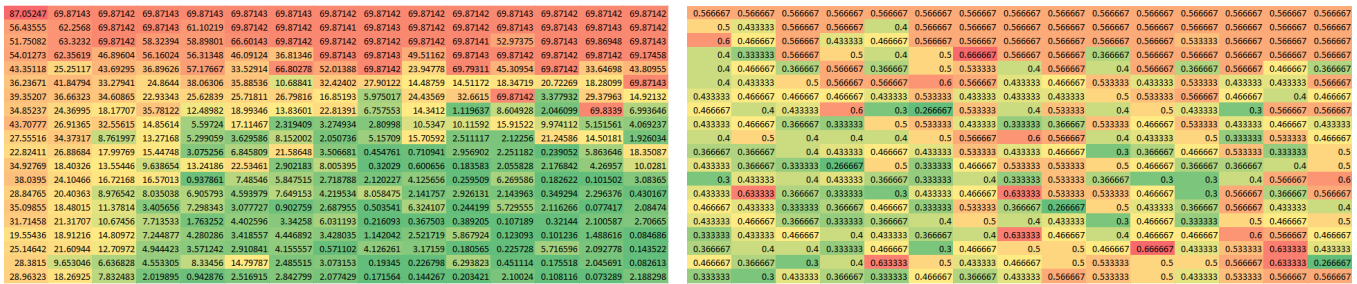


Fig. 14. Adapted learning rate training losses and testing errors

This data is clearly much more demonstrative of any type of second descent and achieved quite good training results. It seems almost that the valley of green possibly corresponds to the first descent. Hopefully extending the picture just a little would be sufficient to see the second descent. The next attempt was exactly this, running the same procedure for larger numbers of hidden layers and larger hidden unit sizes on a slightly simpler dataset. The results for this expanded chart are below:

E. Discussion

There is quite a lot of information in these two rectangles, so let us take this time to have a brief discussion about some of their features. A very prominent feature is the dark bands across the tops of these rectangles. This corresponds to the fact that,

especially in the one hidden unit case, the model is not becoming more complex despite having more parameters being thrown at it. A chain of multiple single nodes is unable to do anything better than a linear separation boundary if its activation is monotone. Consequently, this should not be able to do any better in classifying a training set and the same sort of phenomena holds while the hidden layers are quite small. Generally, as we go down and right we are getting better training losses which is what we expect, but the function is not as smooth as it should be. Optimally, the loss would strictly decrease when any step was made to the right or down. This is clearly not the case in our empirical setting and indicates that our validation dataset will be similarly bumpy. In the testing errors, we see the darkest red spots are at the top left corner (with the worst model) and scattered slightly after the boundary of the 'learning cliff' we can see in the training losses.

Beyond this, there is a mixture of greens and oranges with no immediately obvious patterns. If we recall that this is a continuous function with some sort of noise added upon it, we can start to make out some patterns in this data. There is a green sort of streak coming from the middle left of the data, corresponding to some models which generalized particularly well. We can also see a sort of valley in the top right quadrant near the middle right. It seems this valley could correspond to the first descent of the double descent curve. To see this one has to keep in mind that complexity is not a linear function of layer number and hidden unit size. Consequently, the second descent will not map out a linear curve in this plane and it will be hard to determine exactly what the 'peaks' which correspond to the interpolation maximum will manifest itself as in this plane. Nonetheless, we do start to see the double descent curve on models which are only a few layers deep (from the green section in the bottom left)

IV. CONCLUSION

Overall, this project turned out to be more about how to tune a deep neural network than it did about discovering the double descent curve. This is slightly disappointing, however, the last results produced are clearly the beginning of uncovering the double descent curve in this deep neural network setting. It is likely that other methodologies to help train deep neural networks could be useful in furthering this study so long as they do not inhibit the network's learning by regularization techniques. An example of an unhelpful technique is penalizing the network's weights. The double descent paradigm intends to achieve zero training loss which isn't feasible with a regularizing term on the size of the network's weights. Other schemes like an adapted learning rate which is able to get even lower scores on the training loss function is an example of a technique which would be useful.

As this project was limited by computational resources, it may also be true that for real applications as well, it becomes too difficult to actually go beyond the point of interpolation. Hence, these models would want to fit themselves snugly within the first descent. Ultimately, the double descent curve is a very promising outlook which seems like it will generalize to deeper network models. Unfortunately, the details are not worked out today by this project, but this should only inspire further work on the fascinating topic.

REFERENCES

- [1] Mikhail Belkin et al. "Reconciling modern machine learning and the bias-variance trade-off". In: *arXiv preprint arXiv:1812.11118* (2018).
- [2] Noam Shazeer and Mitchell Stern. "Adafactor: Adaptive Learning Rates with Sublinear Memory Cost". In: *arXiv preprint arXiv:1804.04235* (2018).