

# Fitting and initializing neural networks

Neural networks are almost always fitted with gradient based optimizers, such as variants of Stochastic Gradient Descent<sup>1</sup>. We defer how to compute the gradients to the next note.

## 1 Initialization

How do we set the initial weights before calling an optimizer? *Don't* set all the weights to zero! If different hidden units (adaptable basis functions) start out with the same parameters, they will all compute the same function of the inputs. Each unit will then get the same gradient vector, and be updated in the same way. As each hidden unit remains the same, we can't fit anything much more interesting than logistic regression.

Instead we usually initialize the weights randomly. *Don't* simply set all the weights using `randn()` though! As a concrete example, if all your inputs were  $x_d \in \{-1, +1\}$  the activation  $(\mathbf{w}^{(k)})^\top \mathbf{x}$  to hidden unit  $k$  would have zero mean, but typical size  $\sqrt{D}$  if there are  $D$  inputs. (See the review of random walks on the expectations sheet.) If your units saturate, like the logistic sigmoid, most of the gradients will be close to zero, and it will be hard for the gradient optimizer to update the parameters to useful settings.

Summary: initialize a weight matrix that transforms  $K$  values to small random values, like  $0.1 * \text{randn}() / \text{sqrt}(K)$ , assuming your input features are  $\sim 1$ .

The MLP course points to Glorot and Bengio's (2010) paper Understanding the difficulty of training deep feedforward networks, which suggests a scaling  $\propto 1 / \sqrt{K^{(l)} + K^{(l-1)}}$ , involving the number of hidden units in the layer after the weights, not just before. The argument involves the gradient computations, which I haven't described in detail for neural networks yet, so I will defer the interested reader to the paper or the MLP slides<sup>2</sup>.

Some specialized neural network architectures have particular tricks for initializing them. Do a literature search if you find yourself trying something other than a standard dense feedforward network: e.g., recurrent/recursive architectures, convolutional architectures, transformers, or memory networks. Alternatively, a pragmatic tip: if you are using a neural network toolbox, try to process your data to have similar properties to the standard datasets that are usually used to demonstrate that software. For example, similar dimensionality, means, variances, sparsity (number of non-zero features). Then any initialization tricks that the demonstrations use are more likely to carry over to your setting.

## 2 Local optima

The cost function for neural networks is not unimodal, and so is certainly not convex (a stronger property). We can see why by considering a neural network with two hidden units. Assume we've fitted the network to a (local) optimum of a cost function, so that any small change in parameters will make the network worse. Then we can find another parameter vector that will represent exactly the same function, showing that the optimum is only a local one.

To create the second parameter vector, we simply take all of the parameters associated with hidden unit one, and replace them with the corresponding parameters associated with hidden unit two. Then we take all of the parameters associated with hidden unit two and replace them with the parameters that were associated with hidden unit one. The network is really the same as before, with the hidden units labelled differently, so will have the same cost.

1. Adam (<https://arxiv.org/abs/1412.6980>) has now been popular for some time, although pure SGD is still in use too.

2. <http://www.inf.ed.ac.uk/teaching/courses/mlp/2019-20/lectures/mlp06-enc.pdf>

Models with “hidden” or “latent” representations of data, usually have many equivalent ways to represent the same model. When the goal of a machine learning system is to make predictions, it doesn’t matter whether the parameters are well-specified. However, it’s worth remembering that the values of individual parameters are often completely arbitrary, and can’t be interpreted in isolation.

In practice local optima don’t just correspond to permuting the hidden units. Some local optima will have better cost than others, and some will make predictions that generalize better than others. When I’ve fitted small neural networks, I’ve tried optimizing many times and used the network that cross-validates the best. However, researchers pushing up against available computational resources will find it difficult to optimize a network many times.

One advantage of large neural networks is that fitting far more parameters than necessary tends to work better(!). One intuition is that there are many more ways to set the parameters to get low cost, so it’s less hard to find one good setting.<sup>3</sup> Although it’s difficult to make rigorous statements on this issue. Understanding the difficulties that are faced in really high-dimensional optimization is an open area of research. (For example, <https://arxiv.org/abs/1412.6544>.)

### 3 Regularization by early stopping

We have referred to complex models that generalize poorly as “overfitted”. One idea to avoid “overfitting” is to fit less! That is, stop the optimization routine before it has found a local optimum of the cost function. This heuristic idea is often called “early stopping”.

The most common way to implement early stopping is to periodically monitor performance on a validation set. If the validation score is the best that we have seen so far, we save a copy of the network’s parameters. If the validation score fails to improve upon that cost over some number of future checks (say 20), we stop the optimization and return the weights we’ve saved.

David MacKay’s textbook mentions early stopping (Section 39.4, p479). This book points out that stopping the optimizer prevents the weights from growing too large. Goodfellow et al.’s deep learning textbook (Chapter 7) makes a more detailed link to L2 regularization. MacKay argued that adding a regularization term to the cost function to achieve a similar effect seems more appealing: if we have a well-defined cost function, we’re not tied to a particular optimizer, and it’s probably easier to analyse what we’re doing.

However, I’ve found it hard to argue with early stopping as a pragmatic, sensible procedure. The heuristic directly checks whether continuing to fit is improving predictions for held-out data, which is what we care about. And we might save a lot of computer time by stopping early. Moreover, we can still use a regularized cost function along with early stopping.

### 4 Regularization corrupting the data or model

There are a whole family of methods for regularizing models that involve adding noise to the data or model during training. Like early-stopping, I found this idea unappealing, as it’s hard to understand what objective we are fitting, and it makes the models we obtain depend on which optimizer we are using. However, these methods are often effective. . .

Adding Gaussian noise to the inputs of a linear model during gradient training has the same average effect as L2 regularization<sup>4</sup>. We can also add noise when training neural networks. The procedure will still have a regularization effect, but one that’s harder to understand. We can also add noise to the weights or hidden units of a neural network.

3. The high-level idea is old, but a recent (2018) analysis described the idea that some parts of a large network “get lucky” and identify good features as “The Lottery Ticket Hypothesis”, <https://arxiv.org/abs/1803.03635>

4. Non-examinable: there’s a sketch in these slides: [http://www.cs.toronto.edu/~tijmen/csc321/slides/lecture\\_slides lec9.pdf](http://www.cs.toronto.edu/~tijmen/csc321/slides/lecture_slides lec9.pdf). More detail in Bishop’s (1995) neural network textbook, section 9.3.

In some applications, adding noise has worked better than optimizing easy-to-define cost functions (like L2 regularizers).

Other regularization methods randomly replace some of the weights with zeros (“drop-out”<sup>5</sup>) or features with zeros (such as in “denoising auto-encoders”<sup>6</sup> or a 2006 feature-dropping regularizer). These heuristics prevent the model from fitting delicate combinations of parameters, or fits that depend on careful combinations of features. If used aggressively, “masking noise” makes it hard to fit anything! Often *large* models are needed when using these heuristics.

## 5 Further Reading

Most textbooks are long out-of-date when it comes to recent practical wisdom on fitting neural networks and regularization strategies. However, <http://www.deeplearningbook.org/> is still fairly recent, and is a good starting point. The MLP notes are also more detailed on practical tips for deep nets.

I’ll mention two of the recent tricks that I think are particularly worth knowing that make it easier to fit deep networks. But like most ideas, they don’t always improve things, so experiments are required. And the research landscape and available tools are changing rapidly.

The first trick, *Batch Normalization* (or “batch norm”), is “old” enough to be covered in the deep learning textbook. The discussion in the previous note about initialization pointed out that we don’t want to saturate hidden units. Batch normalization rescales the activations for a unit across a training batch to a target variance, making gradient-based training of neural nets easier in many tasks. In hindsight I’m amazed this trick is so recent: it’s a simple idea that someone could have come up with in a previous decade, but didn’t. Variants are still being actively explored, so I recommend chasing the recent literature if interested.

Another trick is the use of *residual layers*. There are different variants, especially when combined with batch norm in different places. However, the core idea is to take a standard non-linearity  $g$  and transform one hidden layer to another with  $r(\mathbf{h}) = \mathbf{h} + g(W\mathbf{h})$ . The weights  $W$  are used to fit the “residual” difference from a default transformation, that leaves the hidden unit values unchanged. Related, more complicated, layers were previously developed in the recurrent neural network literature, such as the “LSTM” cell, and are also worth knowing about. It is often easier to fit deep stacks of residual layers (or LSTMs, or GRUs), than standard layers.

---

5. <https://www.cs.toronto.edu/~hinton/absps/JMLRdropout.pdf>

6. <http://icml2008.cs.helsinki.fi/papers/592.pdf>