

# Fitting Batch Norm into a neural network

(SPEECH)

>>

(DESCRIPTION)

Text, Batch Normalization. Fitting Batch Norm into a neural network. Website, deep learning, dot, A.I.

(SPEECH)

So you have seen the equations for how to invent Batch Norm for maybe a single hidden layer.

Let's see how it fits into the training of a deep network.

(DESCRIPTION)

New slide, Adding Batch Norm to a network.

(SPEECH)

So, let's say you have a neural network like this, you've seen me say before that you can view each of the unit as computing two things.

First, it computes  $Z$  and then it applies the activation function to compute  $A$ .

And so we can think of each of these circles as representing a two step computation.

And similarly for the next layer, that is  $Z_2$ , and  $A_2$ , and so on.

So, if you were not applying Batch Norm, you would have an input  $X$  fit into the first hidden layer, and then first compute  $Z_1$ , and this is governed by the parameters  $W_1$  and  $B_1$ .

And then ordinarily, you would fit  $Z_1$  into the activation function to compute  $A_1$ .

But what would do in Batch Norm is take this value  $Z_1$ , and apply Batch Norm, sometimes abbreviated BN to it, and that's going to be governed by parameters,  $\beta_1$  and  $\gamma_1$ , and this will give you this new normalized value  $\tilde{Z}_1$ .

And then you fit that to the activation function to get  $A_1$ , which is  $G_1$  applied to  $\tilde{Z}_1$ .

Now, you've done the computation for the first layer, where this Batch Norms that really occurs in between the computation from  $Z$  and  $A$ .

Next, you take this value  $A_1$  and use it to compute  $Z_2$ , and so this is now governed by  $W_2$ ,  $B_2$ .

And similar to what you did for the first layer, you would take  $Z_2$  and apply it through Batch Norm, and we abbreviate it to BN now.

This is governed by Batch Norm parameters specific to the next layer.

So  $\beta_2$ ,  $\gamma_2$ , and now this gives you  $\tilde{Z}_2$ , and you use that to compute  $A_2$  by applying the activation function, and so on.

So once again, the Batch Norms that happens between computing  $Z$  and computing  $A$ .

And the intuition is that, instead of using the un-normalized value  $Z$ , you can use the normalized value  $\tilde{Z}$ , that's the first layer.

The second layer as well, instead of using the un-normalized value  $Z_2$ , you can use the mean and variance normalized values  $\tilde{Z}_2$ .

So the parameters of your network are going to be  $W_1$ ,  $B_1$ .

It turns out we'll get rid of the parameters but we'll see why in the next slide.

But for now, imagine the parameters are the usual  $W_1$ .

$\beta_1$ ,  $\gamma_1$ ,  $\beta_2$ ,  $\gamma_2$ , and so on, for each layer in which you are applying Batch Norm.

For clarity, note that these Betas here, these have nothing to do with the hyperparameter  $\beta$  that we had for momentum over the computing the various exponentially weighted averages.

The authors of the Adam paper use  $\beta$  on their paper to denote that hyperparameter, the authors of the Batch Norm paper had used  $\beta$  to denote this parameter, but these are two completely different Betas.

I decided to stick with  $\beta$  in both cases, in case you read the original papers.

But the  $\beta_1$ ,  $\beta_2$ , and so on, that Batch Norm tries to learn is a different  $\beta$  than the hyperparameter  $\beta$  used in momentum and the Adam and RMSprop algorithms.

So now that these are the new parameters of your algorithm, you would then use whatever optimization you want, such as creating descent in order to implement it.

For example, you might compute  $D\beta_L$  for a given layer, and then update the parameters  $\beta$ , gets updated as  $\beta - \text{learning rate} \times D\beta_L$ . And you can also use Adam or RMSprop or momentum in order to update the parameters  $\beta$  and  $\gamma$ , not just creating descent.

And even though in the previous video, I had explained what the Batch Norm operation does, computes mean and variances and subtracts and divides by them.

If they are using a Deep Learning Programming Framework, usually you won't have to implement the Batch Norm step on Batch Norm layer yourself.

So the probing frameworks, that can be sub one line of code.

So for example, in terms of flow framework, you can implement Batch Normalization with this function.

We'll talk more about probing frameworks later, but in practice you might not end up needing to implement all these details yourself, knowing how it works so that you can get a better understanding of what your code is doing.

But implementing Batch Norm is often one line of code in the deep learning frameworks.

Now, so far, we've talked about Batch Norm as if you were training on your entire training site at the time as if you are using Batch gradient descent.

In practice, Batch Norm is usually applied with mini-batches of your training set.

(DESCRIPTION)

New slide, Working with mini-batches.

(SPEECH)

So the way you actually apply Batch Norm is you take your first mini-batch and compute  $Z_1$ .

Same as we did on the previous slide using the parameters  $\gamma_1$ ,  $\beta_1$  and then you take just this mini-batch and compute mean and variance of the  $Z_1$  on just this mini batch and then Batch Norm would subtract by the mean and divide by the standard deviation and then re-scale by  $\gamma_1$ , to give you  $\tilde{Z}_1$ , and all this is on the first mini-batch, then you apply the activation function to get  $A_1$ , and then you compute  $Z_2$  using  $\gamma_2$ ,  $\beta_2$ , and so on.

So you do all this in order to perform one step of gradient descent on the first mini-batch and then goes to the second mini-batch  $X_2$ , and you do something similar where you will now compute  $Z_1$  on the second mini-batch and then use Batch Norm to compute  $\tilde{Z}_1$ .

And so here in this Batch Norm step, You would be normalizing  $\tilde{Z}$  using just the data in your second mini-batch, so does Batch Norm step here.

Let's look at the examples in your second mini-batch, computing the mean and variances of the  $\tilde{Z}_1$ 's on just that mini-batch and re-scaling by  $\gamma$  and  $\beta$  to get  $\tilde{Z}$ , and so on.

And you do this with a third mini-batch, and keep training.

Now, there's one detail to the parameterization that I want to clean up, which is previously, I said that the parameters was  $W_L$ ,  $B_L$ , for each layer as well as  $\beta_L$ , and  $\gamma_L$ . Now notice that the way  $Z$  was computed is as follows,  $Z_L = W_L \times A_L + B_L$ . But what Batch Norm does, is it is going to look at the mini-batch and normalize  $Z_L$  to first of mean 0 and standard variance, and then a rescale by  $\beta_L$  and  $\gamma_L$ .

But what that means is that, whatever is the value of  $B_L$  is actually going to just get subtracted out, because during that Batch Normalization step, you are going to compute the means of the  $Z_L$ 's and subtract the mean.

And so adding any constant to all of the examples in the mini-batch, it doesn't change anything.

Because any constant you add will get cancelled out by the mean subtractions step.

So, if you're using Batch Norm, you can actually eliminate that parameter, or if you want, think of it as setting it permanently to 0.

So then the parameterization becomes  $Z_L$  is just  $W_L \times A_L$ , And then you compute  $Z_L$  normalized, and we compute  $\tilde{Z}_L = \gamma_L Z_L + \beta_L$ , you end up using this parameter  $\beta_L$  in order to decide what's that mean of  $\tilde{Z}_L$ . Which is why guess post in this layer.

So just to recap, because Batch Norm zeroes out the mean of these  $Z_L$  values in the layer, there's no point having this parameter  $B_L$ , and so you must get rid of it, and instead is sort of replaced by  $\beta_L$ , which is a parameter that controls that ends up affecting the shift or the biased terms.

Finally, remember that the dimension of  $Z_L$ , because if you're doing this on one example, it's going to be  $N_L$  by 1, and so  $B_L$ , a dimension,  $N_L$  by one, if  $N_L$  was the number of hidden units in layer  $L$ . And so the dimension of  $\beta_L$  and  $\gamma_L$  is also going to be  $N_L$  by 1 because that's the number of hidden units you have.

You have  $N_L$  hidden units, and so  $\beta_L$  and  $\gamma_L$  are used to scale the mean and variance of each of the hidden units to whatever the network wants to set them to.

(DESCRIPTION)

New slide, Implementing gradient descent.

(SPEECH)

So, let's pull all together and describe how you can implement gradient descent using Batch Norm.

Assuming you're using mini-batch gradient descent, it rates for  $T = 1$  to the number of many batches.

You would implement forward prop on mini-batch  $X^T$  and doing forward prop in each hidden layer, use Batch Norm to replace  $Z_L$  with  $\tilde{Z}_L$ . And so then it shows that within that mini-batch, the value  $Z$  end up with some normalized mean and variance and the values and the version of the normalized mean that and variance is  $\tilde{Z}_L$ . And then, you use back prop to compute  $DW$ ,  $DB$ , for all the values of  $L$ ,  $D\beta$ ,  $D\gamma$ .

Although, technically, since you have got to get rid of  $B$ , this actually now goes away.

And then finally, you update the parameters.

So,  $W$  gets updated as  $W$  minus the learning rate times, as usual,  $\beta$  gets updated as  $\beta$  minus learning rate times  $DB$ , and similarly for  $\gamma$ .

And if you have computed the gradient as follows, you could use gradient descent.

That's what I've written down here, but this also works with gradient descent with momentum, or RMSprop, or Adam.

Where instead of taking this gradient descent update, mini-batch you could use the updates given by these other algorithms as we discussed in the previous week's videos.

Some of these other optimization algorithms as well can be used to update the parameters  $\beta$  and  $\gamma$  that Batch Norm added to algorithm.

So, I hope that gives you a sense of how you could implement Batch Norm from scratch if you wanted to.

If you're using one of the Deep Learning Programming frameworks which we will talk more about later, hopefully you can just call someone else's implementation in the Programming framework which will make using Batch Norm much easier.

Now, in case Batch Norm still seems a little bit mysterious if you're still not quite sure why it speeds up training so dramatically, let's go to the next video and talk more about why Batch Norm really works and what it is really doing.