

1. Normalization in Deep Learning

Normalization is carried out in the input layer to normalize the raw data. However, for deep neural nets, the activation values of the middle layer nodes can get very large, causing the same problem as our example. This necessitates some type of normalization in the hidden layers as well [1], as shown in Figure 1. Let x be that tensor consisting of a batch of N images. Each of these images has C feature maps or channels with height H and weight W . Therefore, $x \in N \times C \times H \times W$ is a four-dimensional tensor.

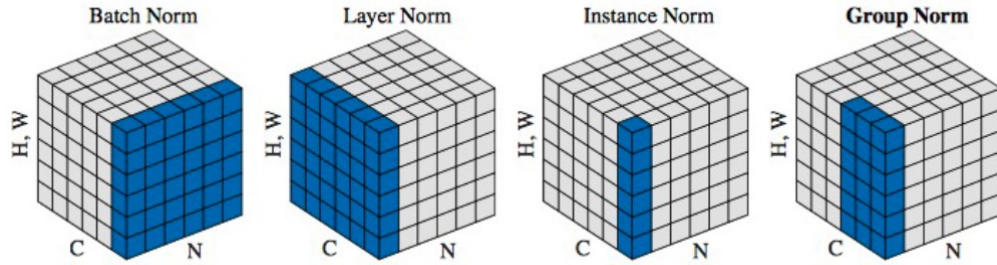


FIGURE 1. An overview of the normalization methods.

1.1. Batch Normalization. BN operates on one channel (one feature map) over all the training samples in the mini-batch. BN normalizes its activation values using the mean (μ_c) and variance (σ_c^2) of the mini-batch and introduces two learnable parameters, γ and β . The network can learn these parameters using a Gradient Descent algorithm.

$$(1) \quad \mu_c = \frac{1}{NHW} \sum_{n=1}^N \sum_{w=1}^W \sum_{h=1}^H x_{ncwh}$$

$$(2) \quad \sigma_c^2 = \frac{1}{NHW} \sum_{n=1}^N \sum_{w=1}^W \sum_{h=1}^H (x_{ncwh} - \mu_c)^2$$

We add the ϵ in computing \hat{x}_{ncwh} to the denominator to stabilize the training in case the variance of the mini-batch is zero.

$$(3) \quad \hat{x}_{ncwh} = \frac{x_{ncwh} - \mu_c}{\sqrt{\sigma_c^2 + \epsilon}}$$

Finally,

$$(4) \quad y_{ncwh} = \gamma \hat{x}_{ncwh} + \beta,$$

where γ and β parameters are the scaling and bias operators. This transformation preserves the network's capacity by ensuring that the normalization layer can produce the identity function. Both normalization methods can accelerate training and make the network converge faster. Plus, BN enables larger training rates and improves the generalization of the neural network. BN depends on the batch size in that, to obtain a statistically more accurate mean and variance, the batch size needs to be large. In addition, since BN is dependent on the batch size, it can't be applied at test time the same way as the training time. The reason is that in the test phase, we usually have one example to process. Therefore, the mean and variance can't be computed like training. Instead, during test time, BN utilizes moving average and variance to perform inference.

1.2. Instance Normalization. Instance normalization was first coined in the StyleNet paper [2]. Instance normalization tells us that it operates on a single sample. To clarify how IN works, let's consider sample feature maps that constitute an input tensor to the IN layer. In instance normalization, we consider one training sample and feature map and take the mean and variance over its spatial locations (W and H).

To perform instance normalization for a single instance x_{nc} , we need to compute the mean and variance:

$$(5) \quad \mu_{nc} = \frac{1}{HW} \sum_{w=1}^W \sum_{h=1}^H x_{ncwh}$$

$$(6) \quad \sigma_{nc}^2 = \frac{1}{HW} \sum_{w=1}^W \sum_{h=1}^H (x_{ncwh} - \mu_{nc})^2.$$

Then, using the obtained mean and variance, each spatial dimension is normalized:

$$(7) \quad \hat{x}_{ncwh} = \frac{x_{ncwh} - \mu_{nc}}{\sqrt{\sigma_{nc}^2 + \epsilon}},$$

where ϵ is a small value added for more stable training. In the end, using two scaling and bias parameters, we linearly transform the result.

While BN transforms the whole mini-batch of samples, IN does that to a single training sample. A small batch size setting is where IN performs better than BN. IN doesn't depend on the batch size, so its implementation is kept the same for both training and testing.

1.3. Layer Normalization. In BN, the statistics are computed across the batch and the spatial dims. In contrast, in Layer Normalization (LN), the statistics (mean and variance) are computed across all channels and spatial dims. Thus, the statistics are independent of the batch. This layer was initially introduced to handle vectors (mostly the RNN outputs).

Nearly nobody spoke about LN until the Transformers paper came out. We can take the mean across the spatial dimension and across all channels, as illustrated below:

$$(8) \quad \mu_n = \frac{1}{CHW} \sum_{c=1}^C \sum_{w=1}^W \sum_{h=1}^H x_{ncwh}$$

$$(9) \quad \sigma_n^2 = \frac{1}{CHW} \sum_{c=1}^C \sum_{w=1}^W \sum_{h=1}^H (x_{ncwh} - \mu_n)^2.$$

Then, using the obtained mean and variance, each spatial dimension is normalized:

$$(10) \quad \hat{x}_{ncwh} = \frac{x_{ncwh} - \mu_n}{\sqrt{\sigma_n^2 + \epsilon}},$$

where ϵ is a small value added for more stable training. In the end, using two scaling and bias parameters, we linearly transform the result.

1.4. Group Normalization. Group normalization (GN) divides the channels into groups and computes the first-order statistics within each group. As a result, GN's computation is independent of batch sizes, and its accuracy is more stable than BN in a wide range of batch sizes.

For *groups* = *number of channels* we get instance normalization, while for *groups* = 1 the method is reduced to layer normalization.

$$(11) \quad \mu_i = \frac{1}{m} \sum_{k \in \mathcal{S}_i} x_k,$$

2

$$(12) \quad \sigma_i^2 = \frac{1}{m} \sum_{k \in \mathcal{S}_i} (x_k - \mu_i)^2,$$

where

$$(13) \quad \mathcal{S}_i = \{k | k_N = i_N, \lfloor \frac{k_C}{C/G} \rfloor = \lfloor \frac{i_C}{C/G} \rfloor\}.$$

Note that G is the number of groups, which is a hyper-parameter. C/G is the number of channels per group. Thus, GN computes μ and σ along the (H, W) axes and along a group of C/G channels.

1.5. Adaptive Instance Normalization. Normalization and style transfer are closely related. In Instance Normalization, what if γ and β are injected from the feature statistics of another image y ? In this way, we will be able to model any arbitrary style by just giving our desired feature image mean as β and variance as γ from style image y .

Adaptive Instance Normalization (AdaIN) receives an input image x (content) and a style input y , and simply aligns the channel-wise mean and variance of x to match those of y . Mathematically:

$$(14) \quad AdaIN(x, y) = \sigma(y) \frac{x - \mu(x)}{\sigma(x)} + \mu(y)$$

References

- [1] Nikolas Adaloglou. In-layer normalization techniques for training very deep neural networks. <https://theaisummer.com/>, 2020.
- [2] Dmitry Ulyanov, Andrea Vedaldi, and Victor Lempitsky. Improved texture networks: Maximizing quality and diversity in feed-forward stylization and texture synthesis. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 6924–6932, 2017.

7.6.1 Injecting Noise at the Input

Some classical regularization techniques can be derived in terms of training on noisy inputs.⁵ Let us consider a regression setting, where we are interested in learning a model $\hat{y}(\mathbf{x})$ that maps a set of features \mathbf{x} to a scalar. The cost function we will use is the least-squares error between the model prediction $\hat{y}(\mathbf{x})$ and the true value y :

$$J = \mathbb{E}_{p(\mathbf{x}, y)} [(\hat{y}(\mathbf{x}) - y)^2], \quad (7.15)$$

where we are given a dataset of m input / output pairs $\{(\mathbf{x}^{(1)}, y^{(1)}), \dots, (\mathbf{x}^{(m)}, y^{(m)})\}$. Our training objective is to minimize the loss function, which in this case is given by the least-squares error between the model prediction $\hat{y}(\mathbf{x})$ and the true label y (where $\mathbf{y}(\mathbf{x}) = [y^{(1)}(\mathbf{x}), \dots, y^{(m)}(\mathbf{x})]$).

Now consider that with each input presentation to the model we also include a random perturbation $\epsilon \sim (\mathbf{0}, \nu \mathbf{I})$, so that the error function becomes

$$\begin{aligned} \tilde{J}_{\mathbf{x}} &= \mathbb{E}_{p(\mathbf{x}, y, \epsilon)} [(\hat{y}(\mathbf{x} + \epsilon) - y)^2] \\ &= \mathbb{E}_{p(\mathbf{x}, y, \epsilon)} [\hat{y}^2(\mathbf{x} + \epsilon) - 2y\hat{y}(\mathbf{x} + \epsilon) + y^2] \\ &= \mathbb{E}_{p(\mathbf{x}, y, \epsilon)} [\hat{y}^2(\mathbf{x} + \epsilon)] - 2\mathbb{E}_{p(\mathbf{x}, y, \epsilon)} [y\hat{y}(\mathbf{x} + \epsilon)] + \mathbb{E}_{p(\mathbf{x}, y, \epsilon)} [y^2] \end{aligned} \quad (7.16)$$

Assuming small noise, we can consider the Taylor series expansion of $\hat{y}(\mathbf{x} + \epsilon)$ around $\hat{y}(\mathbf{x})$.

$$\hat{y}(\mathbf{x} + \epsilon) = \hat{y}(\mathbf{x}) + \epsilon^\top \nabla_{\mathbf{x}} \hat{y}(\mathbf{x}) + \frac{1}{2} \epsilon^\top \nabla_{\mathbf{x}}^2 \hat{y}(\mathbf{x}) \epsilon + O(\epsilon^3) \quad (7.17)$$

Substituting this approximation for $\hat{y}(\mathbf{x} + \epsilon)$ into the objective function (Eq. 7.16) and using the fact that $\mathbb{E}_{p(\epsilon)}[\epsilon] = 0$ and that $\mathbb{E}_{p(\epsilon)}[\epsilon \epsilon^\top] = \nu \mathbf{I}$ to simplify⁶, we get:

$$\begin{aligned} \tilde{J}_{\mathbf{x}} &\approx \mathbb{E}_{p(\mathbf{x}, y, \epsilon)} \left[\left(\hat{y}(\mathbf{x}) + \epsilon^\top \nabla_{\mathbf{x}} \hat{y}(\mathbf{x}) + \frac{1}{2} \epsilon^\top \nabla_{\mathbf{x}}^2 \hat{y}(\mathbf{x}) \epsilon \right)^2 \right] \\ &\quad - 2\mathbb{E}_{p(\mathbf{x}, y, \epsilon)} \left[y\hat{y}(\mathbf{x}) + y\epsilon^\top \nabla_{\mathbf{x}} \hat{y}(\mathbf{x}) + \frac{1}{2} y\epsilon^\top \nabla_{\mathbf{x}}^2 \hat{y}(\mathbf{x}) \epsilon \right] + \mathbb{E}_{p(\mathbf{x}, y, \epsilon)} [y^2] \\ &= \mathbb{E}_{p(\mathbf{x}, y, \epsilon)} [(\hat{y}(\mathbf{x}) - y)^2] + \mathbb{E}_{p(\mathbf{x}, y, \epsilon)} \left[\hat{y}(\mathbf{x}) \epsilon^\top \nabla_{\mathbf{x}}^2 \hat{y}(\mathbf{x}) \epsilon + \left(\epsilon^\top \nabla_{\mathbf{x}} \hat{y}(\mathbf{x}) \right)^2 + O(\epsilon^3) \right] \\ &\quad - 2\mathbb{E}_{p(\mathbf{x}, y, \epsilon)} \left[\frac{1}{2} y \epsilon^\top \nabla_{\mathbf{x}}^2 \hat{y}(\mathbf{x}) \epsilon \right] \\ &= J + \nu \mathbb{E}_{p(\mathbf{x}, y)} [(\hat{y}(\mathbf{x}) - y) \nabla_{\mathbf{x}}^2 \hat{y}(\mathbf{x})] + \nu \mathbb{E}_{p(\mathbf{x}, y)} [\|\nabla_{\mathbf{x}} \hat{y}(\mathbf{x})\|^2] \end{aligned} \quad (7.18)$$

⁵The analysis in this section is mainly based on that in Bishop (1995a,b)

⁶In this derivation we have used two properties of the trace operator: (1) that a scalar is equal to its trace; (2) that, for a square matrix AB , $\text{Tr}(AB) = \text{Tr}(BA)$. These are discussed in Sec. 2.10.

If we consider minimizing this objective function, by taking the functional gradient of $\hat{y}(\mathbf{x})$ and setting the result to zero, we can see that

$$\hat{y}(\mathbf{x}) = \mathbb{E}_{p(y|\mathbf{x})}[y] + O(\nu).$$

This implies that the expectation in the second last term in Eq. 7.18, $\mathbb{E}_{p(\mathbf{x},y)} [(\hat{y}(\mathbf{x}) - y) \nabla_{\mathbf{x}}^2 \hat{y}(\mathbf{x})]$ reduces to $O(\nu)$ because the expectation of the difference $(\hat{y}(\mathbf{x}) - y)$ is reduces to $O(\nu)$.

This leaves us with the objective function of the form

$$\tilde{J}_{\mathbf{x}} = \mathbb{E}_{p(\mathbf{x},y)} [(\hat{y}(\mathbf{x}) - y)^2] + \nu \mathbb{E}_{p(\mathbf{x},y)} [\|\nabla_{\mathbf{x}} \hat{y}(\mathbf{x})\|^2] + O(\nu^2).$$

For small ν , the minimization of J with added noise on the input (with covariance $\nu \mathbf{I}$) is equivalent to minimization of J with an additional *regularization* term given by $\nu \mathbb{E}_{p(\mathbf{x},y)} [\|\nabla_{\mathbf{x}} \hat{y}(\mathbf{x})\|^2]$.

Considering the behavior of this regularization term, we note that it has the effect of penalizing large gradients of the function $\hat{y}(\mathbf{x})$. That is, it has the effect of reducing the *sensitivity* of the output of the network with respect to small variations in its input \mathbf{x} . We can interpret this as attempting to build in some local robustness into the model and thereby promote generalization. We note also that for linear networks, this regularization term reduces to simple weight decay (as discussed in Sec. 7.2.1).

7.6.2 Injecting Noise at the Weights

Rather than injecting noise as part of the input, one could also consider adding noise directly to the model parameters. As we shall see, this can also be interpreted as equivalent (under some assumptions) to a more traditional form of regularization. Adding noise to the weights has been shown to be an effective regularization strategy in the context of recurrent neural networks⁷ Jim *et al.* (1996); Graves (2011b). In the following, we will present an analysis of the effect of weight noise on a standard feedforward neural network (as introduced in Chapter 6).

As we did in the last section, we again consider the regression setting, where we wish to train a function $\hat{y}(\mathbf{x})$ that maps a set of features \mathbf{x} to a scalar using the least-squares cost function between the model predictions $\hat{y}(\mathbf{x})$ and the true values y :

$$J = \mathbb{E}_{p(\mathbf{x},y)} [(\hat{y}(\mathbf{x}) - y)^2]. \quad (7.19)$$

We again assume we are given a dataset of m input / output pairs $\{(\mathbf{x}^{(1)}, y^{(1)}), \dots, (\mathbf{x}^{(m)}, y^{(m)})\}$

We now assume that with each input presentation we also include a random perturbation $\epsilon \mathbf{w} \sim (\mathbf{0}, \eta \mathbf{I})$ of the network weights. Let us imagine that we have

⁷Recurrent neural networks will be discussed in detail in Chapter 10

a standard L -layer MLP, we denote the perturbed model as $\hat{y}_{\epsilon_W}(\mathbf{x})$. Despite the injection of noise, we are still interested in minimizing the squared error of the output of the network. The objective function thus becomes:

$$\begin{aligned}\tilde{J}_W &= \mathbb{E}_{p(\mathbf{x}, y, \epsilon_W)} [(\hat{y}_{\epsilon_W}(\mathbf{x}) - y)^2] \\ &= \mathbb{E}_{p(\mathbf{x}, y, \epsilon_W)} [\hat{y}_{\epsilon_W}^2(\mathbf{x}) - 2y\hat{y}_{\epsilon_W}(\mathbf{x}) + y^2]\end{aligned}\quad (7.20)$$

Assuming small noise, we can consider the Taylor series expansion of $\hat{y}_{\epsilon_W}(\mathbf{x})$ around the unperturbed function $\hat{y}(\mathbf{x})$.

$$\hat{y}_{\epsilon_W}(\mathbf{x}) = \hat{y}(\mathbf{x}) + \epsilon_W^\top \nabla_W \hat{y}(\mathbf{x}) + \frac{1}{2} \epsilon_W^\top \nabla_W^2 \hat{y}(\mathbf{x}) \epsilon_W + O(\epsilon_W^3) \quad (7.21)$$

From here, we follow the same basic strategy that was laid-out in the pervious section in analyzing the effect of adding noise to the input. That is, we substitute the Taylor series expansion of $\hat{y}_{\epsilon_W}(\mathbf{x})$ into the objective function in Eq. 7.20.

$$\begin{aligned}\tilde{J}_W &\approx \mathbb{E}_{p(\mathbf{x}, y, \epsilon_W)} \left[\left(\hat{y}(\mathbf{x}) + \epsilon_W^\top \nabla_W \hat{y}(\mathbf{x}) + \frac{1}{2} \epsilon_W^\top \nabla_W^2 \hat{y}(\mathbf{x}) \epsilon_W \right)^2 \right] \\ &\quad \mathbb{E}_{p(\mathbf{x}, y, \epsilon_W)} \left[-2y \left(\hat{y}(\mathbf{x}) + \epsilon_W^\top \nabla_W \hat{y}(\mathbf{x}) + \frac{1}{2} \epsilon_W^\top \nabla_W^2 \hat{y}(\mathbf{x}) \epsilon_W \right) \right] + \mathbb{E}_{p(\mathbf{x}, y, \epsilon_W)} [y^2] \\ &= \mathbb{E}_{p(\mathbf{x}, y, \epsilon_W)} [(\hat{y}(\mathbf{x}) - y)^2] - 2\mathbb{E}_{p(\mathbf{x}, y, \epsilon_W)} \left[\frac{1}{2} y \epsilon_W^\top \nabla_W^2 \hat{y}(\mathbf{x}) \epsilon_W \right] \\ &\quad + \mathbb{E}_{p(\mathbf{x}, y, \epsilon_W)} \left[\hat{y}(\mathbf{x}) \epsilon_W^\top \nabla_W^2 \hat{y}(\mathbf{x}) \epsilon_W + \left(\epsilon_W^\top \nabla_W \hat{y}(\mathbf{x}) \right)^2 + O(\epsilon_W^3) \right].\end{aligned}\quad (7.22)$$

Where we have used the fact that $\mathbb{E}_{\epsilon_W} \epsilon_W = \mathbf{0}$ to drop terms that are linear in ϵ_W . Incorporating the assumption that $\mathbb{E}_{\epsilon_W} \epsilon_W^2 = \eta \mathbf{I}$, we have:

$$\tilde{J}_W \approx J + \nu \mathbb{E}_{p(\mathbf{x}, y)} [(\hat{y}(\mathbf{x}) - y) \nabla_W^2 \hat{y}(\mathbf{x})] + \nu \mathbb{E}_{p(\mathbf{x}, y)} [\|\nabla_W \hat{y}(\mathbf{x})\|^2] \quad (7.24)$$

Again, if we consider minimizing this objective function, we can see that the optimal value of $\hat{y}(\mathbf{x})$ is:

$$\hat{y}(\mathbf{x}) = \mathbb{E}_{p(y|\mathbf{x})}[y] + O(\eta),$$

implying that the expectation in the middle term in Eq. 7.24, $\mathbb{E}_{p(\mathbf{x}, y)} [(\hat{y}(\mathbf{x}) - y) \nabla_W^2 \hat{y}(\mathbf{x})]$, reduces to $O(\eta)$ because the expectation of the difference $(\hat{y}(\mathbf{x}) - y)$ is reduces to $O(\eta)$.

This leaves us with the objective function of the form

$$\tilde{J}_W = \mathbb{E}_{p(\mathbf{x}, y)} [(\hat{y}(\mathbf{x}) - y)^2] + \eta \mathbb{E}_{p(\mathbf{x}, y)} [\|\nabla_W \hat{y}(\mathbf{x})\|^2] + O(\eta^2).$$

For small η , the minimization of J with added weight noise (with covariance $\eta \mathbf{I}$) is equivalent to minimization of J with an additional *regularization* term: $\eta \mathbb{E}_{p(\mathbf{x}, y)} [\|\nabla_{\mathbf{w}} \hat{y}(\mathbf{x})\|^2]$. This form of regularization encourages the parameters to go to regions of parameter space that have relatively small gradients. In other words, it pushes the model into regions where the model is relatively insensitive to small variations in the weights. Regularization strategies with this kind of behaviour have been considered before TODO restore this broken citation when it is fixed In the simplified case of linear regression (where, for instance, $\hat{y}(\mathbf{x}) = \mathbf{w}^\top \mathbf{x} + b$), this regularization term collapses into $\eta \mathbb{E}_{p(\mathbf{x})} [\|\mathbf{x}\|^2]$, which is not a function of parameters and therefore does not contribute to the gradient of $\tilde{J}_{\mathbf{w}}$ w.r.t the model parameters.

7.7 Early Stopping as a Form of Regularization

When training large models with high capacity, we often observe that training error decreases steadily over time, but validation set error begins to rise again. See Fig. 7.3 for an example of this behavior. This behavior occurs very reliably.

This means we can obtain a model with better validation set error (and thus, hopefully better test set error) by returning to the parameter setting at the point in time with the lowest validation set error. Instead of running our optimization algorithm until we reach a (local) minimum, we run it until the error on the validation set has not improved for some amount of time. Every time the error on the validation set improves, we store a copy of the model parameters. When the training algorithm terminates, we return these parameters, rather than the latest parameters. This procedure is specified more formally in Alg. 7.1.

This strategy is known as *early stopping*. It is probably the most commonly used form of regularization in deep learning. Its popularity is due both to its effectiveness and its simplicity.

One way to think of early stopping is as a very efficient hyperparameter selection algorithm. In this view, the number of training steps is just another hyperparameter. We can see in Fig. 7.3 that this hyperparameter has a U-shaped validation set performance curve, just like most other model capacity control parameters. In this case, we are controlling the effective capacity of the model by determining how many steps it can take to fit the training set precisely. Most of the time, setting hyperparameters requires an expensive guess and check process, where we must set a hyperparameter at the start of training, then run training for several steps to see its effect. The “training time” hyperparameter is unique in that by definition a single run of training tries out many values of the hyperparameter. The only significant cost to choosing this hyperparameter automatically via early stopping is running the validation set evaluation periodically during