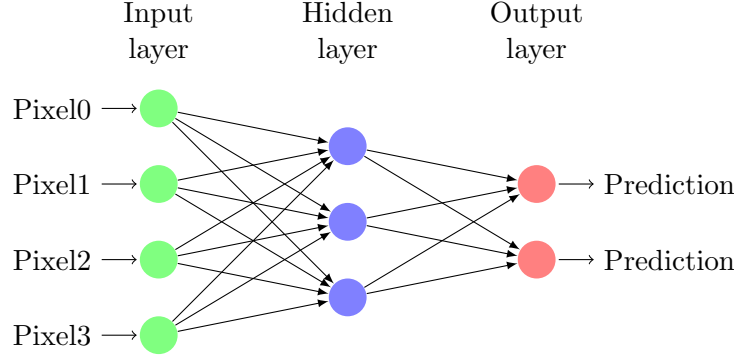
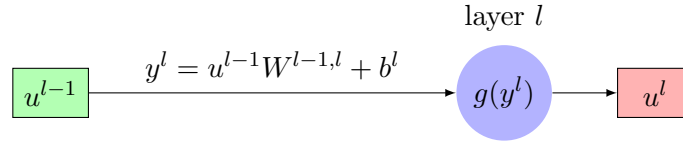


1 Baseline:

Network setup (schematic):



Notations and flows:



u^l is the l th layer's output, $W^{l-1,l}$ is the weight connecting layer $l-1$ and layer l , b^l is the bias at layer l . $g(y)$ is an activation function. The cost function for training can be found in many forms, cross-Entropy, regression etc. In my code, cross-Entropy is used:

$$C(\{\psi_s\}, \{\phi_s\}, \theta) = -\frac{1}{M} \left[\sum_{s=1}^M \phi_s \log[\mathcal{N}_\theta(\psi_s)] + (1 - \phi_s) \log[1 - \mathcal{N}_\theta(\psi_s)] \right]$$

ψ_s is a set of inputs, ϕ_s is the corresponding labels, $\theta = \{W, b, \hat{\theta}\}$ contains the weights W and bias b to be learned, as well as hyperparameters $\hat{\theta}$ (such as learning rate η). \mathcal{N}_θ is a network specified by parameters θ . Given an input ψ_s , its output $\mathcal{N}_\theta(\psi_s) = g(y_s^L) \equiv u_s^L$.

Stochastic Gradient Descent:

The workflow of SGD goes as the following:

- Shuffle the training sets, partition the shuffled set into m batches, each contains M_B samples
- For each batch:
 1. Feedforward, crossing all layers.
 2. BackPropagate, computing the the direction for descent: $\mathbf{v} = \left[\frac{\partial C}{\partial W^{i-1,i}}, \frac{\partial C}{\partial b^i} \right], i = 2 \cdots N_L$, N_L being number of layers.
 3. Update the parameters: $W, b \leftarrow W, b - \eta \cdot \mathbf{v}$, in BaseLine, the learning rate is isotropic.

Back-Propagation:

Introduce error rate/sensitivity at layer l for sample s : $\Delta_s^l = \frac{\partial C}{\partial y_s^l}$. This quantity is crucial for BP algorithm.

At the last layer, the error rate is:

$$\Delta_s^L \equiv \frac{\partial C}{\partial y_s^L} = \frac{\partial C}{\partial u_s^L} g'(y_s^L)$$

If activation function is sigmoid, then the error rate at the last layer becomes:

$$\Delta_s^L = \frac{1}{M_B}(u_s^L - \phi_s),$$

simply proportional to the difference between network's prediction and the labels. For 2 neighboring layers $\{l, l+1\}$,

$$\begin{aligned}\Delta_s^l &= \frac{\partial C}{\partial y_s^l} = \sum_{s'} \frac{\partial C}{\partial y_{s'}^{l+1}} \frac{\partial y_{s'}^{l+1}}{\partial y_s^l} = \sum_{s'} \Delta_{s'}^{l+1} \frac{\partial}{\partial y_s^l} (W^{l,l+1} g(y_{s'}^l) + b^l) \\ &= \sum_{s'} W^{l,l+1} \Delta_{s'}^{l+1} g'(y_s^l) \delta_{s,s'} \\ &= W^{l,l+1} \Delta_s^{l+1} g'(y_s^l)\end{aligned}$$

$\delta_{s,s'} = \{1(s = s'), 0(s \neq s')\}$. For the adjustment of weight:

$$\frac{\partial C}{\partial W^{l-1,l}} = \sum_s \frac{\partial C}{\partial y_s^l} \frac{\partial y_s^l}{\partial W^{l-1,l}} = \sum_s \Delta_s^l u_s^{l-1}$$

For the adjustment of bias:

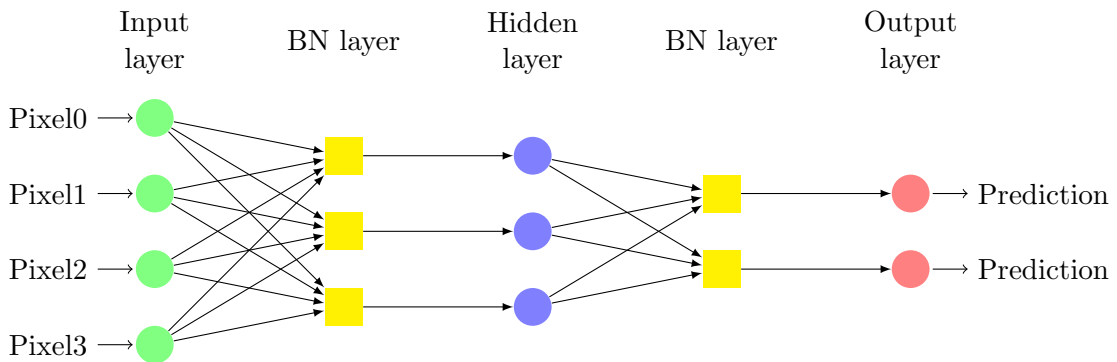
$$\frac{\partial C}{\partial b^l} = \sum_s \frac{\partial C}{\partial y_s^l} \frac{\partial y_s^l}{\partial b^l} = \sum_s \Delta_s^l$$

All in all, the BP algorithm for Baseline is:

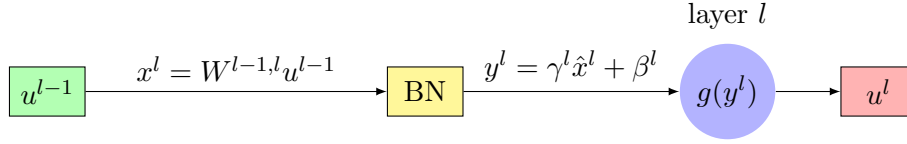
$$\begin{aligned}\Delta_s^l &= W^{l,l+1} \Delta_s^{l+1} g'(y_s^l) \\ \frac{\partial C}{\partial W^{l-1,l}} &= \sum_s \Delta_s^l u_s^{l-1} \\ \frac{\partial C}{\partial b^l} &= \sum_s \Delta_s^l \\ \Delta_s^L &= \frac{1}{M_B}(u_s^L - \phi_s)\end{aligned}$$

2 Batch Normalization:

Network setup (schematic):



Notations and flows:



One important thing about BN algorithm is that it normalize the input right **before** activation:

$$\begin{aligned} y^l &\leftarrow \mathcal{BN}(W^{l-1,l} u^{l-1}) \\ u^l &= g(y^l) \end{aligned}$$

$g(u)$ is activation function. It can be *sigmoid*, *ReLU* etc. $W^{l-1,l}$ is the weights linking layers $\{l-1, l\}$, u^l is the nonlinear activation of layer l . To summarize the normalization in B.N.:

$$\begin{aligned} x^l &= W^{l-1,l} u^{l-1} \\ \hat{x}^l &= \frac{x^l - \mu_B^l}{\sqrt{\sigma_B^{2,l}}} \\ y^l &= \gamma^l \hat{x}^l + \beta^l \\ u^l &= g(y^l) \end{aligned}$$

Back-Propagation:

Cost function $C(\gamma, \beta, W, \hat{\theta}, \psi, \phi)$ is a function of hyper parameter $\hat{\theta}$, training parameter $\{\gamma, \beta, W\}$, input ψ and label ϕ . The whole business is centered at minimizing this beast. Define growth rate at layer l : $\Delta_s^l \equiv \frac{\partial C}{\partial y_s^l}$. Here y_s^l is the linear combination of sample s that waits to be activated. Using chain rule, BP for BN can be derived as below. It differs with the baseline BP slightly.

$$\begin{aligned} \frac{\partial C}{\partial W^{l-1,l}} &= \sum_s \frac{\partial C}{\partial y_s^l} \frac{\partial y_s^l}{\partial W^{l-1,l}} \equiv \sum_s \Delta_s^l \frac{\partial y_s^l}{\partial W^{l-1,l}} \\ \frac{\partial y_s^l}{\partial W^{l-1,l}} &= \frac{\gamma^l}{\sqrt{\sigma_B^{2,l}}} \left[u_s^{l-1} - \langle u^{l-1} \rangle - \hat{x}_s^l \langle u^{l-1} \hat{x}^l \rangle \right] \end{aligned} \quad (1)$$

$$\Delta_s^l = \frac{\gamma^{l+1}}{\sqrt{\sigma_B^{2,l+1}}} W^{l,l+1} g'(y_s^l) \left[\Delta_s^{l+1} - \langle \Delta^{l+1} \rangle - \hat{x}_s^{l+1} \langle \Delta^{l+1} \hat{x}^{l+1} \rangle \right] \quad (2)$$

$$\begin{aligned} \frac{\partial C}{\partial \beta^l} &= \sum_s \frac{\partial C}{\partial y_s^l} \\ \frac{\partial C}{\partial \gamma^l} &= \sum_s \frac{\partial C}{\partial y_s^l} \hat{x}_s^l \end{aligned}$$

$\langle \dots \rangle$ is the sample average. s is the label of sample in a mini-batch. l labels layer. $g'(x)$ is the derivative of activation function. At the boundary:

$$\Delta_s^L = \frac{\partial C}{\partial u_s^L} g'(y_s^L)$$

L labels the last layer in network. Derivation of $\frac{\partial y_s^l}{\partial W^{l-1,l}}$ and Δ_s^l both involves computing quantity $\frac{\partial y_s^l}{\partial x_{s'}^l}$, which exemplifies the mixing of training samples:

$$\begin{aligned}
\frac{\partial}{\partial x_{s'}^l} (\sigma^l)^2 &= \frac{\partial}{\partial x_{s'}^l} \left[\frac{1}{M_B} \sum_s^{M_B} (x_s^l - \mu_B^l)^2 \right] = \frac{2}{M_B} \sum_s (x_s^l - \mu_B^l) (\delta_{s,s'} - \frac{1}{M_B}) \\
&= \frac{2}{M_B} (x_{s'}^l - \mu_B^l - \mu_B^l + \mu_B^l) = \frac{2}{M_B} (x_{s'}^l - \mu_B^l) \\
\frac{\partial y_s^l}{\partial x_{s'}^l} &= \frac{\partial}{\partial x_{s'}^l} \left[\gamma^l \frac{x_s^l - \mu_B^l}{\sigma^l} + \beta_l \right] \\
&= \frac{\gamma^l}{\sigma^l} \left[\delta_{s,s'} - \frac{1}{M_B} \right] - \frac{\gamma^l}{2(\sigma^l)^3} (x_s^l - \mu_B^l) \frac{\partial}{\partial x_{s'}^l} (\sigma^l)^2 \\
&= \frac{\gamma^l}{\sigma^l} \left[\delta_{s,s'} - \frac{1}{M_B} \right] - \frac{\gamma^l}{2(\sigma^l)^3} (x_s^l - \mu_B^l) \frac{2}{M_B} (x_{s'}^l - \mu_B^l) \\
&= \frac{\gamma^l}{\sigma^l} \left[\delta_{s,s'} - \frac{1}{M_B} - \frac{1}{M_B} \frac{x_s^l - \mu_B^l}{\sigma^l} \frac{x_{s'}^l - \mu_B^l}{\sigma^l} \right] \\
&= \frac{\gamma^l}{\sigma^l} \left[\delta_{s,s'} - \frac{1}{M_B} - \frac{1}{M_B} \hat{x}_s^l \hat{x}_{s'}^l \right] \tag{3}
\end{aligned}$$

With the above identity, one can compute the errors Δ and the weight adjustment $\frac{\partial y}{\partial W}$.

We first look into the errors. When the input to activation function at layer l changes: $y^l \rightarrow y^l + \delta y^l$, it will cause chain correction for all the layers $l' \geq l$. E.g. the change of cost δC caused by layer l equals to summing all the changes that happen at layer $l+1$ caused by δy^l :

$$\begin{aligned}
\Delta_s^l &\equiv \frac{\partial C}{\partial y_s^l} = \sum_{s'} \frac{\partial C}{\partial y_{s'}^{l+1}} \frac{\partial y_{s'}^{l+1}}{\partial y_s^l} \\
&= \sum_{s'} \Delta_{s'}^{l+1} \frac{\partial y_{s'}^{l+1}}{\partial x_s^{l+1}} \frac{\partial x_s^{l+1}}{\partial y_s^l} = \sum_{s'} \Delta_{s'}^l \frac{\partial y_{s'}^{l+1}}{\partial x_s^{l+1}} \frac{\partial}{\partial y_s^l} [W^{l,l+1} u_s^l] \\
&= \sum_{s'} \Delta_{s'}^{l+1} W^{l,l+1} g'(y_s^l) \frac{\partial y_{s'}^{l+1}}{\partial x_s^{l+1}}
\end{aligned}$$

Plug the kernel term (3) in,

$$\begin{aligned}
\Delta_s^l &= \sum_{s'} \Delta_{s'}^{l+1} W^{l,l+1} g'(y_s^l) \frac{\gamma^l}{\sigma^l} \left[\delta_{s,s'} - \frac{1}{M_B} - \frac{1}{M_B} \hat{x}_s^l \hat{x}_{s'}^l \right] \\
&= \frac{\gamma^l W^{l,l+1} g'(y_s^l)}{\sigma^l} \left[\Delta_s^{l+1} - \langle \Delta^{l+1} \rangle - \hat{x}_s^l \langle \Delta^{l+1} \hat{x}^l \rangle \right]
\end{aligned}$$

Similarly, one can also obtain the term:

$$\begin{aligned}
\frac{\partial y_s^l}{\partial W^{l-1,l}} &= \sum_{s'} \frac{\partial y_s^l}{\partial x_{s'}^l} \frac{\partial x_{s'}^l}{\partial W^{l-1,l}} \\
&= \sum_{s'} u_{s'}^{l-1} \frac{\gamma^l}{\sigma^l} \left[\delta_{s,s'} - \frac{1}{M_B} - \frac{1}{M_B} \hat{x}_s^l \hat{x}_{s'}^l \right] \\
&= \frac{\gamma^l}{\sigma^l} \left[u_s^{l-1} - \langle u^{l-1} \rangle - \hat{x}_s^l \langle \hat{x}^l u^{l-1} \rangle \right]
\end{aligned}$$

When one layer contains more than one neuron, one can just re-arrange the multiplication and addition/subtraction accordingly through dimension analysis. For example, when layer $l-1$ contains n_{l-1} neuron and layer

l contains n_l neurons, $\frac{\partial y_s^l}{\partial W^{l-1,l}}$ is a $n_{l-1} \times n_l$ matrix, the averaging term $\langle \hat{x}_s^l u^{l-1} \rangle$ can be written as the sample average of outer product $\frac{1}{M_B} \sum_s u_s^{l-1} \hat{x}_s^l$, which satisfies the dimension requirement.

3 Insights:

The key difference between BN and BL is the mixing among different samples in a minibatch. It can be clearly seen from the kernel expression (3). Its effect can be the modification of weights change and (1) and the error change(2). The weights adjustment in BP can be expressed in a matrix form in **sample space**:

$$\begin{aligned} \frac{\partial C}{\partial W^{l-1,l}} &= \frac{\gamma^l}{\sigma^l} \mathbf{\Delta}^l \cdot \mathbf{M} \cdot \mathbf{u}^{l-1} \\ \mathbf{\Delta}^l &= (\Delta_1^l, \Delta_2^l, \dots, \Delta_{M_B}^l) \\ \mathbf{u}^{l-1} &= (u_1^{l-1}, u_2^{l-1}, \dots, u_{M_B}^{l-1})^T \\ \mathbf{M} &= \mathbf{I} - \frac{1}{M_B} [\mathbf{O} + \hat{\mathbf{X}}^l] \end{aligned}$$

\mathbf{I} is an $M_B \times M_B$ identity matrix, \mathbf{O} is an $M_B \times M_B$ matrix with $O_{i,j} = 1, \forall i, j$; and $\hat{\mathbf{X}}_{i,j}^l = \hat{x}_i^l \hat{x}_j^l$.

To see what (3) does, it is helpful to simplify the neural network such that each layer only contains 1 neuron. Then the weights W between each layer becomes a scalar. Under such simplification, $\hat{x}_s^l = \frac{W u_s^{l-1} - \langle W u^{l-1} \rangle}{STD(W u^{l-1})} \equiv \hat{u}_s^{l-1}$. Keep in mind that the STDs for x and u are different: $\sigma^l \equiv STD(W u^{l-1})$, $\tilde{\sigma}^{l-1} \equiv STD(u^{l-1}) = W^{-1} \sigma^l$. Since \hat{x}_s^l is already normalized, $\sum_s \hat{x}_s^l = 0$, $\sum_s (\hat{x}_s^l)^2 = M_B, \forall l$. Now examine the weights update.

$$\begin{aligned} \mathbf{M} \cdot \mathbf{u}^{l-1} &= \mathbf{u}^{l-1} - \frac{1}{M_B} [\mathbf{O} + \hat{\mathbf{X}}^l] \cdot \left(\mathbf{u}^{l-1} - \langle u^{l-1} \rangle \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix} + \langle u^{l-1} \rangle \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix} \right) \frac{\tilde{\sigma}^{l-1}}{\tilde{\sigma}^{l-1}} \\ &= \mathbf{u}^{l-1} - \frac{1}{M_B} [\hat{\mathbf{O}} + \hat{\mathbf{X}}^l] \cdot \left(\hat{\mathbf{u}}^{l-1} + \frac{\langle u^{l-1} \rangle}{\tilde{\sigma}^{l-1}} \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix} \right) \tilde{\sigma}^{l-1} \\ &= \mathbf{u}^{l-1} - \frac{1}{M_B} \left(M_B \frac{\langle u^{l-1} \rangle}{\tilde{\sigma}^{l-1}} \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix} + M_B \frac{\hat{\mathbf{u}}^{l-1} \cdot \hat{\mathbf{u}}^{l-1}}{M_B} \hat{\mathbf{u}}^{l-1} \right) \tilde{\sigma}^{l-1} \\ &= \mathbf{u}^{l-1} - \langle u^{l-1} \rangle \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix} - \frac{\hat{\mathbf{u}}^{l-1} \cdot \hat{\mathbf{u}}^{l-1}}{M_B} \hat{\mathbf{u}}^{l-1} \tilde{\sigma}^{l-1} \\ &= \tilde{\sigma}^{l-1} \left(\hat{\mathbf{u}}^{l-1} - \frac{\hat{\mathbf{u}}^{l-1} \cdot \hat{\mathbf{u}}^{l-1}}{M_B} \hat{\mathbf{u}}^{l-1} \right) \\ &= \mathbf{0} !!! \end{aligned}$$

Weights Update

Now we consider a more general case when each layer has more than one neuron. Use subscripted latin letter i, j, k to denote the neuron. The BN procedure can be written as:

$$\begin{aligned} x_{s,i}^l &= \sum_j W_{i,j}^{l-1,l} u_{s,j}^{l-1} \\ \hat{x}_{s,i}^l &= \frac{x_{s,i}^l - \langle x_i^l \rangle}{\sigma_i^l} \\ y_{s,i}^l &= \gamma_i^l \hat{x}_{s,i}^l + \beta_i^l \\ u_{s,i}^l &= g(y_{s,i}^l) \end{aligned}$$

Since the i th neuron's BN input $y_{s,i}^l$ is only related with the i th neuron's $x_{s,i}^l$, the kernel function for each neuron is then :

$$\frac{\partial y_{s,i}^l}{\partial x_{s',i}^l} = \frac{\gamma_i^l}{\sigma_i^l} \left[\delta_{s,s'} - \frac{1}{M_B} - \frac{1}{M_B} \hat{x}_{s,i}^l \hat{x}_{s',i}^l \right]; \quad \frac{\partial y_{s,i}^l}{\partial x_{s',j}^l} = 0, i \neq j$$

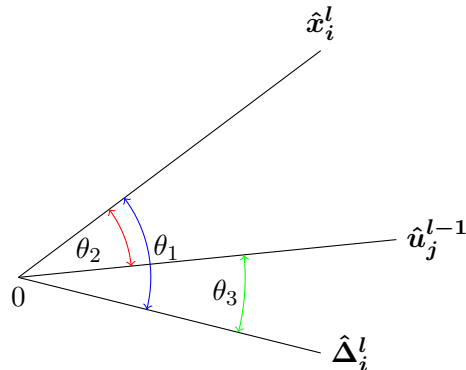
With this setup, the weight update rule in back-propagation becomes:

$$\begin{aligned} \frac{\partial C}{\partial W_{i,j}^{l-1,l}} &= \sum_{s,k} \frac{\partial C}{\partial y_{s,k}^l} \frac{\partial y_{s,k}^l}{\partial W_{i,j}^{l-1,l}} = \sum_{s,k} \Delta_{s,k}^l \frac{\partial y_{s,k}^l}{\partial W_{i,j}^{l-1,l}} \\ \frac{\partial y_{s,k}^l}{\partial W_{i,j}^{l-1,l}} &= \delta_{i,k} \frac{\gamma_i^l}{\sigma_i^l} \left[u_{s,j}^{l-1} - \langle u_j^{l-1} \rangle - \hat{x}_{s,i}^l \langle u_j^{l-1} \hat{x}_i^l \rangle \right] \\ \Rightarrow \frac{\partial C}{\partial W_{i,j}^{l-1,l}} &= M_B \frac{\gamma_i^l}{\sigma_i^l} \left(\langle \Delta_i^l u_j^{l-1} \rangle - \langle \Delta_i^l \rangle \langle u_j^{l-1} \rangle - \langle \Delta_i^l \hat{x}_i^l \rangle \langle u_j^{l-1} \hat{x}_i^l \rangle \right) \end{aligned}$$

Introduce two new variance, the variance of error term $\sigma(\Delta_i^l)$ and the variance of output $\sigma(u_j^{l-1})$. The normalized error $\hat{\Delta}_i^l$ and normalized output \hat{u}_j^{l-1} can be written as vectors in sample space: $\hat{\Delta}_i^l = (\hat{\Delta}_{1,i}^l, \hat{\Delta}_{2,i}^l, \dots, \hat{\Delta}_{M_B,i}^l)$, $\hat{u}_j^{l-1} = (\hat{u}_{1,j}^{l-1}, \hat{u}_{2,j}^{l-1}, \dots, \hat{u}_{M_B,i}^{l-1})$. The above equation can be rewritten as:

$$\begin{aligned} \frac{\partial C}{\partial W_{i,j}^{l-1,l}} &= M_B \gamma_i^l \frac{\sigma(\Delta_i^l) \sigma(u_j^{l-1})}{\sigma(x_i^l)} \left(\text{Corr}(\Delta_i^l, u_j^{l-1}) - \frac{\hat{\Delta}_i^l \cdot \hat{x}_i^l}{M_B} \frac{\hat{u}_j^{l-1} \cdot \hat{x}_i^l}{M_B} \right) \\ \text{Corr}(\Delta_i^l, u_j^{l-1}) &= \left\langle \frac{\Delta_{s,i}^l - \langle \Delta_i^l \rangle}{\sigma(\Delta_i^l)} \frac{u_{s,j}^{l-1} - \langle u_j^{l-1} \rangle}{\sigma(u_j^{l-1})} \right\rangle = \frac{\hat{\Delta}_i^l \cdot \hat{u}_j^{l-1}}{M_B} \end{aligned} \quad (4)$$

Recall that all the normalized quantity has fixed length in sample space $|\hat{\mathbf{a}}| = \sqrt{M_B}$, the above equation has a geometrical interpretation.



$$\begin{aligned}
BN : \quad \frac{\partial C}{\partial W_{i,j}^{l-1,l}} &= M_B \gamma_i^l \frac{\sigma(\Delta_i^l) \sigma(u_j^{l-1})}{\sigma(x_i^l)} [\cos(\theta_3) - \cos(\theta_1) \cos(\theta_2)] \\
BL : \quad \frac{\partial C}{\partial W_{i,j}^{l-1,l}} &= \sum_s \left[\sigma(\Delta_i^l) \hat{\Delta}_{s,i}^l + \langle \Delta_i^l \rangle \right] \left[\sigma(u_j^{l-1}) \hat{u}_{s,j}^{l-1} + \langle u_j^{l-1} \rangle \right] \\
&= M_B \sigma(\Delta_i^l) \sigma(u_j^{l-1}) \cos(\theta_3) + M_B \langle \Delta_i^l \rangle \langle u_j^{l-1} \rangle
\end{aligned} \tag{5}$$

Compare BN and BL, one can see 2 major differences. The first one being the amplitude of the adjustment. BN performs a more precise micro-tuning: the adjustment is proportional to variance. While BL also take the sample average into count, which can cause orders of difference when it comes to adjustment.

Secondly, the stopping condition is different. Assume that the layer's output u_i^l is normalized. For BL, the learning continues until the normalized error of current layer $\hat{\Delta}_i^l$ and normalized output of previous layer \hat{u}_j^{l-1} is not correlated anymore.

For BN, the weights has 2 functionality. (1) it rotates the previous layer's output so the rotated quantity \hat{x}^l can strongly correlate with the error term propagated back. This is reflected by the fact that the learning will stop if $\sin(\theta_1) = 0$. Once $\sin(\theta_1) = 0$, the learning is delegated to γ as shown in the next subsection. (2) Assume $\sin(\theta_1) \neq 0$. When the correlation between \hat{x}_i^l and \hat{u}_j^{l-1} is maximized, e.g. $\sin(\theta_2) = 0$, the update of weight $W_{i,j}^{l-1,l}$ stops. For another neuron in layer $l-1$, say neuron k , if its sample space vector is not aligned with neuron j in layer $l-1$, it will be updated. However, if \hat{u}_k^{l-1} is aligned with \hat{x}_i^l , its update will stop too. The net effect is that if \hat{u}_k^{l-1} and \hat{u}_j^{l-1} are strongly correlated in sample space (e.g. their vectors are aligned), the update of both $W_{i,j}^{l-1,l}$ and $W_{i,k}^{l-1,l}$ stop. In essence, BN could make all neurons of layer $l-1$ connecting to i th neuron in layer l to be synched.

γ Update

The update of γ can be analyzed in the similar way.

$$\frac{\partial C}{\partial \gamma_i^l} = \sigma(\Delta_i^l) \sum_s \hat{\Delta}_i^l \hat{x}_i^l = \sigma(\Delta_i^l) M_B \cos(\theta_1)$$

Since correlation is the same as dot product of two normalized quantity divided by sample number: $Corr(\hat{A}, \hat{B}) = \hat{A} \cdot \hat{B} / M_B$, the above equation can also be expressed as:

$$\frac{\partial C}{\partial \gamma_i^l} = \sigma(\Delta_i^l) M_B Corr(\hat{\Delta}_i^l, \hat{x}_i^l) \tag{6}$$

e.g., when normalized error and input x is highly correlated across the samples, the most probable cause will be γ hence γ will be adjusted most. If the correlation is 0, which indicates the errors have nothing to do with inputs, γ doesn't need to be adjusted.

Error Propagation

For the update of γ, β , the error propagation is crucial. With each layer having more than one neuron, the BP for Δ_s^l is:

$$\begin{aligned}
\Delta_{s,i}^l &= \frac{\partial C}{\partial y_{s,i}^l} = \sum_{s',k} \frac{\partial C}{\partial y_{s',k}^{l+1}} \frac{\partial y_{s',k}^{l+1}}{\partial y_{s,i}^l} = \sum_{s',s'',k} \Delta_{s',k}^{l+1} \frac{\partial y_{s',k}^{l+1}}{\partial x_{s'',k}^{l+1}} \frac{\partial x_{s'',k}^{l+1}}{\partial y_{s,i}^l} = \sum_{s',s'',k} \Delta_{s',k}^{l+1} \frac{\partial y_{s',k}^{l+1}}{\partial x_{s'',k}^{l+1}} W_{k,i}^{l,l+1} g'(y_{s,i}^l) \delta_{s,s''} \\
&= g'(y_{s,i}^l) \sum_{s',k} \Delta_{s',k}^{l+1} W_{k,i}^{l,l+1} \frac{\partial y_{s',k}^{l+1}}{\partial x_{s,k}^{l+1}} \\
&= g'(y_{s,i}^l) \sum_{s',k} \Delta_{s',k}^{l+1} W_{k,i}^{l,l+1} \frac{\gamma_k^{l+1}}{\sigma_k^{l+1}} \left[\delta_{s',s} - \frac{1}{M_B} - \frac{1}{M_B} \hat{x}_{s',k}^{l+1} \hat{x}_{s,k}^{l+1} \right] \\
&= g'(y_{s,i}^l) \sum_k W_{k,i}^{l,l+1} \frac{\gamma_k^{l+1}}{\sigma(x_k^{l+1})} \left[\Delta_{s,k}^{l+1} - \langle \Delta_k^{l+1} \rangle - \hat{x}_{s,k}^{l+1} \langle \Delta_k^{l+1} \hat{x}_k^{l+1} \rangle \right] \\
&= g'(y_{s,i}^l) \sum_k W_{k,i}^{l,l+1} \frac{\gamma_k^{l+1}}{\sigma(x_k^{l+1})} \sigma(\Delta_k^{l+1}) \left[\hat{\Delta}_{s,k}^{l+1} - \hat{x}_{s,k}^{l+1} \langle \hat{\Delta}_k^{l+1} \hat{x}_k^{l+1} \rangle \right]
\end{aligned}$$

Again, in matrix form (Sample Space):

$$\Delta_i^{l-1} = g'(\mathbf{y}_i^{l-1}) \odot \left[\sum_k W_{k,i}^{l-1,l} \frac{\gamma_k^l}{\sigma(x_k^l)} \sigma(\Delta_k^l) \left(\hat{\Delta}_k^l - \cos(\theta_1) \hat{\mathbf{x}}_k^l \right) \right]$$

where θ_1 is the angle between sample space vector $\hat{\Delta}_k^l$ and $\hat{\mathbf{x}}_k^l$. Given all other things equal, if $\hat{\Delta}_k^l$ and $\hat{\mathbf{x}}_k^l$ is aligned in sample space, the error won't propagate from this node backwards. If these 2 vectors are opposite, the error back-propagates with largest magnitude.