

WEEKLY PARTICIPATION 9

Recall that batch normalization layers are used as follows¹:

$$\begin{aligned}\mathbf{a}^\ell &= \text{BN}_{\gamma^\ell, \beta^\ell}(\mathbf{W}^\ell \mathbf{o}^{\ell-1}) = \gamma^\ell \odot \left(\frac{\mathbf{W}^\ell \mathbf{o}^{\ell-1} - \boldsymbol{\mu}_B}{\sqrt{\boldsymbol{\sigma}_B^2 + \epsilon}} \right) + \beta^\ell \\ &= \gamma^\ell \odot \overline{\mathbf{o}^{\ell-1}} + \beta^\ell, \\ \mathbf{o}^\ell &= \sigma(\mathbf{a}^\ell).\end{aligned}$$

For convenience, we have denoted the affinely transformed output of layer $\ell - 1$ by

$$\overline{\mathbf{o}^{\ell-1}} = \frac{\mathbf{W}^\ell \mathbf{o}^{\ell-1} - \boldsymbol{\mu}_B}{\sqrt{\boldsymbol{\sigma}_B^2 + \epsilon}}.$$

In the above expressions, ϵ is a small positive number to guard against division by zero, \odot denotes element-wise multiplication, and $\boldsymbol{\mu}_B, \boldsymbol{\sigma}_B^2 \in \mathbb{R}^{n_{\ell-1}}$ are the vectors of minibatch means and minibatch variances for each neuron in layer $\ell - 1$:

$$\boldsymbol{\mu}_B = \frac{1}{m} \sum_{i=1}^m (\mathbf{o}^{\ell-1})_i, \quad \boldsymbol{\sigma}_B = \frac{1}{m} \sum_{i=1}^m ((\mathbf{o}^{\ell-1})_i - \boldsymbol{\mu}_B)^2.$$

The scale and shift vectors $\gamma^\ell, \beta^\ell \in \mathbb{R}^{n_{\ell-1}}$ are parameters that must be learned during training, using backpropagation. Let f be the training objective.

- (1) Verify that $\mathbf{J}_{\mathbf{o}^\ell}(\mathbf{a}^\ell) = \text{diag}(\sigma'(\mathbf{a}^\ell))$, and use this fact to give an expression for $\nabla_{\mathbf{a}^\ell} f$, assuming that $\nabla_{\mathbf{o}^\ell} f$ is known.
- (2) Verify that $\mathbf{J}_{\mathbf{a}^\ell}(\gamma) = \text{diag}(\overline{\mathbf{o}^{\ell-1}})$, and use this fact to give an expression for $\nabla_{\gamma^\ell} f$ in terms of $\nabla_{\mathbf{a}^\ell} f$.
- (3) Verify that $\mathbf{J}_{\mathbf{a}^\ell}(\beta^\ell) = \mathbf{I}$, and use this fact to give an expression for $\nabla_{\beta^\ell} f$ in terms of $\nabla_{\mathbf{a}^\ell} f$.
- (4) Observe (no need to verify this) that for any $i \in [n_\ell]$,

$$\begin{aligned}[\mathbf{J}_{\mathbf{a}^\ell}(\mathbf{W}^\ell)]_{i,:} &= \left[\frac{\partial (\mathbf{a}^\ell)_i}{\partial (\mathbf{W}^\ell)_{p,q}} \right]_{p,q=1}^{n_\ell, n_{\ell-1}} \\ &= \text{diag} \left(\frac{\gamma^\ell}{\sqrt{\boldsymbol{\sigma}_B^2 + \epsilon}} \right) \mathbf{e}_i (\mathbf{o}^{\ell-1})^T,\end{aligned}$$

so for any vector $\mathbf{v} \in \mathbb{R}^{n_\ell}$,

$$\mathbf{J}_{\mathbf{a}^\ell}(\mathbf{W}^\ell)^T \mathbf{v} = \text{diag} \left(\frac{\gamma^\ell}{\sqrt{\boldsymbol{\sigma}_B^2 + \epsilon}} \right) \mathbf{v} (\mathbf{o}^{\ell-1})^T.$$

Use this fact to give an expression for $\nabla_{\mathbf{W}^\ell} f$ in terms of $\nabla_{\mathbf{a}^\ell} f$.

¹Notice that we do not include a bias \mathbf{b}^ℓ in the call to the BN primitive because β^ℓ is our bias.