

$$\frac{dL}{dW^{(l)}} = \frac{dL}{dZ^{(l)}} \frac{dZ^{(l)}}{dW^{(l)}} = \delta x^T \leftarrow \text{if } x \text{ is "imbalanced", so are gradients.}$$

Therefore, we really want all entries in  $x \rightarrow$  same scale.

Standardization:  $\mu = 0, \sigma = 1$ .

$$\begin{aligned} \bar{x}_i &= x_i - E(x) & E(x) &\approx \frac{1}{N} \sum_{i=1}^N x_i \\ \tilde{x}_i &= \frac{x_i - E(x)}{\sqrt{E[(x_i - E(x))^2]}} & \sigma_i & \end{aligned}$$

Basic idea of Batch Norm: standardize activations at each layer, controlling gradients. computationally expensive idea: so only perform in batch.

$$x^{(l)} \approx \frac{1}{B} \sum_{j=1}^B a_{ij}^{(l)} \quad r^{(l)} \approx \sqrt{\frac{1}{B} \sum_{j=1}^B (a_{ij}^{(l)} - \mu^{(l)})^2} \quad \bar{a}_i^{(l)} = \frac{a_i^{(l)} - \mu^{(l)}}{r^{(l)}} \gamma + \beta \quad \leftarrow \text{learnable scale and bias, same dim as } \bar{a}_i^{(l)}$$

$\rightarrow$  Can be trained w/ backprop since these are differentiable.

$\rightarrow$  Can be placed either before/after nonlinearity.

$\rightarrow$  We can often use a larger learning rate

$\rightarrow$  models can train faster

$\rightarrow$  generally requires less regularization

Basic initialization methods: ensure activations are on reasonable scale that stays constant

More advanced init. methods involve eigenvalues/Jacobians

"Try" to have well-behaved gradients from the get-go.

$$\text{Set } W_{ij} \sim N(0, \sigma^2_w), b_i \approx 0, z_i = \sum_j W_{ij} a_j + b_i \quad \text{assume } a_j \sim N(0, \sigma^2_a)$$

$$\text{so, set std of } W_{ij} = 1/\sqrt{D_a} \quad E[z_i^2] = \sum_j E[W_{ij}^2] E[a_j^2] = D_a \sigma^2_w \sigma^2_a \quad \text{dim of } a$$

"Xavier initialization" If  $D_a \sigma^2_w \gg 1$ , magnitude grows w/ each layer

ReLU issue: zeros out activations. If  $D_a \sigma^2_w < 1$ , magnitude shrinks w/ each layer  $\rightarrow \sigma^2_w = \frac{1}{D_a}$

Ergo, increase std of  $W_{ij} \rightarrow 1/\sqrt{1/2 D_a}$  (proposed in ResNet).

Alternative init.  $\tilde{z}_i = 0.1$  (or small constant) to avoid zero-out in ReLU.

$$\text{Reminder: } \frac{dL}{dW^{(l)}} = J_1 J_2 J_3 \dots J_n \quad \frac{dL}{dZ^{(l)}} \quad \forall x_i, J_i = U_i \Lambda_i V_i \quad (\text{SVD}) \quad \text{diag. matrix of } \lambda\text{'s of } J_i.$$

By ensuring  $\Lambda_i$ 's are scaled, we avoid 0/infinity convergence issue.

Since  $J = W^T$  (derivative w/  $W$ )  $\rightarrow W^{(l)} = U^{(l)} \Lambda^{(l)} V^{(l)}$ , force  $\Lambda^{(l)} = I$ ,  $W^{(l)} \leftarrow U^{(l)} V^{(l)}$

"Measure of last resort": Gradient clipping, because "monster gradients" can occur

$\rightarrow$  per-element clipping:  $\tilde{g}_i = \max(\min(g_i, c_i), -c_i)$

$\rightarrow$  norm clipping:

$$\tilde{g} \leftarrow g \frac{\min(\|g\|, c)}{\|g\|} \quad \leftarrow \text{clips length not dimension, choose } c \text{ thru experimentation}$$

NNS  $\rightarrow$  high-variance (lots of parameters), but with multiple, more agreement on "right"

$$\text{Variance} = E_p(p) [f_p(x) - \bar{f}(x)]^2 \quad \bar{f}(x) = E_p(p) [f_p(x)] \approx \frac{1}{M} \sum_{i=1}^M f_{p_i}(x) \quad \text{bootstrap sample}$$

Select?  $\rightarrow$  principled approach: average  $p(x|x) = \frac{1}{M} \sum_{j=1}^M p_j(x|x)$  (often more robust)

$\rightarrow$  simple approach: majority vote (analogous to "first past the post")

Ensembles in practice: train  $M$  models  $p_j(u|x)$  on the same  $\mathcal{D}$  as training set, maj. vote.

Even faster ensembles? share features for all  $M$  models. Separate ensemble of final classifier "heads" at end (these are task-specific)

Snapshot ensembles: save our parameters as "snapshots" during SGD, use later as a model. The bigger the ensemble, the better, the more expensive.

Dropout: randomly set some activations to 0 in forward pass. Creates a "new" network, helps w/ ensembling out of a single NN.  $2^{k_n}$  diff. models ( $k$  layers,  $n$  nodes/layer)

At test time  $\tilde{W}^{(l)} = \frac{1}{M} W^{(l)}$ , since on average,  $\frac{1}{M}$  of dimensions are forced to 0.

Hyperparameters that affect generalization (validation): ensembling, dropout, arch.

General method to pick hyperparameters: coarse-to-fine, broad sweep before "zeroing in"

$\hookrightarrow$  In practice: random hp. search? grid search.