

Basic Knowledge

$$\underline{C} = \underline{A}\underline{B} \rightarrow \underline{C}^T = \underline{B}^T \underline{A}^T, \quad \underline{C}^{-1} = \underline{B}^{-1} \underline{A}^{-1}$$

$$y \in \mathbb{R}^m, x \in \mathbb{R}^n \quad \frac{\partial y}{\partial x} = \begin{bmatrix} \frac{\partial y_1}{\partial x_1} & \frac{\partial y_1}{\partial x_2} & \dots \\ \vdots & \ddots & \frac{\partial y_m}{\partial x_n} \end{bmatrix} \in \mathbb{R}^{m \times n}$$

$$\underline{y} = \underline{A}\underline{x} \quad \frac{\partial y}{\partial x} = \underline{A} \quad \frac{\partial y}{\partial z} = \underline{A} \frac{\partial x}{\partial z}$$

Key Idea:

General tensor \underline{X} , elements X_i with i : tuple of indexes

$$\underline{Y} = g(\underline{X}), \quad z = f(\underline{Y}) \Rightarrow \nabla_{\underline{X}} z = \sum_j (\nabla_{\underline{X}} Y_j) \frac{\partial z}{\partial Y_j}$$

↳ again, type of indexes

Backprop:

$$\underline{G} = \begin{bmatrix} \leftarrow n^0 \\ \vdots \\ \rightarrow m \end{bmatrix} \begin{matrix} \frac{\partial L}{\partial o_1^e} \\ \vdots \\ \frac{\partial L}{\partial o_m^e} \end{matrix}$$

n^0 : number of outputs (units),
 m : batch size

Every row j is gradient with respect to the j -th sample.

First step: undo activation.

$$\underline{G} \leftarrow \underline{G} \odot \begin{bmatrix} f'(\underline{z}_1^e) \\ \vdots \\ f'(\underline{z}_m^e) \end{bmatrix}$$

$(m \times n^0)$

where $\underline{o}_j^e = f(\underline{z}_j^e)$

↓
activation function
element-wise

Next step: $\nabla_{\underline{W}} J$ and $\nabla_{\underline{b}} J$.

→ sample: $\underline{z}_j^e = \underline{W}^T \underline{o}_j^{e-1} + \underline{b}^e$ (vectors are rows)

$\underline{z}_j^e = \underline{W} \underline{o}_j^{e-1} + \underline{b}^e$ (vectors are columns)

$(\underline{o}^{e-1}) \rightarrow$ whole batch

→ Whole batch:

$$\underline{Z}^e = \begin{bmatrix} \underline{z}_1^e \\ \vdots \\ \underline{z}_m^e \end{bmatrix} = \begin{bmatrix} \underline{o}_1^{e-1} \\ \vdots \\ \underline{o}_m^{e-1} \end{bmatrix} \begin{bmatrix} \underline{W}^e \end{bmatrix} + \begin{bmatrix} \underline{b}^e \end{bmatrix}$$

$(m \times n^0)$ $(m \times n^i)$ $(n^i \times n^0)$ $(1 \times n^0)$

↖ broadcasting

n^i : number of inputs, W_{ij}^e = weight for input i of unit j at layer e

So we have, in

$$\underline{G} = \nabla_{\underline{Z}^l} \mathcal{L} \quad (m \times n^o) \quad G_{ij} = \frac{\partial \mathcal{L}}{\partial Z_{ij}^l} \rightarrow \text{component } j \text{ of sample } i \text{ of layer } l$$

At this point we apply the key idea:

$$\nabla_{\underline{W}^l} \mathcal{L} = \sum_{ij} \frac{\partial Z_{ij}^l}{\partial W_{hk}^l} \frac{\partial \mathcal{L}}{\partial Z_{ij}^l} \Rightarrow (\nabla_{\underline{W}^l} \mathcal{L})_{hk} = \sum_{ij} \frac{\partial Z_{ij}^l}{\partial W_{hk}^l} G_{ij}$$

Considering that $Z_{ij}^l = o_{i1}^{l-1} W_{1j}^l + o_{i2}^{l-1} W_{2j}^l + \dots + o_{in}^{l-1} W_{nj}^l + b_i^l$,

$$\frac{\partial Z_{ij}^l}{\partial W_{hk}^l} = \begin{cases} 0 & \text{if } j \neq k \\ o_{ih}^{l-1} & \text{if } j = k \end{cases}$$

So for every W_{hk} , there are 'm' non-zero values,

$$\frac{\partial Z_{1k}}{\partial W_{hk}^l}, \frac{\partial Z_{2k}}{\partial W_{hk}^l}, \dots, \frac{\partial Z_{mk}}{\partial W_{hk}^l} \quad [i=1, \dots, m; j=k]$$

Which is correct: one for every sample!

$$(\nabla_{\underline{W}^l} \mathcal{L})_{hk} = o_{1h}^{l-1} G_{1k} + o_{2h}^{l-1} G_{2k} + \dots + o_{mh}^{l-1} G_{mk}$$

\hookrightarrow Column h of \underline{O}^{l-1} and k of \underline{G}

So we obtain, equivalently:

$$\nabla_{\underline{W}^l} \mathcal{L} = \underline{O}^{l-1 T} \underline{G}$$

$(n^i \times n^o) \quad (n^i \times m) \quad (m \times n^o)$

Adding regularization,

$\underline{\theta}$ = all parameters

$$\Rightarrow A) \quad \nabla_{\underline{W}^l} \mathcal{J} = \underline{O}^{l-1 T} \underline{G} + \lambda \nabla_{\underline{W}^l} \Omega(\underline{\theta})$$

For bias, write down the correct way
the equation for \underline{z}^l is:

$$\underline{z}^l = \underline{o}^{l-1} \underline{w}^l + \begin{bmatrix} b_1^l \\ b_2^l \\ \vdots \\ b_{n^o}^l \end{bmatrix} \uparrow m \text{ times, without broadcasting}$$

$(m \times n^o)$

So $\underline{z}^l = \underline{o}^{l-1} \underline{w}^l + \begin{bmatrix} 1 \end{bmatrix} [b_1^l \ b_2^l \ \dots \ b_{n^o}^l]$

$(m \times 1) \quad (1 \times n^o)$

Apply key idea

$$\nabla_{\underline{b}^l} \mathcal{L} = \sum_{ij} \frac{\partial \mathcal{L}}{\partial z_{ij}^l} \frac{\partial z_{ij}^l}{\partial b_k^l} = \sum_{ij} \frac{\partial \mathcal{L}}{\partial z_{ij}^l} G_{ij}$$

$$(\nabla_{\underline{b}^l} \mathcal{L})_k = \sum_{ij} \frac{\partial \mathcal{L}}{\partial z_{ij}^l} G_{ij}$$

Considering again that $z_{ij}^l = \dots + b_j^l$,

$$\frac{\partial z_{ij}^l}{\partial b_k^l} = \begin{cases} \emptyset & \text{if } j \neq k \\ 1 & \text{if } j = k \end{cases}$$

So, again, m non-zero values, all ones.

$$(\nabla_{\underline{b}^l} \mathcal{L})_k = \sum_i G_{ik} \Rightarrow \text{sum of } \underline{G} \text{ } k\text{-th column}$$

$$\Rightarrow B) \nabla_{\underline{b}^l} \mathcal{L} = \begin{bmatrix} 1 & \dots & 1 \end{bmatrix} \underline{G}$$

$(1 \times n^o) \quad (1 \times m) \quad (m \times n^o)$

At this point we need to calculate \underline{G} needed for the next step of the backprop.

We need $\nabla_{\underline{o}^{l-1}} \mathcal{L}$.

Applying key idea: $\nabla_{\underline{o}^{l-1}} \mathcal{L} = \sum_{ij} \frac{\partial \mathcal{L}}{\partial z_{ij}^l} \frac{\partial z_{ij}^l}{\partial o_{i-1}^{l-1}}$

$$(\nabla_{\underline{o}^{l-1}} \mathcal{L})_{n^{l-1}} = \sum_{ij} \frac{\partial \mathcal{L}}{\partial z_{ij}^l} \frac{\partial z_{ij}^l}{\partial o_{n^{l-1}}^{l-1}}$$

Considering again that $L_{ij} = 0_{i1} w_{1j} + \dots + 0_{in^0} w_{n^0 j} + 0_i$

$$\frac{\partial Z_{ij}^e}{\partial 0_{hk}^{e-1}} = \begin{cases} 0 & \text{if } i \neq h \\ w_{kj}^e & \text{if } i = h \end{cases}$$

So for every 0_{hk}^{e-1} there are n^0 non-zero values

$$w_{k1}^e, w_{k2}^e, \dots, w_{kn^0}^e \quad [j=1, \dots, n^0; i=h]$$

$$(\nabla_{\underline{0}^{e-1}} L)_{hk} = w_{k1}^e G_{h1} + \dots + w_{kn^0}^e G_{hn^0}$$

↳ Row k of \underline{w}^e and h of \underline{G}

$$\Rightarrow \hookrightarrow \nabla_{\underline{0}^{e-1}} L = \begin{matrix} \underline{G} & \underline{w}^{eT} \\ (m \times n^0) & (n^0 \times n^i) \end{matrix}$$