

Cross-Entropy Loss

$$\mathcal{L}(\hat{y}, \tilde{y}) = - \sum_i y_i \log \hat{y}_i$$

$$\frac{\partial \mathcal{L}}{\partial \hat{y}} = \begin{bmatrix} -\tilde{y} \\ \hat{y} \end{bmatrix} \leftarrow \text{elementwise division.}$$

Simple enough. Now we need the gradient of softmax:

Softmax

$$\hat{y} = \text{softmax}(\vec{a})$$

$$\frac{\partial \hat{y}}{\partial \vec{a}} \text{ is a matrix} = \begin{bmatrix} \frac{\partial y_1}{\partial a_1} & \dots & \frac{\partial y_1}{\partial a_i} & \dots & \frac{\partial y_1}{\partial a_n} \\ \vdots & & \vdots & & \vdots \\ \frac{\partial y_i}{\partial a_1} & \dots & \frac{\partial y_i}{\partial a_i} & \dots & \frac{\partial y_i}{\partial a_n} \\ \vdots & & \vdots & & \vdots \\ \frac{\partial y_n}{\partial a_1} & \dots & \frac{\partial y_n}{\partial a_i} & \dots & \frac{\partial y_n}{\partial a_n} \end{bmatrix}$$

There are really only two things to find here. $\frac{\partial y_i}{\partial a_j}$ when $i \neq j$ and $i = j$.

$$\text{Let } \mathcal{Z} = \sum_i e^{a_i} \text{ so } \text{softmax}(a_i) = \frac{e^{a_i}}{\mathcal{Z}}$$

$$\text{Then: } \frac{\partial y_i}{\partial a_j} = \frac{\partial}{\partial a_j} \frac{e^{a_i}}{\mathcal{Z}} = \frac{0 \cdot \mathcal{Z} - e^{a_i} e^{a_j}}{\mathcal{Z}^2} = \boxed{-\hat{y}_i \hat{y}_j}$$

(recalling $\frac{d}{dx} \frac{f(x)}{g(x)} = \frac{f'(x)g(x) - f(x)g'(x)}{g(x)^2}$)

$$\frac{\partial y_i}{\partial a_i} = \frac{\partial}{\partial a_i} \frac{e^{a_i}}{\mathcal{Z}} = \frac{e^{a_i} \mathcal{Z} - e^{a_i} e^{a_i}}{\mathcal{Z}^2} = \frac{e^{a_i}}{\mathcal{Z}} \left(1 - \frac{e^{a_i}}{\mathcal{Z}}\right) = \boxed{\hat{y}_i (1 - \hat{y}_i)}$$

To backprop through softmax, we need

$$\frac{\partial \mathcal{L}}{\partial \vec{a}} = \frac{\partial \mathcal{L}}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial \vec{a}}$$

Let $\hat{y}, \tilde{y}, \vec{a} \in \mathbb{R}^{1 \times N}$. Then $\frac{\partial \mathcal{L}}{\partial \hat{y}} \in \mathbb{R}^{1 \times N}$ and $\frac{\partial \hat{y}}{\partial \vec{a}} \in \mathbb{R}^{N \times N}$.

$$\frac{\partial \mathcal{L}}{\partial \vec{a}} = \begin{bmatrix} \frac{-y_1}{\hat{y}_1} & \dots & \frac{-y_N}{\hat{y}_N} \end{bmatrix} \begin{bmatrix} \frac{\partial \hat{y}_1}{\partial a_1} & \dots & \frac{\partial \hat{y}_1}{\partial a_n} \\ \vdots & & \vdots \\ \frac{\partial \hat{y}_n}{\partial a_1} & \dots & \frac{\partial \hat{y}_n}{\partial a_n} \end{bmatrix} = \begin{bmatrix} (y_1 - \hat{y}_1) & \dots & (y_n - \hat{y}_n) \end{bmatrix}$$

$$= \tilde{y} - \hat{y}$$

This dot product ~~is~~ the sum

$$\frac{-y_i}{\hat{y}_i} \cdot \frac{\partial \hat{y}_i}{\partial a_i} = \sum_{j \neq i} \frac{y_i}{\hat{y}_j} \frac{\partial \hat{y}_j}{\partial a_i}$$

$$= -y_i(1 - \hat{y}_i) - \sum_j (-y_j \hat{y}_i) = y_i - \hat{y}_i$$

We also want to backprop through ReLU and Linear:

ReLU

$$\vec{a} = \text{ReLU}(\vec{h}) , \quad \text{ReLU}(x_i) = \begin{cases} x_i & \text{if } x_i > 0 \\ 0 & \text{otherwise} \end{cases}$$

$$\text{Then } \frac{\partial \vec{a}}{\partial \vec{h}} = \begin{cases} 1 & \text{if } h_i > 0 \\ 0 & \text{otherwise} \end{cases}$$

Linear

$$\vec{h} = \vec{x}\vec{w} + \vec{b} , \quad \vec{x} \in \mathbb{R}^{n \times M} , \vec{w} \in \mathbb{R}^{M \times L} , \vec{b} \in \mathbb{R}^{1 \times L}$$

$$\frac{\partial \vec{h}}{\partial \vec{x}} = \vec{w}^T , \quad \frac{\partial \vec{h}}{\partial \vec{b}} = \vec{1} , \quad \frac{\partial \vec{h}}{\partial \vec{w}} = \vec{x}$$

Linking these together gives us a way to backprop through our network.