Linéaire

MSELoss

Forward

$$MSE(y, \hat{y}) = ||\hat{y} - y||^2$$

Backward

$$rac{\partial MSE}{\partial \hat{y}} = rac{\partial \sum_{i}^{N} ||\hat{y}_{i} - y_{i}||^{2}}{\partial \hat{y}}$$

On utilise le denominator layout

$$2egin{pmatrix} \hat{y}_{11}-y_{11} & \hat{y}_{12}-y_{12} & \dots & \hat{y}_{1d}-y_{1d} \ \hat{y}_{21}-y_{21} & \hat{y}_{22}-y_{22} & \dots & \hat{y}_{2d}-y_{2d} \ \dots & \dots & \dots & \dots \ \hat{y}_{b1}-y_{b1} & \hat{y}_{b2}-y_{b2} & \dots & \hat{y}_{bd}-y_{bd} \end{pmatrix} = 2egin{pmatrix} \hat{y}_1-y_1 \ \hat{y}_2-y_2 \ \dots \ \hat{y}_b-y_b \end{pmatrix}$$

Donc:

$$rac{\partial MSE}{\partial \hat{y}} = -2(y-\hat{y})$$

Module Linéaire

Backward update gradient

$$rac{\partial Loss}{\partial W^h} = rac{\partial Loss}{\partial z^h} rac{\partial z^h}{\partial W^h} = \delta^h rac{\partial z^h}{\partial W^h}$$

On utilise le denominator layout

$$z^h = \left(z_1^h \quad z_2^h \quad \dots \quad z_{d'}^h
ight) = \left(\sum_i^d z_i^{h-1} w_{i1}^h \quad \sum_i^d z_i^{h-1} w_{i2}^h \quad \dots \quad \sum_i^d z_i^{h-1} w_{id'}^h
ight) \ W^h = egin{pmatrix} w_1^h \ w_2^h \ \dots \ w_{d'}^h \end{pmatrix} = egin{pmatrix} w_{11}^h \quad w_{21}^h \quad \dots \quad w_{d1}^h \ w_{12}^h \quad w_{22}^h \quad \dots \quad w_{d2}^h \ \dots \quad \dots \quad \dots \ w_{1d'}^h \quad w_{2d'}^h \quad \dots \quad w_{dd'}^h \end{pmatrix}$$

$$rac{\partial z^h}{\partial W^h} = egin{pmatrix} z_1^{h-1} & 0 & \dots & 0 & z_2^{h-1} & 0 & \dots & 0 & z_d^{h-1} & 0 & 0 & 0 \ 0 & z_1^{h-1} & \dots & 0 & 0 & z_2^{h-1} & \dots & 0 & 0 & z_d^{h-1} & 0 & 0 \ \dots & \dots \ 0 & 0 & \dots & z_1^{h-1} & 0 & 0 & \dots & z_2^{h-1} & 0 & 0 & 0 & z_d^{h-1} \end{pmatrix}$$

$$rac{\partial z^h}{\partial W^h} = egin{pmatrix} z^{h-1} & 0 & \dots & 0 \ 0 & z^{h-1} & \dots & 0 \ \dots & \dots & \dots & \dots \ 0 & 0 & \dots & z^{h-1} \end{pmatrix}$$

Dimension (d', dd')

$$\delta^h = egin{pmatrix} \delta_1^h \ \delta_2^h \ \dots \ \delta_{J}^h \end{pmatrix}$$

$$rac{\partial Loss}{\partial W^h} = \delta^h rac{\partial z^h}{\partial W^h} = ig(\delta^h_1 \quad \delta^h_2 \quad \dots \quad \delta^h_{d'}ig) egin{pmatrix} z^{h-1} & 0 & \dots & 0 \ 0 & z^{h-1} & \dots & 0 \ \dots & \dots & \dots & \dots \ 0 & 0 & \dots & z^{h-1} \end{pmatrix} = ig(\delta^h)^T z^{h-1}$$

Dimension $(N, d') \times (d', dd') = (N, dd')$

Backward delta

$$rac{\partial Loss}{\partial z^{h-1}} = rac{\partial Loss}{\partial z^h} rac{\partial z^h}{\partial z^{h-1}} = \delta^h rac{\partial z^h}{\partial z^{h-1}}$$

$$z^{h-1} = egin{pmatrix} z_1^{h-1} \ z_2^{h-1} \ \dots \ z_d^{h-1} \end{pmatrix}$$

$$z^h = ig(z_1^h ig| z_2^h ig) \ldots z_{d'}^hig) = ig(\sum_i^d z_i^{h-1} w_{i1}^h ig| \sum_i^d z_i^{h-1} w_{i2}^h ig| \ldots ig| \sum_i^d z_i^{h-1} w_{id'}^hig)$$

$$rac{\partial z^h}{\partial z^{h-1}} = egin{pmatrix} w_{11}^h & w_{12}^h & \dots & w_{1d'}^h \ w_{21}^h & w_{22}^h & \dots & w_{2d'}^h \ \dots & \dots & \dots & \dots \ w_{d1}^h & w_{d2}^h & \dots & w_{dd'}^h \end{pmatrix} = (W^h)^T$$

$$rac{\partial Loss}{\partial z^{h-1}} = \delta^h rac{\partial z^h}{\partial z^{h-1}} = \delta^h (W^h)^T$$

Comme $z=< x.\,w>$, donc on utilise la transition de W. Or:

$$rac{\partial Loss}{\partial z^{h-1}}=\delta^hrac{\partial z^h}{\partial z^{h-1}}=\delta^hW^h$$

Non-linéaire

Tangente hyperbolique

Forward

$$tanH(x)=rac{e^x-e^{-x}}{e^x+e^{-x}}$$

Backward

$$rac{\partial Loss}{\partial a^h} = rac{\partial Loss}{\partial z^h} rac{\partial z^h}{\partial a^h} = \delta^h rac{\partial z^h}{\partial a^h}$$

$$\frac{\partial z^h}{\partial a^h} = \frac{\frac{\partial \frac{e^{a^h} - e^{-a^h}}{e^{a^h} + e^{-a^h}}}{\partial a^h}}{\partial a^h} = \frac{(e^{a^h} - e^{-a^h})'(e^{a^h} + e^{-a^h}) - (e^{a^h} + e^{-a^h})'(e^{a^h} - e^{-a^h})}{(e^{a^h} + e^{-a^h})^2} = \frac{(e^{a^h} + e^{-a^h})^2 - (e^{a^h} - e^{-a^h})^2}{(e^{a^h} + e^{-a^h})^2} = 1 - \frac{(e^{a^h} - e^{-a^h})^2}{(e^{a^h} + e^{-a^h})^2}$$

Or:

$$rac{\partial z^h}{\partial a^h} = 1 - tan H^2(a^h)$$

Donc:

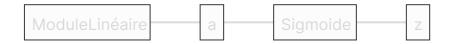
$$rac{\partial Loss}{\partial a^h} = \delta^h rac{\partial z^h}{\partial a^h} = \delta^h (1 - tan H^2(a^h))$$

Sigmoide

Forward

$$\sigma(x) = rac{1}{1+e^{-x}}$$

Backward



$$rac{\partial Loss}{\partial a^h} = rac{\partial Loss}{\partial z^h} rac{\partial z^h}{\partial a^h} = \delta^h rac{\partial z^h}{\partial a^h}$$

$$rac{\partial z^h}{\partial a^h} = rac{\partial rac{1}{1+e^{-a^h}}}{\partial a^h} = rac{-(1+e^{-a^h})'}{(1+e^{-a^h})^2} = rac{e^{-a^h}}{(1+e^{-a^h})^2} = rac{1}{1+e^{-a^h}} rac{e^{-a^h}}{1+e^{-a^h}} = rac{1}{1+e^{-a^h}} rac{1+e^{-a^h}-1}{1+e^{-a^h}} = rac{1}{1+e^{-a^h}} (1-rac{1}{1+e^{-a^h}}) = \sigma(a^h) (1-a^h)$$

Softmax et Cout Entropique

Soft-max

$$Softmax(x) = rac{e^x}{\sum_i^d e^i}$$

Coût Cross-Entropique

Soit y le vecteur supervision codé en one-hot.

Par exemple, si y_i est de 3ème classe.

$$y_i = egin{pmatrix} 0 \ 0 \ 1 \ 0 \end{pmatrix}$$

Et \hat{y}_i est la vecteur de prédiction

$$\hat{y}_i = egin{pmatrix} 0.2 \ 0.4 \ 0.3 \ 0.1 \end{pmatrix}$$

$$CE(y_i, \hat{y}_i) = - < y_i.\,\hat{y}_i> = -(0*0.2+0*0.4+1*0.3+0*0.1) = -0.3$$

Nous allons noter y comme l'indice de la classe à prédire.

$$CE(y,\hat{y}) = -\hat{y}_y$$

Combinaison de Softmax et coût cross-entropique

Afin d'éviter des instabilités numériques, on enchaîne un Softmax passé au logarithme (logSoftMax) et un coût cross entropique comme une combinaison.

Forward

$$CE(y,\hat{y}) = -\lograc{e^{\hat{y}_y}}{\sum_{i=1}^K e^{\hat{y}_i}} = -\hat{y}_y + \log\sum_{i=1}^K e^{\hat{y}_i}$$

Backward

$$rac{\partial Loss}{\partial \hat{y}} = rac{\partial \sum_{i=1}^{N} -\hat{y}_{i,y_i} + log \sum_{j=1}^{K} e^{\hat{y}_{i,j}}}{\partial \hat{y}}$$

On utilise le denominator layout

$$\hat{y} = egin{pmatrix} \hat{y}_1 \ \hat{y}_2 \ \dots \ \hat{y}_N \end{pmatrix} = egin{pmatrix} \hat{y}_{11} & \hat{y}_{12} & \dots & \hat{y}_{1K} \ \hat{y}_{21} & \hat{y}_{22} & \dots & \hat{y}_{2K} \ \dots & & & & \ \hat{y}_{N1} & \hat{y}_{N2} & \dots & \hat{y}_{NK} \end{pmatrix}$$

$$rac{\partial Loss}{\partial \hat{y}_{kf}} = rac{\partial \sum_{i=1}^{N} -\hat{y}_{i,y_i} + \log \sum_{j=1}^{K} e^{\hat{y}_{i,j}}}{\partial \hat{y}_{kf}} = rac{\partial (-\hat{y}_{k,y_k} + \log \sum_{j=1}^{K} e^{\hat{y}_{k,j}})}{\partial \hat{y}_{kf}}$$

Si
$$y_k=f$$
 : $rac{\partial Loss}{\partial \hat{y}_{kf}}=-1+rac{e^{\hat{y}_{kf}}}{\sum_{i=1}^K e^{\hat{y}_{k,j}}}$

Si
$$y_k
eq f$$
 : $rac{\partial Loss}{\partial \hat{y}_{kf}} = rac{e^{\hat{y}_{kf}}}{\sum_{i=1}^{K} e^{\hat{y}_{k,j}}}$

Donc:

$$rac{\partial Loss}{\partial \hat{y}} = egin{pmatrix} -y_{11} + rac{e^{\hat{y}_{11}}}{\sum_{j=1}^{K} e^{\hat{y}_{1,j}}} & -y_{12} + rac{e^{\hat{y}_{12}}}{\sum_{j=1}^{K} e^{\hat{y}_{1,j}}} & \dots & -y_{1K} + rac{e^{\hat{y}_{1K}}}{\sum_{j=1}^{K} e^{\hat{y}_{1,j}}} \ -y_{21} + rac{e^{\hat{y}_{21}}}{\sum_{j=1}^{K} e^{\hat{y}_{2,j}}} & -y_{22} + rac{e^{\hat{y}_{22}}}{\sum_{j=1}^{K} e^{\hat{y}_{2,j}}} & \dots & -y_{2K} + rac{e^{\hat{y}_{2K}}}{\sum_{j=1}^{K} e^{\hat{y}_{2,j}}} \ \dots & \dots & \dots & \dots \ -y_{N1} + rac{e^{\hat{y}_{N1}}}{\sum_{j=1}^{K} e^{\hat{y}_{N,j}}} & -y_{N2} + rac{e^{\hat{y}_{N2}}}{\sum_{j=1}^{K} e^{\hat{y}_{N,j}}} & \dots & -y_{NK} + rac{e^{\hat{y}_{NK}}}{\sum_{j=1}^{K} e^{\hat{y}_{N,j}}} \end{pmatrix}$$

$$rac{\partial Loss}{\partial \hat{y}} = -y + Softmax(\hat{y})$$

Se compresser

Encodage



Décodage



Coût Cross-entropique binaire

Forward

$$BCE(y,\hat{y}) = -(y\log(\hat{y}) + (1-y)\log(1-\hat{y}))$$

Backward

$$\frac{\partial Loss}{\partial \hat{y}} = -(\frac{y}{\hat{y}} + \frac{y-1}{1-\hat{y}})$$