



$$(b) d(\hat{y}_i) = \frac{-y_i}{\hat{y}_i} = (-y \log \hat{y}_i)'$$

$$(c) \frac{\partial \hat{y}_i}{\partial \tilde{y}_j} = \frac{e^{\tilde{y}_i} (\sum_n e^{\tilde{y}_n} - e^{\tilde{y}_j})}{(\sum_n e^{\tilde{y}_n})^2} \quad \text{when } i=j$$

$$= p_i (1 - p_i)$$

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

assume $p_i = \frac{e^{\tilde{y}_i}}{\sum e^{\tilde{y}_n}} = \hat{y}_i$

When $i \neq j$

$$= \frac{-e^{\tilde{y}_i} e^{\tilde{y}_j}}{\sum e^{\tilde{y}_k}} = -p_i p_j = -\hat{y}_i \hat{y}_j$$

$$\frac{\partial \hat{y}}{\partial \tilde{y}} = p_i (\delta_{ij} - p_j) \quad \text{where } \delta_{ij} = \begin{cases} 1 & \text{if } i=j \\ 0 & \text{if } i \neq j \end{cases}$$

$$= \hat{y}_i (\delta_{ij} - \hat{y}_j)$$

chain rule

$$d\tilde{y}_i = \frac{\partial \mathcal{L}}{\partial \tilde{y}} = -\sum_k y_k \frac{\partial \log(\hat{y}_k)}{\partial \tilde{y}} = -\sum_k y_k \frac{1}{\hat{y}_k} + \frac{\partial \hat{y}_k}{\partial \tilde{y}_i}$$

$$= -y_i(1-p_i) + \sum_{k \neq i} p_i y_k = -y_i + p_i(\underbrace{y_i + \sum_{k \neq i} y_k}_{=1})$$

$$d(\tilde{y}) = p_i - y_i = \hat{y}_i - y_i$$

On the right hand side: $\frac{-y_i}{\hat{y}_i} + \hat{y}_i - \sum_j \frac{-y_j}{\hat{y}_j} \hat{y}_j \hat{y}_i$

$$= -y_i + \sum_j y_j \hat{y}_i = \boxed{\hat{y}_i - y_i} = d(\tilde{y})$$

$$(d) \quad d(h)_i = \sum_j \frac{\partial (\tilde{y})_j}{\partial h_i} d(\tilde{y})_j = (\hat{y}_i - y_i) W^T$$

$$(e) \quad d(z)_i = \frac{\partial \mathcal{L}}{\partial h_i} + \frac{\partial h}{\partial z} = \frac{e^{-z_i}}{(1+e^{-z_i})^2} * (\hat{y}_i - y_i) W^T$$

$$= \sigma(z_i)(1-\sigma(z_i)) * (\hat{y}_i - y_i) W^T$$

$$(f) \quad d(W)_{ij} = \frac{\partial \mathcal{L}}{\partial z_i} \frac{\partial z_i}{\partial W_{ij}} = X^T (\hat{y}_i - y_i) W^T + h_i(1-h_i)$$