

Identities	Model
$(A \otimes C)(B \otimes D) = (AB) \otimes (CD)$	<b>1: procedure</b> forward( $\mathbf{x}; \mathbf{W}_{i,\dots,K}, \boldsymbol{\beta}_{i,\dots,K}$ )
$\text{diag}(\mathbf{u})\mathbf{1} = \mathbf{u}$	<b>2:   </b> $\mathbf{z}_0 \leftarrow \mathbf{x}$
$\frac{\partial}{\partial \mathbf{u}} f^T(\mathbf{u}) = \text{diag}[f'(\mathbf{u})]$	<b>3:   for</b> $i = 0, \dots, K - 2$ <b>do</b>
$\frac{\partial}{\partial v} f^T(\mathbf{u}) = \frac{\partial \mathbf{u}^T}{\partial v} \text{diag}[f'(\mathbf{u})]$	<b>4:       </b> $\mathbf{h}_{i+1} \leftarrow \mathbf{W}_i \mathbf{z}_i + \boldsymbol{\beta}_i$
$\frac{\partial}{\partial X^T} (AX) = I_k \otimes A, \quad X \in \mathbb{R}^{n \times k}$	<b>5:       </b> $\mathbf{z}_{i+1} = \sigma(\mathbf{h}_{i+1})$
$\frac{\partial}{\partial X} f(\mathbf{u}) = \frac{\partial \mathbf{u}^T}{\partial X} \left[ \frac{\partial f(\mathbf{u})}{\partial \mathbf{u}} \otimes I \right]$	<b>6:   end for</b>
$f(\mathbf{x}) = \frac{u(\mathbf{x})}{v(\mathbf{x})}$	<b>7:   </b> $\mathbf{h}_K \leftarrow \mathbf{W}_{K-1} \mathbf{z}_{K-1} + \boldsymbol{\beta}_{K-1}$
$\frac{\partial}{\partial \mathbf{x}} f(\mathbf{x}) = \frac{1}{v^2} \left( \frac{\partial \mathbf{u}}{\partial \mathbf{x}} v - \frac{\partial v}{\partial \mathbf{x}} \mathbf{u} \right)$	<b>8:   </b> $\mathbf{z}_K \leftarrow \begin{cases} \sigma(\mathbf{h}_K) & \text{scaler} \\ \frac{\exp(\mathbf{h}_K)}{\mathbf{1}^T \exp(\mathbf{h}_K)} & \text{vector} \end{cases}$

## Loss and gradients

Scaler Output	Vector Output
<ul style="list-style-type: none"> <li>Output</li> </ul>	<ul style="list-style-type: none"> <li>Output</li> </ul>
$z_K = \sigma(h_K), \quad \frac{dz_K}{dh_K} = z_K(1 - z_K)$	$\mathbf{z}_K = \frac{\exp(\mathbf{h}_K)}{\mathbf{1}^T \exp(\mathbf{h}_K)} \quad \frac{\partial \mathbf{z}_K^T}{\partial \mathbf{h}_K} = \text{diag}(\mathbf{z}_K) - \mathbf{z}_K \mathbf{z}_K^T$
<ul style="list-style-type: none"> <li>Loss</li> </ul> $J(z_K; \mathbf{y}) = -y \ln z_K - (1 - y) \ln(1 - z_K)$	<ul style="list-style-type: none"> <li>Loss</li> </ul> $J(\mathbf{z}_K; \mathbf{y}) = -\mathbf{y}^T \ln \mathbf{z}_K$
<ul style="list-style-type: none"> <li>Loss gradients</li> </ul> $\frac{d}{dz_K} J(z_K; \mathbf{y}) = \frac{z_K - y}{z_K(1 - z_K)}$ $\frac{d}{dh_K} J(z_K; \mathbf{y}) = \frac{dz_K}{dh_K} \frac{d}{dz_K} J(z_K; \mathbf{y}) = z_K - y$	<ul style="list-style-type: none"> <li>Loss gradients</li> </ul> $\frac{\partial}{\partial \mathbf{z}_K} J(\mathbf{z}_K; \mathbf{y}) = \text{diag}(\mathbf{z}_K)^{-1} \mathbf{y}$ $\frac{\partial}{\partial \mathbf{h}_K} J(\mathbf{z}_K; \mathbf{y}) = \frac{\partial \mathbf{z}_K^T}{\partial \mathbf{h}_K} \frac{\partial J}{\partial \mathbf{z}_K} = \mathbf{z}_K - \mathbf{y}$