

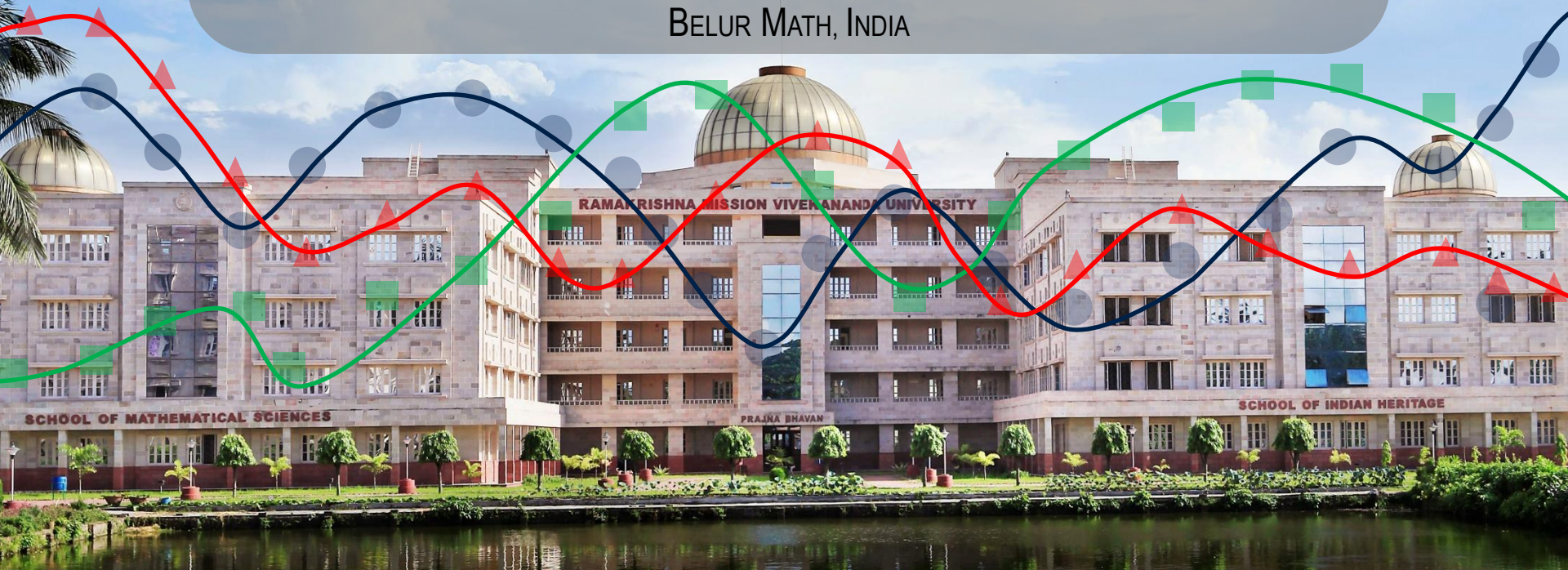
Training Deep Neural Networks: Exploding/vanishing gradients

DRIPTA MJ

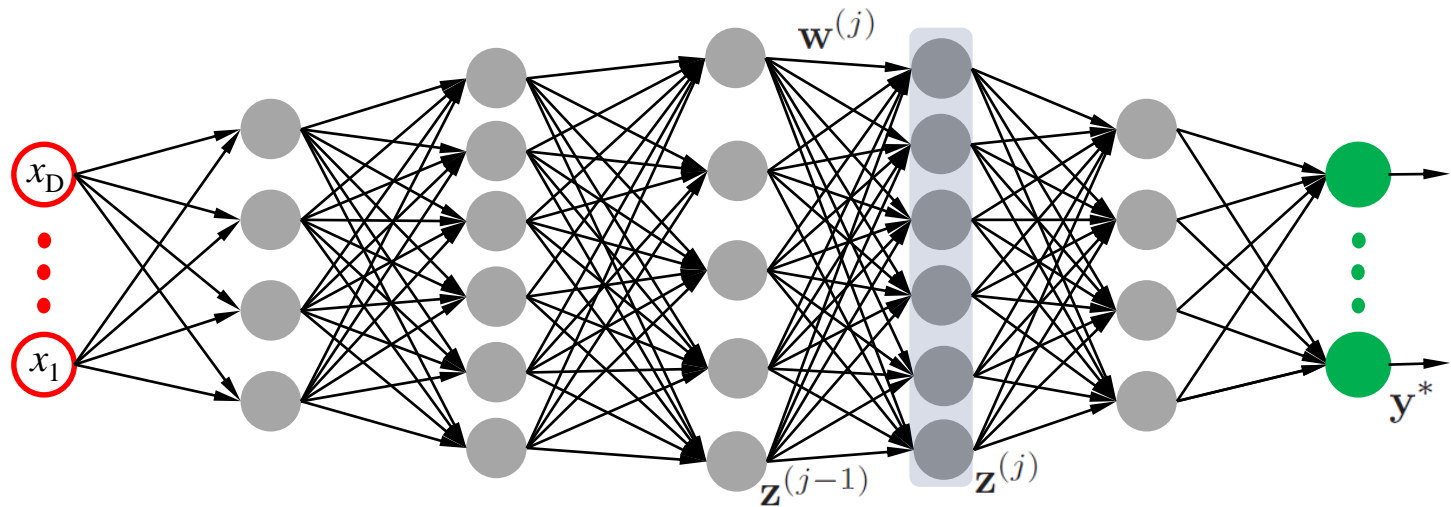
Department of Mathematics

RAMAKRISHNA MISSION VIVEKANANDA EDUCATIONAL AND RESEARCH INSTITUTE

BELUR MATH, INDIA

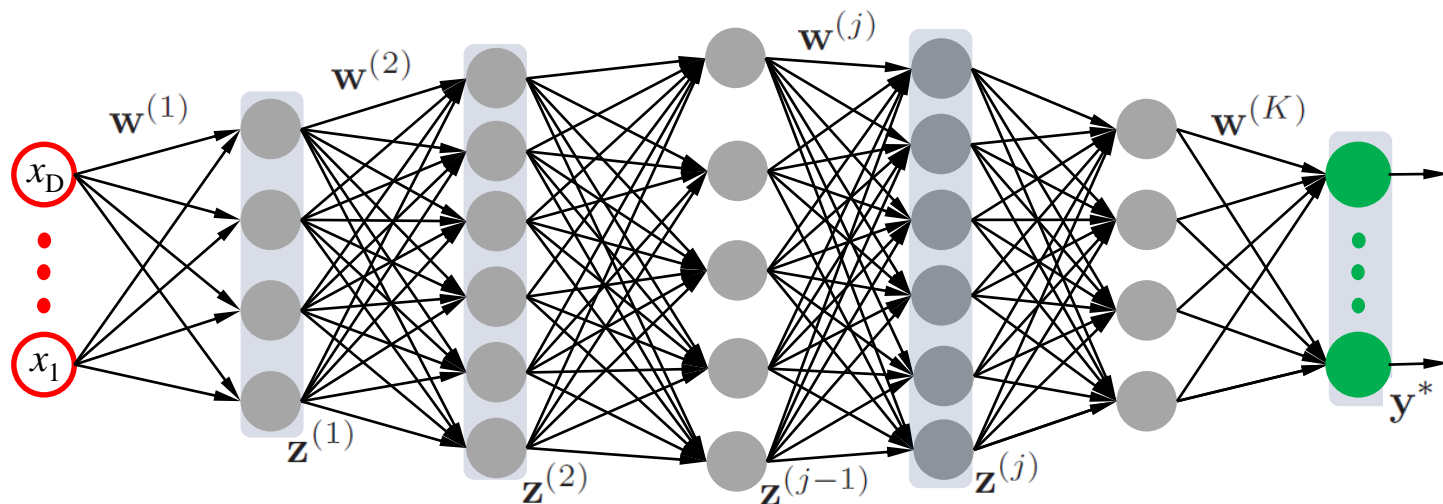


Network setup



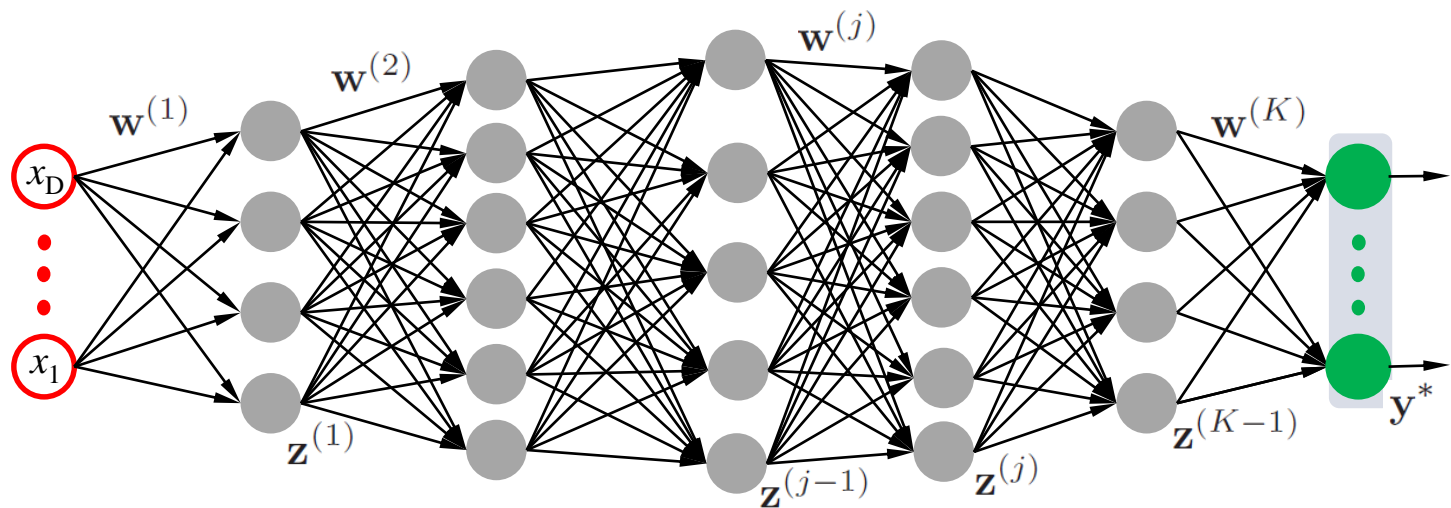
- Total number of layer: K .
- Input (vector) to the j th layer: $\mathbf{z}^{(j-1)}$.
- Output (vector) of the j th layer: $\mathbf{z}^{(j)}$.
- Weight matrix associated with the j th layer: $\mathbf{w}^{(j)}$
- Number of neurons in the j th layer: H_j .

Layer outputs



- Output from layer 1: $\mathbf{z}^{(1)} = \mathcal{A}^{(1)}(\mathbf{w}^{(1)\text{T}}\mathbf{x})$
- Output from layer 2: $\mathbf{z}^{(2)} = \mathcal{A}^{(2)}(\mathbf{w}^{(2)\text{T}}\mathbf{z}^{(1)})$
- Output from layer j : $\mathbf{z}^{(j)} = \mathcal{A}^{(j)}(\mathbf{w}^{(j)\text{T}}\mathbf{z}^{(j-1)})$
- Output from layer K (last layer): $\mathbf{y}^* = \mathcal{A}^{(K)}(\mathbf{w}^{(K)\text{T}}\mathbf{z}^{(K-1)})$

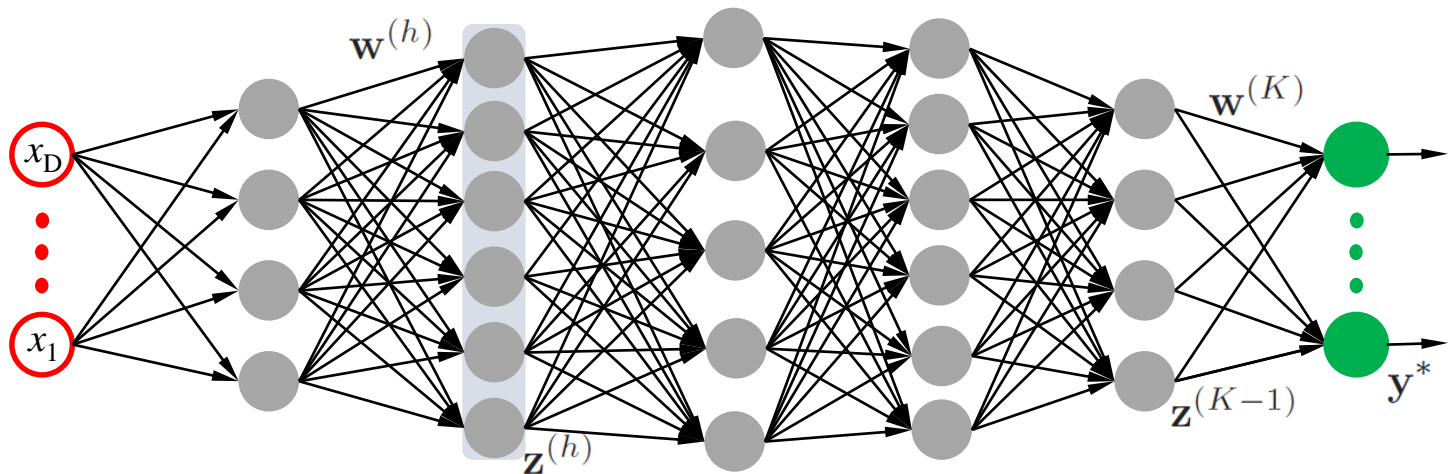
Final output



- Final output can also be written as

$$\begin{aligned}
 y^* &= \mathcal{A}^{(K)} \left(\mathbf{w}^{(K)T} \mathbf{z}^{(K-1)} \right) \\
 &= \mathcal{A}^{(K)} \left(\mathbf{w}^{(K)T} \mathcal{A}^{(K-1)} \left(\mathbf{w}^{(K-1)T} \mathbf{z}^{(K-2)} \right) \right) \\
 &= \mathcal{A}^{(K)} \left(\mathbf{w}^{(K)T} \mathcal{A}^{(K-1)} \left(\mathbf{w}^{(K-1)T} \dots \mathcal{A}^{(j)} \left(\mathbf{w}^{(j)T} \mathbf{z}^{(j-1)} \right) \dots \right) \right) \\
 &= \mathcal{A}^{(K)} \left(\mathbf{w}^{(K)T} \mathcal{A}^{(K-1)} \left(\mathbf{w}^{(K-1)T} \dots \mathcal{A}^{(1)} \left(\mathbf{w}^{(1)T} \mathbf{x} \right) \dots \right) \right)
 \end{aligned}$$

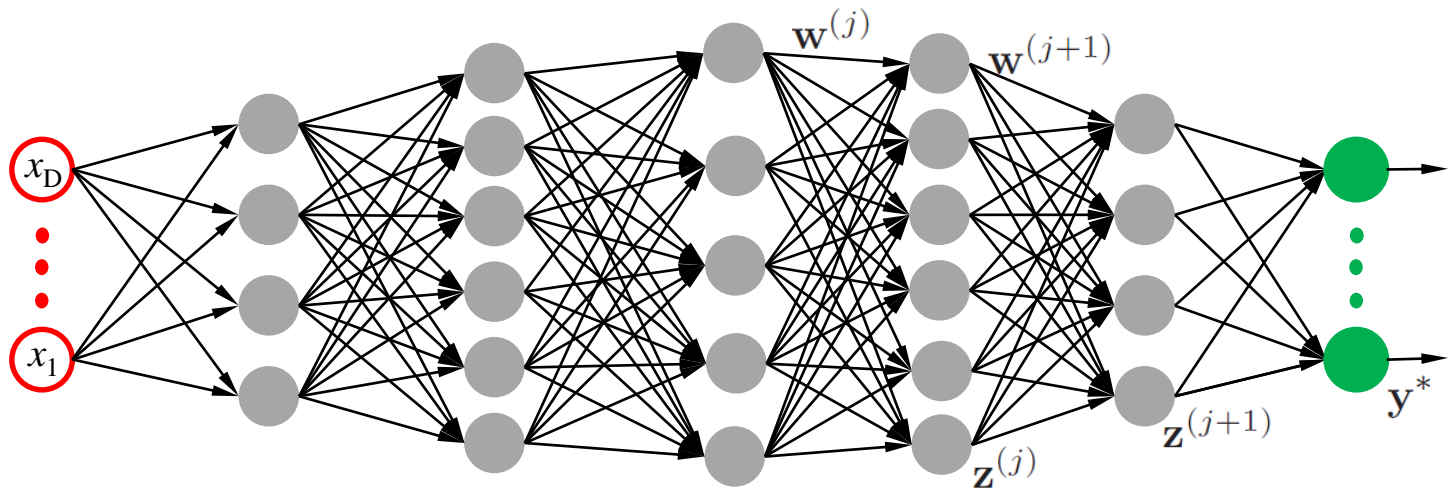
Partial derivative of the loss function



- Partial derivatives of the loss function \mathcal{L} w.r.t. the weights in the h th layer:

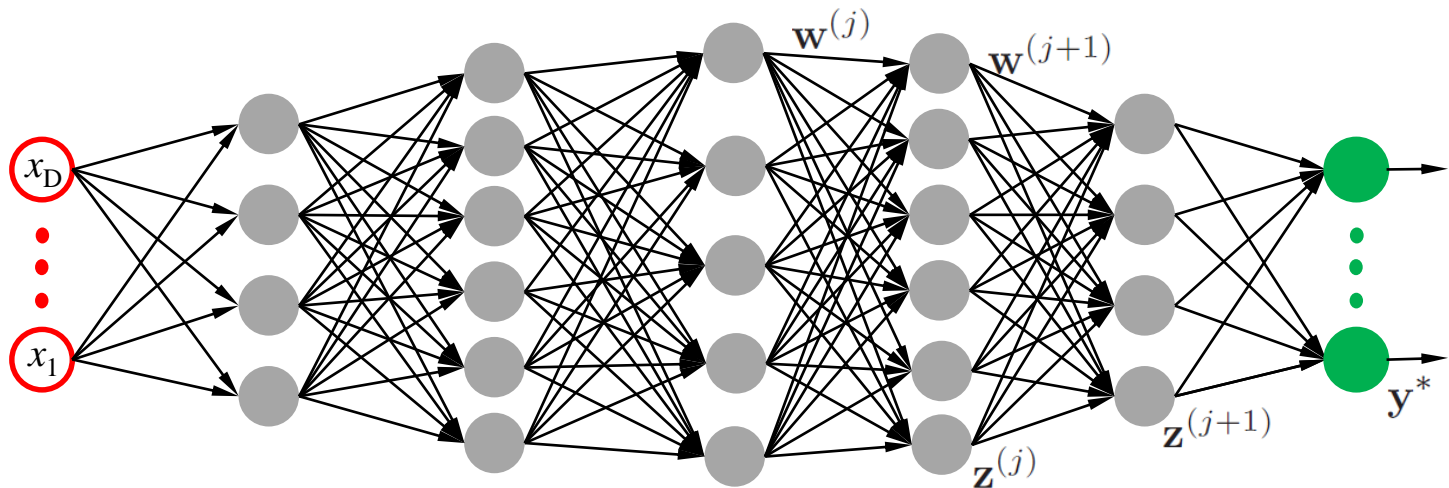
$$\begin{aligned}
 \frac{\partial \mathcal{L}}{\partial \mathbf{w}^{(h)}} &= \frac{\partial \mathcal{L}}{\partial \mathbf{y}^*} \frac{\partial \mathbf{y}^*}{\partial \mathbf{w}^{(h)}} \\
 &= \frac{\partial \mathcal{L}}{\partial \mathbf{y}^*} \frac{\partial \mathbf{y}^*}{\partial \mathbf{z}^{(K-1)}} \frac{\partial \mathbf{z}^{(K-1)}}{\partial \mathbf{w}^{(h)}} \\
 &= \frac{\partial \mathcal{L}}{\partial \mathbf{y}^*} \frac{\partial \mathbf{y}^*}{\partial \mathbf{z}^{(K-1)}} \frac{\partial \mathbf{z}^{(K-1)}}{\partial \mathbf{z}^{(K-2)}} \frac{\partial \mathbf{z}^{(K-2)}}{\partial \mathbf{w}^{(h)}} \\
 &\quad \cdot \\
 &\quad \cdot \\
 &= \frac{\partial \mathcal{L}}{\partial \mathbf{y}^*} \frac{\partial \mathbf{y}^*}{\partial \mathbf{z}^{(K-1)}} \frac{\partial \mathbf{z}^{(K-1)}}{\partial \mathbf{z}^{(K-2)}} \frac{\partial \mathbf{z}^{(K-2)}}{\partial \mathbf{z}^{(K-3)}} \cdots \frac{\partial \mathbf{z}^{(h+1)}}{\partial \mathbf{z}^{(h)}} \frac{\partial \mathbf{z}^{(h)}}{\partial \mathbf{w}^{(h)}}
 \end{aligned}$$

Jacobian matrix



$$\frac{\partial \mathbf{z}^{(j+1)}}{\partial \mathbf{z}^{(j)}} = \begin{bmatrix} \frac{\partial z_1^{(j+1)}}{\partial z_1^{(j)}} & \frac{\partial z_1^{(j+1)}}{\partial z_2^{(j)}} & \cdot & \cdot & \frac{\partial z_1^{(j+1)}}{\partial z_{H_j}^{(j)}} \\ \frac{\partial z_2^{(j+1)}}{\partial z_1^{(j)}} & \frac{\partial z_2^{(j+1)}}{\partial z_2^{(j)}} & \cdot & \cdot & \frac{\partial z_2^{(j+1)}}{\partial z_{H_j}^{(j)}} \\ \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot \\ \frac{\partial z_{H_{j+1}}^{(j+1)}}{\partial z_1^{(j)}} & \frac{\partial z_{H_{j+1}}^{(j+1)}}{\partial z_2^{(j)}} & \cdot & \cdot & \frac{\partial z_{H_{j+1}}^{(j+1)}}{\partial z_{H_j}^{(j)}} \end{bmatrix}$$

Jacobian matrix



- Consider one of the Jacobians $\frac{\partial \mathbf{z}^{(j+1)}}{\partial \mathbf{z}^{(j)}}$:

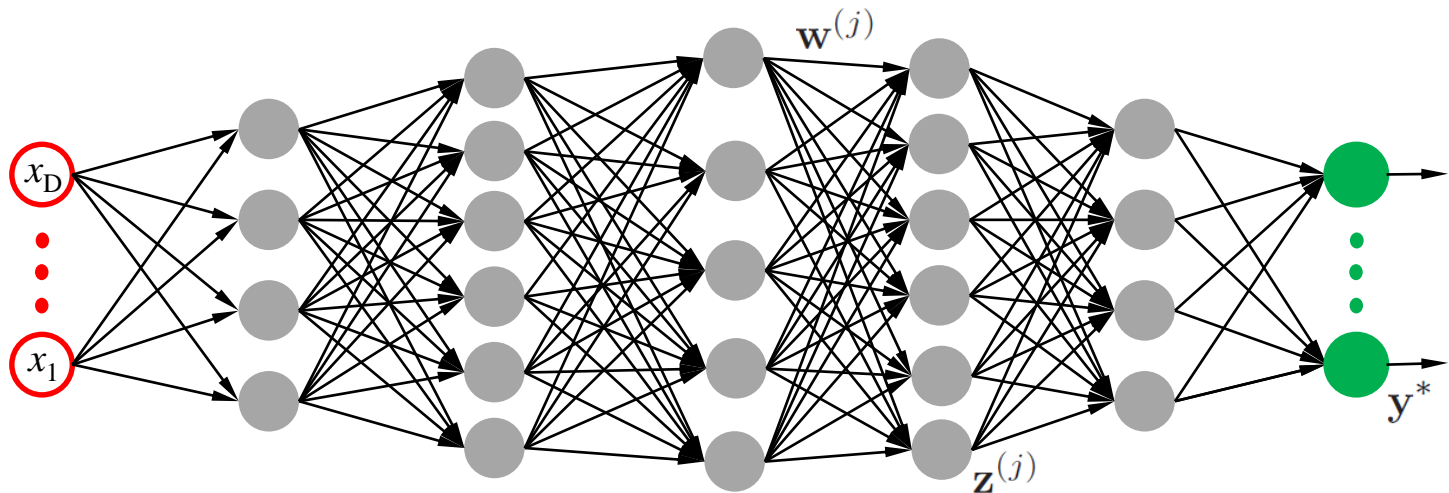
$$\begin{aligned}\frac{\partial \mathbf{z}^{(j+1)}}{\partial \mathbf{z}^{(j)}} &= \frac{\partial \left(\mathcal{A}^{(j+1)} \left(\mathbf{w}^{(j+1)\text{T}} \mathbf{z}^{(j)} \right) \right)}{\partial \mathbf{z}^{(j)}} \\ &= \text{diag} \left[\mathcal{A}^{(j+1)'} \left(\mathbf{w}^{(j+1)\text{T}} \mathbf{z}^{(j)} \right) \right] \mathbf{w}^{(j+1)\text{T}}\end{aligned}$$

Jacobian matrix: components

$$\mathbf{w}^{(j+1)} = \begin{bmatrix} w_{11}^{(j+1)} & w_{12}^{(j+1)} & \cdot & \cdot & w_{1H_{j+1}}^{(j+1)} \\ w_{21}^{(j+1)} & w_{22}^{(j+1)} & \cdot & \cdot & w_{2H_{j+1}}^{(j+1)} \\ \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot \\ w_{H_j 1}^{(j+1)} & w_{H_j 2}^{(j+1)} & \cdot & \cdot & w_{H_j H_{j+1}}^{(j+1)} \end{bmatrix}$$

$$\text{diag} \left[\mathcal{A}^{(j+1)'} (\mathbf{w}^{(j+1)\text{T}} \mathbf{z}^{(j)}) \right] = \begin{bmatrix} \mathcal{A}_1^{(j+1)'} (\mathbf{w}_{\cdot 1}^{(j+1)\text{T}} \mathbf{z}^{(j)}) & 0 & \cdot & \cdot & 0 \\ 0 & \mathcal{A}_2^{(j+1)'} (\mathbf{w}_{\cdot 2}^{(j+1)\text{T}} \mathbf{z}^{(j)}) & \cdot & \cdot & 0 \\ \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot \\ 0 & 0 & \cdot & \cdot & \mathcal{A}_{H_{j+1}}^{(j+1)'} (\mathbf{w}_{\cdot H_{j+1}}^{(j+1)\text{T}} \mathbf{z}^{(j)}) \end{bmatrix}$$

2-norm of the Jacobians

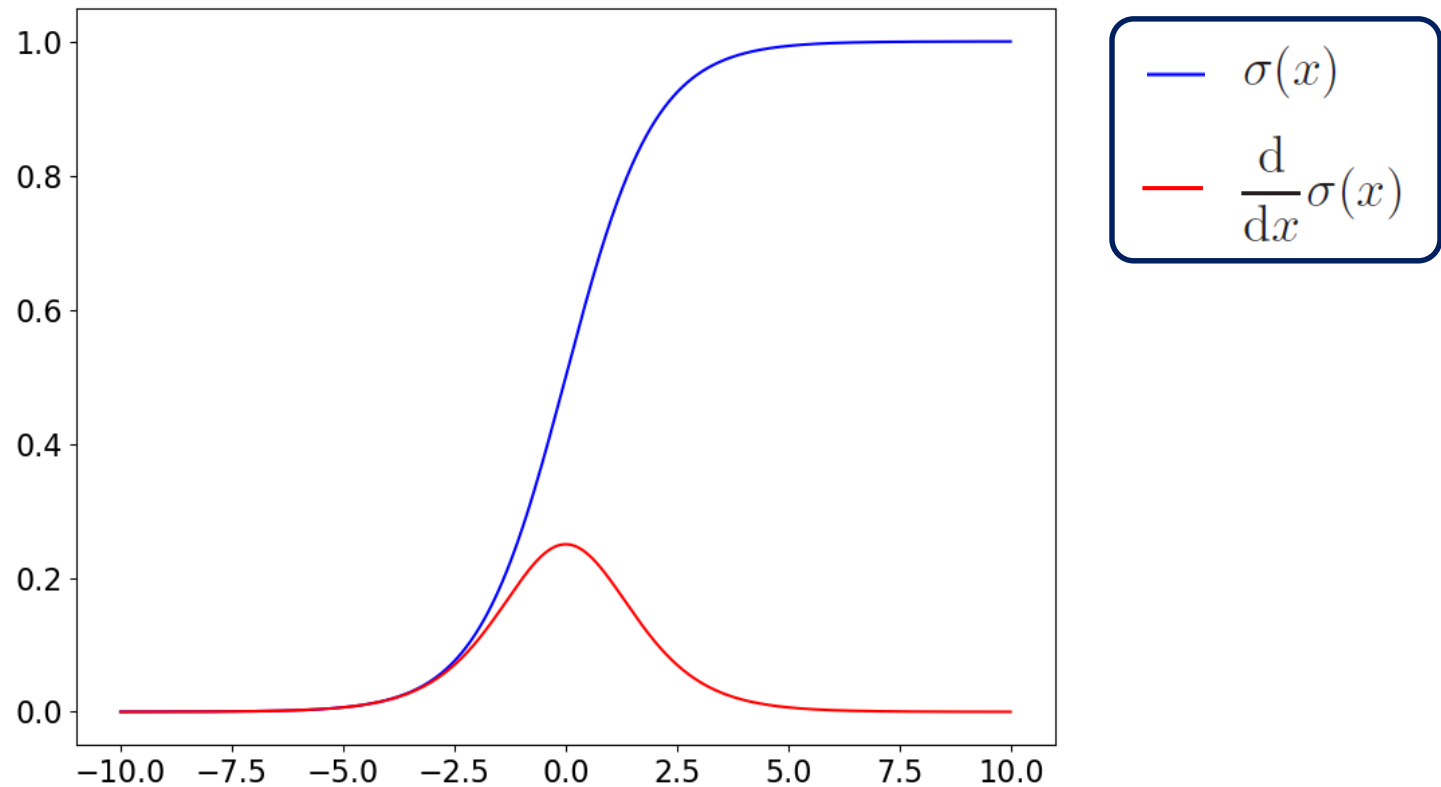


- 2-norm of the Jacobian:

$$\left\| \frac{\partial \mathbf{z}^{(j+1)}}{\partial \mathbf{z}^{(j)}} \right\| = \left\| \text{diag} \left[\mathcal{A}^{(j+1)'} (\mathbf{w}^{(j+1)\text{T}} \mathbf{z}^{(j)}) \right] \mathbf{w}^{(j+1)\text{T}} \right\|$$
$$\leq \left\| \text{diag} \left[\mathcal{A}^{(j+1)'} (\mathbf{w}^{(j+1)\text{T}} \mathbf{z}^{(j)}) \right] \right\| \left\| \mathbf{w}^{(j+1)\text{T}} \right\|$$

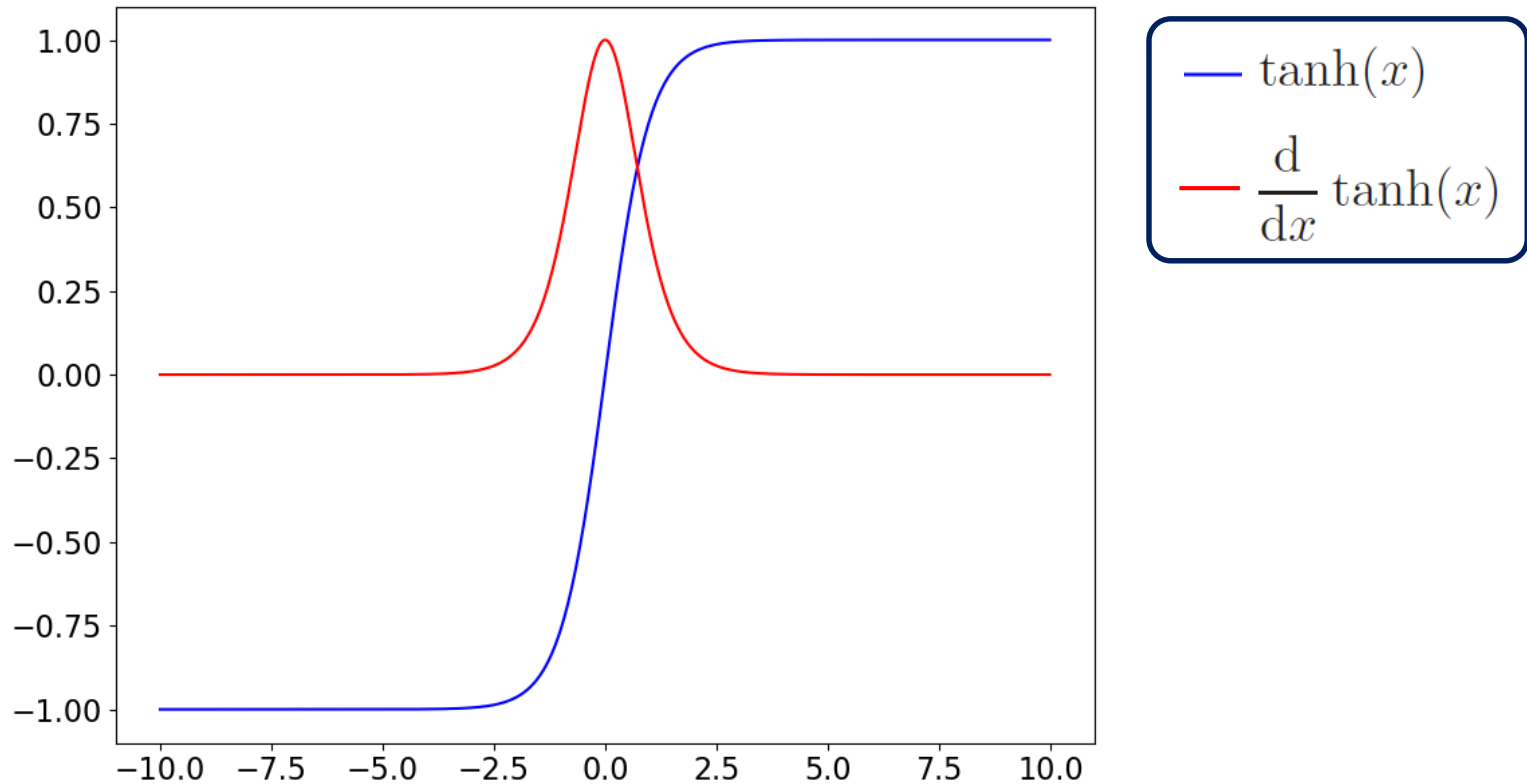
- Nonlinear activation functions $\mathcal{A}(x)$ for which $|\mathcal{A}'(x)|$ is bounded, $\text{diag} \left[\mathcal{A}^{(j+1)'} (\mathbf{w}^{(j+1)\text{T}} \mathbf{z}^{(j)}) \right]$ is also bounded depending on the singular values.
- So there is a limit on how much a vector will be scaled by multiplying the Jacobian with it.

Activation function: sigmoid



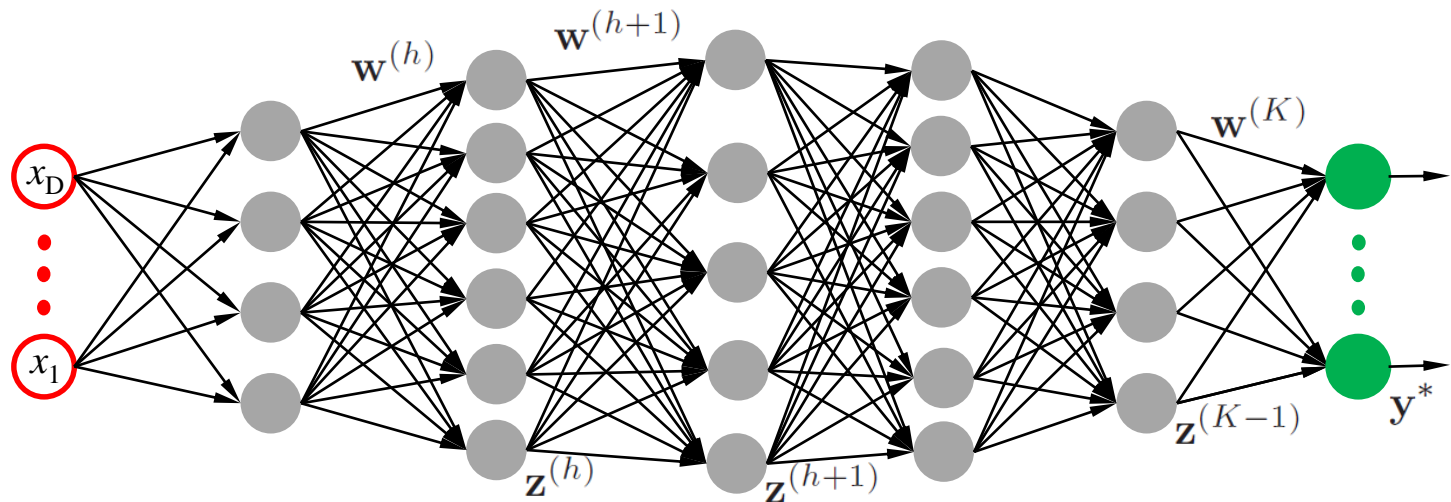
- If $\mathcal{A}^{(j+1)}(\cdot) = \sigma(\cdot)$, then $\left\| \text{diag} \left[\mathcal{A}^{(j+1)'}(\mathbf{w}^{(j+1)\text{T}} \mathbf{z}^{(j)}) \right] \right\| \leq \frac{1}{4}$

Activation function: tanh



- If $\mathcal{A}^{(j+1)}(\cdot) = \tanh(\cdot)$, $\left\| \text{diag} \left[\mathcal{A}^{(j+1)'}(\mathbf{w}^{(j+1)\text{T}} \mathbf{z}^{(j)}) \right] \right\| \leq 1$

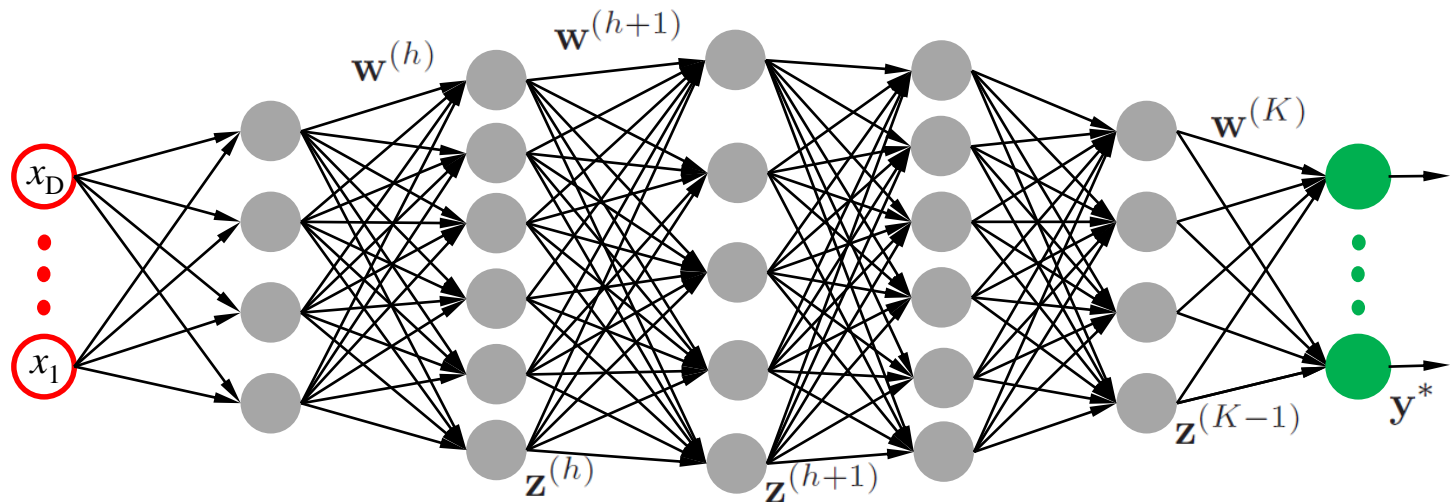
Effect of the weight matrix



$$\frac{\partial \mathcal{L}}{\partial \mathbf{w}^{(h)}} = \frac{\partial \mathcal{L}}{\partial \mathbf{y}^*} \frac{\partial \mathbf{y}^*}{\partial \mathbf{z}^{(K-1)}} \frac{\partial \mathbf{z}^{(K-1)}}{\partial \mathbf{z}^{(K-2)}} \frac{\partial \mathbf{z}^{(K-2)}}{\partial \mathbf{z}^{(K-3)}} \cdots \frac{\partial \mathbf{z}^{(h+1)}}{\partial \mathbf{z}^{(h)}} \frac{\partial \mathbf{z}^{(h)}}{\partial \mathbf{w}^{(h)}}$$

- Effect of the weight matrices on the chain product:
 - Expansion in directions where the singular values of the weight matrices are greater than one.
 - Shrink along directions where the singular values of the weight matrices are less than one.
- Multiplications by the weight matrices can lead to **exploding** or **vanishing** gradients.

Overall effect



$$\frac{\partial \mathcal{L}}{\partial \mathbf{w}^{(h)}} = \frac{\partial \mathcal{L}}{\partial \mathbf{y}^*} \frac{\partial \mathbf{y}^*}{\partial \mathbf{z}^{(K-1)}} \frac{\partial \mathbf{z}^{(K-1)}}{\partial \mathbf{z}^{(K-2)}} \frac{\partial \mathbf{z}^{(K-2)}}{\partial \mathbf{z}^{(K-3)}} \cdots \frac{\partial \mathbf{z}^{(h+1)}}{\partial \mathbf{z}^{(h)}} \frac{\partial \mathbf{z}^{(h)}}{\partial \mathbf{w}^{(h)}}$$

- Overall effect of the Jacobian matrices on the chain product:
 - Expansion in directions where the Jacobian matrices have singular values greater than one.
 - Shrinkage in directions where the Jacobian matrices have singular values less than one.