

# CEE 616: Probabilistic Machine Learning

## M2 Deep Neural Networks: Neural Networks for Structured Data II

**Jimi Oke**

UMassAmherst  

---

College of Engineering

Tue, Oct 21, 2025

# Outline

① Backpropagation

② Training

③ Summary

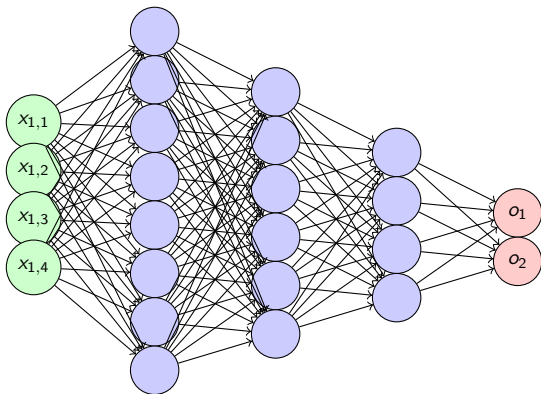
# Forward mode differentiation

# Forward mode differentiation

Let us think of a simple feedforward with inputs  $\mathbf{x} = \mathbf{x}_1 \in \mathbb{R}^{D=4}$ , 3 hidden layers with  $m_1 = 8$ ,  $m_2 = 6$ , and  $m_3 = 4$  neurons, and outputs  $\mathbf{o} \in \mathbb{R}^2$ :

# Forward mode differentiation

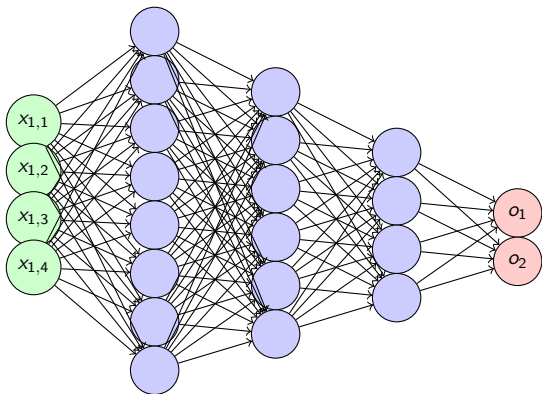
Let us think of a simple feedforward with inputs  $\mathbf{x} = \mathbf{x}_1 \in \mathbb{R}^{D=4}$ , 3 hidden layers with  $m_1 = 8$ ,  $m_2 = 6$ , and  $m_3 = 4$  neurons, and outputs  $\mathbf{o} \in \mathbb{R}^2$ :



Input Layer   Hidden 1 (8)   Hidden 2 (6)   Hidden 3 (4)   Output Layer

# Forward mode differentiation

Let us think of a simple feedforward with inputs  $\mathbf{x} = \mathbf{x}_1 \in \mathbb{R}^{D=4}$ , 3 hidden layers with  $m_1 = 8$ ,  $m_2 = 6$ , and  $m_3 = 4$  neurons, and outputs  $\mathbf{o} \in \mathbb{R}^2$ :

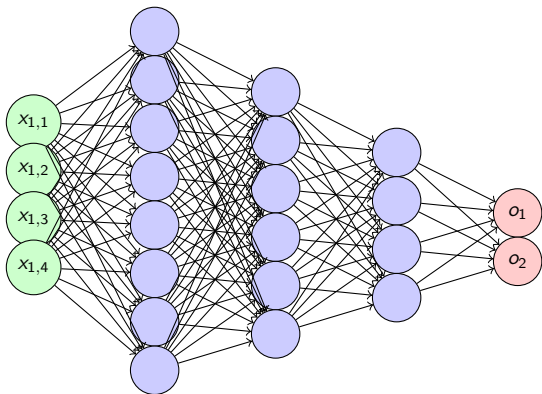


Input Layer Hidden 1 (8) Hidden 2 (6) Hidden 3 (4) Output Layer

We consider each hidden unit as a function  $\mathbf{f}_\ell(\cdot)$ , where  $\ell$  is the layer index, that maps the input  $\mathbf{x}_\ell$  to the output of the hidden unit  $\mathbf{x}_{\ell+1}$ .

# Forward mode differentiation

Let us think of a simple feedforward with inputs  $\mathbf{x} = \mathbf{x}_1 \in \mathbb{R}^{D=4}$ , 3 hidden layers with  $m_1 = 8$ ,  $m_2 = 6$ , and  $m_3 = 4$  neurons, and outputs  $\mathbf{o} \in \mathbb{R}^2$ :



Input Layer Hidden 1 (8) Hidden 2 (6) Hidden 3 (4) Output Layer

We consider each hidden unit as a function  $\mathbf{f}_\ell(\cdot)$ , where  $\ell$  is the layer index, that maps the input  $\mathbf{x}_\ell$  to the output of the hidden unit  $\mathbf{x}_{\ell+1}$ .

# Function composition

We can then express the output of the network as a composition of functions:

$$\boldsymbol{o} = \boldsymbol{f}(\boldsymbol{x}) \quad (1)$$



# Function composition

We can then express the output of the network as a composition of functions:

$$\mathbf{o} = \mathbf{f}(\mathbf{x}) \quad (1)$$

where

$$\mathbf{f} = \mathbf{f}_4 \circ \mathbf{f}_3 \circ \mathbf{f}_2 \circ \mathbf{f}_1 \quad (2)$$

and thus

$$\mathbf{o} = \mathbf{f}_4(\mathbf{f}_3(\mathbf{f}_2(\mathbf{f}_1(\mathbf{x})))) \quad (3)$$

and

$$\begin{aligned} \mathbf{x}_2 &= \mathbf{f}_1(\mathbf{x}_1), & \mathbf{f}_1 : \mathbb{R}^4 &\rightarrow \mathbb{R}^8 \\ \mathbf{x}_3 &= \mathbf{f}_2(\mathbf{x}_2), & \mathbf{f}_2 : \mathbb{R}^8 &\rightarrow \mathbb{R}^6 \\ \mathbf{x}_4 &= \mathbf{f}_3(\mathbf{x}_3), & \mathbf{f}_3 : \mathbb{R}^6 &\rightarrow \mathbb{R}^4 \\ \mathbf{o} &= \mathbf{f}_4(\mathbf{x}_4), & \mathbf{f}_4 : \mathbb{R}^4 &\rightarrow \mathbb{R}^2 \end{aligned}$$

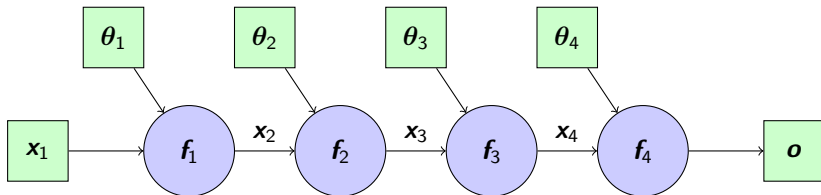
# Computational graph

# Computational graph

We can further visualize the FFNN as a computational graph:

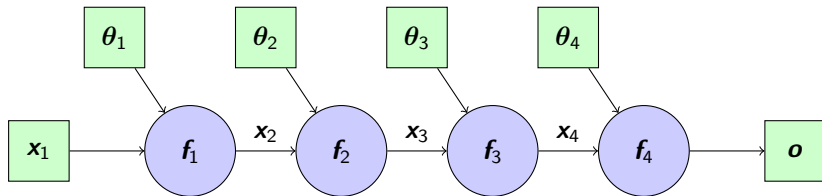
# Computational graph

We can further visualize the FFNN as a computational graph:



# Computational graph

We can further visualize the FFNN as a computational graph:



where  $\theta_\ell$  are the parameters (weights and biases) of layer  $\ell$ .

# Backpropagation

# Backpropagation

To compute the gradient, we need to find:

$$\frac{\partial \mathbf{o}}{\partial \mathbf{x}} = \frac{\partial \mathbf{f}_4(\mathbf{x}_4)}{\partial \mathbf{x}_4} \cdot \frac{\partial \mathbf{f}_3(\mathbf{x}_3)}{\partial \mathbf{x}_3} \cdot \frac{\partial \mathbf{f}_2(\mathbf{x}_2)}{\partial \mathbf{x}_2} \cdot \frac{\partial \mathbf{f}_1(\mathbf{x}_1)}{\partial \mathbf{x}_1} \quad (4)$$

# Backpropagation

To compute the gradient, we need to find:

$$\frac{\partial \mathbf{o}}{\partial \mathbf{x}} = \frac{\partial \mathbf{f}_4(\mathbf{x}_4)}{\partial \mathbf{x}_4} \cdot \frac{\partial \mathbf{f}_3(\mathbf{x}_3)}{\partial \mathbf{x}_3} \cdot \frac{\partial \mathbf{f}_2(\mathbf{x}_2)}{\partial \mathbf{x}_2} \cdot \frac{\partial \mathbf{f}_1(\mathbf{x}_1)}{\partial \mathbf{x}_1} \quad (4)$$

We define the Jacobian matrix of each layer as:

$$\mathbf{J}_f(\mathbf{x}_\ell) = \frac{\partial \mathbf{f}_\ell(\mathbf{x}_\ell)}{\partial \mathbf{x}_\ell}$$



# Backpropagation

To compute the gradient, we need to find:

$$\frac{\partial \mathbf{o}}{\partial \mathbf{x}} = \frac{\partial \mathbf{f}_4(\mathbf{x}_4)}{\partial \mathbf{x}_4} \cdot \frac{\partial \mathbf{f}_3(\mathbf{x}_3)}{\partial \mathbf{x}_3} \cdot \frac{\partial \mathbf{f}_2(\mathbf{x}_2)}{\partial \mathbf{x}_2} \cdot \frac{\partial \mathbf{f}_1(\mathbf{x}_1)}{\partial \mathbf{x}_1} \quad (4)$$

We define the Jacobian matrix of each layer as:

$$\mathbf{J}_f(\mathbf{x}_\ell) = \frac{\partial \mathbf{f}_\ell(\mathbf{x}_\ell)}{\partial \mathbf{x}_\ell} = \begin{pmatrix} \frac{\partial f_{\ell,1}(\mathbf{x}_\ell)}{\partial x_{\ell,1}} & \frac{\partial f_{\ell,1}(\mathbf{x}_\ell)}{\partial x_{\ell,2}} & \cdots & \frac{\partial f_{\ell,1}(\mathbf{x}_\ell)}{\partial x_{\ell,D}} \\ \frac{\partial f_{\ell,2}(\mathbf{x}_\ell)}{\partial x_{\ell,1}} & \frac{\partial f_{\ell,2}(\mathbf{x}_\ell)}{\partial x_{\ell,2}} & \cdots & \frac{\partial f_{\ell,2}(\mathbf{x}_\ell)}{\partial x_{\ell,D}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial f_{\ell,M}(\mathbf{x}_\ell)}{\partial x_{\ell,1}} & \frac{\partial f_{\ell,M}(\mathbf{x}_\ell)}{\partial x_{\ell,2}} & \cdots & \frac{\partial f_{\ell,M}(\mathbf{x}_\ell)}{\partial x_{\ell,D}} \end{pmatrix} \quad (5)$$

# Backpropagation

To compute the gradient, we need to find:

$$\frac{\partial \mathbf{o}}{\partial \mathbf{x}} = \frac{\partial \mathbf{f}_4(\mathbf{x}_4)}{\partial \mathbf{x}_4} \cdot \frac{\partial \mathbf{f}_3(\mathbf{x}_3)}{\partial \mathbf{x}_3} \cdot \frac{\partial \mathbf{f}_2(\mathbf{x}_2)}{\partial \mathbf{x}_2} \cdot \frac{\partial \mathbf{f}_1(\mathbf{x}_1)}{\partial \mathbf{x}_1} \quad (4)$$

We define the Jacobian matrix of each layer as:

$$\mathbf{J}_f(\mathbf{x}_\ell) = \frac{\partial \mathbf{f}_\ell(\mathbf{x}_\ell)}{\partial \mathbf{x}_\ell} = \begin{pmatrix} \frac{\partial f_{\ell,1}(\mathbf{x}_\ell)}{\partial x_{\ell,1}} & \frac{\partial f_{\ell,1}(\mathbf{x}_\ell)}{\partial x_{\ell,2}} & \cdots & \frac{\partial f_{\ell,1}(\mathbf{x}_\ell)}{\partial x_{\ell,D}} \\ \frac{\partial f_{\ell,2}(\mathbf{x}_\ell)}{\partial x_{\ell,1}} & \frac{\partial f_{\ell,2}(\mathbf{x}_\ell)}{\partial x_{\ell,2}} & \cdots & \frac{\partial f_{\ell,2}(\mathbf{x}_\ell)}{\partial x_{\ell,D}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial f_{\ell,M}(\mathbf{x}_\ell)}{\partial x_{\ell,1}} & \frac{\partial f_{\ell,M}(\mathbf{x}_\ell)}{\partial x_{\ell,2}} & \cdots & \frac{\partial f_{\ell,M}(\mathbf{x}_\ell)}{\partial x_{\ell,D}} \end{pmatrix} \quad (5)$$

Thus, the gradient of the output with respect to the input is given by:

# Backpropagation

To compute the gradient, we need to find:

$$\frac{\partial \mathbf{o}}{\partial \mathbf{x}} = \frac{\partial \mathbf{f}_4(\mathbf{x}_4)}{\partial \mathbf{x}_4} \cdot \frac{\partial \mathbf{f}_3(\mathbf{x}_3)}{\partial \mathbf{x}_3} \cdot \frac{\partial \mathbf{f}_2(\mathbf{x}_2)}{\partial \mathbf{x}_2} \cdot \frac{\partial \mathbf{f}_1(\mathbf{x}_1)}{\partial \mathbf{x}_1} \quad (4)$$

We define the Jacobian matrix of each layer as:

$$\mathbf{J}_f(\mathbf{x}_\ell) = \frac{\partial \mathbf{f}_\ell(\mathbf{x}_\ell)}{\partial \mathbf{x}_\ell} = \begin{pmatrix} \frac{\partial f_{\ell,1}(\mathbf{x}_\ell)}{\partial x_{\ell,1}} & \frac{\partial f_{\ell,1}(\mathbf{x}_\ell)}{\partial x_{\ell,2}} & \cdots & \frac{\partial f_{\ell,1}(\mathbf{x}_\ell)}{\partial x_{\ell,D}} \\ \frac{\partial f_{\ell,2}(\mathbf{x}_\ell)}{\partial x_{\ell,1}} & \frac{\partial f_{\ell,2}(\mathbf{x}_\ell)}{\partial x_{\ell,2}} & \cdots & \frac{\partial f_{\ell,2}(\mathbf{x}_\ell)}{\partial x_{\ell,D}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial f_{\ell,M}(\mathbf{x}_\ell)}{\partial x_{\ell,1}} & \frac{\partial f_{\ell,M}(\mathbf{x}_\ell)}{\partial x_{\ell,2}} & \cdots & \frac{\partial f_{\ell,M}(\mathbf{x}_\ell)}{\partial x_{\ell,D}} \end{pmatrix} \quad (5)$$

Thus, the gradient of the output with respect to the input is given by:

$$\frac{\partial \mathbf{o}}{\partial \mathbf{x}} = \mathbf{J}_f(\mathbf{x}_4) \cdot \mathbf{J}_f(\mathbf{x}_3) \cdot \mathbf{J}_f(\mathbf{x}_2) \cdot \mathbf{J}_f(\mathbf{x}_1) \quad (6)$$

# Jacobian matrix (2/2)

# Jacobian matrix (2/2)

We can write the Jacobian as:

# Jacobian matrix (2/2)

We can write the Jacobian as:

$$\mathbf{J}_f(\mathbf{x}) = \frac{\partial \mathbf{f}(\mathbf{x})}{\partial \mathbf{x}} = \begin{pmatrix} \nabla f_1(\mathbf{x})^\top \\ \nabla f_2(\mathbf{x})^\top \\ \vdots \\ \nabla f_M(\mathbf{x})^\top \end{pmatrix}$$

# Jacobian matrix (2/2)

We can write the Jacobian as:

$$\mathbf{J}_f(\mathbf{x}) = \frac{\partial \mathbf{f}(\mathbf{x})}{\partial \mathbf{x}} = \begin{pmatrix} \nabla f_1(\mathbf{x})^\top \\ \nabla f_2(\mathbf{x})^\top \\ \vdots \\ \nabla f_M(\mathbf{x})^\top \end{pmatrix} = \left( \frac{\partial \mathbf{f}}{\partial \mathbf{x}_1}, \frac{\partial \mathbf{f}}{\partial \mathbf{x}_2}, \dots, \frac{\partial \mathbf{f}}{\partial \mathbf{x}_D} \right) \quad (7)$$

# Jacobian matrix (2/2)

We can write the Jacobian as:

$$\mathbf{J}_f(\mathbf{x}) = \frac{\partial \mathbf{f}(\mathbf{x})}{\partial \mathbf{x}} = \begin{pmatrix} \nabla f_1(\mathbf{x})^\top \\ \nabla f_2(\mathbf{x})^\top \\ \vdots \\ \nabla f_M(\mathbf{x})^\top \end{pmatrix} = \left( \frac{\partial \mathbf{f}}{\partial \mathbf{x}_1}, \frac{\partial \mathbf{f}}{\partial \mathbf{x}_2}, \dots, \frac{\partial \mathbf{f}}{\partial \mathbf{x}_D} \right) \quad (7)$$

- $i$ 'th row:  $\nabla f_i(\mathbf{x})^\top \in \mathbb{R}^{1 \times D}$  gradient of the  $i$ 'th output w.r.t. all inputs
- $j$ 'th column:  $\frac{\partial \mathbf{f}}{\partial x_j} \in \mathbb{R}^M$  gradient of all outputs w.r.t. the  $j$ 'th input



# Jacobian matrix (2/2)

We can write the Jacobian as:

$$\mathbf{J}_f(\mathbf{x}) = \frac{\partial \mathbf{f}(\mathbf{x})}{\partial \mathbf{x}} = \begin{pmatrix} \nabla f_1(\mathbf{x})^\top \\ \nabla f_2(\mathbf{x})^\top \\ \vdots \\ \nabla f_M(\mathbf{x})^\top \end{pmatrix} = \left( \frac{\partial \mathbf{f}}{\partial \mathbf{x}_1}, \frac{\partial \mathbf{f}}{\partial \mathbf{x}_2}, \dots, \frac{\partial \mathbf{f}}{\partial \mathbf{x}_D} \right) \quad (7)$$

- $i$ 'th row:  $\nabla f_i(\mathbf{x})^\top \in \mathbb{R}^{1 \times D}$  gradient of the  $i$ 'th output w.r.t. all inputs
- $j$ 'th column:  $\frac{\partial \mathbf{f}}{\partial x_j} \in \mathbb{R}^M$  gradient of all outputs w.r.t. the  $j$ 'th input

Thus:

$$\nabla \mathbf{f}(\mathbf{x}) = \mathbf{J}_f(\mathbf{x})^\top = (\nabla f_1(\mathbf{x}), \nabla f_2(\mathbf{x}), \dots, \nabla f_M(\mathbf{x})) \quad (8)$$

# Vector Jacobian product

# Vector Jacobian product

The vector-Jacobian product (VJP) is the left-multiplication of the Jacobian matrix  $\mathbf{J}_f(\mathbf{x})$  by a vector  $\mathbf{u}$ , which results in a row vector:

# Vector Jacobian product

The vector-Jacobian product (VJP) is the left-multiplication of the Jacobian matrix  $\mathbf{J}_f(\mathbf{x})$  by a vector  $\mathbf{u}$ , which results in a row vector:

$$\mathbf{u}^\top \mathbf{J}_f(\mathbf{x}) \in \mathbb{R}^{1 \times D} \quad (9)$$

# Vector Jacobian product

The vector-Jacobian product (VJP) is the left-multiplication of the Jacobian matrix  $\mathbf{J}_f(\mathbf{x})$  by a vector  $\mathbf{u}$ , which results in a row vector:

$$\mathbf{u}^\top \mathbf{J}_f(\mathbf{x}) \in \mathbb{R}^{1 \times D} \quad (9)$$

- If  $\mathbf{u} \in \mathbb{R}^M$  is a one-hot vector with 1 at index  $i$  and 0 elsewhere, e.g.

$$\mathbf{u} = \begin{pmatrix} 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{pmatrix} = \mathbf{e}_i \quad (10)$$

then the VJP  $\mathbf{e}_i^\top \mathbf{J}_f(\mathbf{x})$  extracts the  $i$ 'th row from  $\mathbf{J}_f(\mathbf{x})$ , which is the gradient of the  $i$ 'th output with respect to all inputs,  $\nabla f_i(\mathbf{x})^\top$ .

# Jacobian-vector product

# Jacobian-vector product

The Jacobian-vector product (JVP) is the right-multiplication of the Jacobian matrix  $\mathbf{J}_f(\mathbf{x})$  by a vector  $\mathbf{v}$ , which results in a column vector:

# Jacobian-vector product

The Jacobian-vector product (JVP) is the right-multiplication of the Jacobian matrix  $\mathbf{J}_f(\mathbf{x})$  by a vector  $\mathbf{v}$ , which results in a column vector:

$$\mathbf{J}_f(\mathbf{x})\mathbf{v} \in \mathbb{R}^M \quad (11)$$



# Jacobian-vector product

The Jacobian-vector product (JVP) is the right-multiplication of the Jacobian matrix  $\mathbf{J}_f(\mathbf{x})$  by a vector  $\mathbf{v}$ , which results in a column vector:

$$\mathbf{J}_f(\mathbf{x})\mathbf{v} \in \mathbb{R}^M \quad (11)$$

- If  $\mathbf{v} \in \mathbb{R}^D$  is a one-hot vector with 1 at index  $j$  and 0 elsewhere, e.g.

$$\mathbf{v} = \begin{pmatrix} 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{pmatrix} = \mathbf{e}_j \quad (12)$$

then the JVP  $\mathbf{J}_f(\mathbf{x})\mathbf{e}_j$  extracts the  $j$ 'th column from  $\mathbf{J}_f(\mathbf{x})$ , which is the gradient of all outputs with respect to the  $j$ 'th input,  $\frac{\partial \mathbf{f}}{\partial x_j}$ .

# Computing the Jacobian

# Computing the Jacobian

- The Jacobian can be computed by performing either  $M$  JVPs (one per input dimension) or  $D$  VJPs (one per output dimension)
- If  $D < M$ , it is more efficient to compute the Jacobian via JVPs for each column  $j$ :

$$\mathbf{J}_f(\mathbf{x})(\mathbf{v}) = \underbrace{\mathbf{J}_f(\mathbf{x})_4}_{M \times M_3} \underbrace{\mathbf{J}_f(\mathbf{x})_3}_{M_3 \times M_2} \underbrace{\mathbf{J}_f(\mathbf{x})_2}_{M_2 \times M_1} \underbrace{\mathbf{J}_f(\mathbf{x})_1}_{M_1 \times D} \underbrace{\mathbf{v}}_{D \times 1} \quad (13)$$

in a right-to-left fashion (**forward mode differentiation**)

- If  $M < D$ , it is more efficient to compute the Jacobian via VJPs for each row  $i$ :

$$\mathbf{u}^\top \mathbf{J}_f(\mathbf{x}) = \underbrace{\mathbf{u}^\top}_{1 \times M} \underbrace{\mathbf{J}_f(\mathbf{x})_4}_{M \times M_3} \underbrace{\mathbf{J}_f(\mathbf{x})_3}_{M_3 \times M_2} \underbrace{\mathbf{J}_f(\mathbf{x})_2}_{M_2 \times M_1} \underbrace{\mathbf{J}_f(\mathbf{x})_1}_{M_1 \times D} \quad (14)$$

in a left-to-right fashion (**reverse mode differentiation**)

- In neural networks, the number of outputs  $M$  is often much smaller than the number of inputs  $D$
- Thus, backpropagation typically employs VJPs to compute the gradients efficiently

# Example: single hidden layer MLP

# Example: single hidden layer MLP

Consider a single hidden layer MLP with an  $\ell_2$  loss for regression (scalar output):

# Example: single hidden layer MLP

Consider a single hidden layer MLP with an  $\ell_2$  loss for regression (scalar output):

$$\mathcal{L}((\mathbf{x}, \mathbf{y}), \boldsymbol{\theta}) = \frac{1}{2} \|\mathbf{y} - \mathbf{W}_2 \varphi(\mathbf{W}_1 \mathbf{x})\|_2^2 \quad (15)$$

# Example: single hidden layer MLP

Consider a single hidden layer MLP with an  $\ell_2$  loss for regression (scalar output):

$$\mathcal{L}((\mathbf{x}, \mathbf{y}), \boldsymbol{\theta}) = \frac{1}{2} \|\mathbf{y} - \mathbf{W}_2 \varphi(\mathbf{W}_1 \mathbf{x})\|_2^2 \quad (15)$$

We can represent the network as follows:

# Example: single hidden layer MLP

Consider a single hidden layer MLP with an  $\ell_2$  loss for regression (scalar output):

$$\mathcal{L}((\mathbf{x}, \mathbf{y}), \boldsymbol{\theta}) = \frac{1}{2} \|\mathbf{y} - \mathbf{W}_2 \varphi(\mathbf{W}_1 \mathbf{x})\|_2^2 \quad (15)$$

We can represent the network as follows:

$$\mathcal{L} = \mathbf{f}_4 \circ \mathbf{f}_3 \circ \mathbf{f}_2 \circ \mathbf{f}_1 \quad (16)$$



# Example: single hidden layer MLP

Consider a single hidden layer MLP with an  $\ell_2$  loss for regression (scalar output):

$$\mathcal{L}((\mathbf{x}, \mathbf{y}), \boldsymbol{\theta}) = \frac{1}{2} \|\mathbf{y} - \mathbf{W}_2 \varphi(\mathbf{W}_1 \mathbf{x})\|_2^2 \quad (15)$$

We can represent the network as follows:

$$\mathcal{L} = \mathbf{f}_4 \circ \mathbf{f}_3 \circ \mathbf{f}_2 \circ \mathbf{f}_1 \quad (16)$$

$$\mathbf{x}_2 = \mathbf{f}_1(\mathbf{x}, \boldsymbol{\theta}_1) = \mathbf{W}_1 \mathbf{x} \quad (17)$$

# Example: single hidden layer MLP

Consider a single hidden layer MLP with an  $\ell_2$  loss for regression (scalar output):

$$\mathcal{L}((\mathbf{x}, \mathbf{y}), \boldsymbol{\theta}) = \frac{1}{2} \|\mathbf{y} - \mathbf{W}_2 \varphi(\mathbf{W}_1 \mathbf{x})\|_2^2 \quad (15)$$

We can represent the network as follows:

$$\mathcal{L} = \mathbf{f}_4 \circ \mathbf{f}_3 \circ \mathbf{f}_2 \circ \mathbf{f}_1 \quad (16)$$

$$\mathbf{x}_2 = \mathbf{f}_1(\mathbf{x}, \boldsymbol{\theta}_1) = \mathbf{W}_1 \mathbf{x} \quad (17)$$

$$\mathbf{x}_3 = \mathbf{f}_2(\mathbf{x}_2, \boldsymbol{\theta}_2) = \varphi(\mathbf{x}_2) \quad (18)$$

# Example: single hidden layer MLP

Consider a single hidden layer MLP with an  $\ell_2$  loss for regression (scalar output):

$$\mathcal{L}((\mathbf{x}, \mathbf{y}), \boldsymbol{\theta}) = \frac{1}{2} \|\mathbf{y} - \mathbf{W}_2 \varphi(\mathbf{W}_1 \mathbf{x})\|_2^2 \quad (15)$$

We can represent the network as follows:

$$\mathcal{L} = \mathbf{f}_4 \circ \mathbf{f}_3 \circ \mathbf{f}_2 \circ \mathbf{f}_1 \quad (16)$$

$$\mathbf{x}_2 = \mathbf{f}_1(\mathbf{x}, \boldsymbol{\theta}_1) = \mathbf{W}_1 \mathbf{x} \quad (17)$$

$$\mathbf{x}_3 = \mathbf{f}_2(\mathbf{x}_2, \boldsymbol{\theta}_2) = \varphi(\mathbf{x}_2) \quad (18)$$

$$\mathbf{x}_4 = \mathbf{f}_3(\mathbf{x}_3, \boldsymbol{\theta}_3) = \mathbf{W}_2 \mathbf{x}_3 \quad (19)$$

# Example: single hidden layer MLP

Consider a single hidden layer MLP with an  $\ell_2$  loss for regression (scalar output):

$$\mathcal{L}((\mathbf{x}, \mathbf{y}), \boldsymbol{\theta}) = \frac{1}{2} \|\mathbf{y} - \mathbf{W}_2 \varphi(\mathbf{W}_1 \mathbf{x})\|_2^2 \quad (15)$$

We can represent the network as follows:

$$\mathcal{L} = \mathbf{f}_4 \circ \mathbf{f}_3 \circ \mathbf{f}_2 \circ \mathbf{f}_1 \quad (16)$$

$$\mathbf{x}_2 = \mathbf{f}_1(\mathbf{x}, \boldsymbol{\theta}_1) = \mathbf{W}_1 \mathbf{x} \quad (17)$$

$$\mathbf{x}_3 = \mathbf{f}_2(\mathbf{x}_2, \boldsymbol{\theta}_2) = \varphi(\mathbf{x}_2) \quad (18)$$

$$\mathbf{x}_4 = \mathbf{f}_3(\mathbf{x}_3, \boldsymbol{\theta}_3) = \mathbf{W}_2 \mathbf{x}_3 \quad (19)$$

$$\mathcal{L} = \mathbf{f}_4(\mathbf{x}_4, \mathbf{y}) = \frac{1}{2} \|\mathbf{y} - \mathbf{x}_4\|_2^2 \quad (20)$$

and  $\mathcal{L} : \mathbb{R}^D \rightarrow \mathbb{R}$ .

## Example: single hidden layer MLP (contd.)

## Example: single hidden layer MLP (contd.)

Using the chain rule, We can express the gradient of the loss wrt the parameters in each layer as:

## Example: single hidden layer MLP (contd.)

Using the chain rule, We can express the gradient of the loss wrt the parameters in each layer as:

$$\frac{\partial \mathcal{L}}{\partial \theta_3} = \frac{\partial \mathcal{L}}{\partial \mathbf{x}_4} \frac{\partial \mathbf{x}_4}{\partial \theta_3} \quad (21)$$

## Example: single hidden layer MLP (contd.)

Using the chain rule, We can express the gradient of the loss wrt the parameters in each layer as:

$$\frac{\partial \mathcal{L}}{\partial \theta_3} = \frac{\partial \mathcal{L}}{\partial \mathbf{x}_4} \frac{\partial \mathbf{x}_4}{\partial \theta_3} \quad (21)$$

$$\frac{\partial \mathcal{L}}{\partial \theta_2} = \frac{\partial \mathcal{L}}{\partial \mathbf{x}_4} \frac{\partial \mathbf{x}_4}{\partial \mathbf{x}_3} \frac{\partial \mathbf{x}_3}{\partial \theta_2} \quad (22)$$



## Example: single hidden layer MLP (contd.)

Using the chain rule, We can express the gradient of the loss wrt the parameters in each layer as:

$$\frac{\partial \mathcal{L}}{\partial \theta_3} = \frac{\partial \mathcal{L}}{\partial \mathbf{x}_4} \frac{\partial \mathbf{x}_4}{\partial \theta_3} \quad (21)$$

$$\frac{\partial \mathcal{L}}{\partial \theta_2} = \frac{\partial \mathcal{L}}{\partial \mathbf{x}_4} \frac{\partial \mathbf{x}_4}{\partial \mathbf{x}_3} \frac{\partial \mathbf{x}_3}{\partial \theta_2} \quad (22)$$

$$\frac{\partial \mathcal{L}}{\partial \theta_1} = \frac{\partial \mathcal{L}}{\partial \mathbf{x}_4} \frac{\partial \mathbf{x}_4}{\partial \mathbf{x}_3} \frac{\partial \mathbf{x}_3}{\partial \mathbf{x}_2} \frac{\partial \mathbf{x}_2}{\partial \theta_1} \quad (23)$$

## Example: single hidden layer MLP (contd.)

Using the chain rule, We can express the gradient of the loss wrt the parameters in each layer as:

$$\frac{\partial \mathcal{L}}{\partial \theta_3} = \frac{\partial \mathcal{L}}{\partial \mathbf{x}_4} \frac{\partial \mathbf{x}_4}{\partial \theta_3} \quad (21)$$

$$\frac{\partial \mathcal{L}}{\partial \theta_2} = \frac{\partial \mathcal{L}}{\partial \mathbf{x}_4} \frac{\partial \mathbf{x}_4}{\partial \mathbf{x}_3} \frac{\partial \mathbf{x}_3}{\partial \theta_2} \quad (22)$$

$$\frac{\partial \mathcal{L}}{\partial \theta_1} = \frac{\partial \mathcal{L}}{\partial \mathbf{x}_4} \frac{\partial \mathbf{x}_4}{\partial \mathbf{x}_3} \frac{\partial \mathbf{x}_3}{\partial \mathbf{x}_2} \frac{\partial \mathbf{x}_2}{\partial \theta_1} \quad (23)$$

Note:

$$\frac{\partial \mathcal{L}}{\partial \theta_k} = (\nabla_{\theta_k} \mathcal{L})^\top \quad (24)$$

is a  $d_k$ -dimensional gradient row vector, where  $d_k$  is the number of parameters in layer  $k$

## Example: single hidden layer MLP (contd.)

Using the chain rule, We can express the gradient of the loss wrt the parameters in each layer as:

$$\frac{\partial \mathcal{L}}{\partial \theta_3} = \frac{\partial \mathcal{L}}{\partial \mathbf{x}_4} \frac{\partial \mathbf{x}_4}{\partial \theta_3} \quad (21)$$

$$\frac{\partial \mathcal{L}}{\partial \theta_2} = \frac{\partial \mathcal{L}}{\partial \mathbf{x}_4} \frac{\partial \mathbf{x}_4}{\partial \mathbf{x}_3} \frac{\partial \mathbf{x}_3}{\partial \theta_2} \quad (22)$$

$$\frac{\partial \mathcal{L}}{\partial \theta_1} = \frac{\partial \mathcal{L}}{\partial \mathbf{x}_4} \frac{\partial \mathbf{x}_4}{\partial \mathbf{x}_3} \frac{\partial \mathbf{x}_3}{\partial \mathbf{x}_2} \frac{\partial \mathbf{x}_2}{\partial \theta_1} \quad (23)$$

Note:

$$\frac{\partial \mathcal{L}}{\partial \theta_k} = (\nabla_{\theta_k} \mathcal{L})^\top \quad (24)$$

is a  $d_k$ -dimensional gradient row vector, where  $d_k$  is the number of parameters in layer  $k$  and can be computed recursively as VJPs:

$$\frac{\partial \mathcal{L}}{\partial \theta_k} = \mathbf{u}_{k+1}^\top \frac{\partial \mathbf{f}_k(\mathbf{x}_k, \theta_k)}{\partial \theta_k} \quad (25)$$

where  $\mathbf{x}_{k+1} = \mathbf{f}_k(\mathbf{x}_k, \theta_k)$ ,  $\mathbf{u}_k^\top = \mathbf{u}_{k+1}^\top \frac{\partial \mathbf{f}_k(\mathbf{x}_k, \theta_k)}{\partial \mathbf{x}_k}$  and  $\mathbf{u}_{K+1} = 1$ .

## Example: single hidden layer MLP (contd.)

## Example: single hidden layer MLP (contd.)

We then need to specify the Jacobians for each layer:

# Example: single hidden layer MLP (contd.)

We then need to specify the Jacobians for each layer:

- Output layer:

$$\frac{\partial \mathcal{L}}{\partial \mathbf{x}_4} = \mathbf{u}_4^\top = \frac{\partial \frac{1}{2} \|\mathbf{y} - \mathbf{x}_4\|_2^2}{\partial \mathbf{x}_4} = (\mathbf{x}_4 - \mathbf{y})^\top \quad (26)$$

# Example: single hidden layer MLP (contd.)

We then need to specify the Jacobians for each layer:

- Output layer:

$$\frac{\partial \mathcal{L}}{\partial \mathbf{x}_4} = \mathbf{u}_4^\top = \frac{\partial \frac{1}{2} \|\mathbf{y} - \mathbf{x}_4\|_2^2}{\partial \mathbf{x}_4} = (\mathbf{x}_4 - \mathbf{y})^\top \quad (26)$$

- Linear layer:

$$\frac{\partial \mathbf{x}_4}{\partial \mathbf{x}_3} = \mathbf{J}_{f_3}(\mathbf{x}_3) = \frac{\partial \mathbf{W}_2 \mathbf{x}_3}{\partial \mathbf{x}_3} = \mathbf{W}_2 \quad (27)$$

$$\frac{\partial \mathbf{x}_4}{\partial \mathbf{W}_2} = \frac{\partial \mathbf{W}_2 \mathbf{x}_3}{\partial \mathbf{W}_2} = \mathbf{x}_3^\top \quad (28)$$

# Example: single hidden layer MLP (contd.)

We then need to specify the Jacobians for each layer:

- Output layer:

$$\frac{\partial \mathcal{L}}{\partial \mathbf{x}_4} = \mathbf{u}_4^\top = \frac{\partial \frac{1}{2} \|\mathbf{y} - \mathbf{x}_4\|_2^2}{\partial \mathbf{x}_4} = (\mathbf{x}_4 - \mathbf{y})^\top \quad (26)$$

- Linear layer:

$$\frac{\partial \mathbf{x}_4}{\partial \mathbf{x}_3} = \mathbf{J}_{f_3}(\mathbf{x}_3) = \frac{\partial \mathbf{W}_2 \mathbf{x}_3}{\partial \mathbf{x}_3} = \mathbf{W}_2 \quad (27)$$

$$\frac{\partial \mathbf{x}_4}{\partial \mathbf{W}_2} = \frac{\partial \mathbf{W}_2 \mathbf{x}_3}{\partial \mathbf{W}_2} = \mathbf{x}_3^\top \quad (28)$$



# Batch learning

The gradient descent method was described without referencing the observations.

# Batch learning

The gradient descent method was described without referencing the observations. In reality, the forward and backward passes are performed for each observation in the training set, with parameter updates obtained by **averaging** the gradients:

# Batch learning

The gradient descent method was described without referencing the observations. In reality, the forward and backward passes are performed for each observation in the training set, with parameter updates obtained by **averaging** the gradients:

$$\theta_{\ell,t+1} =$$

# Batch learning

The gradient descent method was described without referencing the observations. In reality, the forward and backward passes are performed for each observation in the training set, with parameter updates obtained by **averaging** the gradients:

$$\theta_{\ell,t+1} = \theta_{\ell,t} -$$

# Batch learning

The gradient descent method was described without referencing the observations. In reality, the forward and backward passes are performed for each observation in the training set, with parameter updates obtained by **averaging** the gradients:

$$\boldsymbol{\theta}_{\ell,t+1} = \boldsymbol{\theta}_{\ell,t} - \rho \frac{1}{N} \sum_{n=1}^N \frac{\partial \mathcal{L}_n(\boldsymbol{\theta}_{\ell,t})}{\partial \boldsymbol{\theta}_{\ell}} \quad (29)$$

# Batch learning

The gradient descent method was described without referencing the observations. In reality, the forward and backward passes are performed for each observation in the training set, with parameter updates obtained by **averaging** the gradients:

$$\boldsymbol{\theta}_{\ell,t+1} = \boldsymbol{\theta}_{\ell,t} - \rho \frac{1}{N} \sum_{n=1}^N \frac{\partial \mathcal{L}_n(\boldsymbol{\theta}_{\ell,t})}{\partial \boldsymbol{\theta}_{\ell}} \quad (29)$$

This approach is called **batch learning**.

# Batch learning

The gradient descent method was described without referencing the observations. In reality, the forward and backward passes are performed for each observation in the training set, with parameter updates obtained by **averaging** the gradients:

$$\boldsymbol{\theta}_{\ell,t+1} = \boldsymbol{\theta}_{\ell,t} - \rho \frac{1}{N} \sum_{n=1}^N \frac{\partial \mathcal{L}_n(\boldsymbol{\theta}_{\ell,t})}{\partial \boldsymbol{\theta}_{\ell}} \quad (29)$$

This approach is called **batch learning**.

- One sweep through all the training observations is called an *epoch*

# Batch learning

The gradient descent method was described without referencing the observations. In reality, the forward and backward passes are performed for each observation in the training set, with parameter updates obtained by **averaging** the gradients:

$$\boldsymbol{\theta}_{\ell,t+1} = \boldsymbol{\theta}_{\ell,t} - \rho \frac{1}{N} \sum_{n=1}^N \frac{\partial \mathcal{L}_n(\boldsymbol{\theta}_{\ell,t})}{\partial \boldsymbol{\theta}_{\ell}} \quad (29)$$

This approach is called **batch learning**.

- One sweep through all the training observations is called an *epoch*
- In batch learning, only one update results from a training epoch



# Batch learning

The gradient descent method was described without referencing the observations. In reality, the forward and backward passes are performed for each observation in the training set, with parameter updates obtained by **averaging** the gradients:

$$\boldsymbol{\theta}_{\ell,t+1} = \boldsymbol{\theta}_{\ell,t} - \rho \frac{1}{N} \sum_{n=1}^N \frac{\partial \mathcal{L}_n(\boldsymbol{\theta}_{\ell,t})}{\partial \boldsymbol{\theta}_{\ell}} \quad (29)$$

This approach is called **batch learning**.

- One sweep through all the training observations is called an *epoch*
- In batch learning, only one update results from a training epoch
- Thus, several epochs are required for convergence

# Stochastic gradient descent

# Stochastic gradient descent

In standard gradient descent (batch learning, the updates are performed only after the gradient is computed for *all* training observations.

# Stochastic gradient descent

In standard gradient descent (batch learning, the updates are performed only after the gradient is computed for *all* training observations.

The stochastic gradient descent approach approximates the gradient in each iteration using a **randomly selected** observation  $n$ :

# Stochastic gradient descent

In standard gradient descent (batch learning, the updates are performed only after the gradient is computed for *all* training observations.

The stochastic gradient descent approach approximates the gradient in each iteration using a **randomly selected** observation  $n$ :

$$\theta_{\ell,t+1} =$$

# Stochastic gradient descent

In standard gradient descent (batch learning, the updates are performed only after the gradient is computed for *all* training observations.

The stochastic gradient descent approach approximates the gradient in each iteration using a **randomly selected** observation  $n$ :

$$\theta_{\ell,t+1} = \theta_{\ell,t} -$$

# Stochastic gradient descent

In standard gradient descent (batch learning, the updates are performed only after the gradient is computed for *all* training observations.

The stochastic gradient descent approach approximates the gradient in each iteration using a **randomly selected** observation  $n$ :

$$\boldsymbol{\theta}_{\ell,t+1} = \boldsymbol{\theta}_{\ell,t} - \rho \frac{\partial \mathcal{L}_{\textcolor{red}{n}}(\boldsymbol{\theta}_{\ell,t})}{\partial \boldsymbol{\theta}_{\ell}} \quad (30)$$

# Stochastic gradient descent

In standard gradient descent (batch learning, the updates are performed only after the gradient is computed for *all* training observations.

The stochastic gradient descent approach approximates the gradient in each iteration using a **randomly selected** observation  $n$ :

$$\boldsymbol{\theta}_{\ell,t+1} = \boldsymbol{\theta}_{\ell,t} - \rho \frac{\partial \mathcal{L}_n(\boldsymbol{\theta}_{\ell,t})}{\partial \boldsymbol{\theta}_{\ell}} \quad (30)$$

(31)



# Stochastic gradient descent

In standard gradient descent (batch learning, the updates are performed only after the gradient is computed for *all* training observations.

The stochastic gradient descent approach approximates the gradient in each iteration using a **randomly selected** observation  $n$ :

$$\theta_{\ell,t+1} = \theta_{\ell,t} - \rho \frac{\partial \mathcal{L}_n(\theta_{\ell,t})}{\partial \theta_{\ell}} \quad (30)$$

(31)

- Each iteration over observation  $n$  results in a weight/bias update

# Stochastic gradient descent

In standard gradient descent (batch learning, the updates are performed only after the gradient is computed for *all* training observations.

The stochastic gradient descent approach approximates the gradient in each iteration using a **randomly selected** observation  $n$ :

$$\theta_{\ell,t+1} = \theta_{\ell,t} - \rho \frac{\partial \mathcal{L}_n(\theta_{\ell,t})}{\partial \theta_{\ell}} \quad (30)$$

(31)

- Each iteration over observation  $n$  results in a weight/bias update
- Thus, one sweep through the entire training set (an epoch) produces  $N$  updates

# Stochastic gradient descent

In standard gradient descent (batch learning, the updates are performed only after the gradient is computed for *all* training observations.

The stochastic gradient descent approach approximates the gradient in each iteration using a **randomly selected** observation  $n$ :

$$\theta_{\ell,t+1} = \theta_{\ell,t} - \rho \frac{\partial \mathcal{L}_n(\theta_{\ell,t})}{\partial \theta_{\ell}} \quad (30)$$

$$(31)$$

- Each iteration over observation  $n$  results in a weight/bias update
- Thus, one sweep through the entire training set (an epoch) produces  $N$  updates

# Stochastic gradient descent

In standard gradient descent (batch learning, the updates are performed only after the gradient is computed for *all* training observations.

The stochastic gradient descent approach approximates the gradient in each iteration using a **randomly selected** observation  $n$ :

$$\theta_{\ell,t+1} = \theta_{\ell,t} - \rho \frac{\partial \mathcal{L}_n(\theta_{\ell,t})}{\partial \theta_{\ell}} \quad (30)$$

$$(31)$$

- Each iteration over observation  $n$  results in a weight/bias update
- Thus, one sweep through the entire training set (an epoch) produces  $N$  updates

This procedure is also known as **online learning**

# Mini-batch learning

# Mini-batch learning

To introduce stability, we can compute weight updates over a *subset* of training observations

- ① Randomly sample a mini-batch  $\mathcal{B}_t$  of  $B$  samples (this means that  $|\mathcal{B}_t| = B$ ), where  $B \ll N$ .

# Mini-batch learning

To introduce stability, we can compute weight updates over a *subset* of training observations

- ① Randomly sample a mini-batch  $\mathcal{B}_t$  of  $B$  samples (this means that  $|\mathcal{B}_t| = B$ ), where  $B \ll N$ .
- ② Perform a forward and backward pass through each mini-batch to update the weights:

# Mini-batch learning

To introduce stability, we can compute weight updates over a *subset* of training observations

- ① Randomly sample a mini-batch  $\mathcal{B}_t$  of  $B$  samples (this means that  $|\mathcal{B}_t| = B$ ), where  $B \ll N$ .
- ② Perform a forward and backward pass through each mini-batch to update the weights:



# Mini-batch learning

To introduce stability, we can compute weight updates over a *subset* of training observations

- ① Randomly sample a mini-batch  $\mathcal{B}_t$  of  $B$  samples (this means that  $|\mathcal{B}_t| = B$ ), where  $B \ll N$ .
- ② Perform a forward and backward pass through each mini-batch to update the weights:

$$\theta_{\ell,t+1} =$$

# Mini-batch learning

To introduce stability, we can compute weight updates over a *subset* of training observations

- ① Randomly sample a mini-batch  $\mathcal{B}_t$  of  $B$  samples (this means that  $|\mathcal{B}_t| = B$ ), where  $B \ll N$ .
- ② Perform a forward and backward pass through each mini-batch to update the weights:

$$\theta_{\ell,t+1} = \theta_{\ell,t} -$$

# Mini-batch learning

To introduce stability, we can compute weight updates over a *subset* of training observations

- ① Randomly sample a mini-batch  $\mathcal{B}_t$  of  $B$  samples (this means that  $|\mathcal{B}_t| = B$ ), where  $B \ll N$ .
- ② Perform a forward and backward pass through each mini-batch to update the weights:

$$\boldsymbol{\theta}_{\ell,t+1} = \boldsymbol{\theta}_{\ell,t} - \rho \frac{1}{|\mathcal{B}_t|} \sum_{n \in \mathcal{B}_t} \frac{\partial \mathcal{L}_m(\boldsymbol{\theta}_{\ell,t})}{\partial \boldsymbol{\theta}_{\ell}} \quad (32)$$

# Mini-batch learning

To introduce stability, we can compute weight updates over a *subset* of training observations

- ① Randomly sample a mini-batch  $\mathcal{B}_t$  of  $B$  samples (this means that  $|\mathcal{B}_t| = B$ ), where  $B \ll N$ .
- ② Perform a forward and backward pass through each mini-batch to update the weights:

$$\theta_{\ell,t+1} = \theta_{\ell,t} - \rho \frac{1}{|\mathcal{B}_t|} \sum_{n \in \mathcal{B}_t} \frac{\partial \mathcal{L}_m(\theta_{\ell,t})}{\partial \theta_{\ell}} \quad (32)$$

Thus, for each iteration, average the gradient over a *mini-batch* of  $B$  randomly selected observations

- ③ Repeat step 2 until convergence

# Considerations

# Considerations

- Online learning is more efficient for very large datasets

# Considerations

- Online learning is more efficient for very large datasets
- In practice, mini-batch learning is employed, as batch sizes (usually,  $M = 16$  or  $M = 32$ ) can be chosen to take advantage of parallel computing architectures

# Considerations

- Online learning is more efficient for very large datasets
- In practice, mini-batch learning is employed, as batch sizes (usually,  $M = 16$  or  $M = 32$ ) can be chosen to take advantage of parallel computing architectures
- An **adaptive** learning rate  $\rho$  can guarantee convergence, e.g.  $\rho_r = \frac{1}{r}$



# Considerations

- Online learning is more efficient for very large datasets
- In practice, mini-batch learning is employed, as batch sizes (usually,  $M = 16$  or  $M = 32$ ) can be chosen to take advantage of parallel computing architectures
- An **adaptive** learning rate  $\rho$  can guarantee convergence, e.g.  $\rho_r = \frac{1}{r}$
- Input **standardization** (mean zero, SD 1) is recommended for consistent weight initialization and regularization

# Considerations

- Online learning is more efficient for very large datasets
- In practice, mini-batch learning is employed, as batch sizes (usually,  $M = 16$  or  $M = 32$ ) can be chosen to take advantage of parallel computing architectures
- An **adaptive** learning rate  $\rho$  can guarantee convergence, e.g.  $\rho_r = \frac{1}{r}$
- Input **standardization** (mean zero, SD 1) is recommended for consistent weight initialization and regularization
- Weights/biases are **initialized** to *small* values near 0 for better performance

# Considerations

- Online learning is more efficient for very large datasets
- In practice, mini-batch learning is employed, as batch sizes (usually,  $M = 16$  or  $M = 32$ ) can be chosen to take advantage of parallel computing architectures
- An **adaptive** learning rate  $\rho$  can guarantee convergence, e.g.  $\rho_r = \frac{1}{r}$
- Input **standardization** (mean zero, SD 1) is recommended for consistent weight initialization and regularization
- Weights/biases are **initialized** to *small* values near 0 for better performance
- Cost function  $C$  is nonlinear and nonconvex; other optimization approaches (e.g. conjugate gradient) can provide faster convergence compared to stochastic gradient descent

# Regularization

# Regularization

DNNs are prone to overfitting.

# Regularization

DNNs are prone to overfitting. To mitigate this, we can regularize them by:

# Regularization

DNNs are prone to overfitting. To mitigate this, we can regularize them by:

- Early stopping: terminating training if objective does not improve after a specified number of epochs (patience)

# Regularization

DNNs are prone to overfitting. To mitigate this, we can regularize them by:

- Early stopping: terminating training if objective does not improve after a specified number of epochs (patience)
- Weight decay:



# Regularization

DNNs are prone to overfitting. To mitigate this, we can regularize them by:

- Early stopping: terminating training if objective does not improve after a specified number of epochs (patience)
- Weight decay:

$C \leftarrow$

# Regularization

DNNs are prone to overfitting. To mitigate this, we can regularize them by:

- Early stopping: terminating training if objective does not improve after a specified number of epochs (patience)
- Weight decay:

$$C \leftarrow C + \lambda J =$$

# Regularization

DNNs are prone to overfitting. To mitigate this, we can regularize them by:

- Early stopping: terminating training if objective does not improve after a specified number of epochs (patience)
- Weight decay:

$$C \leftarrow C + \lambda J = C + \lambda \left( \sum w^2 + \sum b^2 \right) \quad (33)$$

# Regularization

DNNs are prone to overfitting. To mitigate this, we can regularize them by:

- Early stopping: terminating training if objective does not improve after a specified number of epochs (patience)
- Weight decay:

$$C \leftarrow C + \lambda J = C + \lambda \left( \sum w^2 + \sum b^2 \right) \quad (33)$$

where  $\lambda$  is tuning parameter estimated via cross-validation

# Regularization

DNNs are prone to overfitting. To mitigate this, we can regularize them by:

- Early stopping: terminating training if objective does not improve after a specified number of epochs (patience)
- Weight decay:

$$C \leftarrow C + \lambda J = C + \lambda \left( \sum w^2 + \sum b^2 \right) \quad (33)$$

where  $\lambda$  is tuning parameter estimated via cross-validation

- Model compression via  $\ell_1$  regularization of weights (sparse DNNs)

# Regularization

DNNs are prone to overfitting. To mitigate this, we can regularize them by:

- Early stopping: terminating training if objective does not improve after a specified number of epochs (patience)
- Weight decay:

$$C \leftarrow C + \lambda J = C + \lambda \left( \sum w^2 + \sum b^2 \right) \quad (33)$$

where  $\lambda$  is tuning parameter estimated via cross-validation

- Model compression via  $\ell_1$  regularization of weights (sparse DNNs)
- Dropout:

# Regularization

DNNs are prone to overfitting. To mitigate this, we can regularize them by:

- Early stopping: terminating training if objective does not improve after a specified number of epochs (patience)
- Weight decay:

$$C \leftarrow C + \lambda J = C + \lambda \left( \sum w^2 + \sum b^2 \right) \quad (33)$$

where  $\lambda$  is tuning parameter estimated via cross-validation

- Model compression via  $\ell_1$  regularization of weights (sparse DNNs)
- Dropout: randomly switching off connections from each neuron with probability  $p$

# Regression MLP architecture



# Regression MLP architecture

Typical hyperparameter values are:

# Regression MLP architecture

Typical hyperparameter values are:

Hyperparameter	Value
# input neurons	1 per input feature

# Regression MLP architecture

Typical hyperparameter values are:

Hyperparameter	Value
# input neurons	1 per input feature
# hidden layers	Usually 1 – 5

# Regression MLP architecture

Typical hyperparameter values are:

Hyperparameter	Value
# input neurons	1 per input feature
# hidden layers	Usually 1 – 5
# neurons per hidden layer	Usually 10 – 100

# Regression MLP architecture

Typical hyperparameter values are:

Hyperparameter	Value
# input neurons	1 per input feature
# hidden layers	Usually 1 – 5
# neurons per hidden layer	Usually 10 – 100
# output neurons	1 per prediction dimension

# Regression MLP architecture

Typical hyperparameter values are:

Hyperparameter	Value
# input neurons	1 per input feature
# hidden layers	Usually 1 – 5
# neurons per hidden layer	Usually 10 – 100
# output neurons	1 per prediction dimension
hidden layer activation	ReLU

# Regression MLP architecture

Typical hyperparameter values are:

Hyperparameter	Value
# input neurons	1 per input feature
# hidden layers	Usually 1 – 5
# neurons per hidden layer	Usually 10 – 100
# output neurons	1 per prediction dimension
hidden layer activation	ReLU
output activation	None (if unbounded)

# Regression MLP architecture

Typical hyperparameter values are:

Hyperparameter	Value
# input neurons	1 per input feature
# hidden layers	Usually 1 – 5
# neurons per hidden layer	Usually 10 – 100
# output neurons	1 per prediction dimension
hidden layer activation	ReLU
output activation	None (if unbounded)
loss function	MSE or MAE/Huber



# Classification MLP architecture

# Classification MLP architecture

- For classification, input and hidden layers are chosen in similar fashion to the regression case

# Classification MLP architecture

- For classification, input and hidden layers are chosen in similar fashion to the regression case
- However, the number of output neurons is given by the name of classes/labels
- The output layer activation is typically the softmax function:

$$o_c = \mathcal{S}(\mathbf{a}_c) = \frac{e^{\mathbf{a}_c}}{\sum_{c'=1}^C e^{\mathbf{a}_{c'}}} \quad (34)$$

where  $\mathbf{a}_c$  is the unnormalized log probability of each class  $c$

# Classification MLP architecture

- For classification, input and hidden layers are chosen in similar fashion to the regression case
- However, the number of output neurons is given by the name of classes/labels
- The output layer activation is typically the softmax function:

$$o_c = \mathcal{S}(\mathbf{a}_c) = \frac{e^{\mathbf{a}_c}}{\sum_{c'=1}^C e^{\mathbf{a}_{c'}}} \quad (34)$$

where  $\mathbf{a}_c$  is the unnormalized log probability of each class  $c$

- The loss function is taken as the **categorical cross-entropy**:

# Classification MLP architecture

- For classification, input and hidden layers are chosen in similar fashion to the regression case
- However, the number of output neurons is given by the name of classes/labels
- The output layer activation is typically the softmax function:

$$o_c = \mathcal{S}(\mathbf{a}_c) = \frac{e^{\mathbf{a}_c}}{\sum_{c'=1}^C e^{\mathbf{a}_{c'}}} \quad (34)$$

where  $\mathbf{a}_c$  is the unnormalized log probability of each class  $c$

- The loss function is taken as the **categorical cross-entropy**:

$$\mathcal{L} = - \sum_{c=1}^C y_c \log p_c = - \sum_{c=1}^C y_c \log(\mathcal{S}(\mathbf{a}_c)) \quad (35)$$

# Other types of neural networks

# Other types of neural networks

The standard ANN architecture (MLP) we have studied is also called the feed-forward network.

Other architectures have been shown to give better performance for various applications:

# Other types of neural networks

The standard ANN architecture (MLP) we have studied is also called the feed-forward network.

Other architectures have been shown to give better performance for various applications:

- Recurrent neural networks (RNNs): time-series forecasting



# Other types of neural networks

The standard ANN architecture (MLP) we have studied is also called the feed-forward network.

Other architectures have been shown to give better performance for various applications:

- Recurrent neural networks (RNNs): time-series forecasting
- Convolutional neural networks (CNNs): image classification

# Other types of neural networks

The standard ANN architecture (MLP) we have studied is also called the feed-forward network.

Other architectures have been shown to give better performance for various applications:

- Recurrent neural networks (RNNs): time-series forecasting
- Convolutional neural networks (CNNs): image classification
- Long short-term memory networks (LSTMs): time-series, pattern identification, etc.

# Other types of neural networks

The standard ANN architecture (MLP) we have studied is also called the feed-forward network.

Other architectures have been shown to give better performance for various applications:

- Recurrent neural networks (RNNs): time-series forecasting
- Convolutional neural networks (CNNs): image classification
- Long short-term memory networks (LSTMs): time-series, pattern identification, etc.

# Other types of neural networks

The standard ANN architecture (MLP) we have studied is also called the feed-forward network.

Other architectures have been shown to give better performance for various applications:

- Recurrent neural networks (RNNs): time-series forecasting
- Convolutional neural networks (CNNs): image classification
- Long short-term memory networks (LSTMs): time-series, pattern identification, etc.