

## Lecture 3

### MLP

#### What is an MLP (and why it works)

- **MLP = feed-forward network:** layers of neurons, each layer fully connected to the previous one.
- **Key ingredients:** linear transforms + **nonlinear** activations. Without nonlinearity, a stack of layers collapses to one linear map.
- **Power:** With enough hidden units, an MLP can approximate any continuous function on a compact set (Universal Approximation Theorem).

#### Anatomy of a layer

For layer  $l$ :

$$\mathbf{z}^{[l]} = W^{[l]}\mathbf{a}^{[l-1]} + \mathbf{b}^{[l]}, \quad \mathbf{a}^{[l]} = f^{[l]}(\mathbf{z}^{[l]})$$

- $W^{[l]} \in \mathbb{R}^{m_l \times m_{l-1}}$ ,  $\mathbf{b}^{[l]} \in \mathbb{R}^{m_l}$ .
- $\mathbf{a}^{[0]} = \mathbf{x}$  (input),  $\mathbf{a}^{[L]}$  is output.

#### Activations (and when to use them)

- **ReLU:**  $f(z) = \max(0, z)$ . Fast, sparse gradients; watch “dead” neurons (use good init).
- **Leaky/Parametric ReLU:** fixes dead ReLU ( $f(z) = \max(\alpha z, z)$ ).
- **Tanh:** zero-centered; can saturate (small gradients at large  $|z|$ ).
- **Sigmoid:** for probabilities in **binary** output; avoid in deep hidden layers (saturation).
- **Softmax:** multi-class probabilities:  $\hat{y}_k = \frac{e^{z_k}}{\sum_j e^{z_j}}$ .

#### Losses (likelihood view)

- **Binary classification** (target  $y \in \{0, 1\}$ , output  $\hat{y} = \sigma(z)$ ):

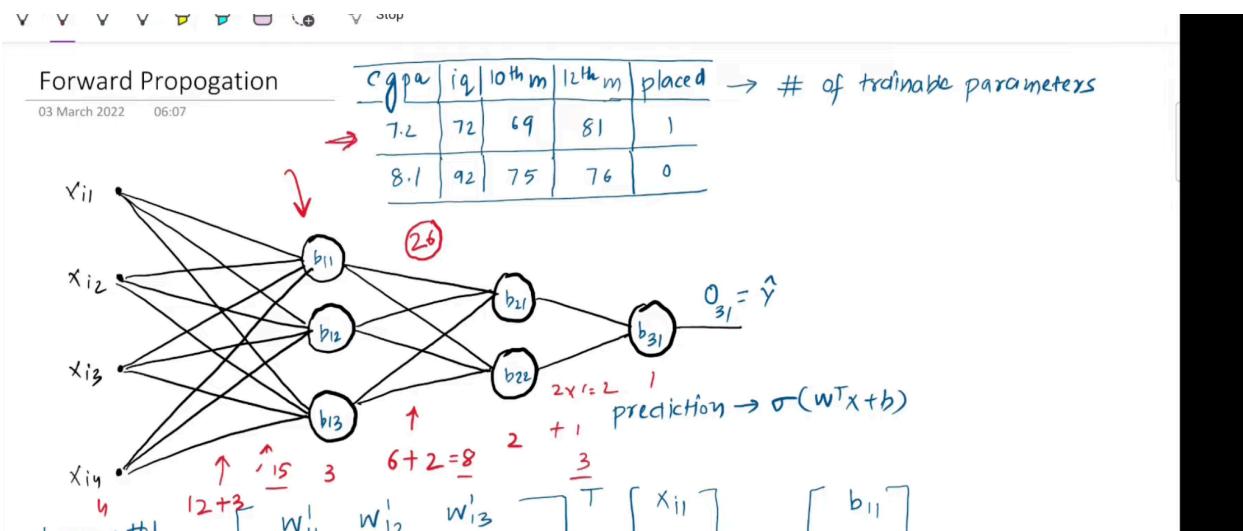
$$\mathcal{L}(\hat{y}, y) = -[y \log \hat{y} + (1 - y) \log(1 - \hat{y})]$$

- **Multi-class** (one-hot  $y$ , softmax  $\hat{\mathbf{y}}$ ):

$$\mathcal{L}(\hat{\mathbf{y}}, \mathbf{y}) = - \sum_k y_k \log \hat{y}_k$$

- **Regression:**  $\frac{1}{2} \|\hat{\mathbf{y}} - \mathbf{y}\|^2$  or MAE/Huber.

## Forward Propagation



Layer #1

$$\begin{aligned}
 & \text{Inputs: } x_{11}, x_{12}, x_{13}, x_{14} \\
 & \text{Weights: } w_{11}, w_{12}, w_{13}, w_{21}, w_{22}, w_{23}, w_{31}, w_{32}, w_{33} \\
 & \text{Biases: } b_{11}, b_{12}, b_{13} \\
 & \text{Dimensions: } 4 \times 3 \xrightarrow{T} 3 \times 4 \quad 4 \times 1 \quad 3 \times 1
 \end{aligned}$$

$$\begin{aligned}
 & = \begin{bmatrix} w_{11}' x_{11} + w_{21}' x_{12} + w_{31}' x_{13} + w_{41}' x_{14} \\ w_{12}' x_{11} + w_{22}' x_{12} + w_{32}' x_{13} + w_{42}' x_{14} \\ w_{13}' x_{11} + w_{23}' x_{12} + w_{33}' x_{13} + w_{43}' x_{14} \end{bmatrix} + \begin{bmatrix} b_{11} \\ b_{12} \\ b_{13} \end{bmatrix}
 \end{aligned}$$

$$= \begin{bmatrix} O_{11} \\ O_{12} \\ O_{13} \end{bmatrix}$$

Layer #2

$$\begin{bmatrix} w_{11}^2 & w_{12}^2 \\ w_{21}^2 & w_{22}^2 \\ w_{31}^2 & w_{32}^2 \end{bmatrix}^T \begin{bmatrix} O_{11} \\ O_{12} \\ O_{13} \end{bmatrix} + \begin{bmatrix} b_{21} \\ b_{22} \end{bmatrix}_{2 \times 1}$$

3x2  $\xrightarrow{T}$  2x3      3x1

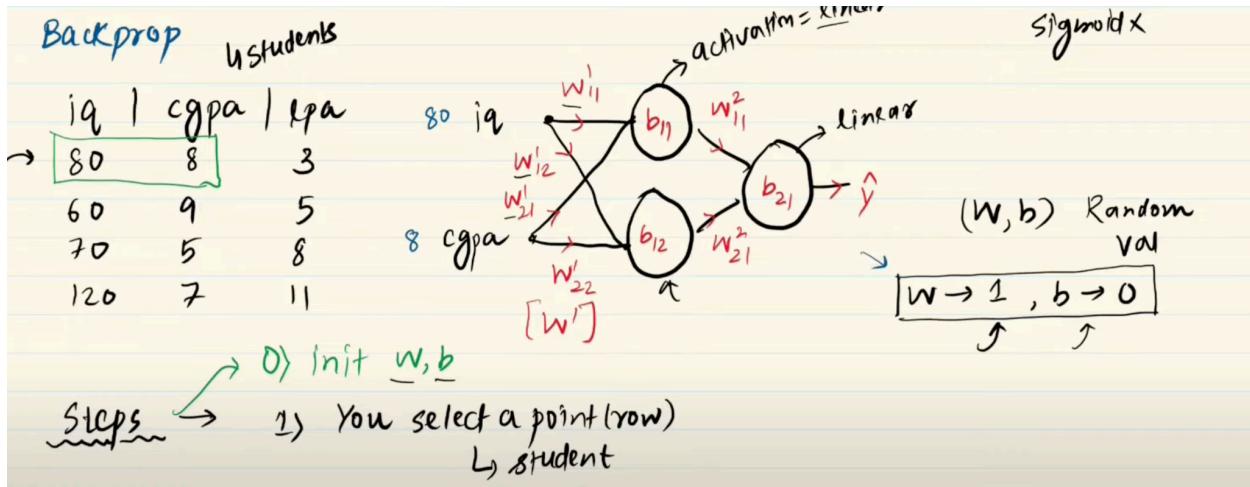
Layer #3

$$\begin{bmatrix} w_{11}^3 \\ w_{21}^3 \end{bmatrix}^T \begin{bmatrix} O_{21} \\ O_{22} \end{bmatrix} + \begin{bmatrix} b_{31} \end{bmatrix}_{1 \times 1}$$

2x1  $\xrightarrow{T}$  1x2      2x1  
1x1

$$= \sigma \left( \begin{bmatrix} w_{11}^3 O_{21} + w_{21}^3 O_{22} + b_{31} \end{bmatrix} \right) = \hat{y}_i$$

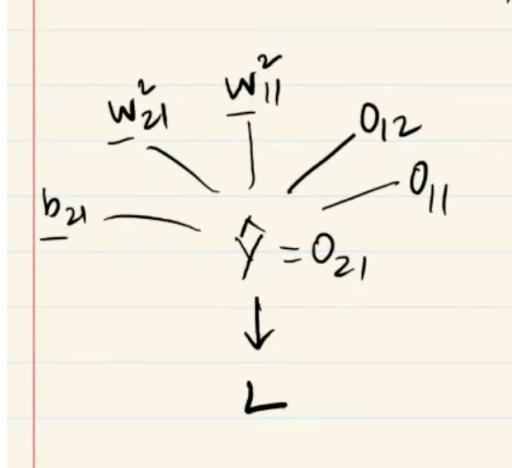
**Backpropagation**

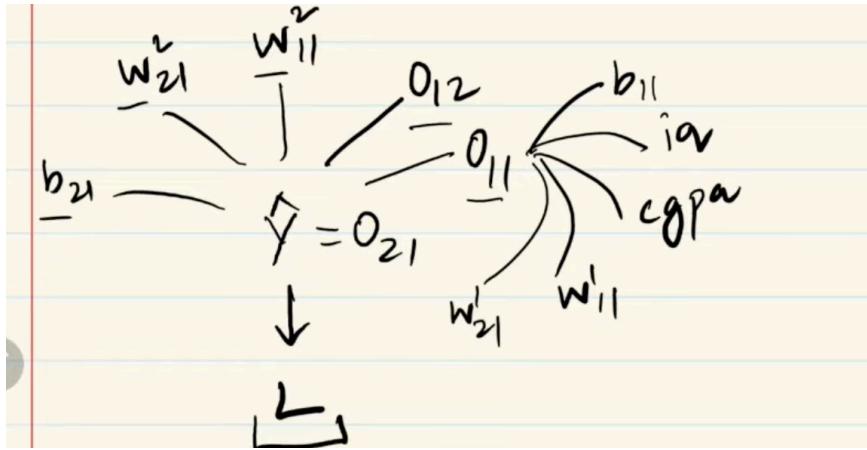


2) Predict (lpa)  $\rightarrow$  forward prop [dot product]

3) Choose a loss function

$$O_{21} = W_{11}^2 O_{11} + W_{21}^2 O_{12} + b_{21}$$





4) Weights and bias update  $\curvearrowleft$

$\curvearrowleft$  Gradient descent

$$w_{\text{new}} = w_{\text{old}} - \eta \frac{\partial L}{\partial w_{\text{old}}}$$

$$b_{\text{new}} = b_{\text{old}} - \eta \frac{\partial L}{\partial b_{\text{old}}}$$

$$w_{11}^2_{\text{new}} = w_{11}^2_{\text{old}} - \eta \frac{\partial L}{\partial w_{11}^2}$$

$$b_{\text{new}} = b_{\text{old}} - \eta \frac{\partial L}{\partial b_{\text{old}}}$$

$$w_{21}^2_{\text{new}} = w_{21}^2_{\text{old}} - \eta \frac{\partial L}{\partial w_{21}^2}$$

$$b_{21}^2_{\text{new}} = b_{21}^2_{\text{old}} - \eta \frac{\partial L}{\partial b_{21}}$$

$$\frac{\partial L}{\partial w_{11}^2} = \boxed{\frac{\partial L}{\partial \hat{y}}} \times \boxed{\frac{\partial \hat{y}}{\partial w_{11}^2}} \rightarrow \text{Chain of differentials}$$

$$\boxed{\frac{\partial L}{\partial \hat{y}} = \frac{\partial}{\partial \hat{y}} (y - \hat{y})^2 = [-2(y - \hat{y})]}$$

$$\frac{\partial \hat{y}}{\partial w_{11}^2} = \frac{\partial}{\partial w_{11}^2} [D_{11} w_{11}^2 + D_{12} w_{21}^2 + b_{21}] = D_{11}$$

$$\boxed{\frac{\partial L}{\partial w_{11}^2} = -2(y - \hat{y})D_{11}}$$

Note

The **Universal Approximation Theorem (UAT)** is one of the most important results in neural network theory. It formally states that a feedforward neural network with at least one hidden layer, containing a finite number of neurons, can approximate any continuous function to any desired degree of accuracy—provided the activation function satisfies certain conditions (nonlinear, bounded, and continuous, like sigmoid, tanh, or even ReLU in some formulations).

---

## Intuition

- Imagine you want a neural network to "learn" any curve or surface.
- The theorem says: **as long as you allow enough hidden neurons, the network can approximate the function arbitrarily well.**
- It doesn't say how many neurons are needed, nor how to train the network efficiently—it just guarantees that such a network exists.

