

ENGR 3321: Introduction to Deep Learning for Robotics

Neural Network NNN:
Multi-Layer Perceptron Model

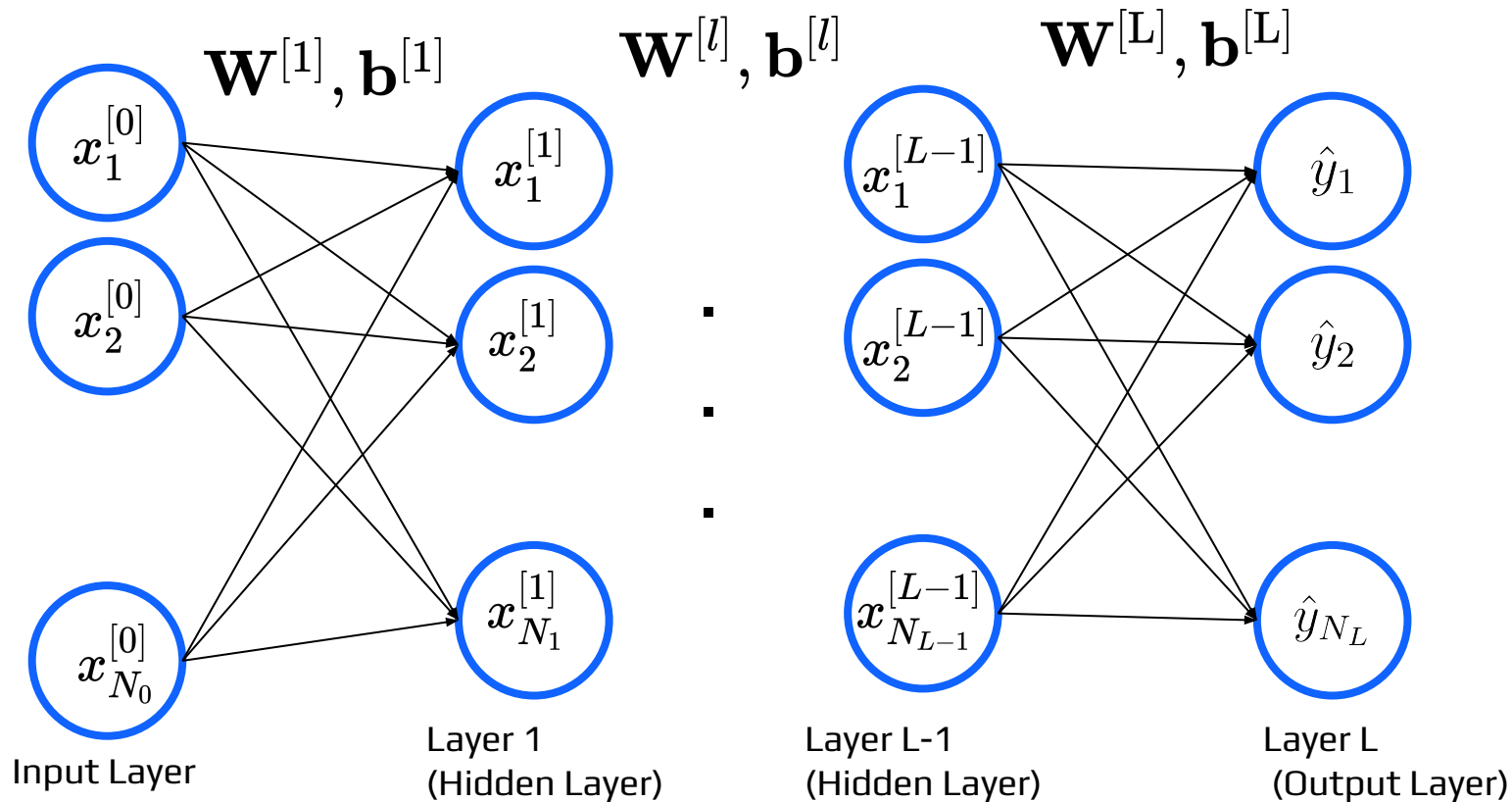
10/21/2024



Outline

- Representations
- Multi-Class Classification

Multi-Layer Perceptron Model



Individual Representation

$$x_n^{[l]} = a(w_{1n}^{[l]}x_1^{[l-1]} + w_{2n}^{[l]}x_2^{[l-1]} + \dots + w_{N_{l-1}n}^{[l]}x_{N_{l-1}}^{[l-1]} + b_n^{[l]})$$

Matrix Form

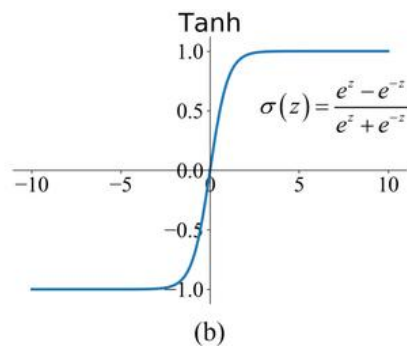
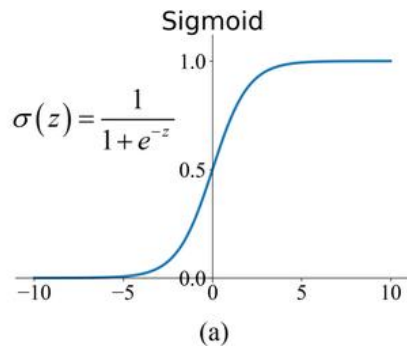
$$\mathbf{X}^{[l]} = a(\mathbf{Z}^{[l]}) = a(\mathbf{X}^{[l-1]} \cdot \mathbf{W}^{[l]T} + \mathbf{b}^{[l]})$$

(M, N_l) (M, N_l) (M, N_{l-1}) (N_{l-1}, N_l) $(1, N_l)$

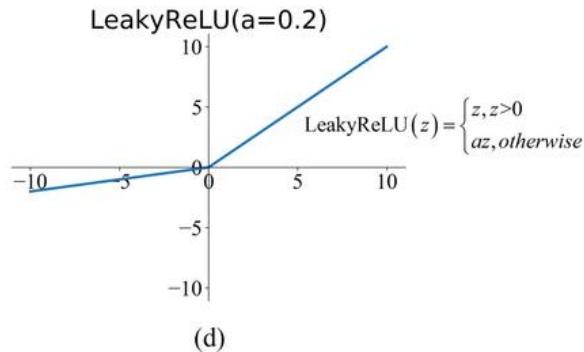
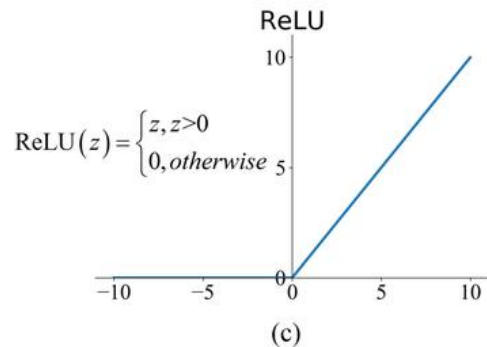
$a(\cdot)$ activation function

Activation Functions

$$\sigma'(z) = \sigma(z)(1 - \sigma(z))$$



$$\sigma'(z) = 1 - \sigma^2(z)$$



$$\text{ReLU}'(z) = \begin{cases} 1, & z > 0 \\ 0, & \text{otherwise} \end{cases}$$

$$\text{LeakyReLU}'(z) = \begin{cases} 1, & z > 0 \\ a, & \text{otherwise} \end{cases}$$

Feature (Input) Matrix

$$\mathbf{X}^{[0]} = \begin{bmatrix} (1)x_1^{[0]} & (1)x_2^{[0]} & \dots & (1)x_{N_0}^{[0]} \\ (2)x_1^{[0]} & (2)x_2^{[0]} & \dots & (2)x_{N_0}^{[0]} \\ \dots & \dots & \dots & \dots \\ (M)x_1^{[0]} & (1)x_2^{[0]} & \dots & (M)x_{N_0}^{[0]} \end{bmatrix} (M, N_0)$$

Trainable Parameters

$$\mathbf{W}^{[l]} = \begin{bmatrix} w_{11}^{[l]} & w_{21}^{[l]} & \dots & w_{N_{l-1}1}^{[l]} \\ w_{12}^{[l]} & w_{22}^{[l]} & \dots & w_{N_{l-1}2}^{[l]} \\ & & \dots & \\ w_{1N_l}^{[l]} & w_{2N_l}^{[l]} & \dots & w_{N_{l-1}N_l}^{[l]} \end{bmatrix} (N_l, N_{l-1})$$

$$\mathbf{b}^{[l]} = \begin{bmatrix} b_1^{[l]} & b_2^{[l]} & \dots & b_{N_l}^{[l]} \end{bmatrix} (1, N_l)$$

Forward Propagation

$$\mathbf{z}^{[l]} = \mathbf{X}^{[l-1]} \cdot \mathbf{W}^{[l]T} + \mathbf{b}^{[l]}$$

$$\mathbf{z}^{[l]} = \begin{bmatrix} (1) \mathbf{x}_1^{[l-1]} & (1) \mathbf{x}_2^{[l-1]} & \dots & (1) \mathbf{x}_{N_{l-1}}^{[l-1]} \\ (2) \mathbf{x}_1^{[l-1]} & (2) \mathbf{x}_2^{[l-1]} & \dots & (2) \mathbf{x}_{N_{l-1}}^{[l-1]} \\ \dots & \dots & \dots & \dots \\ (M) \mathbf{x}_1^{[l-1]} & (M) \mathbf{x}_2^{[l-1]} & \dots & (M) \mathbf{x}_{N_{l-1}}^{[l-1]} \end{bmatrix} \cdot \begin{bmatrix} w_{11}^{[l]} & w_{12}^{[l]} & \dots & w_{1N_l}^{[l]} \\ w_{21}^{[l]} & w_{22}^{[l]} & \dots & w_{2N_l}^{[l]} \\ \dots & \dots & \dots & \dots \\ w_{N_{l-1}1}^{[l]} & w_{N_{l-1}2}^{[l]} & \dots & w_{N_{l-1}N_l}^{[l]} \end{bmatrix} + \begin{bmatrix} b_1^{[l]} & b_2^{[l]} & \dots & b_{N_l}^{[l]} \\ b_1^{[l]} & b_2^{[l]} & \dots & b_{N_l}^{[l]} \\ \dots & \dots & \dots & \dots \\ b_1^{[l]} & b_2^{[l]} & \dots & b_{N_l}^{[l]} \end{bmatrix}$$

$$\mathbf{X}^{[l]} = a(\mathbf{z}^{[l]})$$

Special Case:

$$\hat{\mathbf{Y}} = a(\mathbf{X}^{[L-1]} \mathbf{W}^{[L]T} + \mathbf{b}^{[L]}) = a(\mathbf{z}^{[L]}) = \mathbf{X}^{[L]}$$

Prediction (output) Matrix

$$\hat{\mathbf{Y}} = \begin{bmatrix} {}^{(1)}y_1 & {}^{(1)}y_2 & \dots & {}^{(1)}y_{N_L} \\ {}^{(2)}y_1 & {}^{(2)}y_2 & \dots & {}^{(2)}y_{N_L} \\ \dots & \dots & \dots & \dots \\ {}^{(M)}y_1 & {}^{(M)}y_2 & \dots & {}^{(M)}y_{N_L} \end{bmatrix}_{(M, N_L)}$$

Gradient Descent Optimization

Given dataset: $\left\{ \left({}^{(1)}\mathbf{x}, {}^{(1)}\mathbf{y} \right), \left({}^{(2)}\mathbf{x}, {}^{(2)}\mathbf{y} \right), \dots, \left({}^{(M)}\mathbf{x}, {}^{(M)}\mathbf{y} \right) \right\}$

Initialize $\mathbf{W}^{[l]}$, $\mathbf{b}^{[l]}$

Repeat until converge {

 compute $\mathcal{L}(\hat{\mathbf{Y}}, \mathbf{Y})$

 compute $\nabla \mathcal{L}$

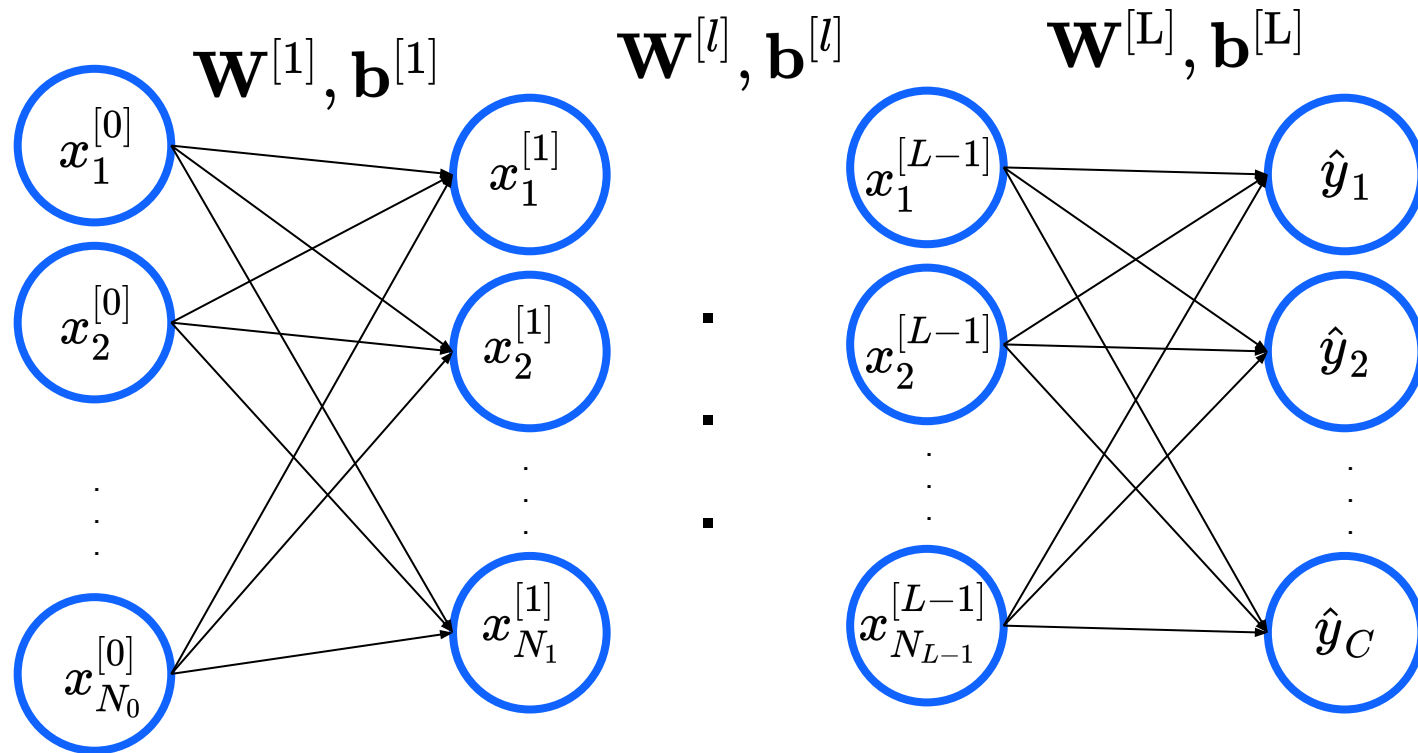
$\mathbf{W}^{[l]} := \mathbf{W}^{[l]} - \alpha \cdot d\mathbf{W}^{[l]}$

$\mathbf{b}^{[l]} := \mathbf{b}^{[l]} - \alpha \cdot d\mathbf{b}^{[l]}$

}

where α is learning rate

Multi-Class Classification

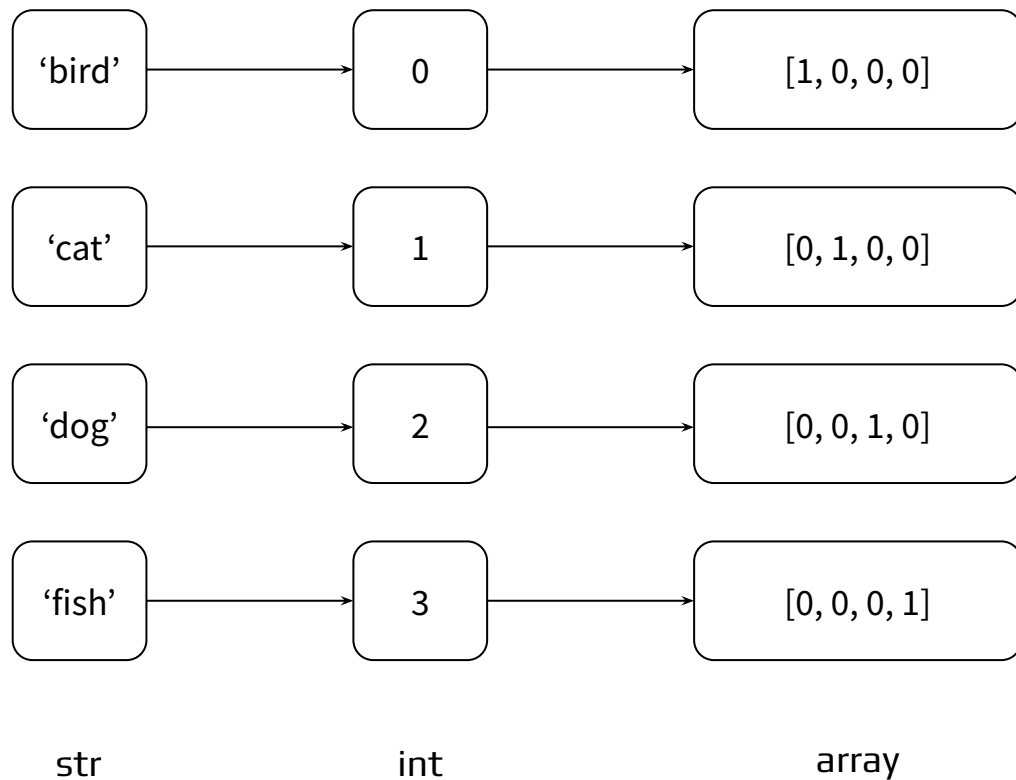


Multi-Class Classification

$$\mathbf{y} = \begin{bmatrix} {}^{(1)}y \\ {}^{(2)}y \\ \cdot \\ \cdot \\ \cdot \\ {}^{(M)}y \end{bmatrix} (M,1)$$

$$\hat{\mathbf{y}} = \begin{bmatrix} {}^{(1)}\hat{y} \\ {}^{(2)}\hat{y} \\ \cdot \\ \cdot \\ \cdot \\ {}^{(M)}\hat{y} \end{bmatrix} (M,1)$$

One-Hot Encoding on Targets



Softmax Activation on Predictions

$$\hat{y}_c = \frac{e^{z_c^{[L]}}}{\sum_{c=1}^C e^{z_c^{[L]}}}, \forall c = 1, \dots, C$$

$$\sum \left[{}^{(m)}\hat{y}_1 \quad {}^{(m)}\hat{y}_2 \quad \dots \quad {}^{(m)}\hat{y}_C \right] = 1$$

Probability of the m -th sample being predicted as a member in class 1

Multi-Class Classification

$$\mathbf{Y} = \begin{bmatrix} {}^{(1)}y_1 & {}^{(1)}y_2 & \dots & {}^{(1)}y_C \\ {}^{(2)}y_1 & {}^{(2)}y_2 & \dots & {}^{(2)}y_C \\ \dots & \dots & \dots & \dots \\ {}^{(M)}y_1 & {}^{(M)}y_2 & \dots & {}^{(M)}y_C \end{bmatrix}_{(M,C)}$$

$$\hat{\mathbf{Y}} = \begin{bmatrix} {}^{(1)}\hat{y}_1 & {}^{(1)}\hat{y}_2 & \dots & {}^{(1)}\hat{y}_C \\ {}^{(2)}\hat{y}_1 & {}^{(2)}\hat{y}_2 & \dots & {}^{(2)}\hat{y}_C \\ \dots & \dots & \dots & \dots \\ {}^{(M)}\hat{y}_1 & {}^{(M)}\hat{y}_2 & \dots & {}^{(M)}\hat{y}_C \end{bmatrix}_{(M,C)}$$

Multi-Class Cross Entropy Loss

$$\mathcal{L}(\hat{\mathbf{Y}}, \mathbf{Y}) = \frac{1}{M} \sum_{m=1}^M \left[\sum_{c=1}^C \left(-^{(m)}y_c \ln^{(m)} \hat{y}_c \right) \right]$$

Back-Propagation

$$\nabla \mathcal{L} = \left[\cdots \quad \frac{\partial \mathcal{L}}{\partial w_{l-1,l}^{[l]}} \quad \cdots \quad \frac{\partial \mathcal{L}}{\partial b_l^{[l]}} \quad \cdots \right]$$

$$d\mathbf{Z}^{[L]} = \frac{\partial \mathcal{L}}{\partial \hat{\mathbf{Y}}} \cdot \frac{\partial \hat{\mathbf{Y}}}{\partial \mathbf{Z}^{[L]}} = \hat{\mathbf{Y}} - \mathbf{Y}$$

NOTE: Only valid if

1. last layer is softmax activated;
2. Prediction is evaluated by cross entropy loss

For l from L to 1

$$d\mathbf{W}^{[l]} = d\mathbf{Z}^{[l]} \cdot \frac{\partial \mathbf{Z}^{[l]}}{\partial \mathbf{W}^{[l]}} = d\mathbf{Z}^{[l]T} \cdot \mathbf{X}^{[l-1]}$$

$$d\mathbf{b}^{[l]} = d\mathbf{Z}^{[l]} \cdot \frac{\partial \mathbf{Z}^{[l]}}{\partial \mathbf{b}^{[l]}} = \text{mean}(d\mathbf{Z}^{[l]}, \text{axis}=0, \text{keepdims}=\text{True})$$

$$d\mathbf{X}^{[l-1]} = d\mathbf{Z}^{[l]} \cdot \frac{\partial \mathbf{Z}^{[l]}}{\partial \mathbf{X}^{[l-1]}} = d\mathbf{Z}^{[l]} \cdot \mathbf{W}^{[l]}$$

$$d\mathbf{Z}^{[l-1]} = d\mathbf{X}^{[l-1]} * a'(\mathbf{Z}^{[l-1]})$$