

Introduction to Deep Learning

2- Shallow Versus Deep Neural Networks

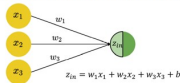
Prof. Monir EL ANNAS

Shallow Neural Network

The Perceptron

- The basic computational unit for neural networks is the perceptron. It is a weighted sum of input values plus bias term, transformed by a non-linear activation function, resulting in an output value (single neuron).

- ❶ **Affine Transformation:** weighted sum of inputs plus bias.



- ❷ **Non-linear Activation:** a non-linear transformation applied to the weighted sum.



- A suitable choice of the activation function τ leads to known functions $f(x)$.
- The identity function gives us the simple linear regression:

$$y = \tau(\mathbf{w}^\top \mathbf{x}) = \mathbf{w}^\top \mathbf{x}$$

- The logistic function (sigmoid function) gives us the logistic regression:

$$y = \tau(\mathbf{w}^\top \mathbf{x}) = \frac{1}{1 + \exp(-\mathbf{w}^\top \mathbf{x})}$$

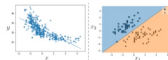
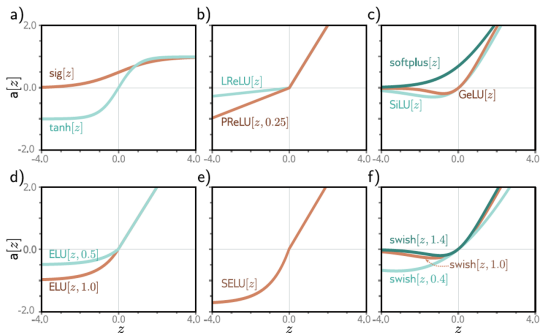


Figure: Left: A regression line learned by a single neuron. Right: A decision-boundary learned by a single neuron in a binary classification task.

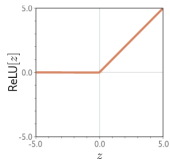
Shallow Neural Network

Activation Functions



$$a[z] = \text{ReLU}[z] = \begin{cases} 0 & z < 0 \\ z & z \geq 0 \end{cases}.$$

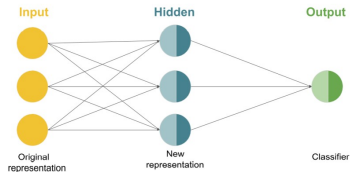
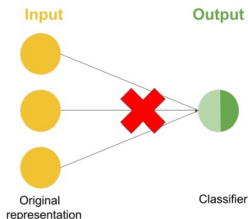
Rectified Linear Unit
(particular kind of activation function)



Shallow Neural Network

Shallow Neural Network

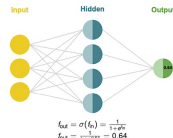
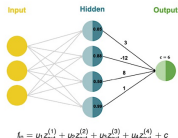
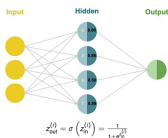
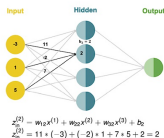
- Instead of a single neuron, we use more complex networks.
- Fully-connected neural networks describe the simplest neural network architecture where all neurons from one layer are connected to the neurons of the next layer.
- Fully-connected neural networks with a single layer are known as shallow neural networks.



Shallow Neural Network

Forward Pass: Binary Classification Example

- Following the computation from left to right is called a **forward pass**. It is how Neural Networks make their predictions.
- Each neuron in the hidden layer performs an affine transformation on the inputs.
- Each hidden neuron performs a non-linear activation transformation on the weighted sum.
- The output neuron performs an affine transformation on its inputs.
- The output neuron performs a non-linear activation transformation on the weighted sum.

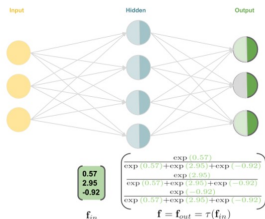


- Binary Classification:** one neuron in the output layer.
- Activation function:** sigmoid activation function used in the output layer.

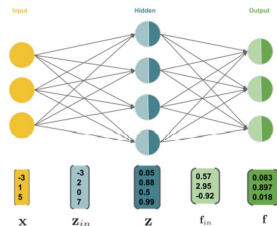
Shallow Neural Network

Forward Pass: Multi-Class Classification Example

Forward pass (Hidden: Sigmoid, Output: Softmax).



Forward pass (Hidden: Sigmoid, Output: Softmax).

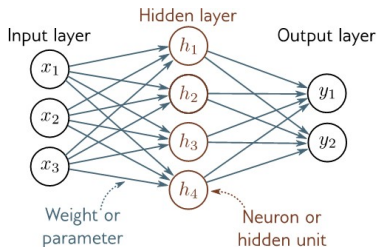


- **Multi-Class Classification:** Add additional neurons to the output layer. Each neuron will represent a specific class.
- **Activation function:** softmax activation function used in the output layer $f_{out,k} = \tau(f_{in,k}) =$

$$\frac{\exp(f_{in,k})}{\sum_{k'=1}^G \exp(f_{in,k'})}$$

Shallow Neural Network

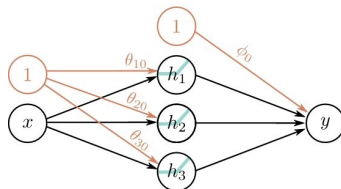
Formulation: Terminology



- Y-offsets = biases
- Slopes = weights
- Everything in one layer connected to everything in the next = fully connected network
- No loops = feedforward network
- Values after activation functions = activations
- Values before activation functions = pre-activations
- One hidden layer = shallow neural network
- More than one hidden layer = deep neural network
- Number of hidden units = capacity

Shallow Neural Network

Formulation: 1 input 1 output



- $y = f[\mathbf{x}, \phi] = \phi_0 + \phi_1 a[\theta_{10} + \theta_{11}\mathbf{x}] + \phi_2 a[\theta_{20} + \theta_{21}\mathbf{x}] + \phi_3 a[\theta_{30} + \theta_{31}\mathbf{x}]$
- Break down into two parts:

$$y = \phi_0 + \phi_1 h_1 + \phi_2 h_2 + \phi_3 h_3$$

- where we refer to h_1 , h_2 , and h_3 as hidden units :

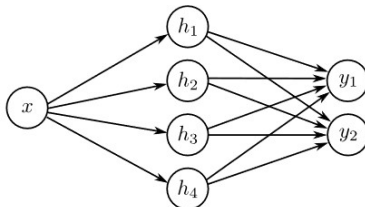
$$h_1 = a[\theta_{10} + \theta_{11}\mathbf{x}]$$

$$h_2 = a[\theta_{20} + \theta_{21}\mathbf{x}]$$

$$h_3 = a[\theta_{30} + \theta_{31}\mathbf{x}]$$

Shallow Neural Network

Formulation: More than one output



- Network with one input, four hidden units, and two outputs

- $y_1 = \phi_{10} + \phi_{11}h_1 + \phi_{12}h_2 + \phi_{13}h_3 + \phi_{14}h_4$

- $y_2 = \phi_{20} + \phi_{21}h_1 + \phi_{22}h_2 + \phi_{23}h_3 + \phi_{24}h_4.$

$$h_1 = a [\theta_{10} + \theta_{11}x]$$

$$h_2 = a [\theta_{20} + \theta_{21}x]$$

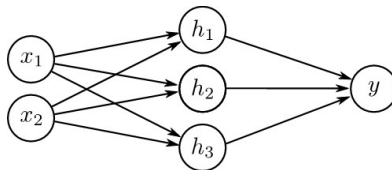
- where :

$$h_3 = a [\theta_{30} + \theta_{31}x]$$

$$h_4 = a [\theta_{40} + \theta_{41}x],$$

Shallow Neural Network

Formulation: More than one input



- Neural network with 2D multivariate input $x = [x_1, x_2]^T$ and scalar output y .
- $y = \phi_0 + \phi_1 h_1 + \phi_2 h_2 + \phi_3 h_3$.

$$h_1 = a [\theta_{10} + \theta_{11}x_1 + \theta_{12}x_2]$$

- where : $h_2 = a [\theta_{20} + \theta_{21}x_1 + \theta_{22}x_2]$

$$h_3 = a [\theta_{30} + \theta_{31}x_1 + \theta_{32}x_2],$$

Shallow Neural Network

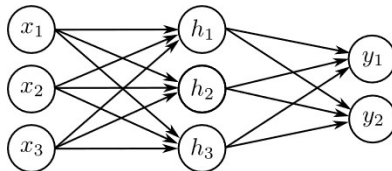
Formulation: General case arbitrary inputs hidden units outputs

- D_o Outputs, D hidden units, and D_i inputs

$$h_d = a \left(\theta_{d0} + \sum_{i=1}^{D_i} \theta_{di} x_i \right)$$

$$y_j = \phi_{j0} + \sum_{d=1}^D \phi_{jd} h_d$$

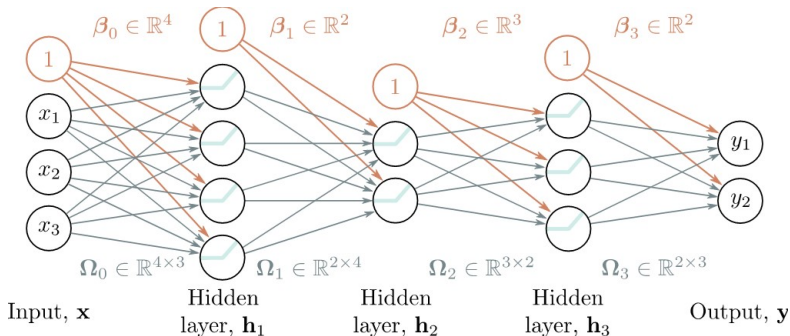
- e.g., Three inputs, three hidden units, two outputs



Deep neural networks

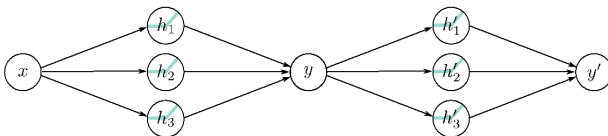
Deep neural networks

- Networks with more than one hidden layer
- Intuition becomes more difficult



Deep neural networks

Formulation: Deep network with two hidden layers



- Networks 1:

$$h_1 = a(\theta_{10} + \theta_{11}x)$$

$$h_2 = a(\theta_{20} + \theta_{21}x)$$

$$h_3 = a(\theta_{30} + \theta_{31}x)$$

$$y = \phi_0 + \phi_1 h_1 + \phi_2 h_2 + \phi_3 h_3$$

- Networks 2:

$$h'_1 = a(\theta'_{10} + \theta'_{11}y)$$

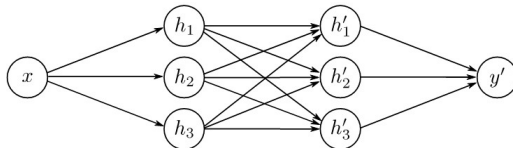
$$h'_2 = a(\theta'_{20} + \theta'_{21}y)$$

$$h'_3 = a(\theta'_{30} + \theta'_{31}y)$$

$$y' = \phi'_0 + \phi'_1 h'_1 + \phi'_2 h'_2 + \phi'_3 h'_3$$

Deep neural networks

Formulation: Deep network with two hidden layers



- Substituting the expression for y gives:

$$h'_1 = a [\theta'_{10} + \theta'_{11}y] = a [\theta'_{10} + \theta'_{11}\phi_0 + \theta'_{11}\phi_1h_1 + \theta'_{11}\phi_2h_2 + \theta'_{11}\phi_3h_3]$$

$$h'_2 = a [\theta'_{20} + \theta'_{21}y] = a [\theta'_{20} + \theta'_{21}\phi_0 + \theta'_{21}\phi_1h_1 + \theta'_{21}\phi_2h_2 + \theta'_{21}\phi_3h_3]$$

$$h'_3 = a [\theta'_{30} + \theta'_{31}y] = a [\theta'_{30} + \theta'_{31}\phi_0 + \theta'_{31}\phi_1h_1 + \theta'_{31}\phi_2h_2 + \theta'_{31}\phi_3h_3]$$

- Which we can rewrite as:

$$h'_1 = a [\psi_{10} + \psi_{11}h_1 + \psi_{12}h_2 + \psi_{13}h_3]$$

$$h'_2 = a [\psi_{20} + \psi_{21}h_1 + \psi_{22}h_2 + \psi_{23}h_3]$$

$$h'_3 = a [\psi_{30} + \psi_{31}h_1 + \psi_{32}h_2 + \psi_{33}h_3]$$

- We can combine the equations to get one expression:

$$\begin{aligned} y' = & \phi'_0 + \phi'_1a [\theta_{10} + \theta_{11}x] + \psi_{12}a [\theta_{20} + \theta_{21}x] + \psi_{13}a [\theta_{30} + \theta_{31}x] \\ & + \phi'_2a [\psi_{20} + \psi_{21}a [\theta_{10} + \theta_{11}x] + \psi_{22}a [\theta_{20} + \theta_{21}x] + \psi_{23}a [\theta_{30} + \theta_{31}x]] \\ & + \phi'_3a [\psi_{30} + \psi_{31}a [\theta_{10} + \theta_{11}x] + \psi_{32}a [\theta_{20} + \theta_{21}x] + \psi_{33}a [\theta_{30} + \theta_{31}x]] \end{aligned}$$

Deep neural networks

Formulation: Hyperparameters

- K layers = depth of network
- D_k hidden units per layer = width of network
- These are called hyperparameters – chosen before training the network.
- Can try retraining with different hyperparameters – hyperparameter optimization or hyperparameter search.

Deep neural networks

Formulation: The compact form

- **Notation:**

$$h_1 = a(\theta_{10} + \theta_{11}x)$$

$$h_2 = a(\theta_{20} + \theta_{21}x)$$

$$h_3 = a(\theta_{30} + \theta_{31}x)$$

$$h'_1 = a(\psi_{10} + \psi_{11}h_1 + \psi_{12}h_2 + \psi_{13}h_3)$$

$$h'_2 = a(\psi_{20} + \psi_{21}h_1 + \psi_{22}h_2 + \psi_{23}h_3)$$

$$h'_3 = a(\psi_{30} + \psi_{31}h_1 + \psi_{32}h_2 + \psi_{33}h_3)$$

$$y' = \phi'_0 + \phi'_1 h'_1 + \phi'_2 h'_2 + \phi'_3 h'_3$$

- **The compact form:**

$$\begin{bmatrix} h_1 \\ h_2 \\ h_3 \end{bmatrix} = \mathbf{a} \left[\begin{bmatrix} \theta_{10} \\ \theta_{20} \\ \theta_{30} \end{bmatrix} + \begin{bmatrix} \theta_{11} \\ \theta_{21} \\ \theta_{31} \end{bmatrix} x \right]$$

$$\begin{bmatrix} h'_1 \\ h'_2 \\ h'_3 \end{bmatrix} = \mathbf{a} \left(\begin{bmatrix} \psi_{10} \\ \psi_{20} \\ \psi_{30} \end{bmatrix} + \begin{bmatrix} \psi_{11} & \psi_{12} & \psi_{13} \\ \psi_{21} & \psi_{22} & \psi_{23} \\ \psi_{31} & \psi_{32} & \psi_{33} \end{bmatrix} \begin{bmatrix} h_1 \\ h_2 \\ h_3 \end{bmatrix} \right)$$

$$y' = \phi'_0 + \begin{bmatrix} \phi'_1 & \phi'_2 & \phi'_3 \end{bmatrix} \begin{bmatrix} h'_1 \\ h'_2 \\ h'_3 \end{bmatrix}$$

Deep neural networks

Formulation: The compact form

• Notation:

$$\begin{bmatrix} h_1 \\ h_2 \\ h_3 \end{bmatrix} = \mathbf{a} \begin{bmatrix} \theta_{10} \\ \theta_{20} \\ \theta_{30} \end{bmatrix} + \begin{bmatrix} \theta_{11} \\ \theta_{21} \\ \theta_{31} \end{bmatrix} \mathbf{x}$$

$$\begin{bmatrix} h'_1 \\ h'_2 \\ h'_3 \end{bmatrix} = \mathbf{a} \left(\begin{bmatrix} \psi_{10} \\ \psi_{20} \\ \psi_{30} \end{bmatrix} + \begin{bmatrix} \psi_{11} & \psi_{12} & \psi_{13} \\ \psi_{21} & \psi_{22} & \psi_{23} \\ \psi_{31} & \psi_{32} & \psi_{33} \end{bmatrix} \begin{bmatrix} h_1 \\ h_2 \\ h_3 \end{bmatrix} \right)$$

$$y' = \phi'_0 + \begin{bmatrix} \phi'_1 & \phi'_2 & \phi'_3 \end{bmatrix} \begin{bmatrix} h'_1 \\ h'_2 \\ h'_3 \end{bmatrix}$$

• The compact form:

$$h = \mathbf{a}(\boldsymbol{\theta}_0 + \boldsymbol{\theta}\mathbf{x})$$

$$\mathbf{h}' = \mathbf{a}(\boldsymbol{\psi}_0 + \boldsymbol{\Psi}\mathbf{h})$$

$$y = \phi'_0 + \boldsymbol{\phi}'\mathbf{h}'$$

Deep neural networks

Formulation: The compact form

- Notation:

$$h = a(\theta_0 + \theta x)$$

$$h' = a(\psi_0 + \psi h)$$

$$y = \phi'_0 + \phi' h'$$

- The compact form:

$$h_1 = a(\beta_0 + \Omega_0 x)$$

$$h_2 = a(\beta_1 + \Omega_1 h_1)$$

$$y = \beta_2 + \Omega_2 h_2$$

- A general deep network $y = f[x, \phi]$ with K layers can now be written as:

$$h_1 = a[\beta_0 + \Omega_0 x]$$

$$h_2 = a[\beta_1 + \Omega_1 h_1]$$

$$h_3 = a[\beta_2 + \Omega_2 h_2]$$

$$\vdots$$

$$h_K = a[\beta_{K-1} + \Omega_{K-1} h_{K-1}]$$

$$y = \beta_K + \Omega_K h_K.$$

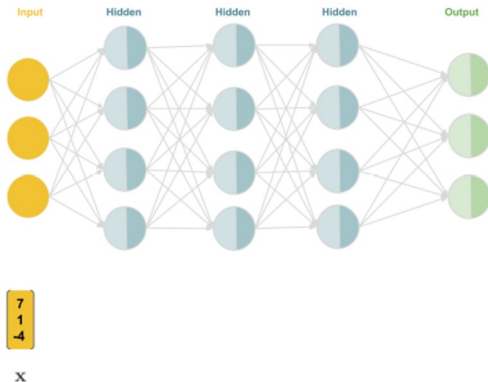
- We can equivalently write the network as a single function:

$$y = \beta_K + \Omega_K a \left[\beta_{K-1} + \Omega_{K-1} a \left[\dots \beta_2 + \Omega_2 a \left[\beta_1 + \Omega_1 a \left[\beta_0 + \Omega_0 x \right] \dots \right] \right] \right].$$

Deep neural networks

Forward Pass: Multi-Class Classification Example

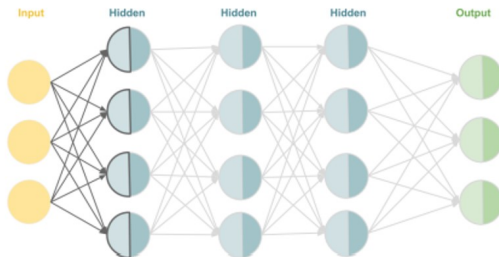
- Forward pass (Hidden: ReLU, Output: Softmax)



Deep neural networks

Forward Pass: Multi-Class Classification Example

- Forward pass (Hidden: ReLU, Output: Softmax)

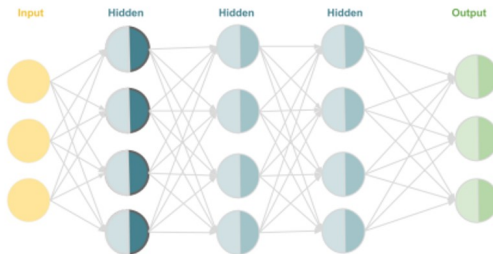


$$\begin{bmatrix} 7 \\ 1 \\ -4 \end{bmatrix} \begin{pmatrix} 7*13 + 1*(-9) + (-4)*2 + 5 \\ 7*(-8) + 1*0 + (-4)*3 + (-2) \\ 7*4 + 1*(-1) + (-4)*5 + 2 \\ 7*(-3) + 1*12 + (-4)*7 + 11 \end{pmatrix}$$
$$\mathbf{x} \quad \mathbf{z}_{in}^{(1)} = \mathbf{W}^{(1)T} \mathbf{x} + \mathbf{b}^{(1)}$$

Deep neural networks

Forward Pass: Multi-Class Classification Example

- Forward pass (Hidden: ReLU, Output: Softmax)



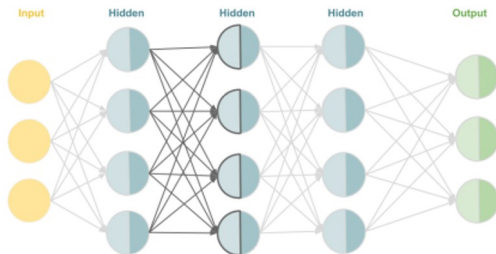
<div style="background-color: yellow; padding: 5px; display: inline-block;">7 1 4</div>	<div style="background-color: lightblue; padding: 5px; display: inline-block;">79 -70 9 -26</div>	<div style="border: 1px solid black; padding: 5px; display: inline-block;">$\begin{pmatrix} \max(0, 79) \\ \max(0, -70) \\ \max(0, 9) \\ \max(0, -26) \end{pmatrix}$</div>
---	---	---

$$\mathbf{x} \quad \mathbf{z}_{in}^{(1)} \quad \mathbf{z}^{(1)} = \mathbf{z}_{out}^{(1)} = \sigma(\mathbf{z}_{in}^{(1)})$$

Deep neural networks

Forward Pass: Multi-Class Classification Example

- Forward pass (Hidden: ReLU, Output: Softmax)



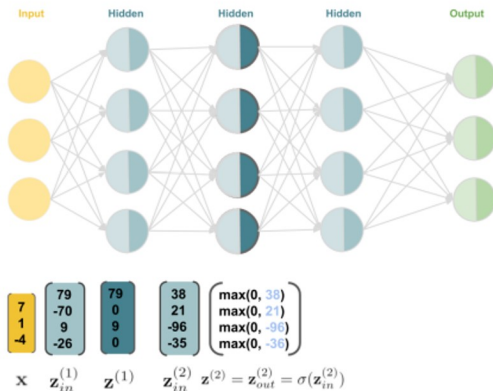
$$\begin{bmatrix} 7 \\ 1 \\ -4 \end{bmatrix} \quad \begin{bmatrix} 79 \\ -70 \\ 9 \\ -26 \end{bmatrix} \quad \begin{bmatrix} 79 \\ 0 \\ 9 \\ 0 \end{bmatrix} \quad \left(\begin{array}{l} 79 \cdot 1 + 0 \cdot 0 + 9 \cdot (-4) + 0 \cdot 1 + (-5) \\ 79 \cdot 0 + 0 \cdot 11 + 9 \cdot 2 + 0 \cdot (-14) + 3 \\ 79 \cdot (-1) + 0 \cdot 5 + 9 \cdot (-2) + 0 \cdot 16 + 1 \\ 79 \cdot 0 + 0 \cdot (-9) + 9 \cdot (-3) + 0 \cdot 4 + (-8) \end{array} \right)$$

$\mathbf{x} \quad \mathbf{z}_{in}^{(1)} \quad \mathbf{z}^{(1)} \quad \mathbf{z}_{in}^{(2)} = \mathbf{W}^{(2)T} \mathbf{z}^{(1)} + \mathbf{b}^{(2)}$

Deep neural networks

Forward Pass: Multi-Class Classification Example

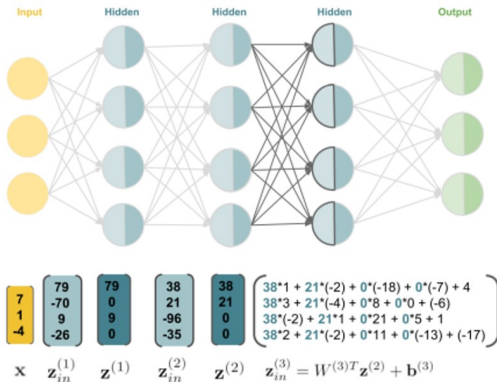
- Forward pass (Hidden: ReLU, Output: Softmax)



Deep neural networks

Forward Pass: Multi-Class Classification Example

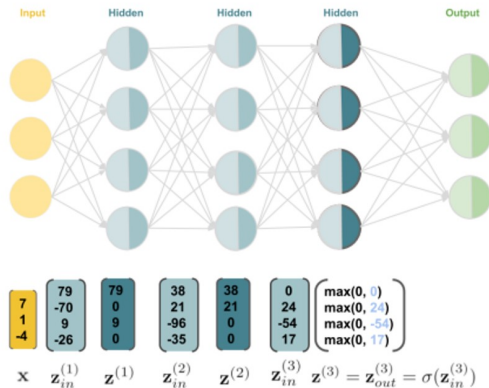
- Forward pass (Hidden: ReLU, Output: Softmax)



Deep neural networks

Forward Pass: Multi-Class Classification Example

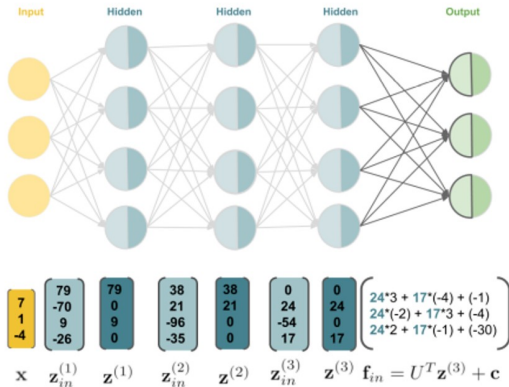
- Forward pass (Hidden: ReLU, Output: Softmax)



Deep neural networks

Forward Pass: Multi-Class Classification Example

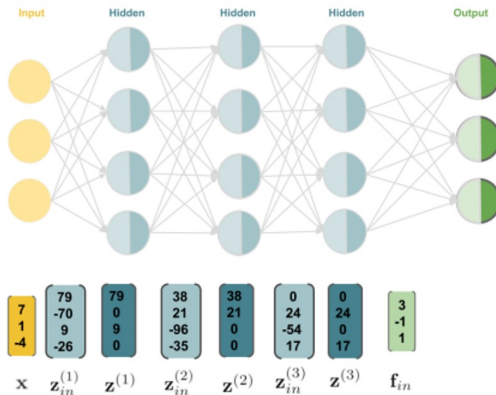
- Forward pass (Hidden: ReLU, Output: Softmax)



Deep neural networks

Forward Pass: Multi-Class Classification Example

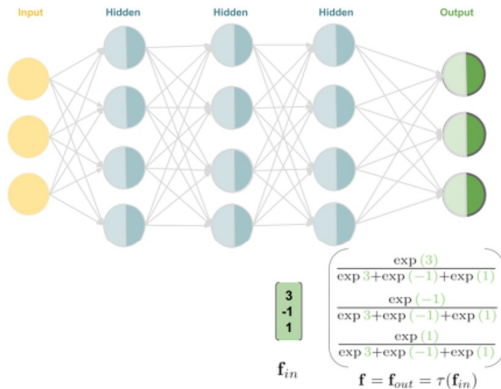
- Forward pass (Hidden: ReLU, Output: Softmax)



Deep neural networks

Forward Pass: Multi-Class Classification Example

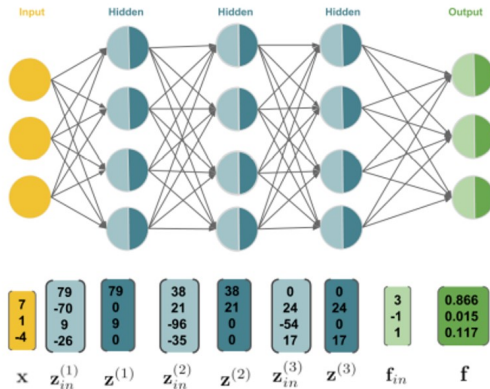
- Forward pass (Hidden: ReLU, Output: Softmax)



Deep neural networks

Forward Pass: Multi-Class Classification Example

- Forward pass (Hidden: ReLU, Output: Softmax)



Where are we going?

- We have defined families of very flexible networks that map multiple inputs to multiple outputs
- Now we need to train them:
 - How to choose loss functions
 - How to find minima of the loss function
 - How to do this in particular for deep networks
- Then we need to test them