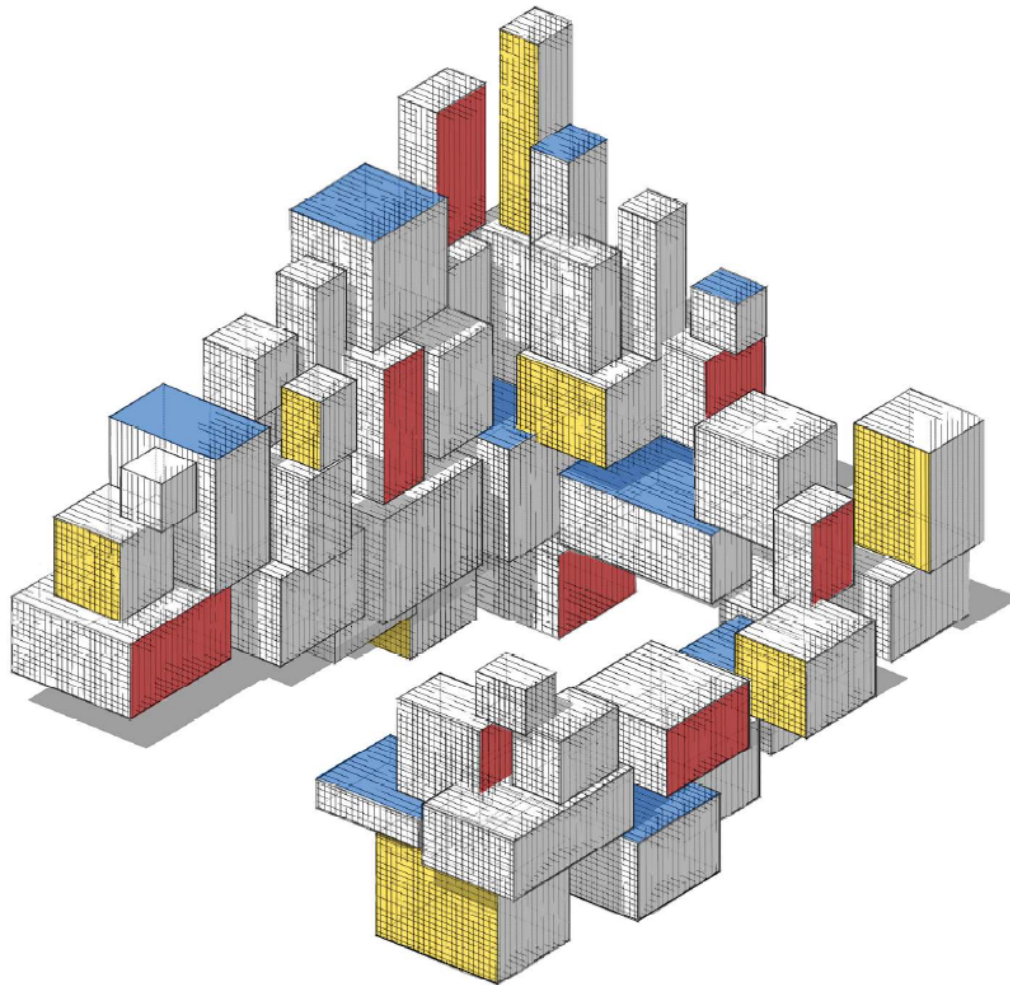


# Deep Learning



# Purpose

In this lecture we introduce the rich class of approximating functions called **neural networks**.

The learners belonging to the neural-network class of functions have attractive properties that have made them ubiquitous in modern machine learning applications — their training is computationally feasible and their complexity is easy to control and fine-tune.

Topics include:

- Activation functions
- Feed-forward neural networks
- Example applications:
  - Nonlinear multi-output regression
  - Multi-logit classification
  - Density estimation

# Supervised Learning

Recall the basic supervised learning task: predict a random output  $Y$  from a random input  $X$ , using a prediction function  $g : \mathbf{x} \mapsto y$  that belongs to a suitably chosen class of approximating functions  $\mathcal{G}$ .

More generally, we may wish to predict a **vector-valued** output  $\mathbf{y}$  using a prediction function  $g : \mathbf{x} \mapsto \mathbf{y}$  from class  $\mathcal{G}$ .

In what follows,  $\mathbf{y}$  denotes the vector-valued output for a given input  $\mathbf{x}$ . This differs from our previous use, where  $y$  denotes a vector of scalar outputs.

The class  $\mathcal{G}$  is sometimes referred to as the **hypothesis space** or the **universe of possible models**, and the **representational capacity** of a hypothesis space  $\mathcal{G}$  is simply its complexity.

# Approximation–Estimation Tradeoff

Suppose that we have a class of functions  $\mathcal{G}_L$ , indexed by a parameter  $L$  such that  $\mathcal{G}_L \subset \mathcal{G}_{L+1} \subset \mathcal{G}_{L+2} \subset \dots$ .

In selecting a suitable class of functions, we have a **approximation–estimation tradeoff**:

- The class  $\mathcal{G}_L$  must be complex (rich) enough to accurately represent the optimal unknown prediction function  $g^*$ , which may require a very large  $L$ .
- The learners in the class  $\mathcal{G}_L$  must be simple enough to train with small estimation error and with minimal demands on computer memory, which may necessitate a small  $L$ .

A class of functions that permits such a natural hierarchical construction is the class of **neural networks**.

# Neural Network

Conceptually, a neural network with  $L$  layers is a nonlinear parametric regression model whose representational capacity is controlled by  $L$ .

Alternatively, we will define the output of a neural network as the repeated composition of linear and (componentwise) nonlinear functions.

This provides a flexible class of nonlinear output functions that can be easily differentiated.

As a result, the training of learners via gradient optimization methods involves mostly standard matrix operations that can be performed very efficiently.

Neural networks were originally intended to mimic the workings of the human brain, with the network nodes modeling neurons and the network links modeling the axons connecting neurons.

# Kolmogorov's Approximation

Many effective machine algorithms have been inspired by mathematical ideas for function approximation.

## Theorem: Kolmogorov (1957)

kolmogorov Every continuous function  $g^* : [0, 1]^p \mapsto \mathbb{R}$  with  $p \geq 2$  can be written as

$$g^*(\mathbf{x}) = \sum_{j=1}^{2p+1} h_j \left( \sum_{i=1}^p h_{ij}(x_i) \right),$$

Handwritten annotations in blue ink:

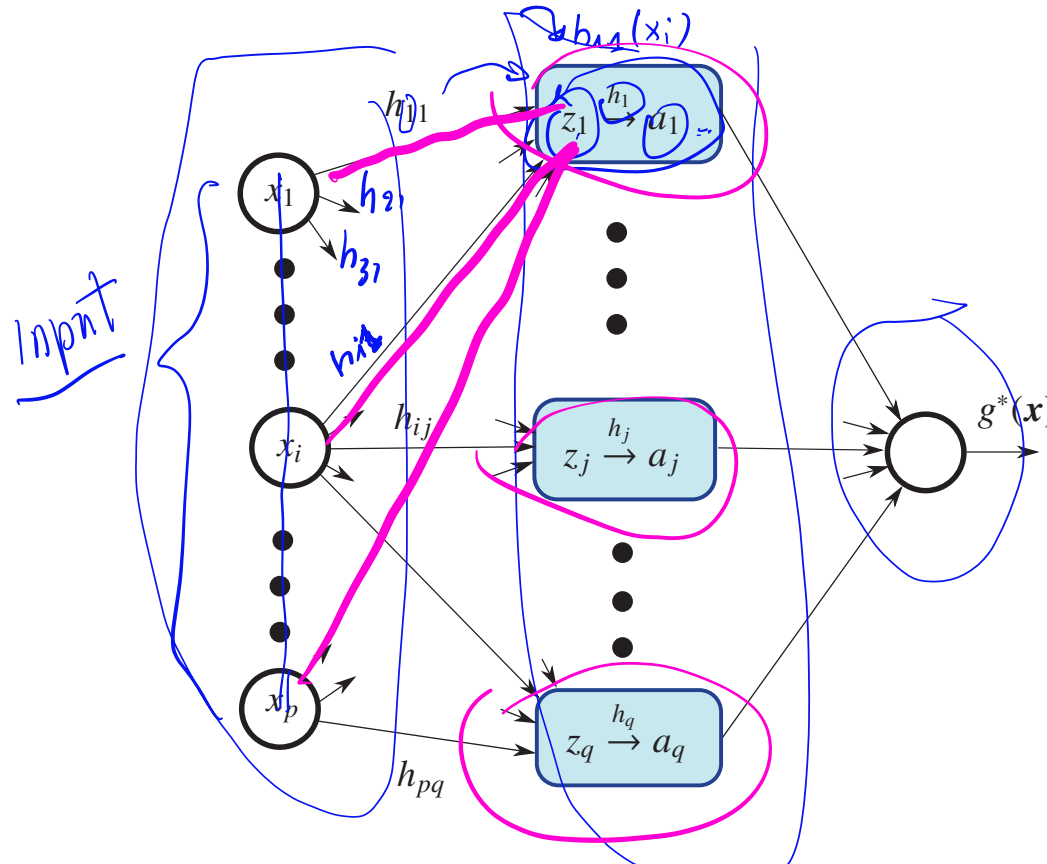
- Above the inner sum:  $(\sin x_1 + x_2^3)^2 + \sqrt{x_1^3 - x_2^2}$
- To the right:  $\sin x_1 + x_2^3 = y_1$
- Below that:  $h_1(y_1) = y_1^2$
- Below the inner sum:  $y_j$

where  $\{h_j, h_{ij}\}$  is a set of univariate continuous functions that depend on  $g^*$ .

This result tells us that any continuous high-dimensional map can be represented as the function composition of much simpler (1D) maps.

# Representation of Kolmogorov's Approximation

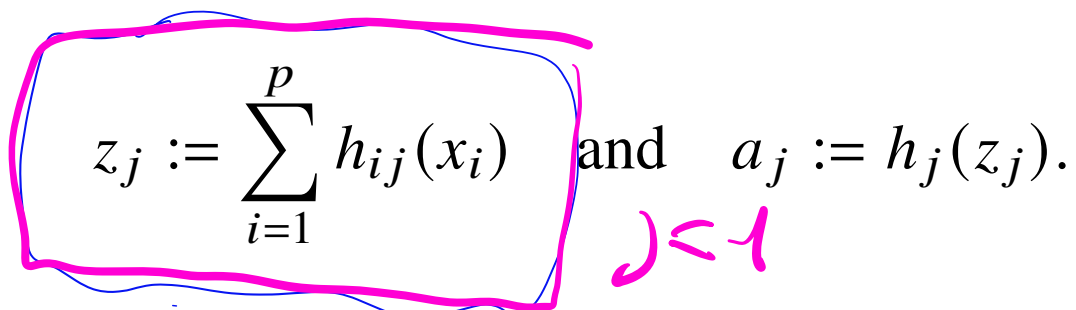
A **neural network** representation of the composition of the maps needed to compute the output  $g^*(\mathbf{x})$  for a given input  $\mathbf{x} \in \mathbb{R}^P$ .



**Figure:** Every continuous function  $g^* : [0, 1]^P \mapsto \mathbb{R}$  can be represented by a neural network with one hidden layer ( $l = 1$ ), an input layer ( $l = 0$ ), and an output layer ( $l = 2$ ).

# In/Hidden/Output Layers

- Each of the  $p$  components of the input  $\mathbf{x}$  is represented as a node in the **input layer** ( $l = 0$ ).
- In the **hidden layer** ( $l = 1$ ) there are  $q := 2p + 1$  nodes, each of which is associated with a pair of variables  $(z, a)$  with values


$$z_j := \sum_{i=1}^p h_{ij}(x_i) \quad \text{and} \quad a_j := h_j(z_j).$$

A link between nodes  $(z_j, a_j)$  and  $x_i$  with weight  $h_{ij}$  signifies that the value of  $z_j$  depends on the value of  $x_i$  via function  $h_{ij}$ .

- The **output layer** ( $l = 2$ ) represents the value  $g^*(\mathbf{x}) = \sum_{j=1}^q a_j$ .



# Activation Functions

In practice, we do not know the collection of (generally nonlinear) functions  $\{h_i, h_{ij}\}$ , because they depend on the unknown  $g^*$ .

Kolmogorov's theorem only asserts the existence of  $\{h_j, h_{ij}\}$ , and does not tell us how to construct these nonlinear functions.

One way out of this predicament is to replace these  $(2p + 1)(p + 1)$  unknown functions with a much larger number of *known* nonlinear functions called **activation functions**.

An example is the logistic activation function:

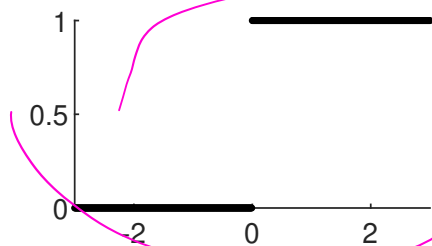
$$S(z) = (1 + \exp(-z))^{-1}.$$

We hope that a network built from a sufficiently large number of activation functions will have similar representational capacity as Kolmogorov's approximation network with  $(2p + 1)(p + 1)$  functions.

# Activation Functions

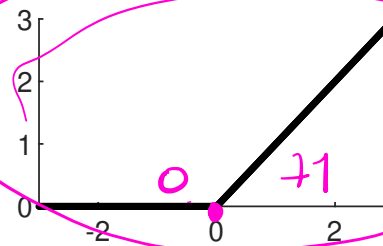
In general, we wish to use the simplest activation functions that will allow us to build a learner with large representational capacity and low training cost. There are many choices of activation functions.

Heaviside or unit step



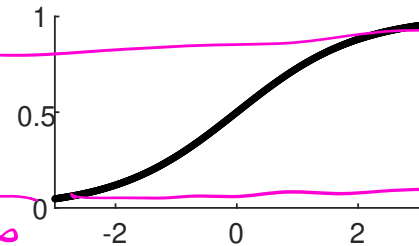
$$\mathbb{I}\{z \geq 0\}$$

rectified linear unit (ReLU)



$$z \times \mathbb{I}\{z \geq 0\}$$

logistic



$$(1 + \exp(-z))^{-1}$$

Another way to improve the representational capacity of the network is to introduce more hidden layers.

# Feed-Forward Neural Networks

In a neural network with  $L + 1$  layers, the input layer ( $l = 0$ ) encodes the input feature vector  $\mathbf{x}$ , and the output layer ( $l = L$ ) encodes the output function  $\mathbf{g}(\mathbf{x})$ . The remaining layers are **hidden layers**.

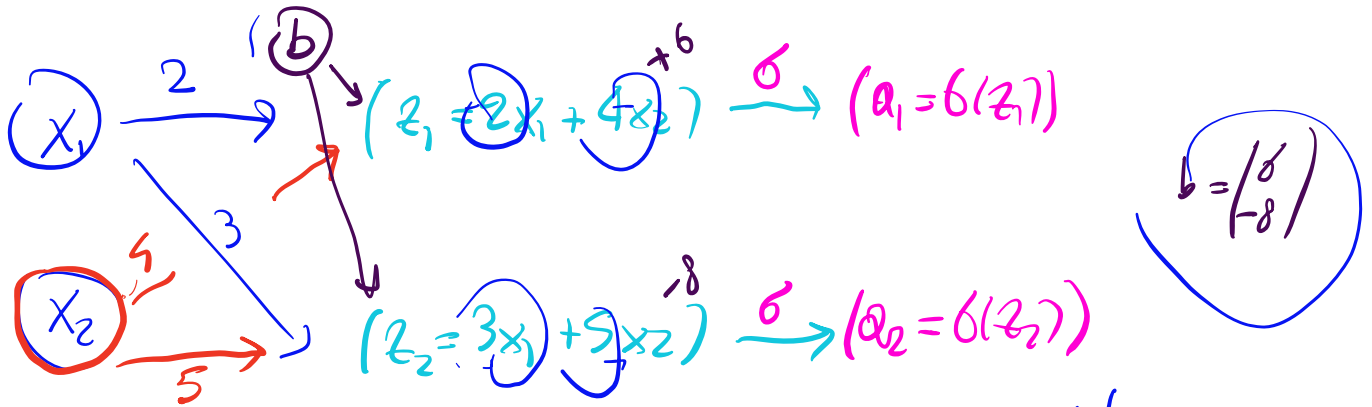
Each layer  $l$  has  $p_l$  nodes.

All nodes in the hidden layers are associated with a pair of variables  $(z, a)$ , gathered into  $p_l$ -dimensional column vectors  $\mathbf{z}_l$  and  $\mathbf{a}_l$ .

In **feed-forward** networks, the variables in any layer  $l$  are simple functions of the variables in the preceding layer  $l - 1$ . In particular,  $\mathbf{z}_l$  and  $\mathbf{a}_{l-1}$  are related via the linear relation

$$\mathbf{z}_l = \mathbf{W}_l \mathbf{a}_{l-1} + \mathbf{b}_l,$$

for some **weight matrix**  $\mathbf{W}_l$  and **bias vector**  $\mathbf{b}_l$ .



$$W_0 = \begin{pmatrix} 2 & 4 \\ 3 & 5 \end{pmatrix}$$

$$W_0 \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} + \begin{pmatrix} b_1 \\ b_2 \end{pmatrix}$$

$$\bar{\sigma} \left( \begin{pmatrix} 2 & 4 \\ 3 & 5 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} + \begin{pmatrix} b_1 \\ b_2 \end{pmatrix} \right) = \begin{pmatrix} \sigma(2x_1 + 4x_2 + 6) \\ \sigma(3x_1 + 5x_2 - 8) \end{pmatrix}$$

$$\bar{\sigma}(a_1, \dots, a_m) := (\sigma(a_1), \dots, \sigma(a_m))$$

$$\begin{matrix} \downarrow & \uparrow \\ (\alpha, \beta) \end{matrix}$$

$$\begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \rightarrow \sigma \left( \boxed{2x_1} + \boxed{4x_2} + \boxed{6} \right)$$

$$\sigma(2x_1 + 6) + 4x_2$$

$$\downarrow \quad h_{11} \quad h_{21}$$

$$h \left( h_{11}^{(x_1)} + h_{21}^{(x_2)} \right)$$

$$\boxed{= \alpha \sigma(u) + \beta \sigma(v)} \\ \in \mathbb{R}$$

$$y = \alpha x + b$$

$$\boxed{X \rightarrow WX + b} \quad \text{Affine}$$

$X \in \mathbb{R}^n \quad W \in \mathbb{R}^{m \times n} \quad b \in \mathbb{R}^m$

# Feed-Forward Neural Networks

Within any hidden layer  $l = 1, \dots, L - 1$ , the components of the vectors  $\mathbf{z}_l$  and  $\mathbf{a}_l$  are related via

$$\mathbf{a}_l = \mathbf{S}_l(\mathbf{z}_l),$$

where  $\mathbf{S}_l : \mathbb{R}^{p_l} \mapsto \mathbb{R}^{p_l}$  is a nonlinear multivalued function.

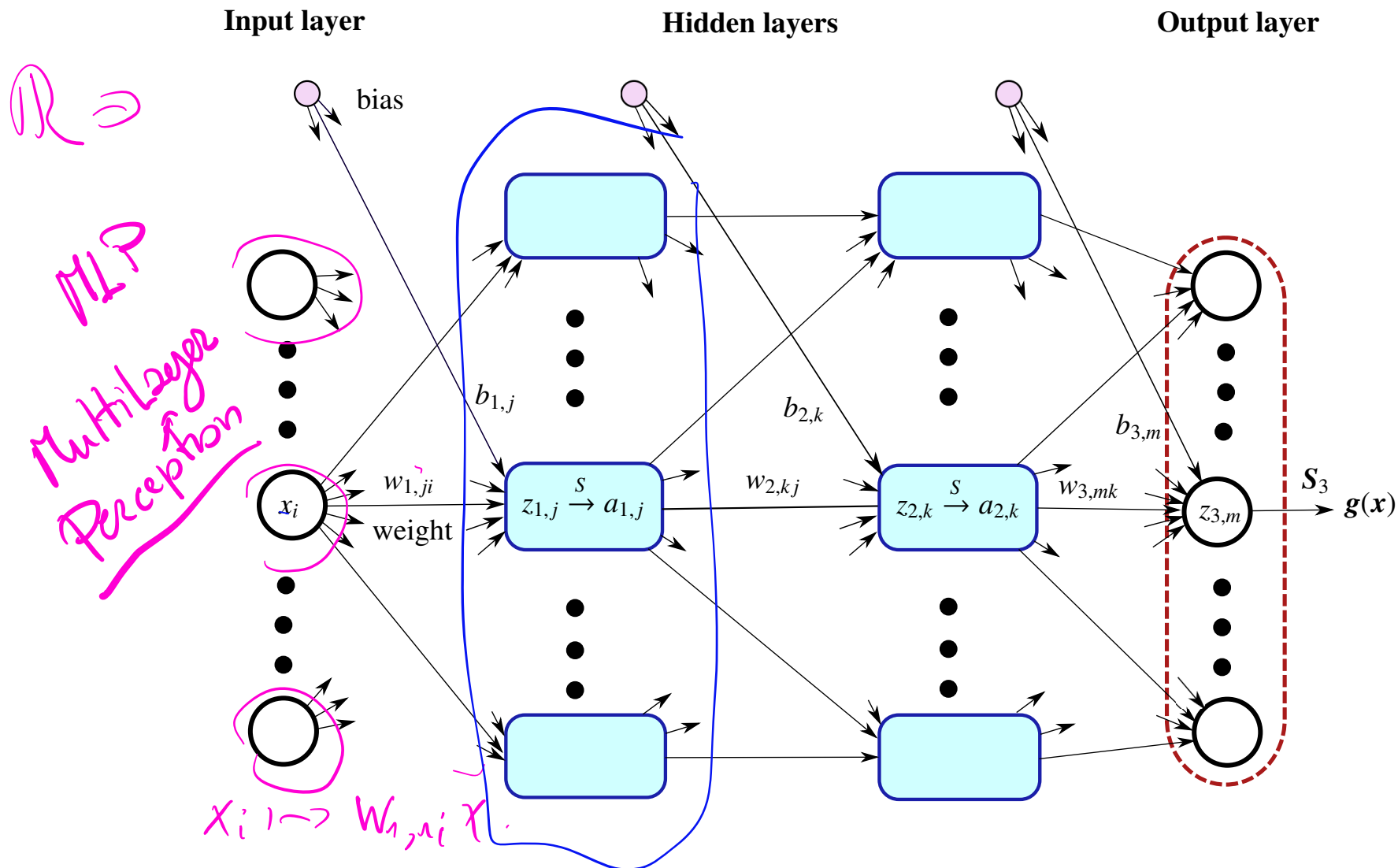
All of these multivalued functions are typically of the form

$$\mathbf{S}_l(\mathbf{z}) = [S(z_1), \dots, S(z_{\dim(\mathbf{z})})]^\top, \quad l = 1, \dots, L - 1, \quad (1)$$

where  $S$  is an **activation** function common to all hidden layers.

The function  $\mathbf{S}_L : \mathbb{R}^{p_{L-1}} \mapsto \mathbb{R}^{p_L}$  in the output layer is more general and its specification depends, for example, on whether the network is used for classification or for the prediction of a continuous output  $Y$ .

# Feed-Forward Neural Networks



**Figure:** A neural network with  $L = 3$ : the  $l = 0$  layer is the input layer, followed by two hidden layers, and the output layer. Hidden layers may have different numbers of nodes.

# Output of a Feed-Forward Neural Network

The output of a feed-forward neural network is thus determined by the input vector  $\mathbf{x}$ , (nonlinear) functions  $\{S_l\}$ , as well as weight matrices  $\mathbf{W}_l = [w_{l,ij}]$  and bias vectors  $\mathbf{b}_l = [b_{l,j}]$  for  $l = 1, 2, 3$ .

The  $(i, j)$ -th element of the weight matrix  $\mathbf{W}_l = [w_{l,ij}]$  is the weight that connects the  $j$ -th node in the  $l$ -th layer with the  $i$ -th node in the  $(l + 1)$ -st layer.

The name given to  $L$  (the number of layers without the input layer) is the **network depth** and  $\max_l p_l$  is called the **network width**.

While we mostly study networks that have an equal number of nodes in the hidden layers ( $p_1 = \dots = p_{L-1}$ ), in general there can be different numbers of nodes in each hidden layer.

# Function Composition

The output  $\mathbf{g}(\mathbf{x})$  of a multiple-layer neural network is obtained from the input  $\mathbf{x}$  via the following sequence of computations:

$$\underbrace{\mathbf{x}}_{\mathbf{a}_0} \rightarrow \underbrace{\mathbf{W}_1 \mathbf{a}_0 + \mathbf{b}_1}_{\mathbf{z}_1} \rightarrow \underbrace{S_1(\mathbf{z}_1)}_{\mathbf{a}_1} \rightarrow \underbrace{\mathbf{W}_2 \mathbf{a}_1 + \mathbf{b}_2}_{\mathbf{z}_2} \rightarrow \underbrace{S_2(\mathbf{z}_2)}_{\mathbf{a}_2} \rightarrow \cdots$$
$$\rightarrow \underbrace{\mathbf{W}_L \mathbf{a}_{L-1} + \mathbf{b}_L}_{\mathbf{z}_L} \rightarrow \underbrace{S_L(\mathbf{z}_L)}_{\mathbf{a}_L} = \mathbf{g}(\mathbf{x}). \in \mathbb{R} \quad (2)$$

Denoting the function  $\mathbf{z} \mapsto \mathbf{W}_l \mathbf{z} + \mathbf{b}_l$  by  $\mathbf{M}_l$ , the output  $\mathbf{g}(\mathbf{x})$  can thus be written as the function composition

$$\mathbf{g}(\mathbf{x}) = S_L \circ \mathbf{M}_L \circ \cdots \circ S_2 \circ \mathbf{M}_2 \circ S_1 \circ \mathbf{M}_1(\mathbf{x}). \quad (3)$$

The algorithm for computing the output  $\mathbf{g}(\mathbf{x})$  for an input  $\mathbf{x}$  is called **feed-forward propagation**.



---

**Algorithm 1:** Feed-Forward Propagation for a Neural Network

---

**Input:** Feature vector  $\mathbf{x}$ ; weights  $\{w_{l,ij}\}$ , biases  $\{b_{l,i}\}$  for each layer  $l = 1, \dots, L$ .

**Output:** The value of the prediction function  $g(\mathbf{x})$ .

```
1  $\mathbf{a}_0 \leftarrow \mathbf{x}$            // the zero or input layer
2 for  $l = 1$  to  $L$  do
3   Compute the hidden variable  $z_{l,i}$  for each node  $i$  in layer  $l$ :

           
$$\mathbf{z}_l \leftarrow \mathbf{W}_l \mathbf{a}_{l-1} + \mathbf{b}_l$$


4   Compute the activation function  $a_{l,i}$  for each node  $i$  in layer  $l$ :

           
$$\mathbf{a}_l \leftarrow S_l(\mathbf{z}_l)$$


5 return  $g(\mathbf{x}) \leftarrow \mathbf{a}_L$            // the output layer
```

---

# Example: Nonlinear Multi-Output Regression

Given the input  $\mathbf{x} \in \mathbb{R}^{p_0}$  and an activation function  $S : \mathbb{R} \mapsto \mathbb{R}$ , the output  $\mathbf{g}(\mathbf{x}) := [g_1(\mathbf{x}), \dots, g_{p_2}(\mathbf{x})]^\top$  of a **nonlinear multi-output regression** model can be computed via a neural network with:

$$\begin{aligned} \mathbf{z}_1 &= \mathbf{W}_1 \mathbf{x} + \mathbf{b}_1, \quad \text{where } \mathbf{W}_1 \in \mathbb{R}^{p_1 \times p_0}, \mathbf{b}_1 \in \mathbb{R}^{p_1}, \\ a_{1,k} &= S(z_{1,k}), \quad k = 1, \dots, p_1, \\ \mathbf{g}(\mathbf{x}) &= \mathbf{W}_2 \mathbf{a}_1 + \mathbf{b}_2, \quad \text{where } \mathbf{W}_2 \in \mathbb{R}^{p_2 \times p_1}, \mathbf{b}_2 \in \mathbb{R}^{p_2}, \end{aligned}$$

which is a neural network with one hidden layer and output function  $S_2(\mathbf{z}) = \mathbf{z}$ .

In the special case where  $p_1 = p_2 = 1$ ,  $\mathbf{b}_2 = 0$ ,  $\mathbf{W}_2 = 1$ , and we collect all parameters into the vector  $\boldsymbol{\theta}^\top = [\mathbf{b}_1, \mathbf{W}_1] \in \mathbb{R}^{p_0+1}$ , the neural network can be interpreted as a **generalized linear model** with  $\mathbb{E}[Y \mid \mathbf{X} = \mathbf{x}] = h([1, \mathbf{x}^\top] \boldsymbol{\theta})$  for some activation function  $h$ .

# Example: Multi-Logit Classification

Suppose that, for a classification problem, an input  $\mathbf{x}$  has to be classified into one of  $c$  classes, labeled  $0, \dots, c - 1$ . We can perform the classification via a neural network with one hidden layer, with  $p_1 = c$  nodes. In particular, we have

$$\mathbf{z}_1 = \mathbf{W}_1 \mathbf{x} + \mathbf{b}_1, \quad \mathbf{a}_1 = \mathbf{S}_1(\mathbf{z}_1) = \begin{pmatrix} \exp(z_1) / \sum \exp(z_k) \\ \exp(z_2) / \sum \exp(z_k) \\ \vdots \end{pmatrix}$$

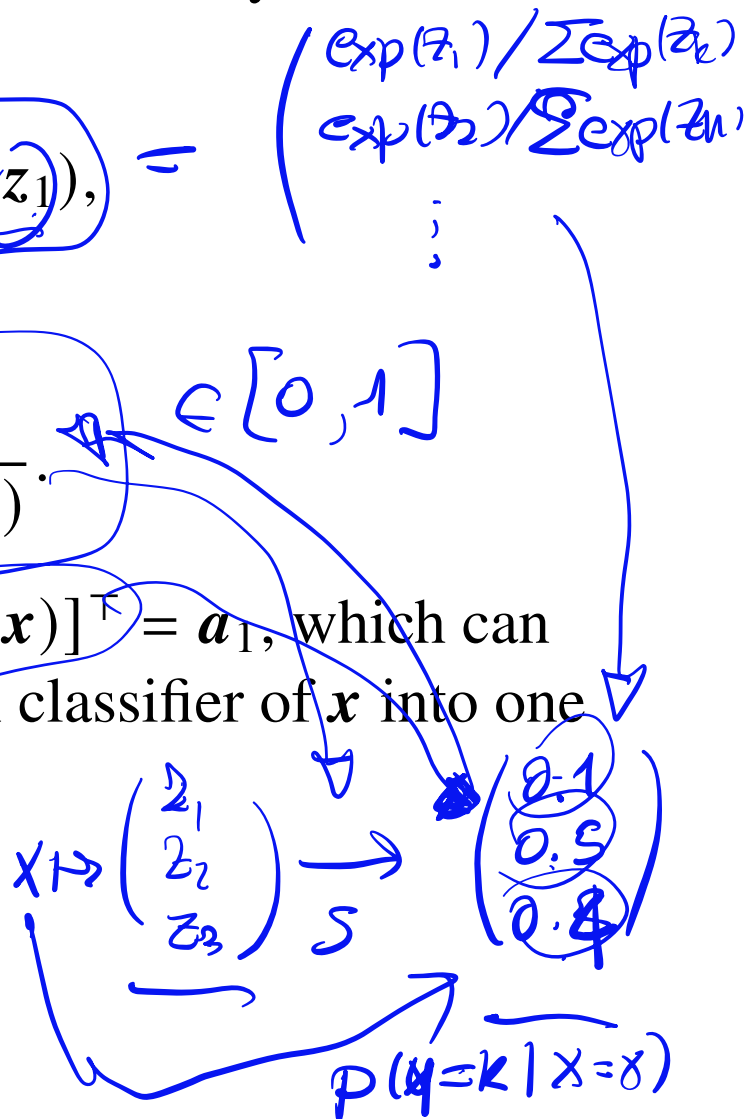
where  $\mathbf{S}_1$  is the **softmax** function:

$$\text{softmax} : \mathbf{z} \mapsto \frac{\exp(z)}{\sum_k \exp(z_k)}$$

For the output, we take  $\mathbf{g}(\mathbf{x}) = [g_1(\mathbf{x}), \dots, g_c(\mathbf{x})]^T = \mathbf{a}_1$ , which can then be used as a **pre-classifier** of  $\mathbf{x}$ . The actual classifier of  $\mathbf{x}$  into one of the categories  $0, 1, \dots, c - 1$  is then

$$\underset{k \in \{0, \dots, c-1\}}{\operatorname{argmax}} g_{k+1}(\mathbf{x}).$$

This is equivalent to the multi-logit classifier.



Regression

$$f(x) = y \in \mathbb{R}$$

$$X = \begin{pmatrix} x_1 \\ \vdots \\ 1 \\ \vdots \\ x_d \end{pmatrix} \rightarrow W_0 x + b \rightarrow S_0(W_0 x + b) \rightarrow \dots \rightarrow S_\ell(-) \in \mathbb{R}$$

Regression

Extend  $\rightarrow f: \mathbb{R}^d \rightarrow \mathbb{R}^e$

$$S_\ell(-) \in \mathbb{R}^e$$

Classification  $X \rightarrow \text{same} \rightarrow S_\ell(-) = \text{softmax}$

$$f(z) = \begin{pmatrix} \exp(z_1) / \sum \exp(z_k) \\ \vdots \end{pmatrix} \leftarrow \text{Prob.} \Rightarrow X \in \text{Class}(k) \text{ s.t. } \max_k \text{pr}(y=k | x=x)$$

$\hookrightarrow$  dimensione = # Class =  $C$

$$C \begin{cases} \vdots \\ \vdots \end{cases} \xrightarrow{\text{Soft}} \begin{cases} \vdots \\ \vdots \end{cases} \rightarrow \underline{\text{CLASSE PU' Prob.}}$$

$\underbrace{\quad}_{e-1} \quad \quad \quad e \quad e$

# Example: Density Estimation

A Gaussian mixture density with  $p_1$  components and a common scale parameter  $\sigma > 0$  can be viewed as an NN with two hidden layers.

Let the activation function in the first hidden layer,  $S_1$ , be of the form (1) with  $S(z) := \exp(-z^2/(2\sigma^2))/\sqrt{2\pi\sigma^2}$ . Then the density value  $g(x)$  is computed via:

No

$$\begin{aligned}z_1 &= \mathbf{W}_1 x + \mathbf{b}_1, & \mathbf{a}_1 &= S_1(z_1), \\z_2 &= \mathbf{W}_2 \mathbf{a}_1 + \mathbf{b}_2, & \mathbf{a}_2 &= S_2(z_2), \\g(x) &= \mathbf{a}_1^\top \mathbf{a}_2,\end{aligned}$$

where  $\mathbf{W}_1 = \mathbf{1}$  is a  $p_1 \times 1$  column vector of ones,  $\mathbf{W}_2 = \mathbf{O}$  is a  $p_1 \times p_1$  matrix of zeros, and  $S_2$  is the softmax function.

We identify the column vector  $\mathbf{b}_1$  with the  $p_1$  location parameters,  $[\mu_1, \dots, \mu_{p_1}]^\top$  of the Gaussian mixture and  $\mathbf{b}_2 \in \mathbb{R}^{p_1}$  with the  $p_1$  weights of the mixture.

# Network Architecture

The **network architecture** is an important design consideration in neural networks.

For example, if the connectivity from one layer to the next is sparse and the links share the same weight values  $\{w_{l,ij}\}$  (called **parameter sharing**) for all  $\{(i, j) : |i - j| = 0, 1, \dots\}$ , then the weight matrices will be sparse and **Toeplitz**, allowing fast matrix-vector product calculations, e.g., via the **fast Fourier transform**.

An important example is the **convolution neural network** (CNN), in which the network layers encode the **convolution** operation:

$$\mathbf{W}_l \mathbf{a}_{l-1} = \mathbf{w}_l * \mathbf{a}_{l-1},$$

where  $[\mathbf{x} * \mathbf{y}]_i := \sum_k x_k y_{i-k+1}$ .

# Convolution Neural Networks

CNNs are particularly suited to image processing problems.

Suppose that the input image is given by an  $m_1 \times m_2$  matrix of pixels.

Define a  $k \times k$  matrix (sometimes called a **kernel**, where  $k$  is generally taken to be 3 or 5).

Then, the convolution layer output can be calculated using the discrete convolution of all possible  $k \times k$  input matrix regions and the kernel matrix; the convolution layer output size is  $(m_1 - k + 1) \times (m_2 - k + 1)$ .

The figure shows a  $5 \times 5$  input image and a  $2 \times 2$  kernel with a  $4 \times 4$  output matrix.

