

STAT 535

10/31/2023

# Lecture 10

# Lecture Notes III – Neural Networks

Marina Meilă  
`mmp@stat.washington.edu`

Department of Statistics  
University of Washington

October, 2023

# Multi-layer/Deep neural networks

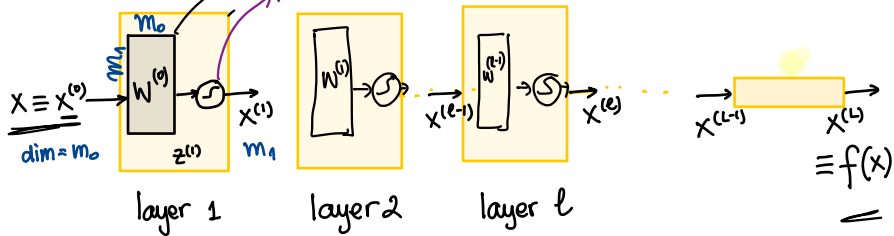
The construction can be generalized recursively to arbitrary numbers of layers.

Each layer is a linear combination of the outputs from a previous layer (a multivariate operation), followed by a non-linear transformation via the logistic function  $\phi$ . Let  $x \equiv x^{(0)}, y \equiv x^{(L)}, n_0 = n, n_L = 1$  and define the recursion:

$$x_j^{(l)} = \phi \left( \underbrace{(w_j^{(l)})^T x^{(l-1)}}_{\text{linear}} \right), \text{ for } j = 1 : n_l \quad (9)$$

layers  $l = 1, 2, 3, \dots, L$

The vector variable  $x^{(l)} \in \mathbb{R}^{n_l}$  is the output of layer  $l$  of the network. As before, the sigmoid of the last layer may be omitted.



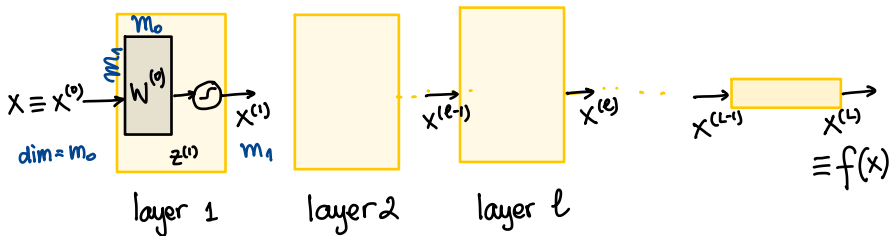
# Multi-layer/Deep neural networks

The construction can be generalized recursively to arbitrary numbers of layers. Each layer is a linear combination of the outputs from a previous layer (a multivariate operation), followed by a non-linear transformation via the logistic function  $\phi$ . Let  $x \equiv x^{(0)}, y \equiv x^{(L)}, n_0 = n, n_L = 1$  and define the recursion:

$$x_j^{(l)} = \phi \left( (w_j^{(l)})^T x^{(l-1)} \right), \text{ for } j = 1 : n_l \quad (9)$$

*Handwritten notes:*  $w_j^{(l)} = w_j^{(l-1)}$  (with a squiggle),  $x_j^{(l-1)} = \phi(\dots)$  (with a squiggle), and  $z$  (with a squiggle).

The vector variable  $x^{(l)} \in \mathbb{R}^{n_l}$  is the output of layer  $l$  of the network. As before, the sigmoid of the last layer may be omitted.



Forward propagation:

$$x \mapsto f(x)$$

$$\downarrow$$

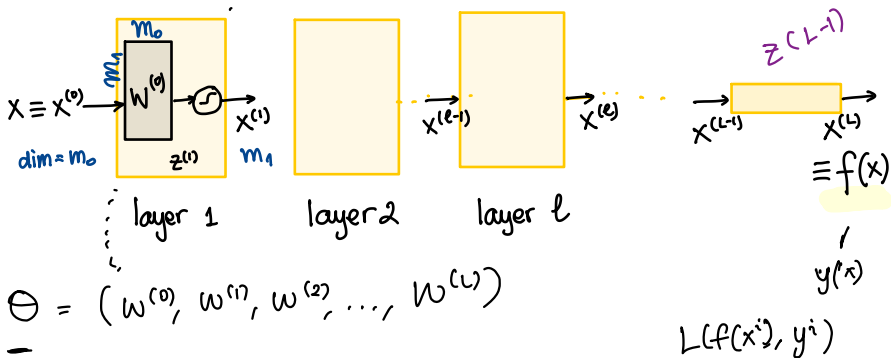
$$x^{(0)} \mapsto x^{(L)}$$

## Multi-layer/Deep neural networks

The construction can be generalized recursively to arbitrary numbers of layers. Each layer is a linear combination of the outputs from a previous layer (a multivariate operation), followed by a non-linear transformation via the logistic function  $\phi$ . Let  $x \equiv x^{(0)}, y \equiv x^{(L)}, n_0 = n, n_L = 1$  and define the recursion:

$$x_j^{(l)} = \phi \left( \underbrace{(w_j^{(l)})^T x^{(l-1)}}_{z^{(l)}} \right), \text{ for } j = 1 : n_l \quad (9)$$

The vector variable  $x^{(l)} \in \mathbb{R}^{n_l}$  is the output of layer  $l$  of the network. As before, the sigmoid of the last layer may be omitted.



Training data  $i$

$$(x^{(i)}, y^{(i)}) \rightarrow f(x^{(i)})$$

$$L(f(x^{(i)}), y^{(i)})$$

$$\nabla_{\theta} L(f(x^{(i)}), y^{(i)}) = \frac{\partial L}{\partial \theta}$$

for example:

$$L(f(x^{(i)}), y^{(i)}) = (f(x^{(i)}) - y^{(i)})^2$$

$$\frac{\partial L}{\partial w^{(L-1)}} \stackrel{\text{chain rule}}{=} \frac{\partial L}{\partial f(x_i)} \cdot \frac{\partial f(x_i)}{\partial z^{(L-1)}} \cdot \frac{\partial z^{(L-1)}}{\partial w^{(L-1)}}$$

$$\downarrow$$
$$x^{(L-2)}$$

$$z^{(L)} = w^{(L)} \top x^{(L-1)}$$

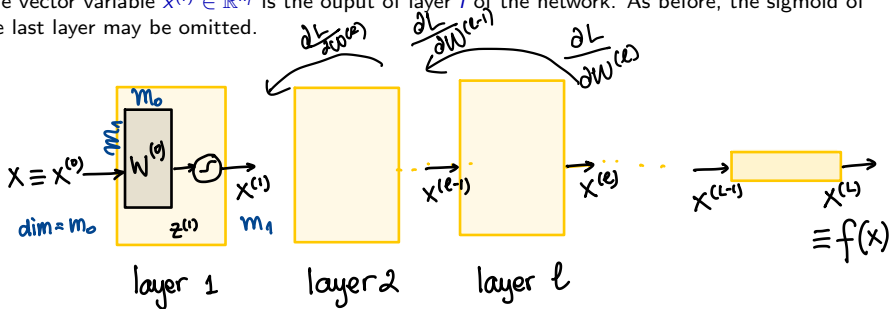
$$f(x_i) = \phi(z_i^{L-1})$$

## Multi-layer/Deep neural networks

The construction can be generalized recursively to arbitrary numbers of layers. Each layer is a linear combination of the outputs from a previous layer (a multivariate operation), followed by a non-linear transformation via the logistic function  $\phi$ . Let  $x \equiv x^{(0)}, y \equiv x^{(L)}, n_0 = n, n_L = 1$  and define the recursion:

$$x_j^{(l)} = \phi \left( (w_j^{(l)})^T x^{(l-1)} \right), \text{ for } j = 1 : n_l \quad (9)$$

The vector variable  $x^{(l)} \in \mathbb{R}^{n_l}$  is the output of layer  $l$  of the network. As before, the sigmoid of the last layer may be omitted.



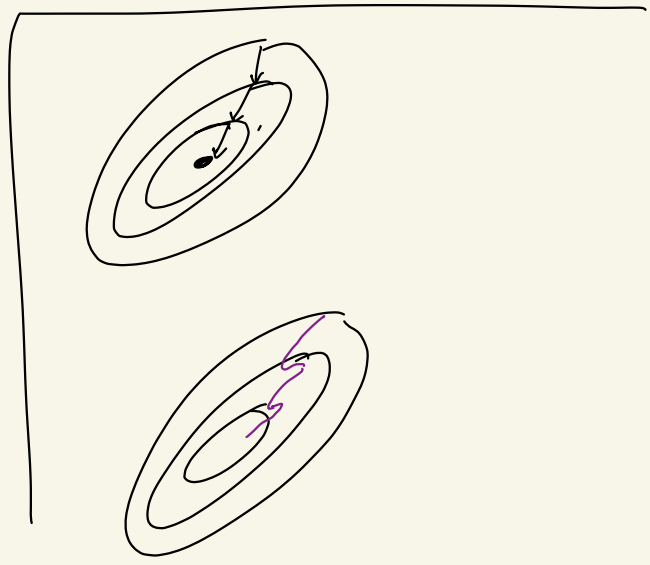
# Backpropagation.

$$z^{(l-1)} = w^{(l-1)} \phi(z^{(l-2)})$$

$$z^{(l-2)} = w^{(l-2)} x^{(l-3)}$$

$$\frac{\partial L}{\partial w^{(l-2)}} = \frac{\partial L}{\partial z^{(l-2)}} \cdot \frac{\partial z^{(l-2)}}{\partial w^{(l-2)}} \rightarrow x^{(l-3)}$$

$$\frac{\partial L}{\partial z^{(l-2)}} = \frac{\partial L}{\partial z^{(l-1)}} \cdot \frac{\partial z^{(l-1)}}{\partial z^{(l-2)}}$$



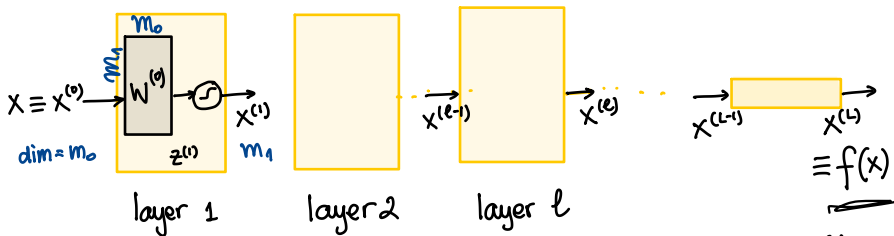


## Multi-layer/Deep neural networks

The construction can be generalized **recursively** to arbitrary numbers of layers. Each layer is a linear combination of the outputs from a previous layer (a multivariate operation), followed by a non-linear transformation via the logistic function  $\phi$ . Let  $x \equiv x^{(0)}$ ,  $y \equiv x^{(L)}$ ,  $n_0 = n$ ,  $n_L = 1$  and define the recursion:

$$x_j^{(l)} = \phi \left( (w_j^{(l)})^T x^{(l-1)} \right), \text{ for } j = 1 : n_l \quad (9)$$

The vector variable  $x^{(l)} \in \mathbb{R}^{n_l}$  is the output of layer  $l$  of the network. As before, the sigmoid of the last layer may be omitted.



Forward propagation :  $x_i(0) \rightarrow x_i(1) \rightarrow \dots \rightarrow x_i(l) \rightarrow \dots \rightarrow x_i(L) = f(x_i)$

Backward propagation :  $\frac{\partial L}{\partial z^{L-1}}, \frac{\partial L}{\partial w^{L-1}}, \frac{\partial L}{\partial z^{L-2}}, \frac{\partial L}{\partial w^{L-2}}, \dots, \frac{\partial L}{\partial z^{(1)}}, \frac{\partial L}{\partial w^{(1)}}$

## Are multiple layers necessary?

- ▶ 1990's: NO
  - ▶ 2000's: YES
  - ▶ 2020's: The more the better!
- 
- ▶ A theoretical result

### Theorem (Cybenko, $\approx$ 1986)

Any continuous function from  $[0, 1]^d$  to  $\mathbb{R}$  can be approximated arbitrarily closely by a linear output, **two layer neural network** defined in (2) with a sufficiently large number of hidden units  $m$ .

- ▶ A practical result



## Deep Learning

- ▶ **Deep learning** = multi-layer neural net
- ▶ So, what is new?
  - ▶ small variations in the “units”, e.g. switch stochastically w.p.  $\phi(w^T x^{in})$  (Restricted Boltzmann Machine), Rectified Linear units
  - ▶ training method stochastic gradient, auto-encoders vs. back-propagation (we will return to this when we talk about training predictors)
  - ▶ lots of data
  - ▶ double descent

## Resnets – Residual networks

Idea What is the "simplest" input-output function?  $f_0(x) = x$

(identity function)

► Hence, a NN layer should learn the difference w.r.t. identity  $f_0$

$$x_{l+1} = B_l \phi(W_l x_l) + x_l$$

non-linear

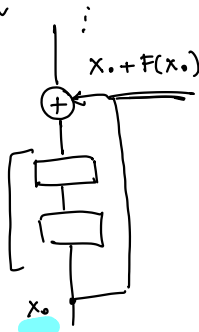
(10)

Generalization DenseNet

► Layer  $l$  gets inputs from  $l-1, l-2, \dots$

"residual  
function"

$F(x_0)$



$f(x) = mx + b$

Goal:  $f_0(x) = x$

$F(x) = f_0(x) - y$

# ConvNets – Convolutional Networks

- ▶ **discrete convolution** let  $f, g : \mathbb{Z} \rightarrow \mathbb{R}$

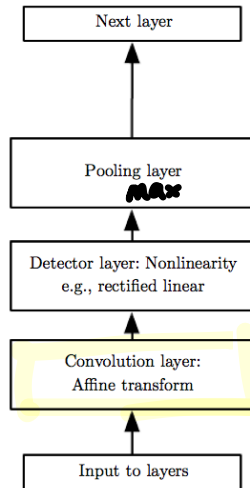
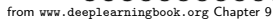
$\mathbb{Z}$  = all integers

$$(f * g)(t) = \sum_{i \in \mathbb{Z}} f(t-i) \underbrace{g(i)}_{w(i), w_i} \quad (11)$$

- ▶ convolution as **Toeplitz** matrix vector multiplication
- ▶ in ConvNets,  $\mathbb{Z}$  is replaced by  $1 : m$ ,  $f$  is **padded with 0's**
  - ▶  $g$  is a (smoothing) kernel
  - ▶ i.e.  $g(i) = g(-i) > 0$  and  $|\text{supp } g| = 2s + 1 \ll m$ ,  $\sum_i g(i) = 1$
- ▶ Convolutional layer  $f \leftarrow x$  input,  $g \leftarrow w$  weights,  $s$  output

$$s(t) = \sum_{i=t-s}^{t+s} w_i s(t-i) = (s * w)(t) \quad (12)$$

- ▶ Pooling



# Autoencoders

← unsupervised learning

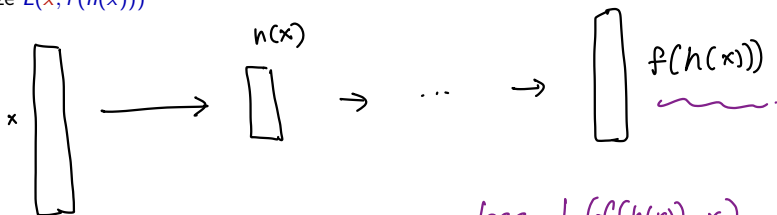
**Question** How to learn from data without outputs  $y$ ?

This is **unsupervised learning**, not prediction

**Idea** Learn a **low dimensional/sparse** representation  $h(x)$  of data  $x \in \mathbb{R}^d$

$$h(x) \in \mathbb{R}^{(m)} \text{ with } \underline{m} < d \quad f(h(x)) \approx x! \quad (13)$$

► Optimize  $L(x, f(h(x)))$



- compressing input data

$$\text{loss: } L(f(h(x)), x)$$

↓  
optimize this

# Autoencoders

**Question** How to learn from data without outputs  $y$ ?

This is **unsupervised learning**, not prediction

**Idea** Learn a **low dimensional/sparse** representation  $h(x)$  of data  $x \in \mathbb{R}^d$

$$h(x) \in \mathbb{R}^m, \text{ with } m < d \quad f(h(x)) \approx x! \quad (13)$$

► Optimize  $L(x, f(h(x)))$