

STAT 535

10/27/22

# Lecture 9

Neural Networks

Q&A, Q1 back  
HW 3 posted

# Lecture Notes III – Neural Networks

Marina Meilă  
mmp@stat.washington.edu

Department of Statistics  
University of Washington

October, 2020

Two-layer Neural Networks ←

Multi-layer neural networks

A zoo of multilayer networks

**Reading** HTF Ch.: 11.3 Neural networks, Murphy Ch.: (16.5 neural nets) and Dive Into Deep Learning 4.1-4.3

## Two-layer Neural Networks

- ▶ The **activation function** (a term borrowed from neuroscience) is any continuous, bounded and strictly increasing function on  $\mathbb{R}$ . Almost universally, the activation function is the **logistic** (or **sigmoid**)

$$\phi(z) = \frac{1}{1 + e^{-z}} \quad (1)$$

because of its nice additional computational and statistical properties.

- ▶ We build a **two-layer neural network** in the following way:

Inputs	$x_k$	$k = 1 : d$
Bottom layer <sup>1</sup>	$z_j = \phi(w_j^T x)$	$j = 1 : m, w_j \in \mathbb{R}^n$
Top layer	$f = \phi(\beta^T z)$	$\beta \in \mathbb{R}^m$
Output	$f$	$\in [0, 1]$

In other words, the neural network implements the function

$$f(x) = \sum_{j=1}^m \beta_j z_j = \sum_{j=1}^m \beta_j \phi\left(\sum_{k=1}^n w_{kj} x_k\right) \in (-\infty, \infty) \quad (2)$$

Note that this is just a linear combination of logistic functions.

---

<sup>1</sup>In neural net terminology, each variable  $z_j$  is a **unit**, the bottom layer is **hidden**, while top one is **visible**, and the units in this layer are called hidden/visible units as well. Sometimes the inputs are called **input units**; imagine neurons or individual circuits in place of each  $x, y, z$  variable.

$$x \in \mathbb{R}^d$$

$$y \in \mathcal{Y} \begin{cases} \mathbb{R} \\ \pm 1 \\ 1, 2, \dots, r \end{cases}$$

output [layer]

$$f(x) = \beta^T z(x) \Rightarrow \text{parameters}$$

$$z_j(x) = \varphi(W_j x)$$

$z \in \mathbb{R}^m$  → activation function

↳ hidden units variable

$$\beta \in \mathbb{R}^m$$

$$W = [W_{jk}]_{\substack{j=1:m \\ k=1:d}}$$

( $x_k$  = input unit)  
( $f$  = output unit)

Hidden layer

$\varphi = \begin{cases} \text{sigmoid} \\ \text{(logistic!)} \\ \text{ReLU} \end{cases}$

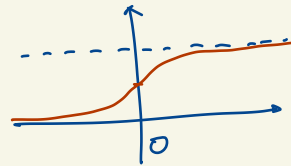
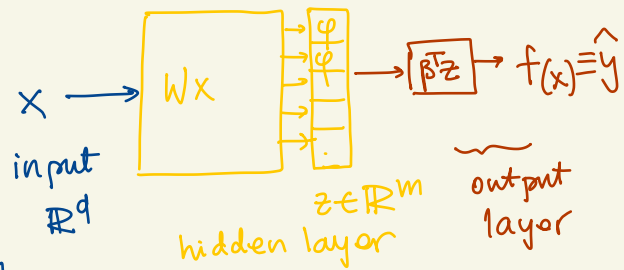
$$\varphi(u) = \frac{1}{1 + e^{-u}} = \frac{e^u}{1 + e^u}$$

bounded, smooth  
(Lipschitz)

$$\varphi(u) = \begin{cases} u & \text{for } u \geq 0 \\ 0 & \text{for } u < 0 \end{cases}$$

(Rectified Linear Unit)

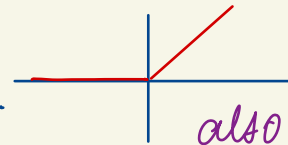
unbounded  
non-smooth at 0  
Lipschitz  
 $L = 1$



$f(u)$  is  $L$ -Lipschitz  $\Leftrightarrow$

$$|f(u) - f(u')| \leq L \|u - u'\|$$

$L > 0$  a parameter



also called 'hinge'

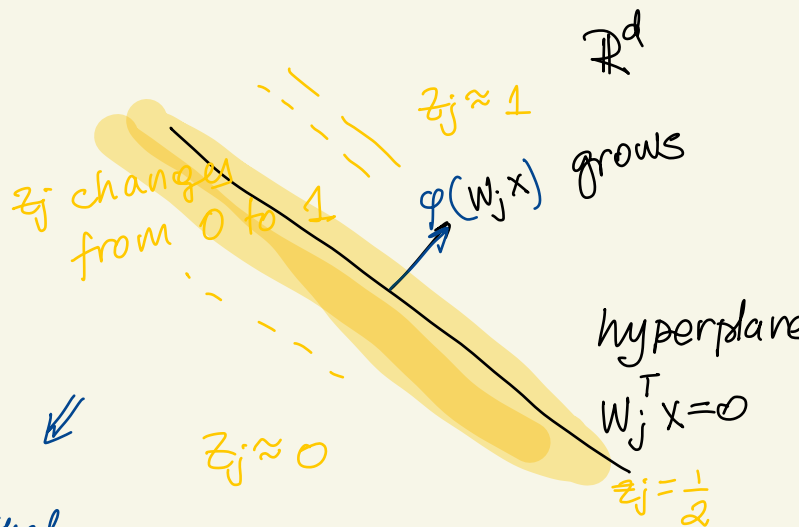
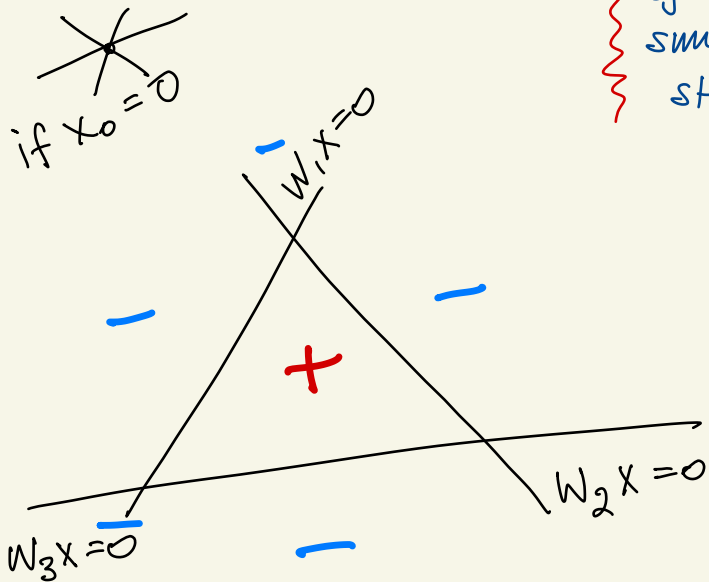
What does 2 layer NN represent?

$\varphi = \text{sigmoid}$

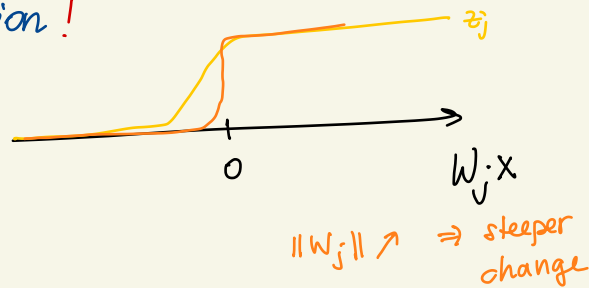
$$f = \sum \beta_j \varphi(W_j x)$$

$W_j^T \in \mathbb{R}^d$  or  $\mathbb{R}^{d+1}$

$$\text{augmented } x : x = \begin{bmatrix} x_0=1 \\ x_1 \\ \vdots \\ x_d \end{bmatrix}$$



$z_j$  is smoothed step function!



$$f = \beta_1 z_1 + \beta_2 z_2 + \beta_3 z_3 + \beta_4 z_4$$

classification: sign  $f$

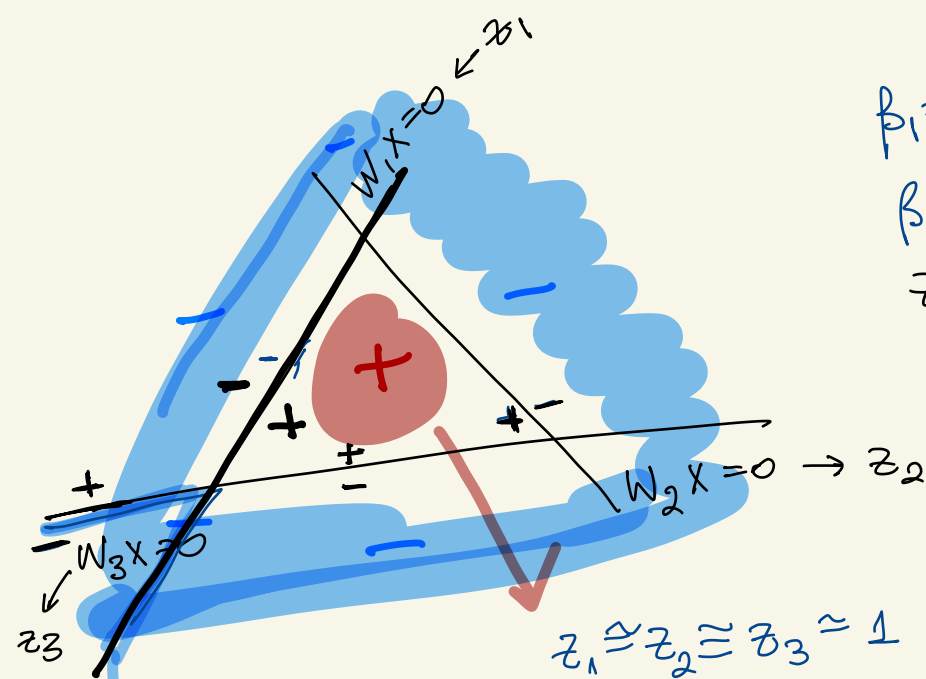
$x = [1 \ x_1 \ x_2]$

Exercise  
 $W_4 = ?$

$$\beta_1 = \beta_2 = \beta_3 = 1$$

$$\beta_4 = -2.5 \leftarrow \text{threshold}$$

$$z_4 \equiv 1$$



$$z_1 \approx z_2 \approx z_3 \approx 1 \Rightarrow \text{sgn } f > 0$$

$$\text{sgn } f < 0$$

Increasing accuracy

- sharper transition  $\leftarrow$  increase  $\|W_j\|$
- more complex decision boundary

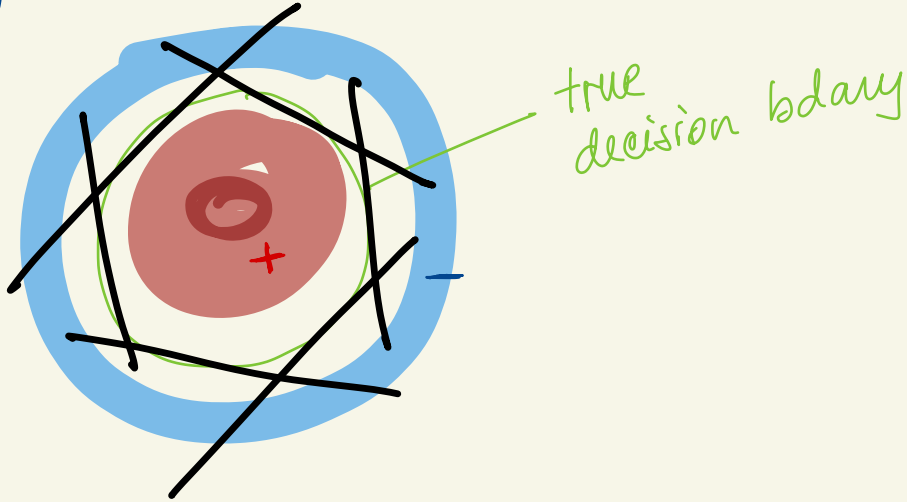
$$z_1 < \frac{1}{2} \approx 0$$

$$z_2 > \frac{1}{2} \approx 1$$

$$z_3 < \frac{1}{2} \approx 0$$

## Increasing accuracy

- sharper transition  $\Leftarrow$  increase  $\|W_j\|$
- more complex decision boundary  $\Leftarrow$  increase  $m$





$$f = \beta^T z$$

$$\rightarrow \beta^T z + \beta_0$$

$$\text{OR } \tilde{z} = \begin{bmatrix} 1 \\ z_1 \\ \vdots \\ z_m \end{bmatrix} \} \rightarrow \tilde{\beta}^T \tilde{z}$$

$$\tilde{\beta} = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_m \end{bmatrix}$$

This page contains  
answers to some after  
class questions

What if no  $\beta_0$ ?  $\Rightarrow m$  larger !!

$x \in \mathbb{R}^d$   
data

$x^{(1)}, \dots, x^{(n)} \in \mathbb{R}^d$   $d \times n$

$m, q \rightarrow \mathcal{F}_{m,q}$  model  
class

$W, \beta \rightarrow$  parameters  
or weights

## Output layer options

- ▶ **linear** layer as in (2)  $f = \sum_j \beta_j z_j$
- ▶ **logistic** layer: in **classification**  $f(x) \in [0, 1]$  is interpreted as the probability of the + class.

$$f(x) = \phi \left( \sum_{j=1}^m \beta_j z_j \right) = \phi \left( \sum_{j=1}^m \beta_j \phi \left( \sum_j w_{kj} x_k \right) \right) \quad (3)$$

- ▶ **softmax** layer in multiway classification

The **softmax** function  $\phi(z) : \mathbb{R}^m \rightarrow (0, 1)^m$

$$\phi_k(z) = \frac{e^{z_k}}{\sum_{j=1}^m e^{z_j}} \quad (4)$$

- ▶ Properties

- ▶  $\sum_{j=1}^m \phi_j(z) = 1$  for all  $z$
- ▶ for  $z_k \gg z_j, j \neq k$   $\phi_k(z) \rightarrow 1$ .
- ▶ derivatives  $\frac{\partial \phi_j}{\partial z_k} = \phi_k \delta_{jk} - \phi_j \phi_k$

## Output layer

$$z \in \mathbb{R}^m \rightarrow y \in \mathcal{Y}$$

- $y \in \mathbb{R}$

$$f(x) = \beta^T z$$

regression

$$\text{Loss} = L_{LS}$$

- $y \in \{\pm 1\}$

$$f(x) = \varphi(\beta^T z) - \frac{1}{2}$$

binary  
classif

$$\Rightarrow \hat{y} = \text{sgn } f = \arg \max \{ \beta^T z, -\beta^T z \}$$

parameters

- $y \in \{1, 2, \dots, r\}$

multiway  
classification

$$\text{Loss} = L_{\text{logit}} = -\ln P[y = +1|x]$$

$$\hat{y} = \arg \max_k \{ \beta_1^T z, \beta_2^T z, \dots, \beta_r^T z \}$$

$$W$$
$$[\beta_1 \dots \beta_r]$$

$$\text{Loss} = -\ln P_{y=\text{true}} | x$$

$$f(x) = \underbrace{\text{softmax}}_{\varphi} (\beta_1^T z, \dots, \beta_r^T z) \in (0, 1)^r$$

(should be soft arg max !!)

$$\varphi_k(u) = \frac{e^{u_k}}{\sum_{k'=0}^r e^{u_{k'}}}$$

$$\Rightarrow \varphi_k(u) \in (0, 1), k = 1:r$$
$$\sum_k \varphi_k(u) = 1 \text{ for all } u \in \mathbb{R}^r$$

$$k = 1:r$$

Ex:  $r=2 \Rightarrow \varphi_1(u) = \text{sigmoid (logistic)}$

$$\varphi_2(u) = 1 - \varphi_1$$

- special cases of GLIM

Interpret  $\varphi_k(u) = \Pr[u=k]$

$$f_k(x) = \Pr[y=k]$$

↑  
model's  
confidence

# Generalized Linear Models (GLM)

A GLM is a regression where the "noise" distribution is in the exponential family.

- ▶  $y \in \mathbb{R}$ ,  $y \sim P_\theta$  with  
*single variable  $y$*

$$P_\theta(y) = e^{\theta y - \ln \psi(\theta)} = \frac{1}{\psi(\theta)} e^{\theta y} \quad (5)$$

*normalizing constant*

- ▶ the parameter  $\theta$  is a linear function of  $x \in \mathbb{R}^d$

$$P_{y|x} \Rightarrow \theta(x) \quad \theta = \beta^T x \quad \text{this } W \text{ in NN} \quad (6)$$

- ▶ We denote  $E_\theta[y] = \mu$ . The function  $g(\mu) = \theta$  that relates the mean parameter to the natural parameter is called the **link function**.

The log-likelihood (w.r.t.  $\beta$ ) is

**FACT 1 obvious:**  $l(\beta) = \ln P_\theta(y|x) = \theta y - \psi(\theta)$  where  $\theta = \beta^T x$  **FACT 2 interesting:** (7)

and the gradient w.r.t.  $\beta$  is therefore

$$\nabla_\beta l = \nabla_\theta l \nabla_\beta (\beta^T x) = (y - \mu) x \quad \frac{\partial \ln P}{\partial \theta} = -\mu(\theta) + y \quad (8)$$

This simple expression for the gradient is the generalization of the gradient expression you obtained for the two layer neural network in the homework. [Exercise: This means that the sigmoid function is the *inverse link function* defined above. Find what is the link function that corresponds to the neural network.]

$$\frac{\partial}{\partial \beta} (\beta^T x) = \frac{\partial \theta}{\partial \beta} = x$$

Proof of  
FACT 2

$$\ln P(y) = \theta y - \ln \psi(\theta)$$

$$\psi(\theta) = \int_{-\infty}^{\infty} e^{\theta y} dy$$

$$\frac{\partial}{\partial \theta} \ln P_{\theta}(y) = y - \frac{\int y e^{\theta y} dy}{\underbrace{\psi(\theta)}_{P_{\theta}(y)}} = y - \overbrace{\int y \frac{e^{\theta y}}{\psi(\theta)} dy}^{E_{\theta}[y]} = y - \mu(\theta)$$