# Lecture 9

HW 3 due 2/4
Sol 1 , Sol 2
available
[HW 4 : TB posted]
OH tomorrow
        9:30, 10:30
L$_{11}$ corrections

# Lecture Notes IV – Neural Networks, Part 1

Marina Meilă

`mmp@uwaterloo.ca`

With Thanks to Pascal Poupart & Gautam Kamath
Cheriton School of Computer Science
University of Waterloo

February 2, 2026

A little history ✔

The single "neuron" ✔

Two-layer Neural Networks ←
    Hidden layer options
    Output layer options

Multi-layer neural networks ←

Training a neural network by backpropagation

**Reading** HTF Ch.: 11.3 Neural networks, Murphy Ch.: (16.5 neural nets), Bach Ch.: –, Deep Learning Book (Goodfellow, Bengio, Courville) 6.1-4, ResNet 7.6, ConvNet 9 , Autoencoders 14.1, Dive Into Deep Learning 4.1-4.3.

# (Artificial) Neural Network (nn) unit

▶ For each **unit** $i$

$$y_i \equiv f_i(x) = \phi(\sum_j w_{ij} x_j + w_{i0}) \tag{1}$$

▶ **Weigth vector** $w_i$
  ▶ $w_{ij}$ = strength of the link from unit $i$ to input $j$
  ▶ $w_{ij} = 0$: no link
  ▶ $w_{ij}$ can be positive or negative
  ▶ Sometimes we call the input vector $x = [x_{1:d}]$ input units
▶ **activation function** $\phi()$
  ▶ must be non-linear (otherwise the unit is a linear transformation)
  ▶ wanted: monotonically increasing, differentiable, gradient non-zero [1]

Notation $\phi$ is overloaded

▶ When we talk about nn in general: $\phi$ is any activation function
▶ When we do calculations with nn: $\phi$ is by default the logistic function (unless specified otherwise)

$$\textbf{logistic or sigmoid function} \quad \phi(u) = \frac{1}{1 + e^{-u}} \tag{2}$$

▶ When we do statistics or ML (but not nn): $\phi$ is the logistic function

▶ Exercise: compare $f_i(x)$ from (1) with $p(x) = Pr[y = 1 \mid x]$ from logistic regression.

---

[1] More technically: $\phi$ can be any continuous, bounded and strictly increasing function on $\mathbb{R}$.
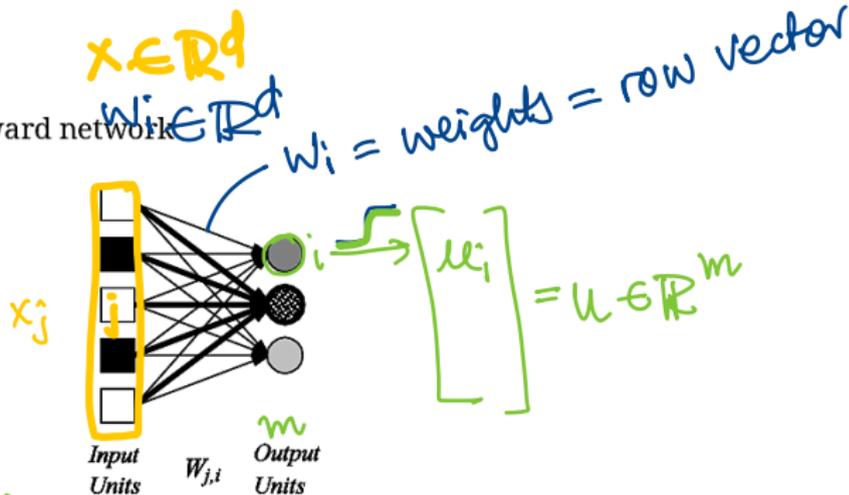
Marina Meila | CS480/680 Winter 2026: Lecture IV – Neural Networks 1

6

# Perceptron

- Single layer feed-forward network

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

$W_i$ = weights = row vector



$$W_i^T x = z_i$$

$$u_i = \varphi(z_i)$$

activation fction

$$u_i = u \in \mathbb{R}^m$$

*Input Units*    $W_{j,i}$    *Output Units*

$x_j$

UNIVERSITY OF
WATERLOO

$$W = \begin{bmatrix} W_1 \\ W_2 \\ \cdots \\ W_m \end{bmatrix} \in \mathbb{R}^{m \times d}$$

# Perceptron weights for OR, AND, XOR

## 1 unit implementation of OR



$(x_1, x_2) \in$

$f(x) = 0$

$\varphi = \dfrac{sgn(z) + 1}{2}$

$w_1 = w_2 = 1$

$w_0 = -\dfrac{1}{2}$

$f^{true} = x_1 \text{ OR } x_2 \qquad x_{1,2} \in \{0,1\}$

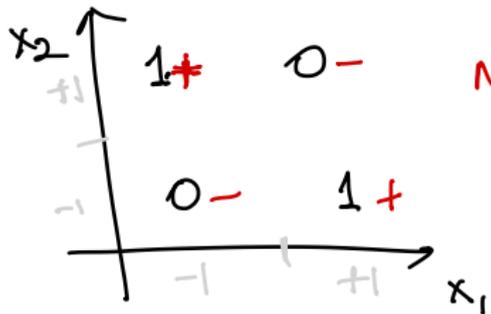$$f(x_1, x_2) = \varphi \left( \underbrace{w_1 x_1 + w_2 x_2 + w_0}_{z} \right)$$

activation

Intercept as slope "trick"

$\tilde{x} = \begin{bmatrix} 1 \\ x_1 \\ x_2 \end{bmatrix} \qquad \tilde{W} = \begin{bmatrix} w_0 \\ w_1 \\ w_2 \end{bmatrix} \Rightarrow z = \tilde{W}^T \tilde{x}$

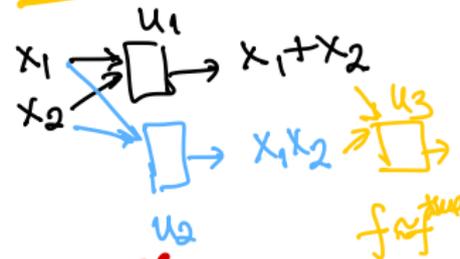# Perceptron weights for OR, AND, XOR

$x_2$

$1$ ✱     $0$ −

$0$ −     $1$ +

$x_1$

NOT
linearly
separable!

$f^{true} = x_1 \text{ XOR } x_2$

add layers

**Neural Nets**

**2)**

$x_1 \to \boxed{} \to x_1 + x_2$ $u_1$

$x_2 \to \boxed{} \to x_1 x_2$ $u_2$

$\to \boxed{} \to f \approx f^{true}$ $u_3$

$\Rightarrow \exists w_{0,1,2}$

so that
$f = f^{true}$ **on** $\{0,1\}^2$

$$f(x_1, x_2) = \varphi(w_1 x_1 + w_2 x_2 + w_0)$$

$\hookrightarrow$ monotonic ↗

**Ways around:**

**1)** add "features"

$x_3 = x_1 x_2$

$\tilde{x} = [1 \quad x_1 \quad x_2 \quad x_1 x_2]$

$\tilde{w} = [w_0 \quad w_1 \quad w_2 \quad w_3]$ → SVM

# Two-layer Neural Networks



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

| | Inputs | $x_j$ | $j = 1 : d$ |
|---|---|---|---|
| Hidden layer | Bottom layer[2] | $u_i = \phi(w_i^T x)$ | $i = 1 : m,\ w_i \in \mathbb{R}^d$ |
| | Top layer | $f = \beta^T u$ | $\beta \in \mathbb{R}^m$ |
| | Output | $f$ | $\in \mathbb{R}$ |

In other words, the neural network implements the function

$$f(x) = \sum_{i=1}^{m} \beta_i u_i = \sum_{i=1}^{m} \beta_i \phi(\sum_{j=1}^{d} w_{ij} x_j) \in (-\infty, \infty) \tag{3}$$

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

▶ As we will see shortly, in general, $f(x)$ can also be non-linear

---

[2] In neural net terminology, each variable $u_i$ 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 $x, u, y$ variable.

# Activation functions for the hidden layer

For the hidden layer, we have to choose

- number of units $m$
- activation function $\varphi$

Common activation functions

- **Functions that approximate a step function**
  - threshold function (or step function) $1$ for $u \geq 0$, and $0$ otherwise (not used)
  - logistic $\phi$
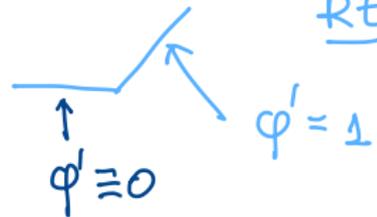  - hyperbolic tangent $tanh$, arctangent $\tan^{-1}$
- **Hinge** functions
  - RELU $= max(u, 0)$
  - softplus $= \ln(1 + e^u)$
  
    in practical implementations, these unbounded functions are bounded at a large value $M$
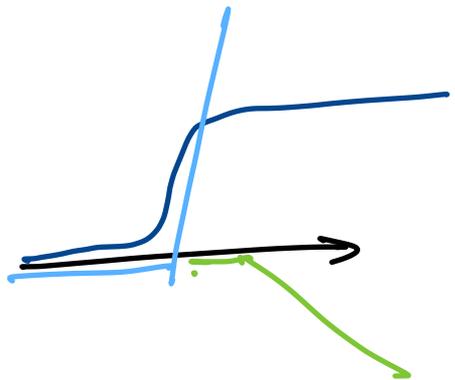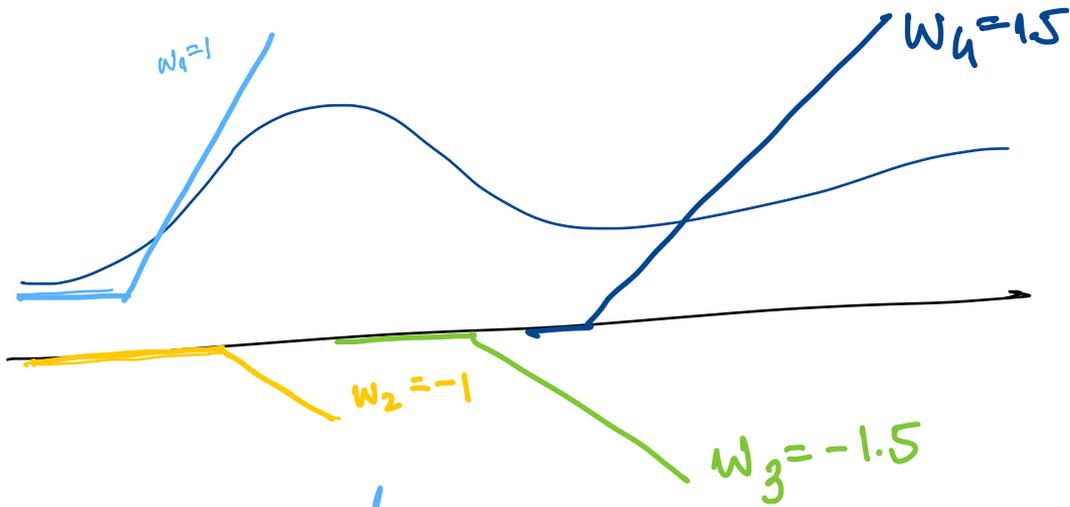  - Why hinge functions? Gradient is 1 or 0 (approximately), faster computation!!, and no saturation

THEORY $\varphi$ continuous
$\varphi$ monotonic ↗
$\varphi$ bounded

PRACTICAL $\varphi' > 0$
$\varphi'$ computed fast

REctified Linear Unit

$$\varphi'(u) = \begin{cases} 1 & \text{if } u > 0 \\ 0 & u \leq 0 \end{cases}$$

$\varphi' = 1$

$\varphi' \cong 0$

$u_i \rightarrow \acute{u}_i \rightarrow w^T u = \sum w_i u_i$
$= \sum w_i$
$u_i > 0$

$W_1=1$

$W_4=1.5$

$W_2=-1$

$W_3=-1.5$

$\wedge \int = \int + \frown$

# Output layer

$\mathbb{R}^m$  — hidden layer $\in \mathbb{R}^m$

$z \in \mathbb{R}^r$
$y \in \{1, 2, \ldots r\}$

Multi Layer

▶ Let $\boxed{z = \beta^T u}$ be the linear output of the nn.

▶ For problems other than regression, it unifies the analysis to apply a non-linear function $\phi_{\text{out}}$ to $z$. This is why, theoretically but not in practice, we will write $f(z(x)) = \phi_{\text{out}}(z(x))$

▶ Why? $\phi_{\text{out}}(z)$ matches the loss $\mathcal{L}$, which in turn matches the prediction problem

| Prediction problem | Predict |
|---|---|
| regression | $\hat{y} = z$ |
| binary classification | $\hat{y} = \mathrm{sgn} z$ |
| multiway classification | $\hat{y} = \underset{1:r}{\arg\max} z_k$ |

| Output layer $\phi_{\text{out}}$ | Range of $f$ |
|---|---|
| linear $\phi_{\text{out}} = z$ | $\in \mathbb{R}$ |
| logistic $\phi_{\text{out}} = \phi(\beta^T u)$ | $\in [0, 1]$ |
| softmax $\phi_{\text{out} k} = \phi_k(\beta_k^T u)$ | $\in [0, 1]^r$ |

USE    n.n. for prediction                          TRAINING

▶ Regression: **linear** layer as in (3) $f = \sum_i \beta_i u_i$

▶ Classification (binary): **logistic** layer $f(x) \in [0, 1]$ is interpreted as the probability of the $+$ class.

$$f(x) = \phi\left(\sum_{j=1}^m \beta_i u_i\right) = \phi\left(\sum_{i=1}^m \beta_j \phi(\sum_j w_{ij} x_j)\right) \quad (4)$$

▶ Multiway classification with $r$ classes
  ▶ Output is vector of $r$ functions $f_1, \ldots f_r$
  ▶ $f_k$ is the probability of $y = k$
  ▶ (sometimes $f_k$ can be a "confidence")
  ▶ This is done with a **softmax** layer (next page)

## OPTIONAL - Generalized Linear Models (GLM)

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

▶ $y \in \mathbb{R}$, $y \sim P_\theta$ with

$$P_\theta(y) = e^{\theta y - \psi(\theta)} \tag{8}$$

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

$$\theta = \beta^T x \tag{9}$$

▶ 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

$$l(\beta) = \ln P_\theta(y|x) = \theta y - \psi(\theta) \quad \text{where } \theta = \beta^T x \tag{10}$$

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

$$\nabla_\beta l = \nabla_\theta l \nabla_\beta (\beta^T x) = (y - \mu)x \tag{11}$$

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.]
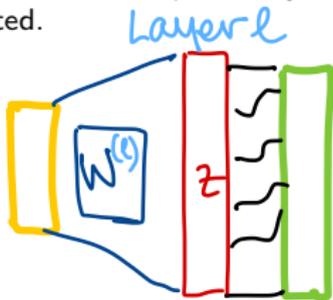
# Multi-layer/Deep neural networks

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
$x \equiv x^{(0)}, y \equiv x^{(L)}$, $m_0 = d, m_L = \dim y$ (typically 1) and define the recursion:

$$x_j^{(l)} = \phi \left( \underbrace{(w_j^{(l)})^T x^{(l-1)}}_{z^{(l)}} \right), \text{ for } j = 1 : m_l, \, l = 1 : L \tag{12}$$

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



$l = 0$
$x^{(0)} = x$
$m_0 = d$
input
layer

Layer $l$

$\ell = 1 : L$

$W^{(\ell)}$

$z$

$l = L$
$m_L = 1$ for example
$\phi$out $(z^{(L)} = \frac{1}{z^{(L)}})$

$x^{(\ell-1)}$ for ex $\in \mathbb{R}^{m_{\ell-1}}$

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

$u \in \mathbb{R}^m \rightarrow x^{(\ell)} \in \mathbb{R}^{m_\ell}$

$z = w^T x \rightarrow z^{(\ell)} \in \mathbb{R}^{m_\ell}$

$u = \varphi(z) \Longleftrightarrow u_i = \varphi(z_i) \quad i = 1 : m$

# 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* MIT Technologies from Google *can be approximated arbitrarily closely by a linear output, two layer neural* **Technology** *n (3) with a sufficiently large number of hidden units* **Review** *m.*
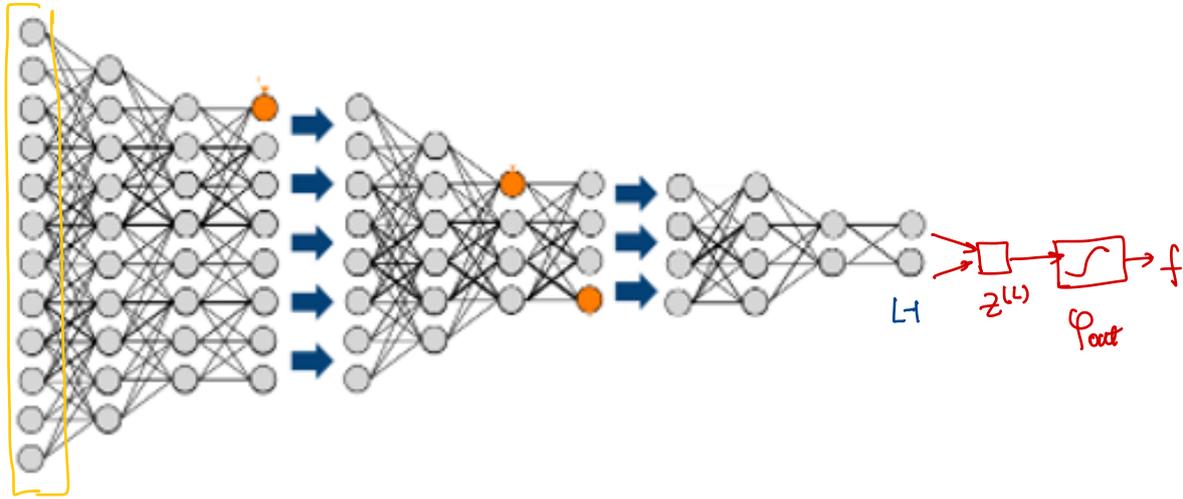
- A practical result

**10 BREAKTHROUGH TECHNOLOGIES 2013**

## 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 Bolzmann 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**

$X \equiv x^{(0)}$
input

# #Parameters

1 unit $\qquad W = \{ w \} \qquad d+1$ params

1 hidden layer $\qquad W = \{ W \in \mathbb{R}^{m \times d}, \phi \in \mathbb{R}^m \} \qquad m + md$ param

L layers $\qquad W = \{ W^{(1)}, \ldots W^{(L)} \} \qquad \sum_{l=1}^{L} m^{(l-1)} m^{(l)}$ params