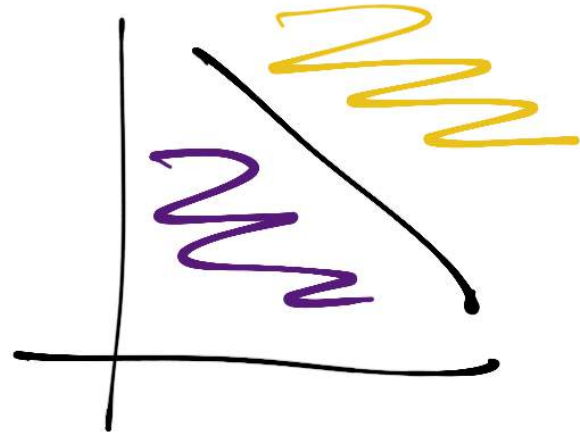# Advanced Machine Learning, 20 Aug 2019

## Feedforward Neural Networks

### Perceptron model

#### Linear separability



$$\bar{X} \longrightarrow Y$$

Introduce nonlinearity

$$\bar{X} \longrightarrow \phi(\bar{x}) \longrightarrow Y$$

$$x \rightarrow \phi(x)$$ Geometric transformation

$\hookrightarrow$ Use "kernel trick"

Find $K$ s.t. $K(x,y) = \phi(x) \cdot \phi(y)$

How to discover $\phi$? (Or, equivalently, $K$)?
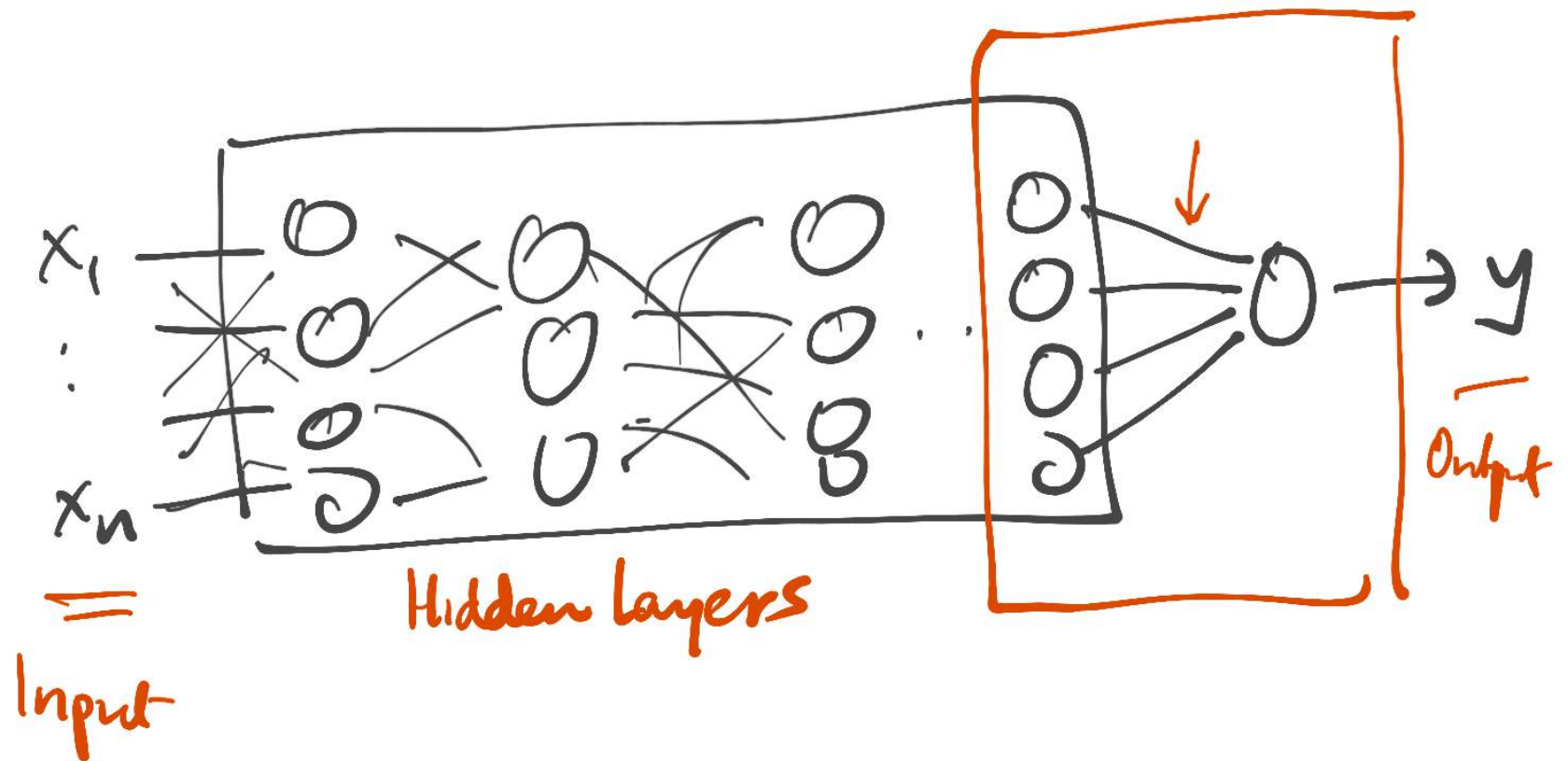
Brute Force
Very high dimension (infinite)
RBF (radial basis)

Hand craft
by experience

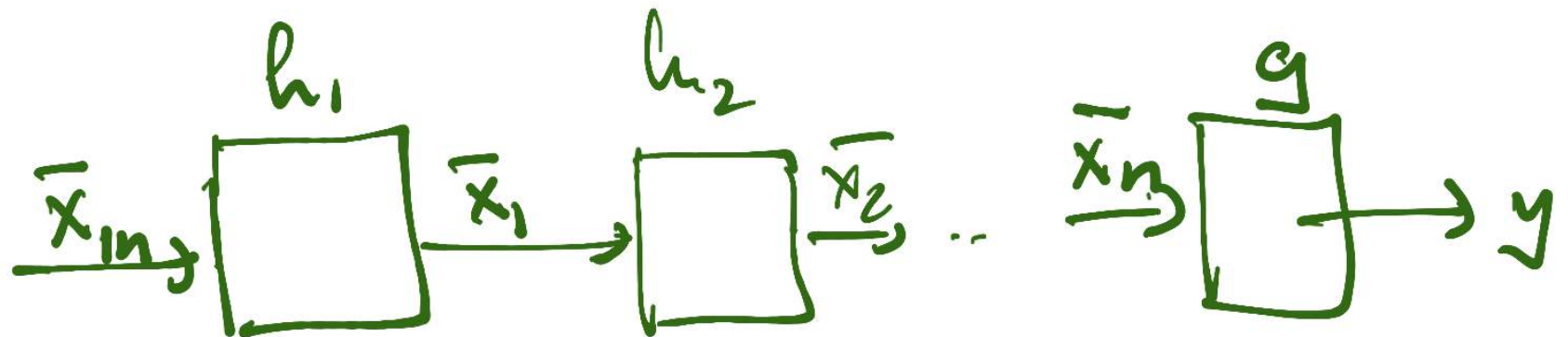# Feedforward NN — network of perceptrons

## Directed acyclic graph



$x_1$
$\vdots$
$x_n$

Input

Hidden layers

$y$

Output

Hidden layers compute $\phi(x)$

NN $\rightarrow$ learn $\phi$ from training data

Weights of edges in the hidden layers
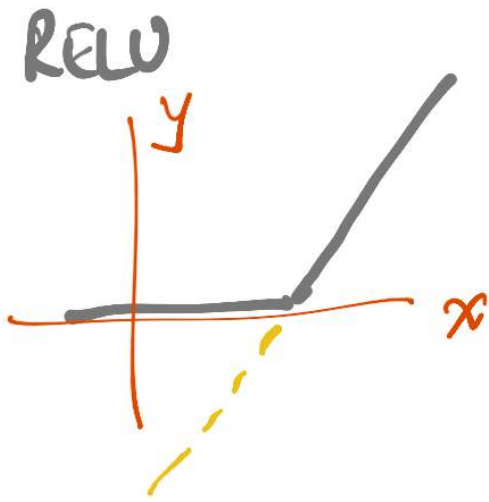


$$y = W^T x + b$$

$g(x$

Cascaded computation

$$g\left( h_n\left( \cdots h_2\left( h_1\left( \bar{x}_{in} \right) \right) \right) \right)$$

Composition of linear fns is a linear fn

Need to introduce non-linearity (in the hidden layers)
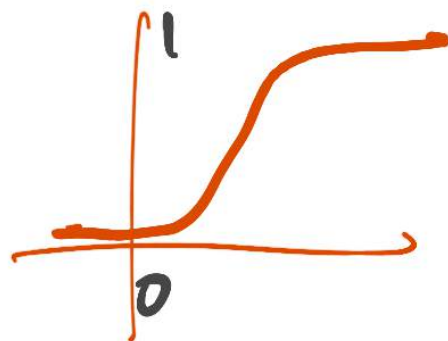
Rectified Linear Unit (RELU)

$$\max\left(0, \underbrace{W^T x + c}_{g(x)}\right) + b$$

RELU
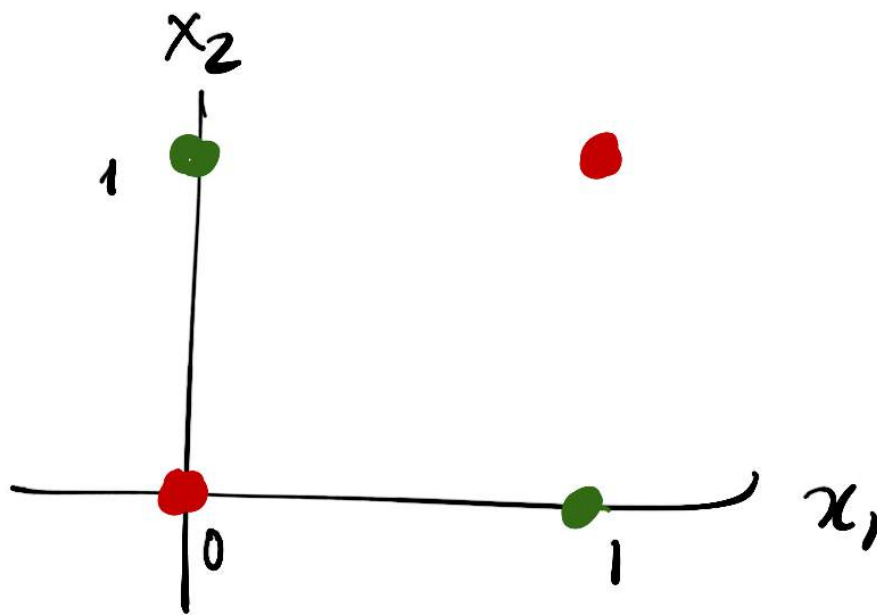
# Sigmoid

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

$$z = w^T x + c$$

Classical example $XOR(x_1, x_2)$

$$x_1, x_2 \in \{0, 1\}$$

# One hidden layer

$$\bar{x} \xrightarrow{\quad} \boxed{h} \xrightarrow{h(\bar{x})} \boxed{g} \xrightarrow{g \cdot h(\bar{x})}$$

$$Wx + c \qquad\qquad\qquad Wx + b$$

## Solution for XOR

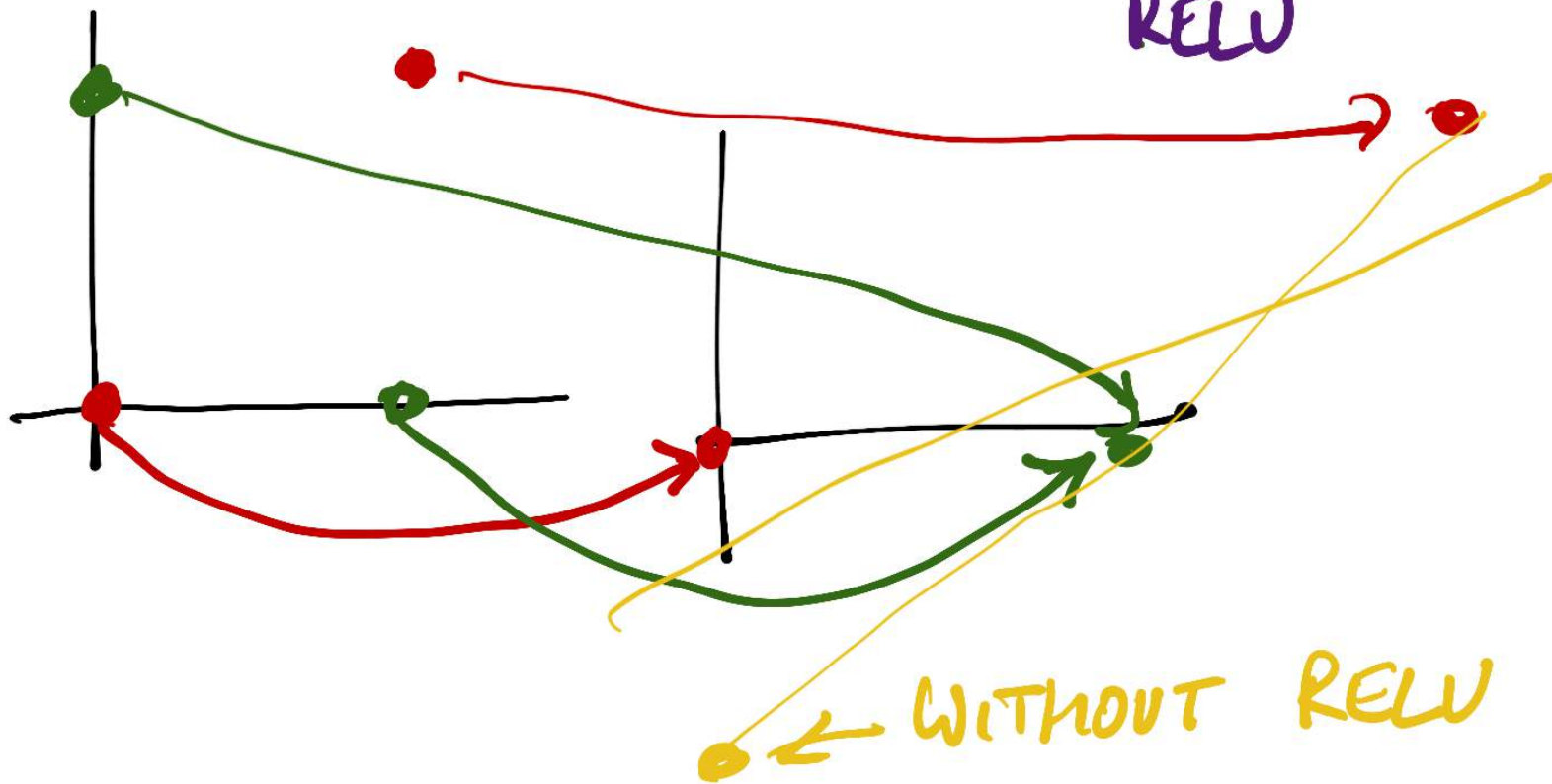$$W = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix} \qquad\qquad w = \begin{bmatrix} 1 \\ -2 \end{bmatrix}$$

$$c = \begin{bmatrix} 0 \\ -1 \end{bmatrix} \qquad\qquad b = 0$$

$$\begin{bmatrix} 0 & 0 \\ 0 & 1 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix} + \begin{bmatrix} 0 & -1 \end{bmatrix} = \begin{bmatrix} 0 & -1 \\ 1 & 0 \\ 1 & 0 \\ 2 & 1 \end{bmatrix} \begin{bmatrix} 0 & 0 \\ 1 & 0 \\ 1 & 0 \\ 2 & 1 \end{bmatrix}$$

$\xrightarrow{\text{RELU}}$



← WITHOUT RELU

Goal — find W through learning

1. Gradient descent

Loss function

Back propagation → Calculating gradients for hidden weights

Stochastic gradient descent

2. Choice of "activation functions"

3. Architecture — layers, connectivity etc

# Gradient descent

## Loss function / Cost

Shift weights according to gradient wrt loss function

## Mean Squared Error (MSE)

## Cross entropy

## Why these functions ?

# Maximum likelihood estimators (MLE)

Assuming we are computing probabilistic output. $P(y | \bar{x})$

Sequence of Heads & Tail $\Rightarrow$ Observation

$$\Rightarrow p^* = \frac{\#\text{Heads}}{\text{Total}} \quad \nearrow^{0}$$

$$n \text{ heads } t \text{ tail}$$

Hypotheses $h$ for the coin probability

$$h^n (1-h)^t$$

$$\arg\max_h \; p(O|h) = p^*$$

NN computes MLE for $p(y|\bar{x})$

$p^*(y|\bar{x})$ — What NN reports

"Real" $p(y|\bar{x})$

Want $p^*$ to be as close to $p$ as possible

# Information theory

Entropy $\quad H = -\sum p \log(p)$

$$= \mathbb{E}_{x \sim p} \log(x)$$

Minimizes encoding of given distribution

Optimal encoding for $Q$, apply it to $P$

KL distance

$$D_{\mathrm{KL}}(P \| Q) = \mathbb{E}_{x \sim P}\left[\log \frac{P(x)}{Q(x)}\right] = \mathbb{E}_{x \sim P}[\log P(x) - \log Q(x)].$$

# Cross Entropy

$$H(P, Q) = H(P) + D_{\mathrm{KL}}(P \| Q)$$

$$-\sum_{p} \log(P) + \sum_{p} \log(P) - \sum_{p} \log(Q)$$

$$-\sum_{p} \log(Q)$$

NN learns $p^*(y|x)$

Distance from $p(Y|x)$

$$H(p, p^{\neq}) = -E_{x \sim p} \log(p^{\neq})$$

Minimy $-E_{x \sim p} \cdots \Rightarrow$ Maximize $E_{x \sim p} \log(p^{\neq})$

Back to MLE

Observations $O_1, O_2, \ldots, O_n$

hypothesis $h$

$$\underset{h}{\text{argmax}} \prod_{i=1}^{n} P(O_i | h)$$

$$\underset{h}{\text{argmax}} \prod_{i=1}^{n} P(O_i | h)$$

$$= \underset{h}{\text{argmax}} \log \left( \prod_{i=1}^{n} P(O_i | h) \right)$$

$$= \underset{h}{\text{argmax}} \frac{1}{n} \sum_{i=1}^{n} \log \left( P(O_i | h) \right)$$

$$= \mathbb{E}_{x \sim p} \log(p^*)$$

# Cross Entropy is equivalent to MLE
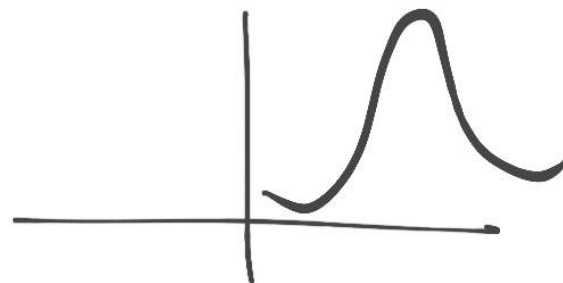
Specific cases

Linear regression

$P(y|x)$

$x \rightarrow y^*$

$\downarrow$

$N(\text{---}, \sigma)$

Normal distribution centered at computed $y^*$

$\downarrow$

MSE

# 2 way sigmoid

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

$$P(y=1) = y$$

$$P(y=0) = 1-y$$

$$P^*(x) = \sigma(z) = a$$

$$\sum p(x) \underset{a}{\overset{\log}{p^*}}(x) = y \log(a) + (1-y)\log(1-a)$$