

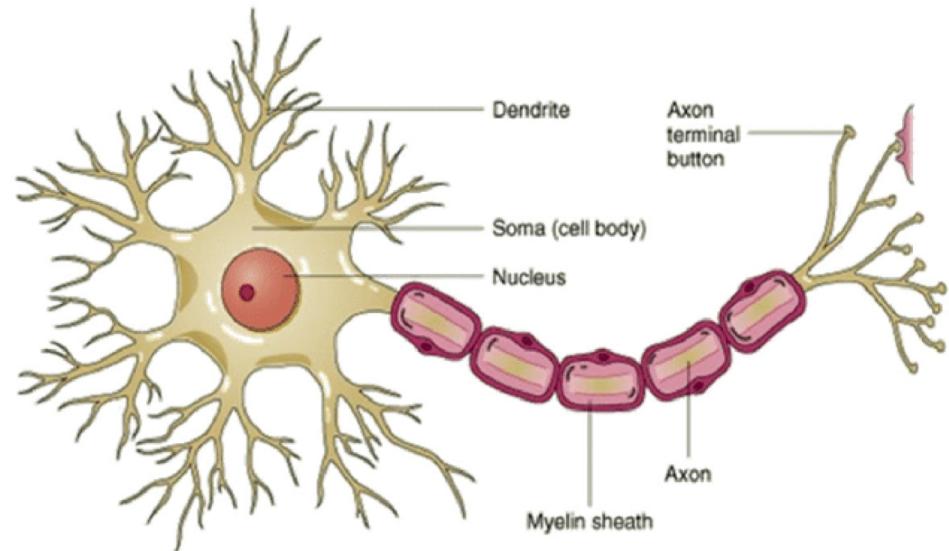
Deep feed-forward networks

Goodfellow et al. 2016 (sec. 6)

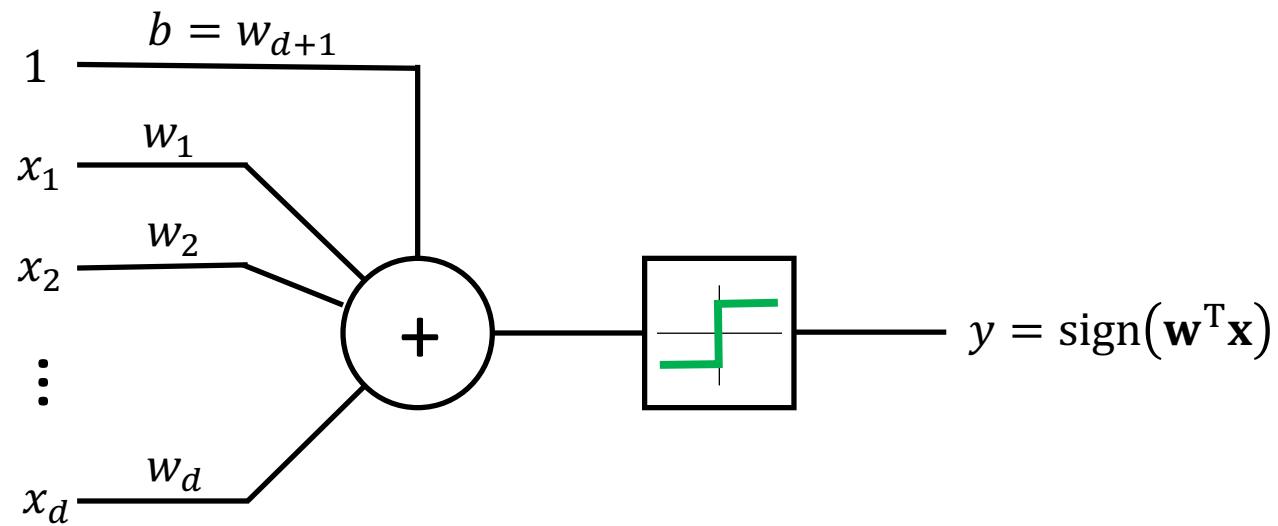
Perceptron



Rosenblatt 1957



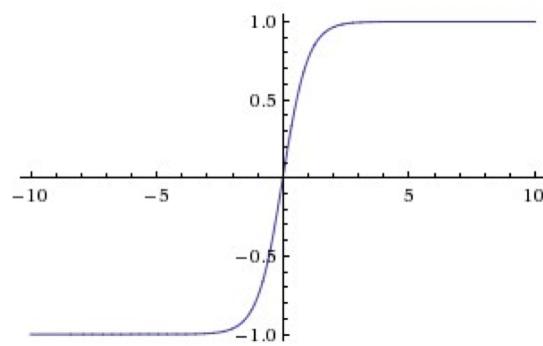
Perceptron



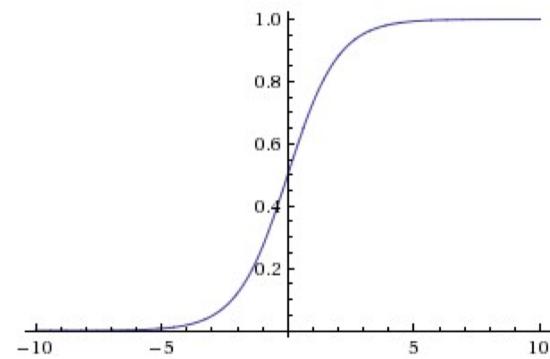
* Note the class labels here are +1 and -1 instead of +1 and 0 as before
Rosenblatt 1957

Activation functions

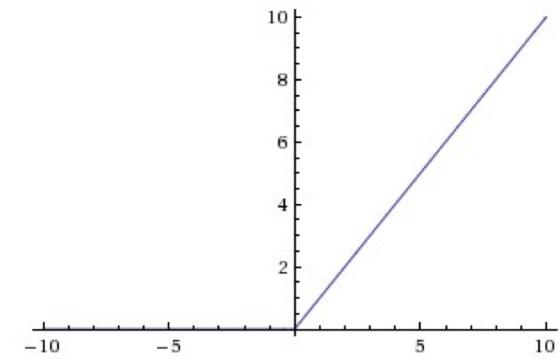
tanh



sigmoid



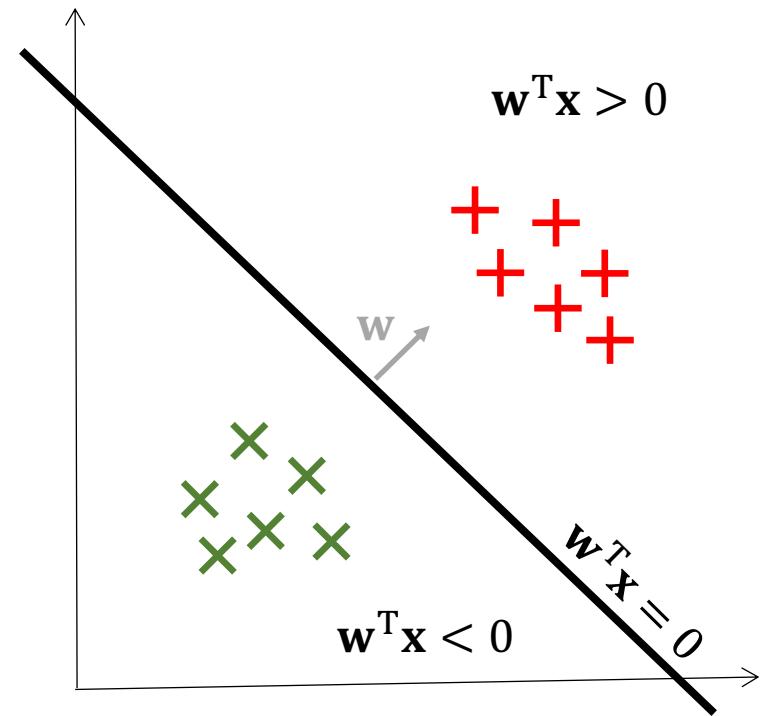
ReLU



Original Perceptron learning algorithm

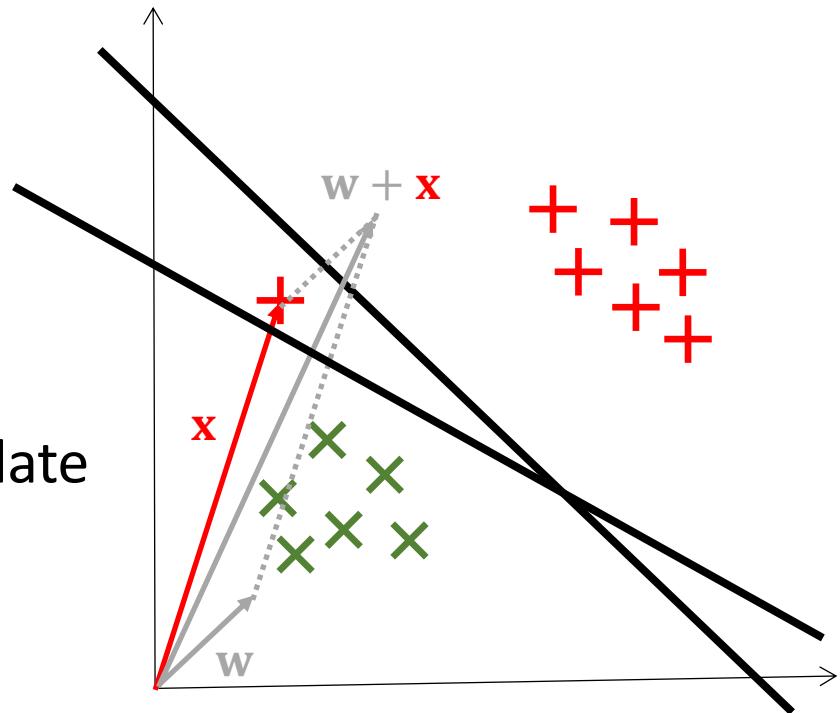
- Initialize $\mathbf{w} = \mathbf{0}$
- Given example \mathbf{x} , output

$$\hat{y} = \text{sign}(\mathbf{w}^T \mathbf{x})$$



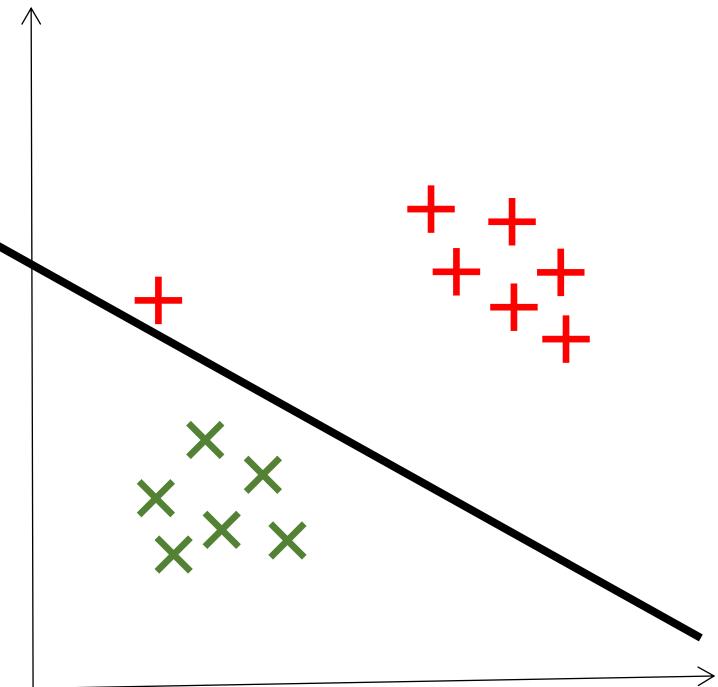
Original Perceptron learning algorithm

- Initialize $\mathbf{w} = \mathbf{0}$
- Given example \mathbf{x} , output
$$\hat{y} = \text{sign}(\mathbf{w}^T \mathbf{x})$$
- If predicted wrong ($\hat{y} \neq y$) update
 - If $y = 1$: $\mathbf{w} \leftarrow \mathbf{w} + \mathbf{x}$



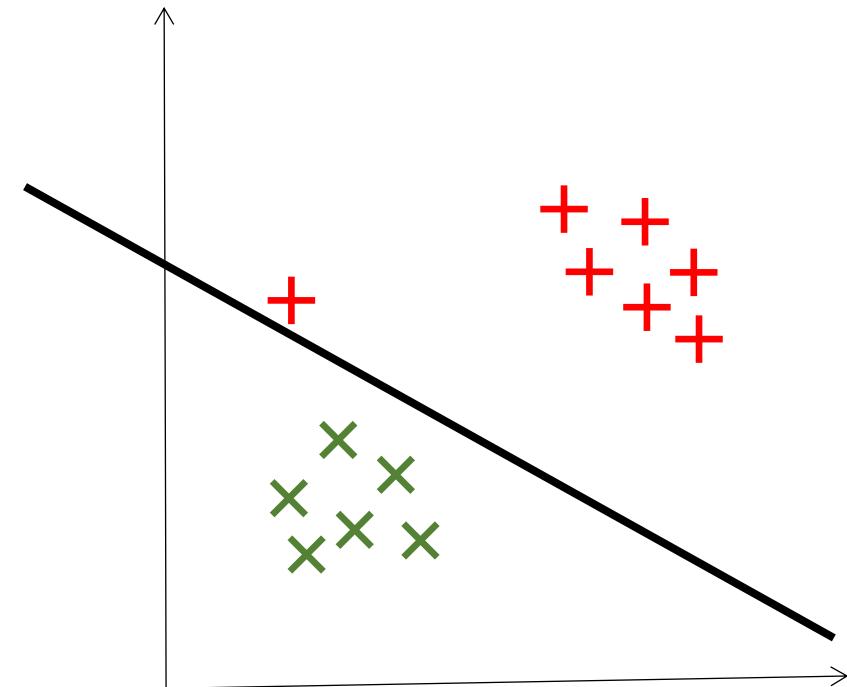
Original Perceptron learning algorithm

- Initialize $\mathbf{w} = \mathbf{0}$
- Given example \mathbf{x} , output
$$\hat{y} = \text{sign}(\mathbf{w}^T \mathbf{x})$$
- If predicted wrong ($\hat{y} \neq y$) update
 - If $y = 1$: $\mathbf{w} \leftarrow \mathbf{w} + \mathbf{x}$
 - If $y = -1$: $\mathbf{w} \leftarrow \mathbf{w} - \mathbf{x}$



Original Perceptron learning algorithm

- Initialize $\mathbf{w} = \mathbf{0}$
- Given example \mathbf{x} , output
$$\hat{y} = \text{sign}(\mathbf{w}^T \mathbf{x})$$
- Update
$$\mathbf{w} \leftarrow \mathbf{w} + y\mathbf{x} \mathbb{I}_{\text{sign}(\mathbf{w}^T \mathbf{x}) \neq y}$$



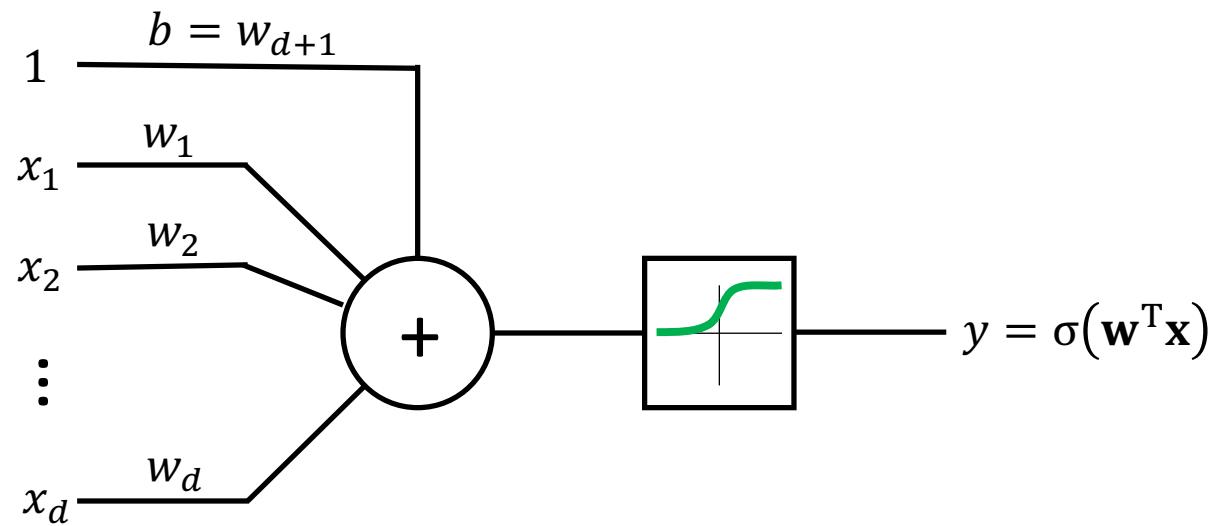
Perceptron training as SGD

$$\hat{L}(\mathbf{w}) = - \sum_{i=1}^n y_i \mathbf{w}^T \mathbf{x}_i \mathbb{I}_{\text{sign}(\mathbf{w}^T \mathbf{x}_i) \neq y_i} \rightarrow \min$$

- SGD with step size $\alpha^{(t)} = 1$

$$\begin{aligned}\mathbf{w}^{(t+1)} &= \mathbf{w}^{(t)} - \nabla_{\mathbf{w}} \left(-y_t \mathbf{w}^T \mathbf{x}_t \mathbb{I}_{\text{sign}(\mathbf{w}^T \mathbf{x}_t) \neq y_t} \right) \\ &= \mathbf{w}^{(t)} + y_t \mathbf{x}_t \mathbb{I}_{\text{sign}(\mathbf{w}^T \mathbf{x}_t) \neq y_t}\end{aligned}$$

Logistic regression



Logistic regression with SGD

$$\hat{L}(\mathbf{w}) = -\frac{1}{n} \sum_{i=1}^n \log \sigma(y_i \mathbf{w}^T \mathbf{x}_i) \rightarrow \min$$

- Computing the gradients

$$\begin{aligned}\nabla_{\mathbf{w}} \log \sigma(y_i \mathbf{w}^T \mathbf{x}_i) &= \frac{\nabla_{\mathbf{w}} \sigma(y_i \mathbf{w}^T \mathbf{x}_i)}{\sigma(y_i \mathbf{w}^T \mathbf{x}_i)} \\ &= \frac{\sigma(y_i \mathbf{w}^T \mathbf{x}_i)(1 - \sigma(y_i \mathbf{w}^T \mathbf{x}_i))}{\sigma(y_i \mathbf{w}^T \mathbf{x}_i)} y_i \mathbf{x}_i\end{aligned}$$

* Note the class labels here are +1 and -1. Since $\sigma(-x) = 1 - \sigma(x)$, we can write the log-likelihood function as

$$\hat{L}(\mathbf{w}) = -\frac{1}{n} \sum_{i:y_i=1} \log \sigma(\mathbf{w}^T \mathbf{x}_i) - \frac{1}{n} \sum_{i:y_i=-1} \log(1 - \sigma(\mathbf{w}^T \mathbf{x}_i)) = -\frac{1}{n} \sum_{i=1}^n \log \sigma(y_i \mathbf{w}^T \mathbf{x}_i)$$

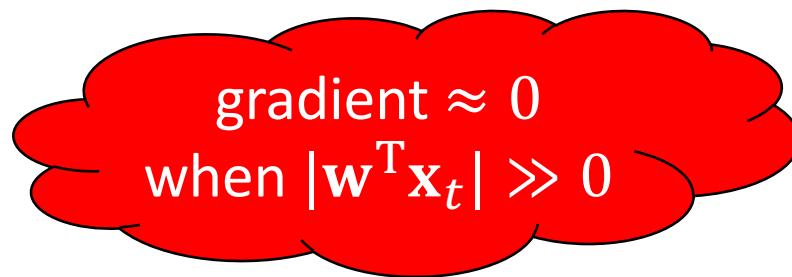
** Derivative of sigmoid $\sigma'(x) = \sigma(x)(1 - \sigma(x))$

Logistic regression with SGD

$$\hat{L}(\mathbf{w}) = -\frac{1}{n} \sum_{i=1}^n \log \sigma(y_i \mathbf{w}^T \mathbf{x}_i) \rightarrow \min$$

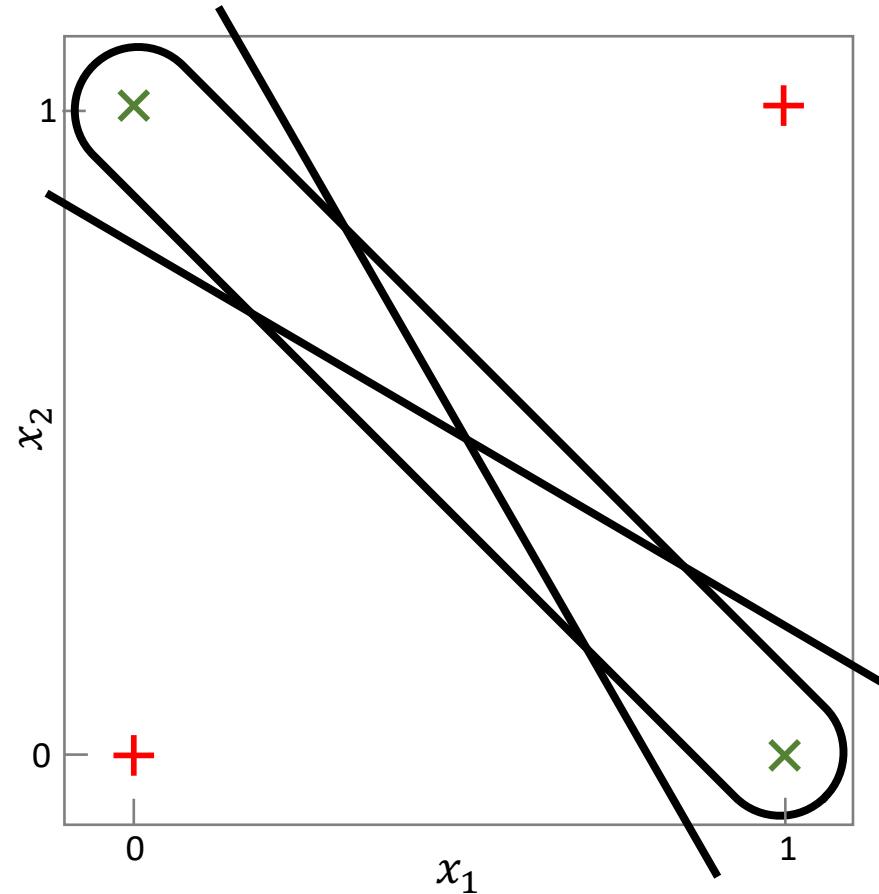
- SGD step

$$\mathbf{w}^{(t+1)} = \mathbf{w}^{(t)} + \frac{\alpha^{(t)}}{n} \frac{\sigma(y_t \mathbf{w}^T \mathbf{x}_t)(1 - \sigma(y_t \mathbf{w}^T \mathbf{x}_t))}{\sigma(y_t \mathbf{w}^T \mathbf{x}_t)} y_t \mathbf{x}_t$$



gradient ≈ 0
when $|\mathbf{w}^T \mathbf{x}_t| \gg 0$

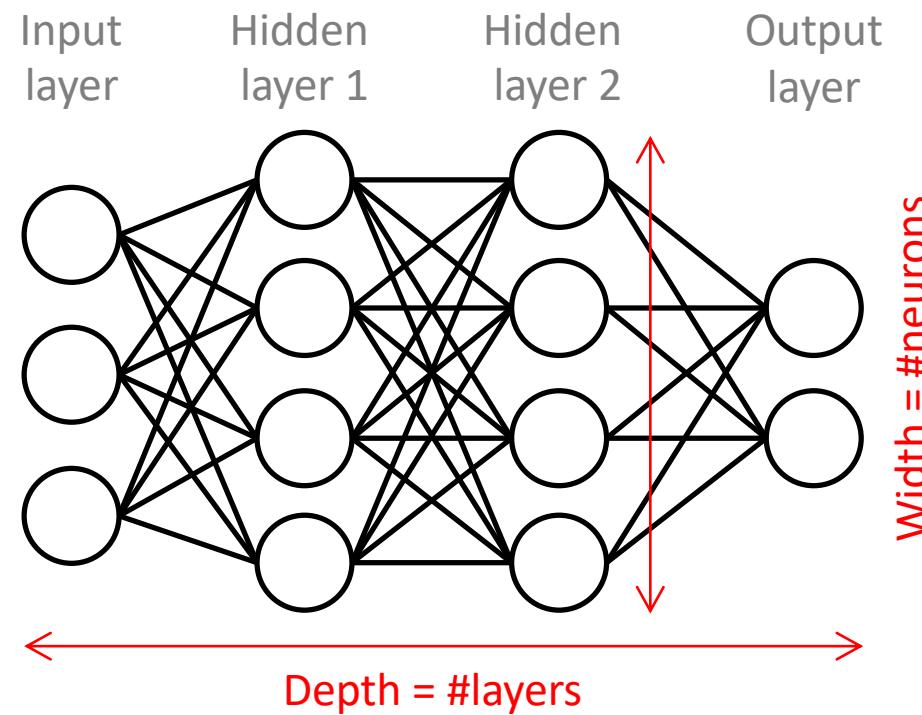
The XOR problem



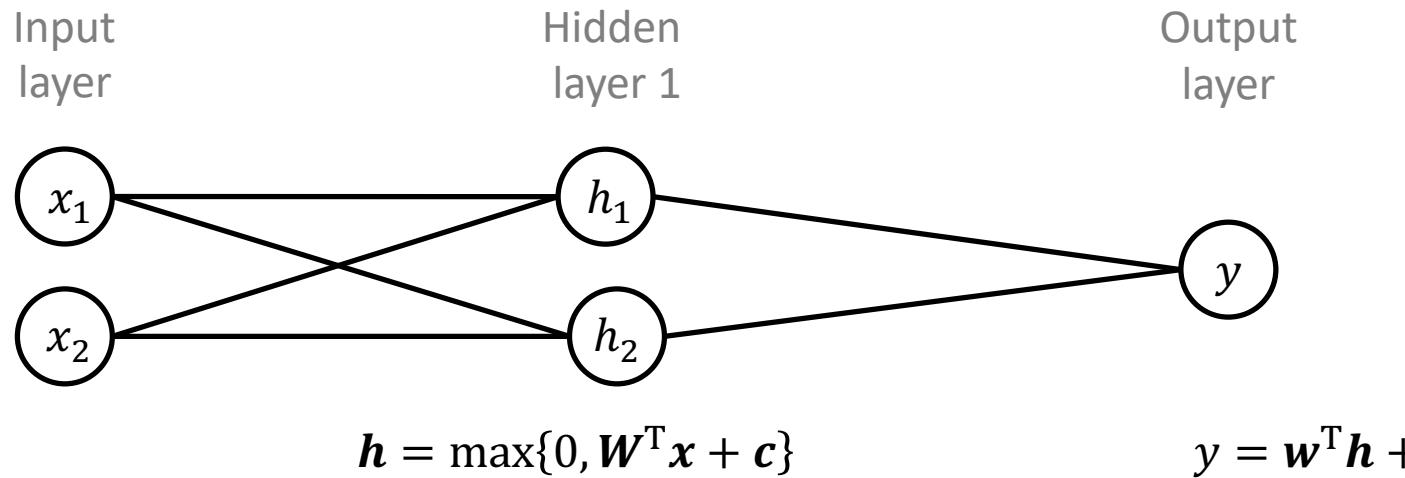
Not linearly separable!

Minsky, Papert 1969

Multi-layer perceptron



Solving XOR with one hidden layer



Input:

$$\mathbf{X} = \begin{bmatrix} 0 & 0 \\ 0 & 1 \\ 1 & 0 \\ 1 & 1 \end{bmatrix}$$

Hidden layer:

$$\mathbf{w} = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix} \quad \mathbf{c} = \begin{bmatrix} 0 \\ -1 \end{bmatrix}$$

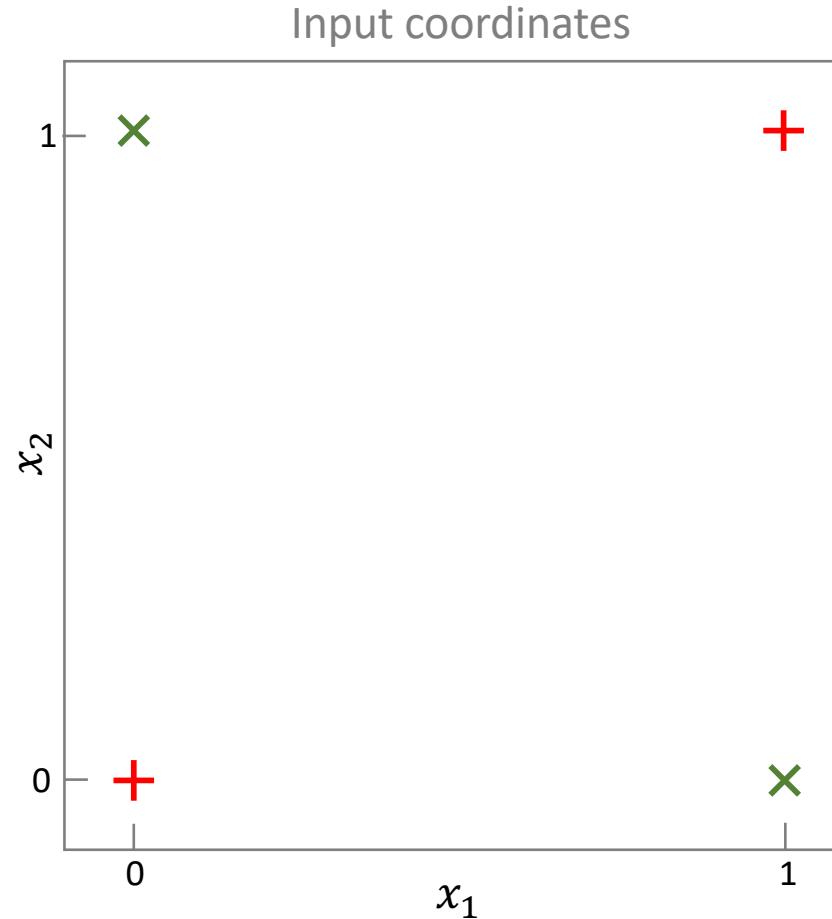
$$\mathbf{H} = \max \left\{ 0, \begin{bmatrix} 0 & -1 \\ 1 & 0 \\ 1 & 0 \\ 2 & 1 \end{bmatrix} \right\} = \begin{bmatrix} 0 & 0 \\ 1 & 0 \\ 1 & 0 \\ 2 & 1 \end{bmatrix}$$

Output layer:

$$\mathbf{w} = \begin{bmatrix} 1 \\ -2 \end{bmatrix} \quad b = 0$$

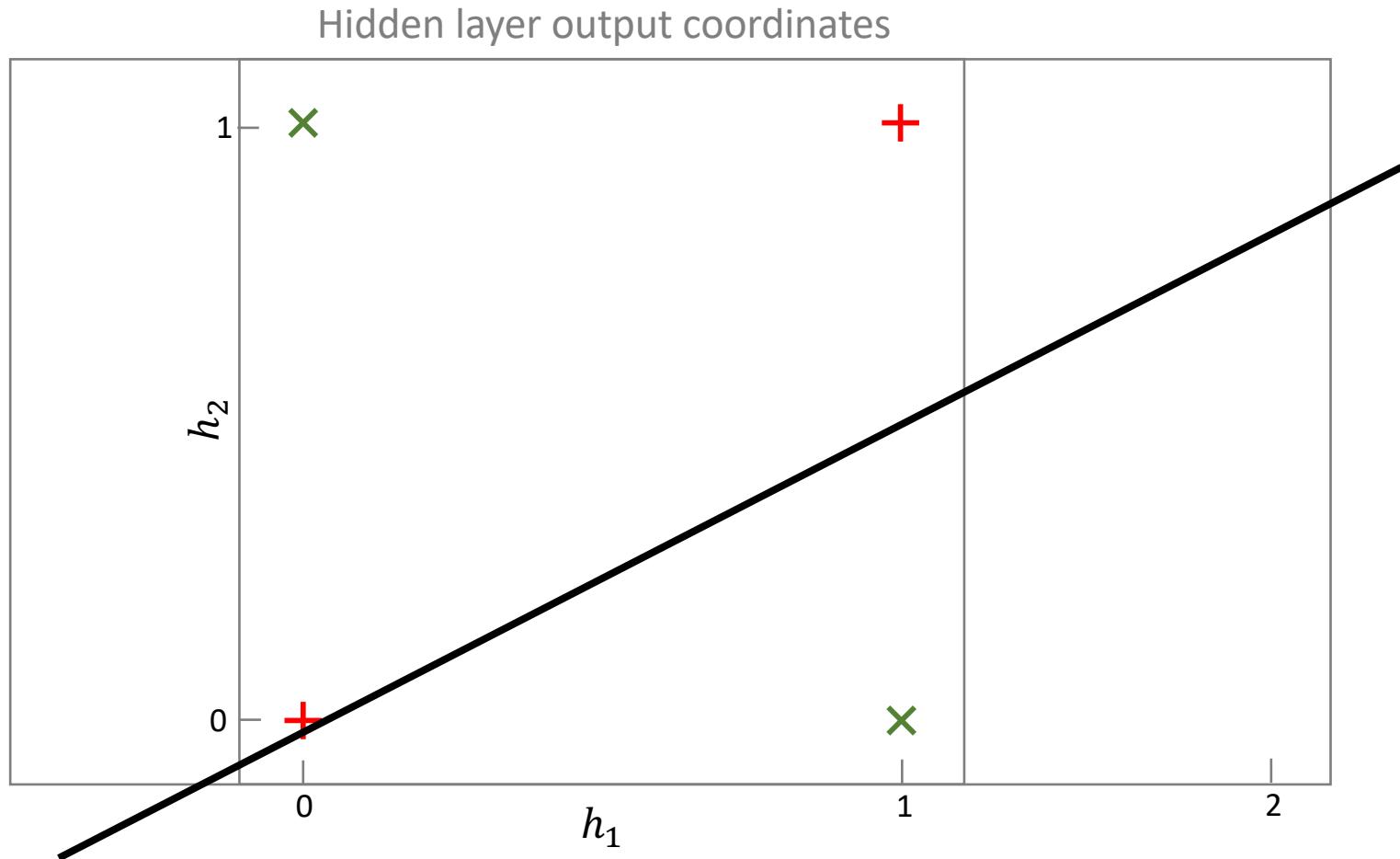
$$\mathbf{Y} = \begin{bmatrix} 0 \\ 1 \\ 1 \\ 0 \end{bmatrix} \quad (\text{Learned XOR})$$

Solving XOR with one hidden layer



Example: Goodfellow et al.

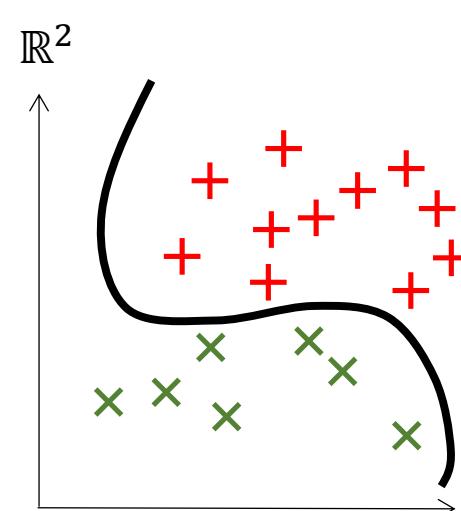
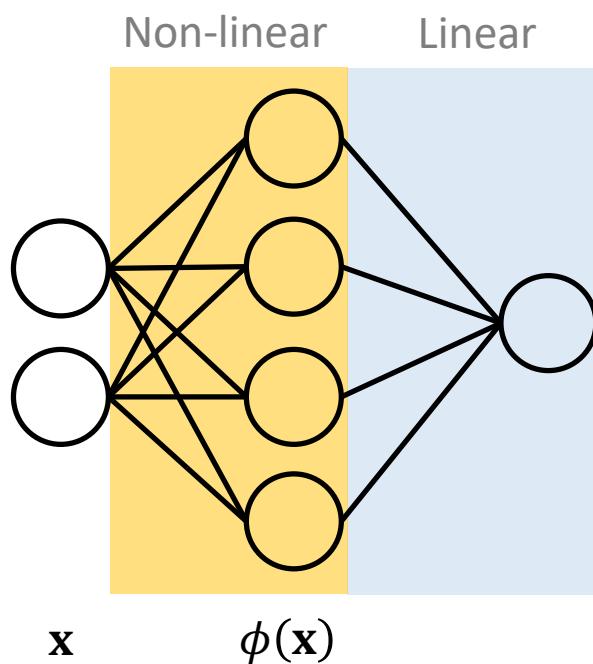
Solving XOR with one hidden layer



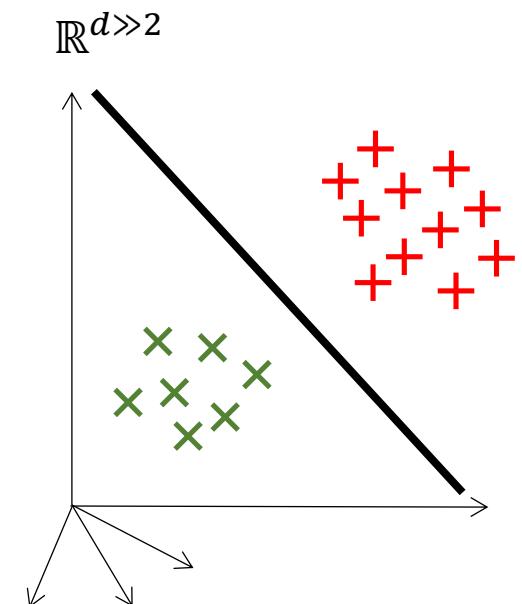
linearly separable!

Example: Goodfellow et al.

Representation learning



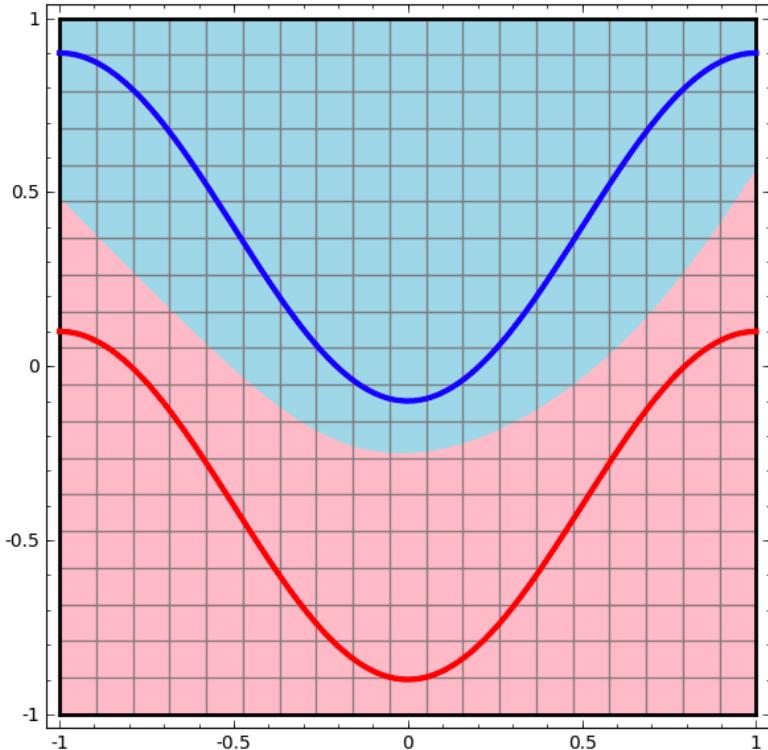
Original data \mathbf{x}



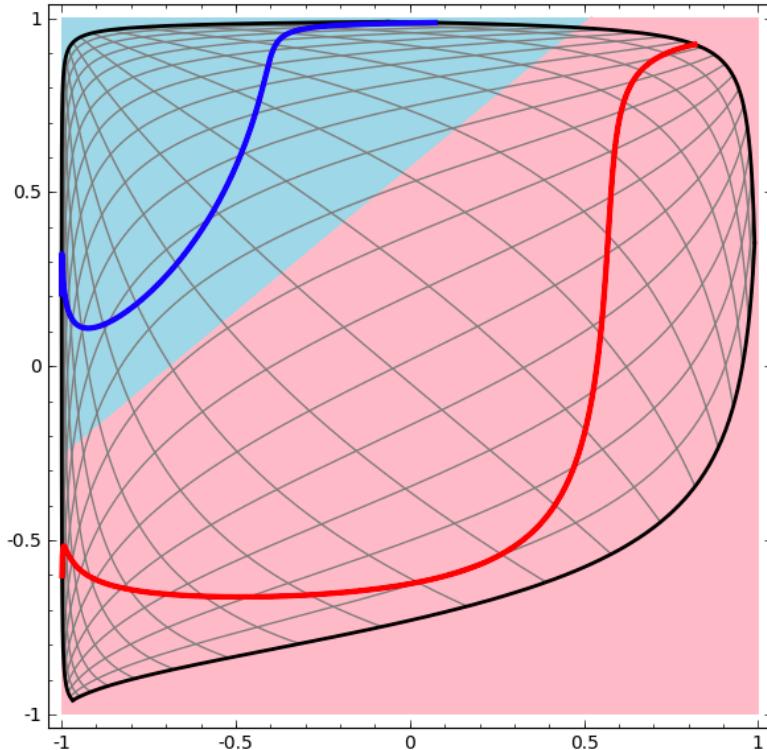
Non-linear features $\phi(\mathbf{x})$

Representation learning

Input coordinates



Hidden layer coordinates



Hidden layer learns a representation
that is linearly separable

Image credit: Christopher Olah

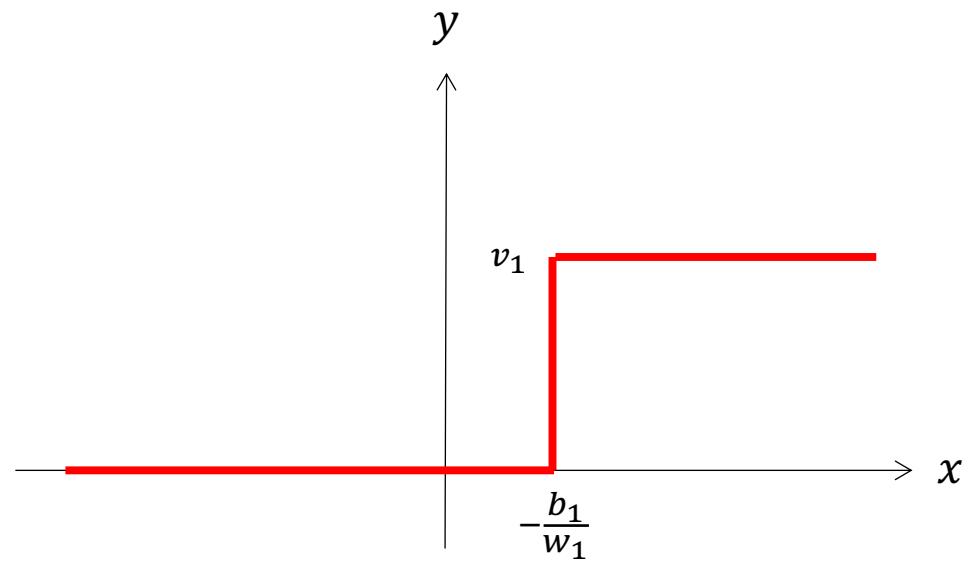
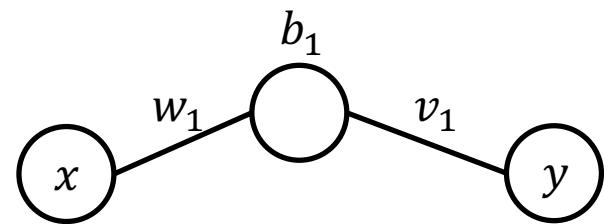
Universal approximation

A perceptron with one hidden layer of finite width can **arbitrarily accurately approximate any continuous function** on a compact subset of \mathbb{R}^n (under mild assumptions on the activation function)

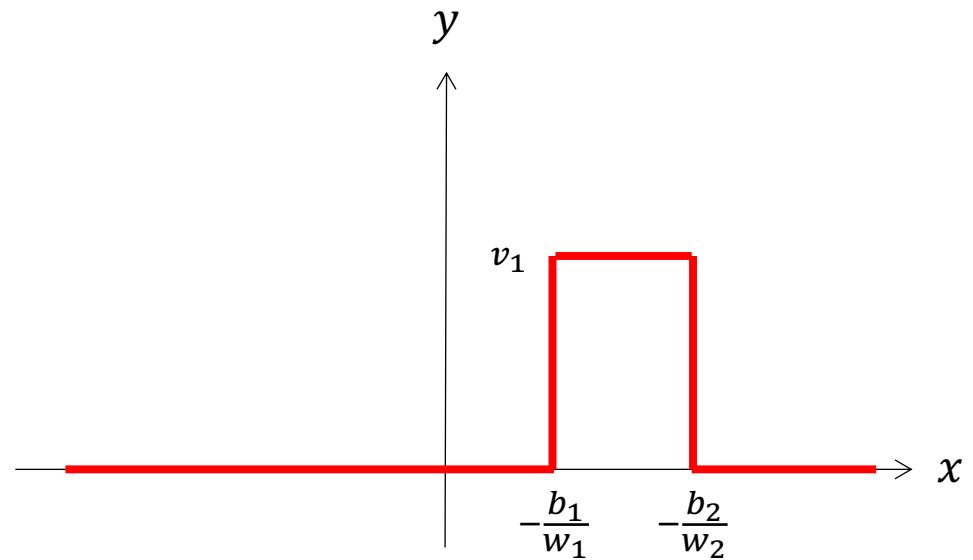
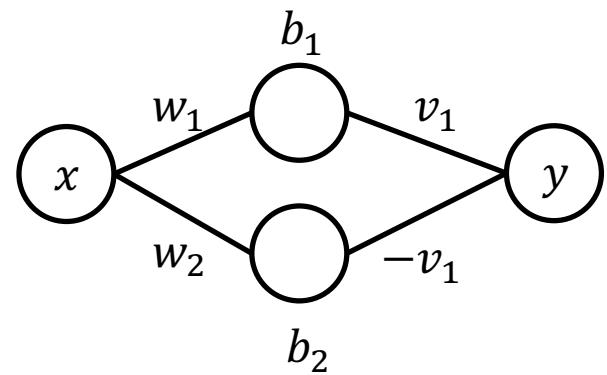
- 😊 one hidden layer is enough to represent a reasonable function
- 😢 how to find the optimal weights?
- 😢 how wide should the layer be?
- 😡 width grows exponentially with dimensions

*More generally, a Borel-measurable function
Cybenko 1989; Hornik 1991

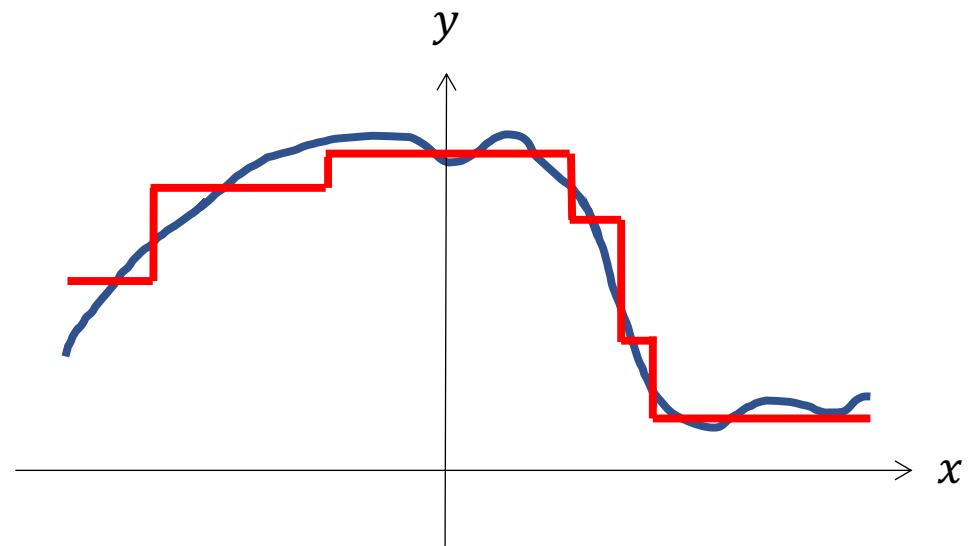
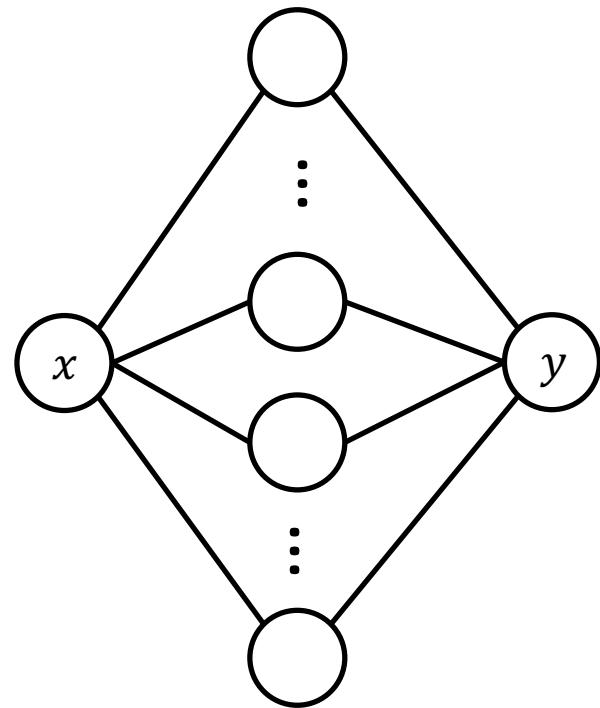
Intuition



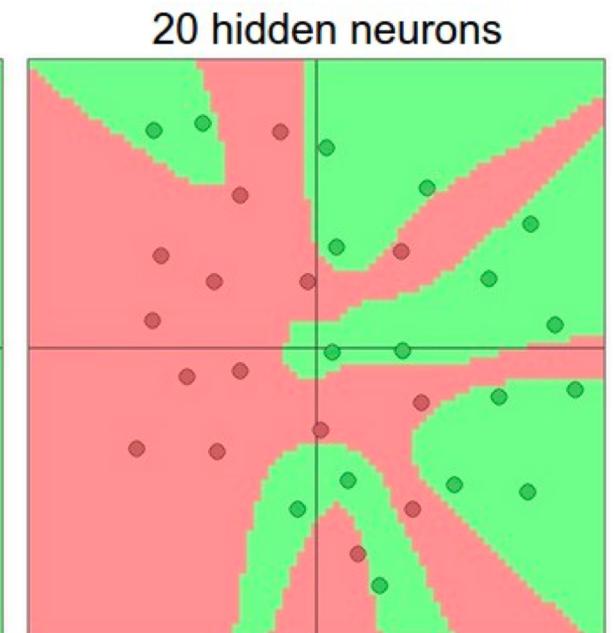
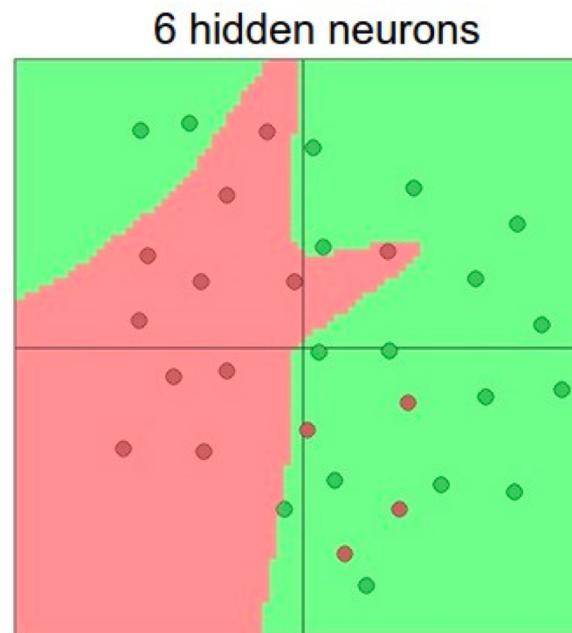
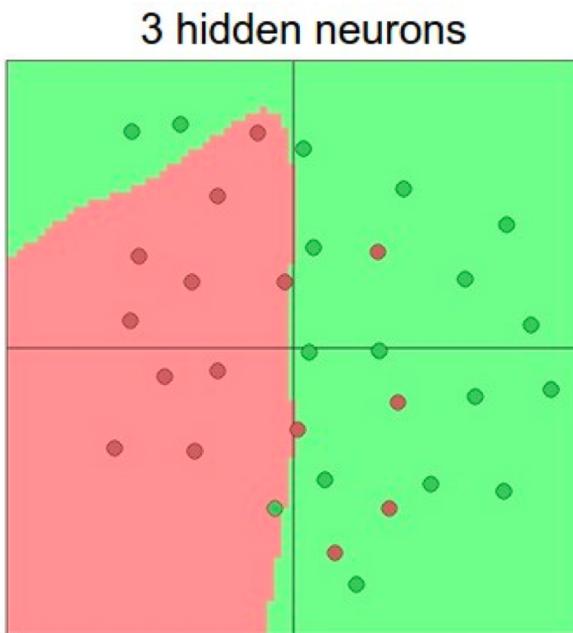
Intuition



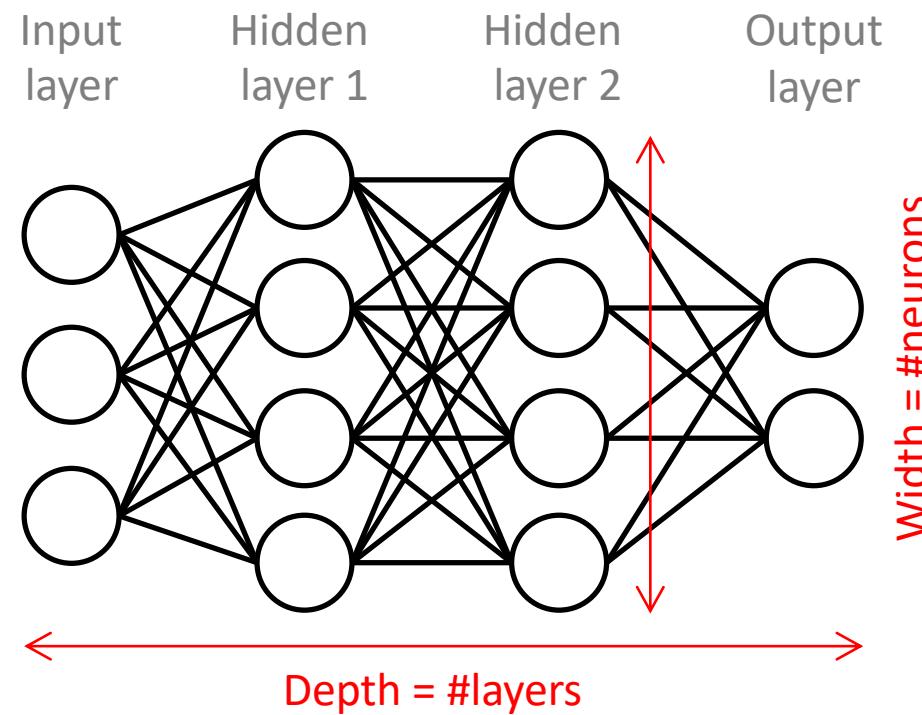
Intuition



Universal approximation

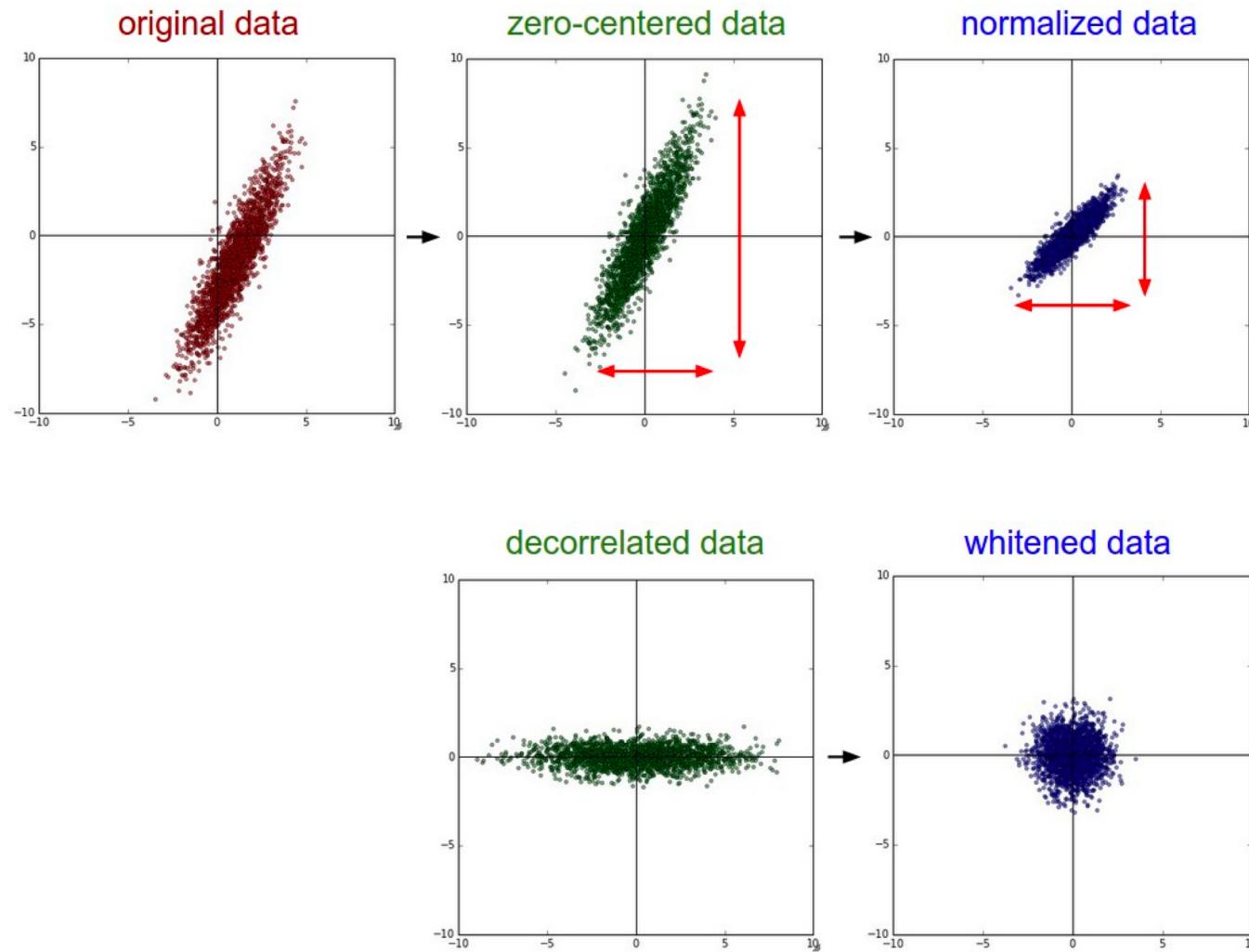


Depth vs Width



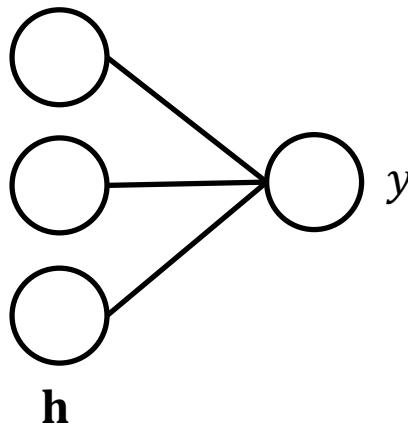
What is better: depth or width?

Input pre-processing



Popular output layer types

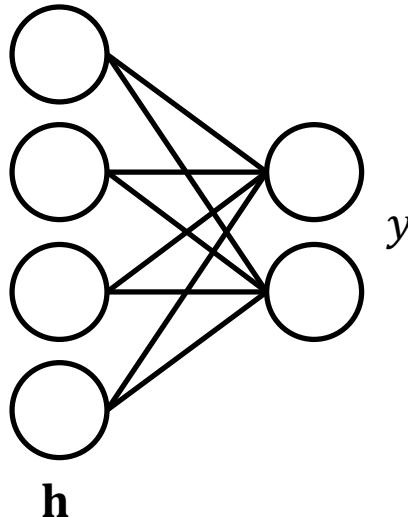
- 1D regression
- Binary classification



$$y = \mathbf{w}^T \mathbf{h}$$

$$y = \sigma(\mathbf{w}^T \mathbf{h})$$

- nD regression
- Softmax (multi-class classification)



$$\mathbf{y} = \mathbf{W}^T \mathbf{h}$$

$$\mathbf{y} = \frac{e^{\mathbf{w}_i^T \mathbf{h}}}{\sum_i e^{\mathbf{w}_i^T \mathbf{h}}}$$