

ECE 884 Deep Learning

Lecture 10: Neural Networks

02/18/2021

Review of last lecture

- Training
 - GD
 - SGD
 - Momentum, Adam, AdamW

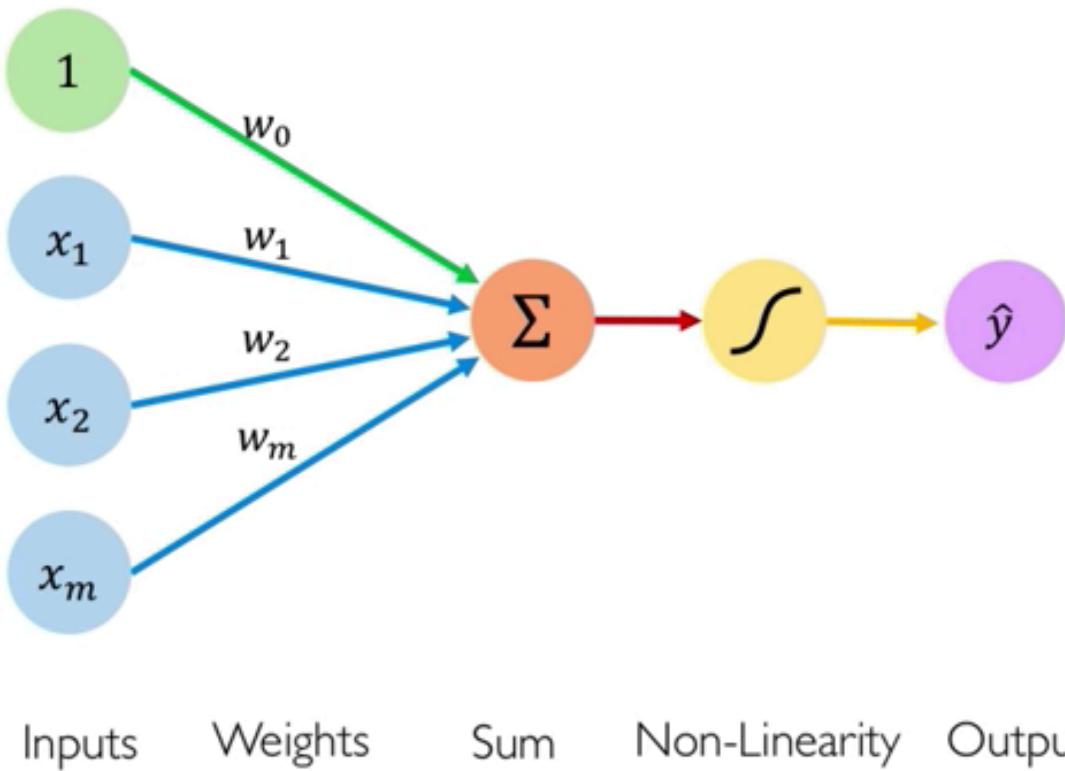
Today's lecture

- (Deep) Neural Networks
 - Perceptron
 - Multilayer Perceptron (MLP)
 - Activation Functions
 - Why Neural Networks are powerful

Perceptron

The structural building block of deep learning

Perceptron



Linear combination of inputs

Output

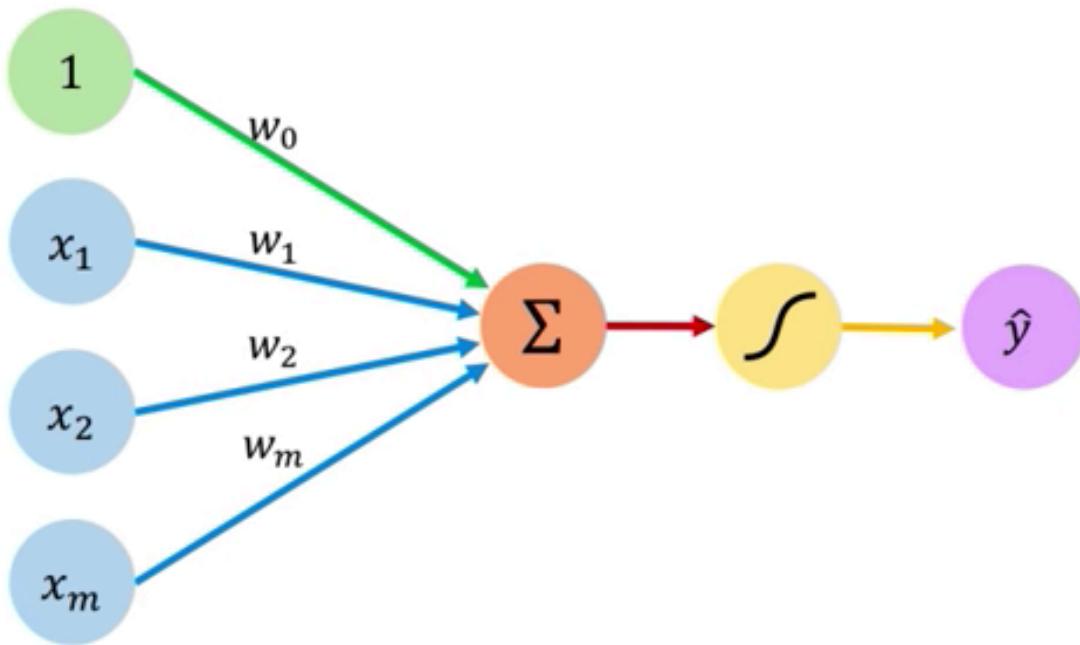
$\hat{y} = g \left(w_0 + \sum_{i=1}^m x_i w_i \right)$

Non-linear activation function

Bias

Diagram illustrating the mathematical formula for a Perceptron's output. The output \hat{y} is the result of applying a non-linear activation function g to the linear combination of inputs. The linear combination is the sum of the bias w_0 and the weighted sum of inputs x_i and their corresponding weights w_i .

Perceptron



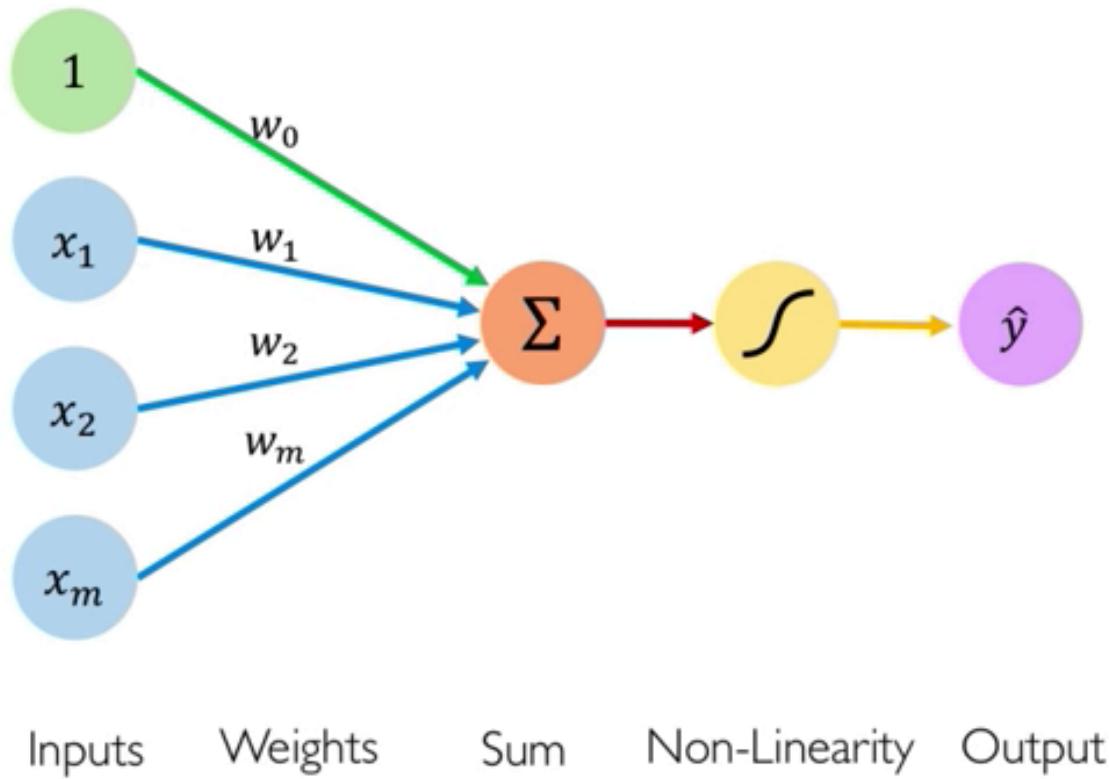
Inputs Weights Sum Non-Linearity Output

$$\hat{y} = g \left(w_0 + \sum_{i=1}^m x_i w_i \right)$$

$$\hat{y} = g (w_0 + \mathbf{X}^T \mathbf{W})$$

where: $\mathbf{X} = \begin{bmatrix} x_1 \\ \vdots \\ x_m \end{bmatrix}$ and $\mathbf{W} = \begin{bmatrix} w_1 \\ \vdots \\ w_m \end{bmatrix}$

Perceptron

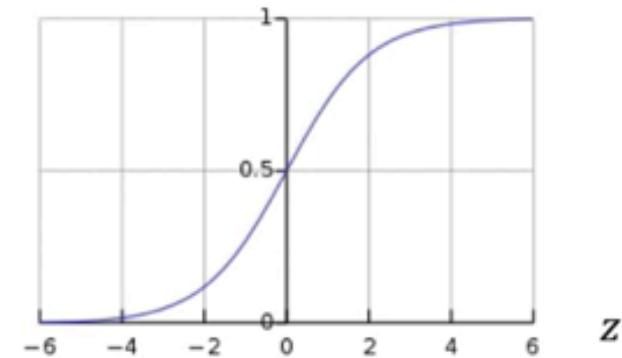


Activation Functions

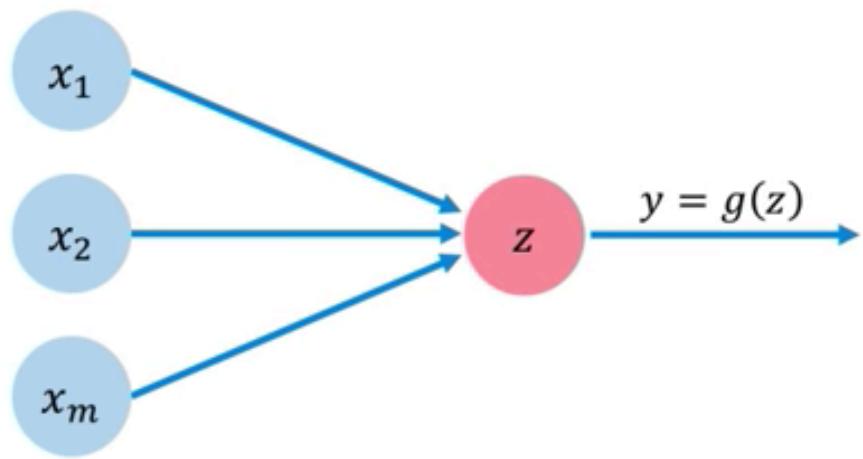
$$\hat{y} = g(w_0 + X^T \mathbf{w})$$

- Example: sigmoid function

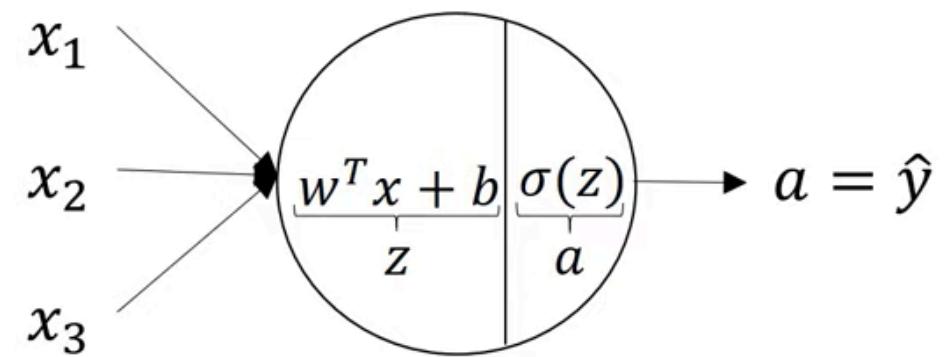
$$g(z) = \sigma(z) = \frac{1}{1 + e^{-z}}$$



Building Neural Networks with Perceptron



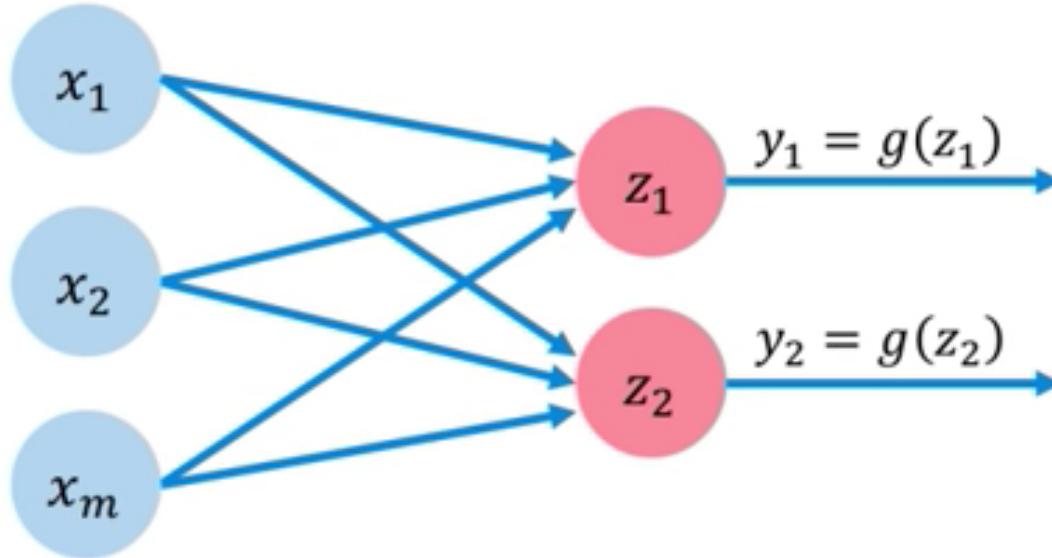
$$z = w_0 + \sum_{j=1}^m x_j w_j$$



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

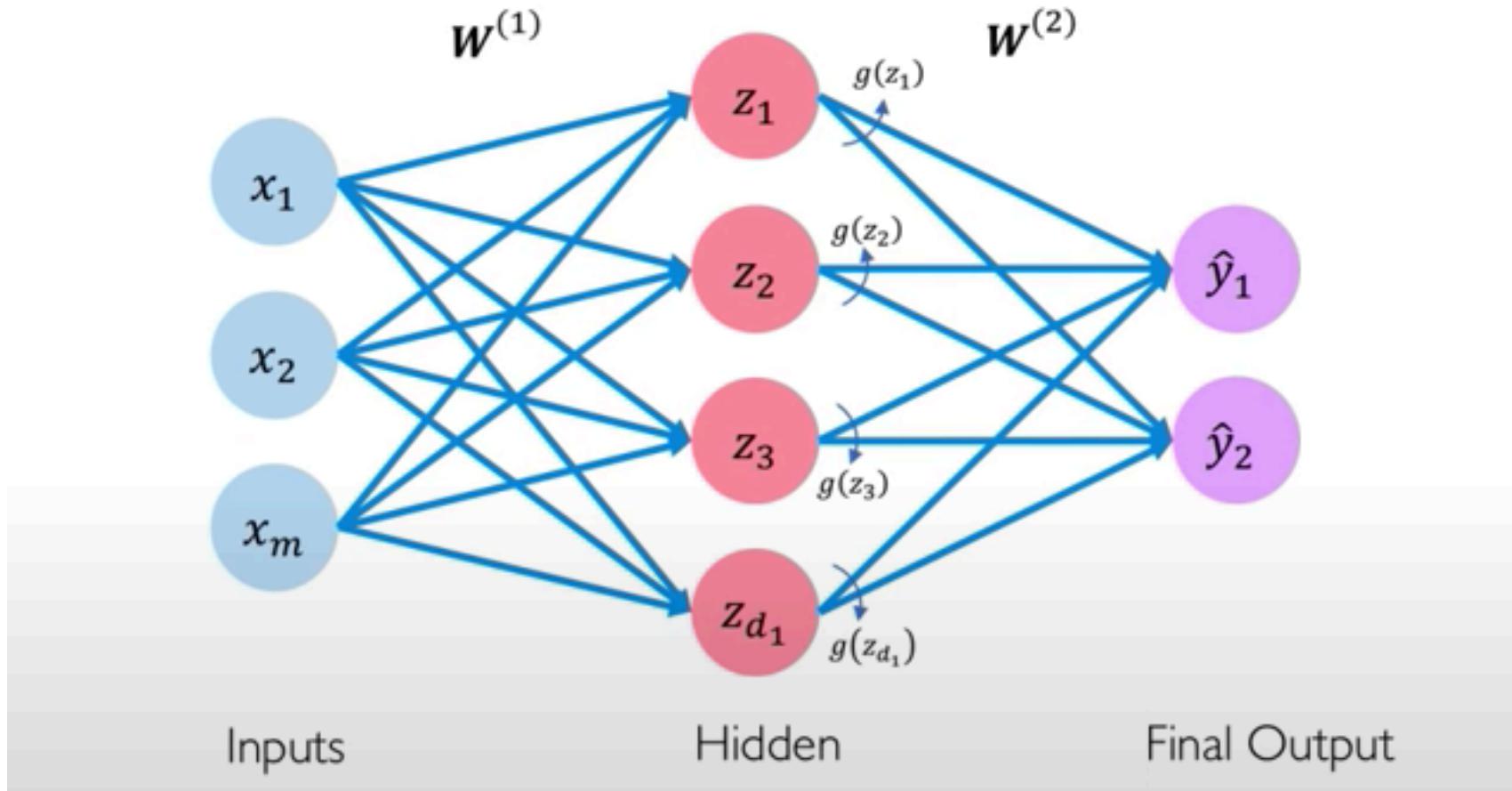
$$a = \sigma(z)$$

Building Neural Networks with Perceptron



$$z_i = w_{0,i} + \sum_{j=1}^m x_j w_{j,i}$$

2-layer Neural Networks



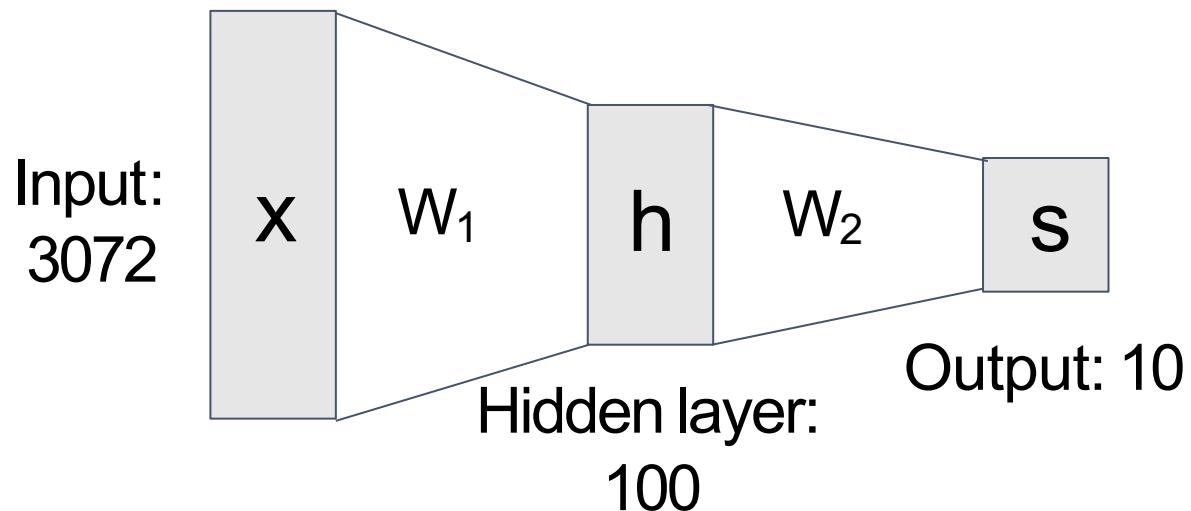
2-layer Neural Networks

Before: Linear classifier

$$f(x) = Wx + b$$

Now: 2-layer Neural Network

$$f(x) = W_2 g(W_1 x + b_1) + b_2$$



2-layer Neural Networks

Before: Linear classifier

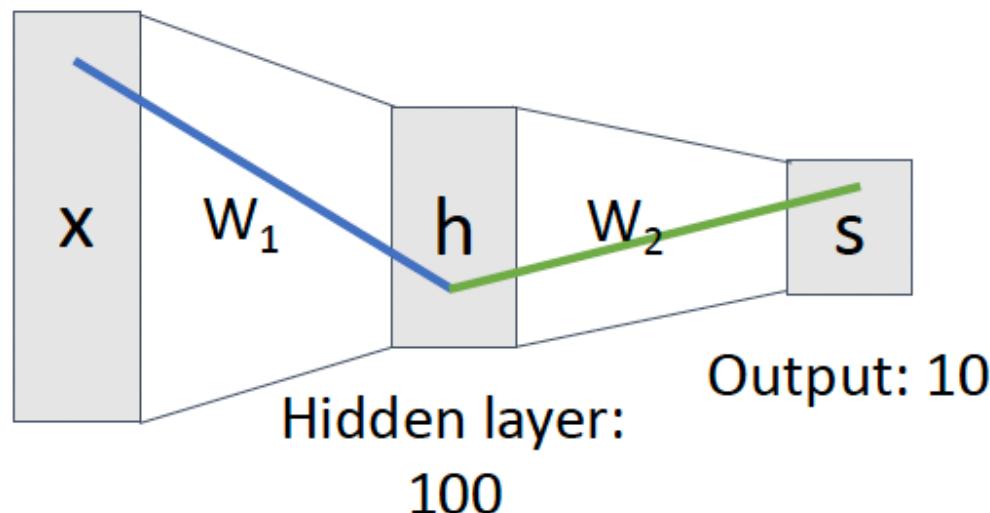
$$f(x) = Wx + b$$

Now: 2-layer Neural Network

$$f(x) = W_2 g(W_1 x + b_1) + b_2$$

Element (i, j)
of W_1 gives
the effect on
 h_i from x_j

Input:
3072



Element (i, j)
of W_2 gives
the effect on
 s_i from h_j

2-layer Neural Networks

Before: Linear classifier

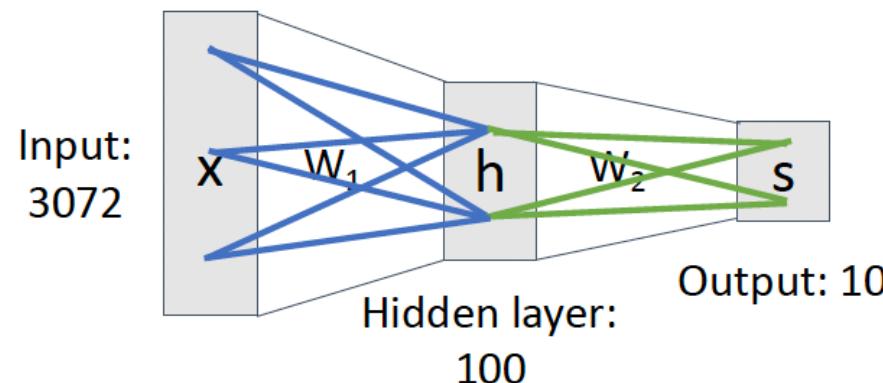
$$f(x) = Wx + b$$

Now: 2-layer Neural Network

$$f(x) = W_2 g(W_1 x + b_1) + b_2$$

Element (i, j) of W_1
gives the effect on
 h_i from x_j

All elements
of x affect all
elements of h



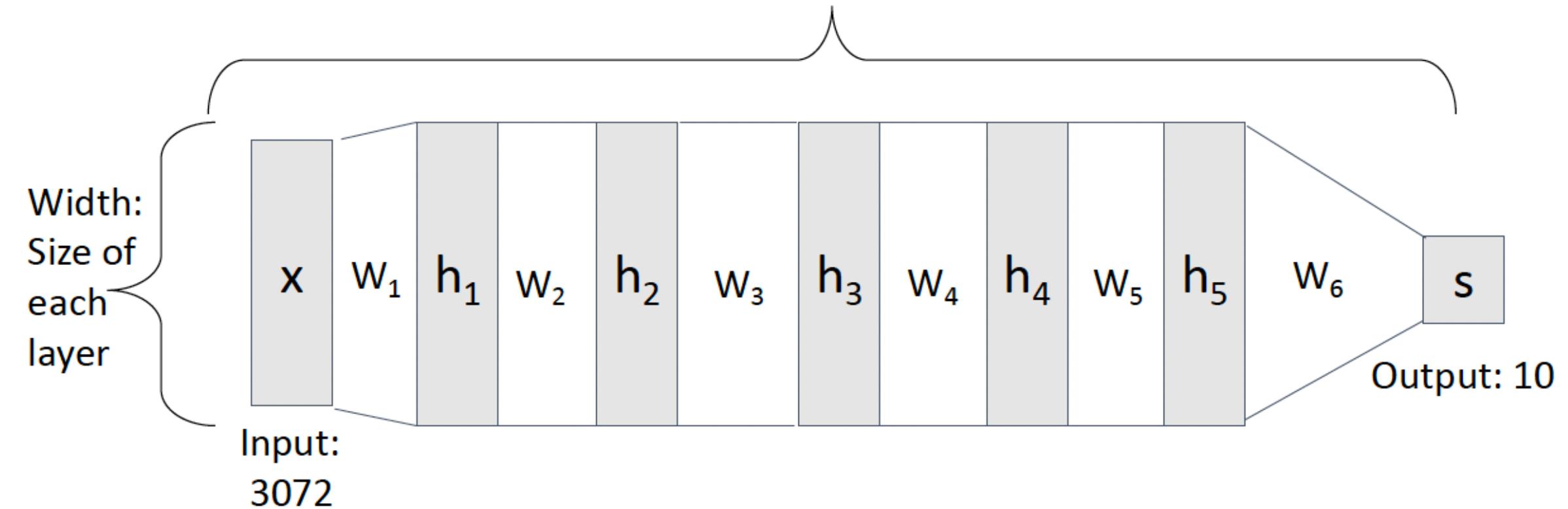
Element (i, j) of W_2
gives the effect on
 s_i from h_j

All elements
of h affect all
elements of s

Fully-connected neural network
Also “Multi-Layer Perceptron” (MLP)

Deep Neural Networks

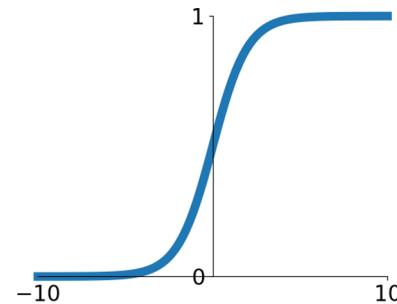
Depth = number of layers



Activation Functions

Sigmoid

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

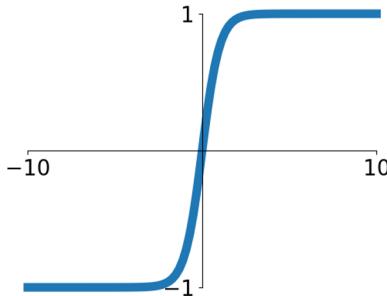


Probability interpretation

Activation Functions

tanh

$$\tanh(x) = \frac{e^{2x} - 1}{e^{2x} + 1}$$



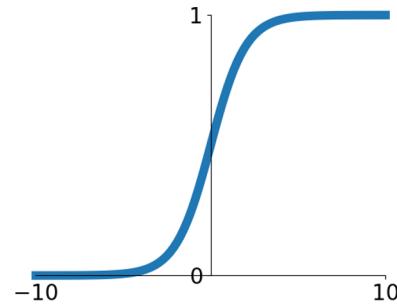
Zero mean / Zero centered

tanh is always superior than sigmoid

Activation Functions

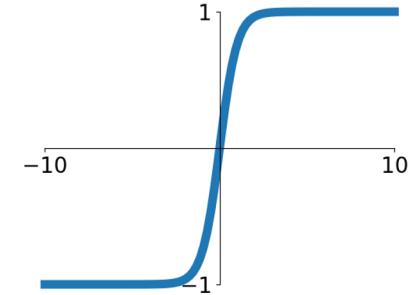
Sigmoid

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



tanh

$$\tanh(x) = \frac{e^{2x} - 1}{e^{2x} + 1}$$

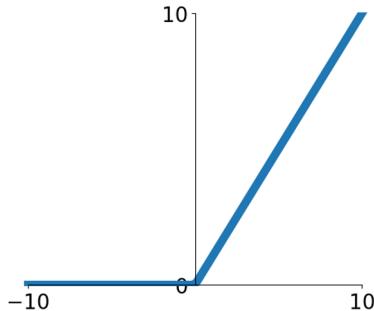


Drawbacks:

- gradient is small, and thus training (SGD) is slow
 - $\text{Exp}()$ is computationally expensive

Activation Functions

ReLU (Rectified Linear Unit)
 $\max(0, x)$



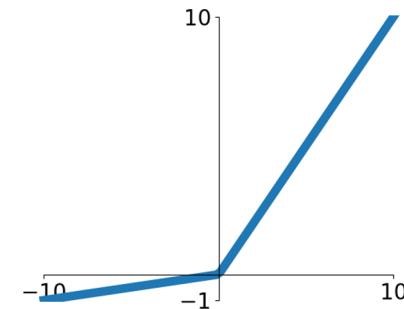
Default choice of activation function, ReLU is a good default choice for most problems

Drawbacks:

- Not zero centered.
- Dead ReLU when $x < 0$

Activation Functions

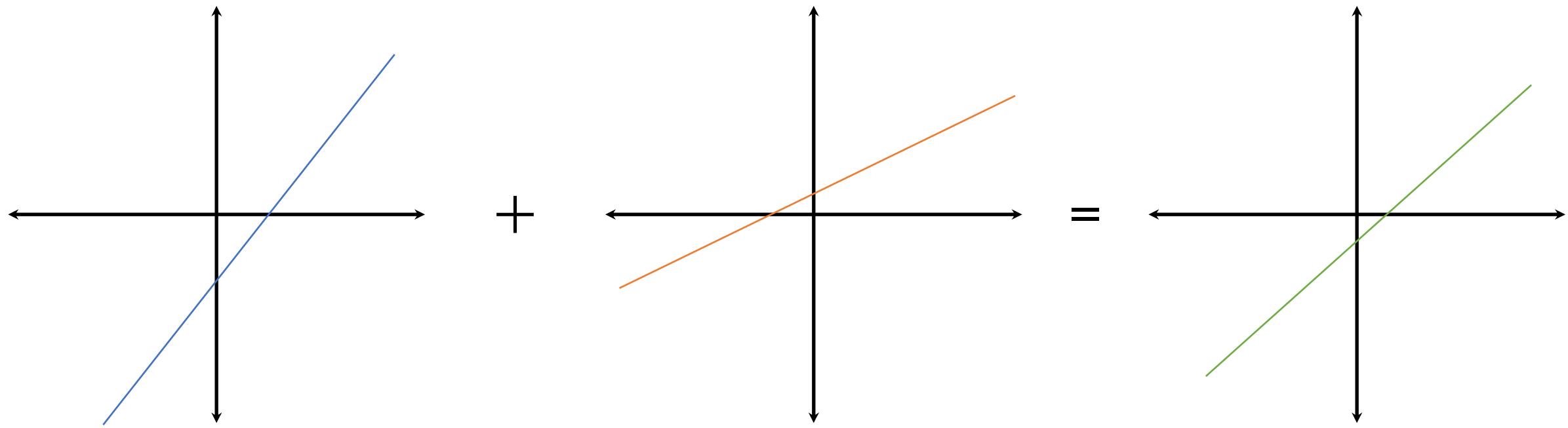
Leaky ReLU
 $\max(0.2x, x)$



Less used

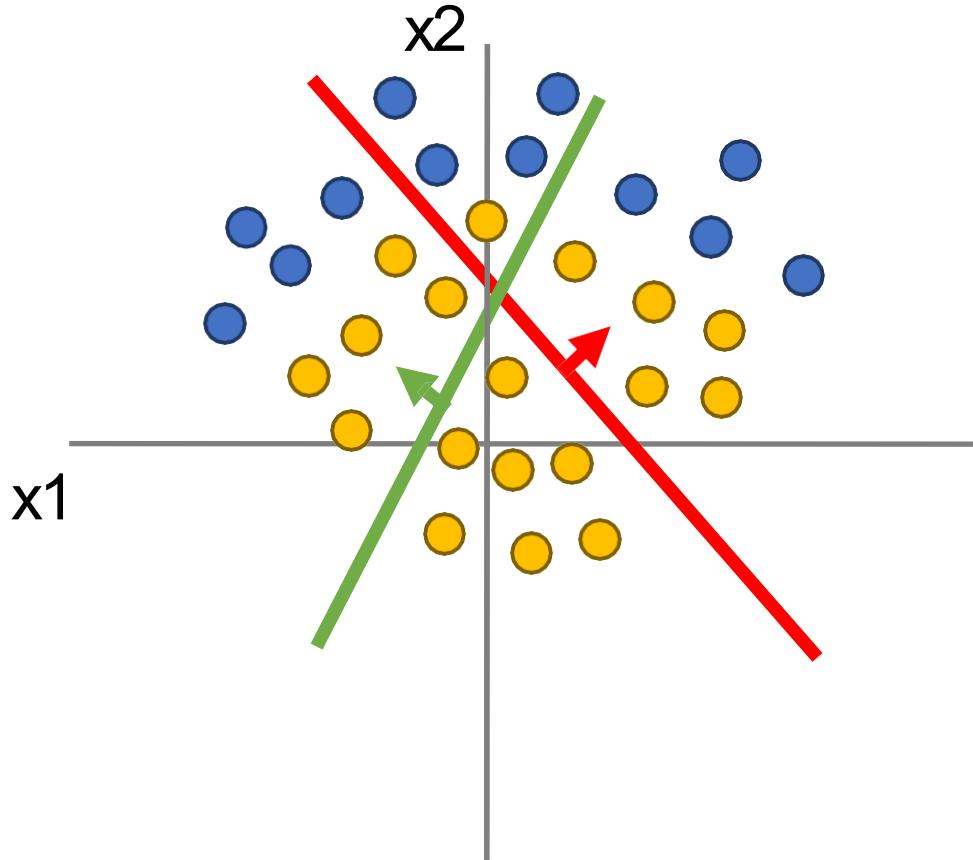
What happens if we build a neural network
with no non-linear activation function?

N linear classifiers = 1 linear classifier



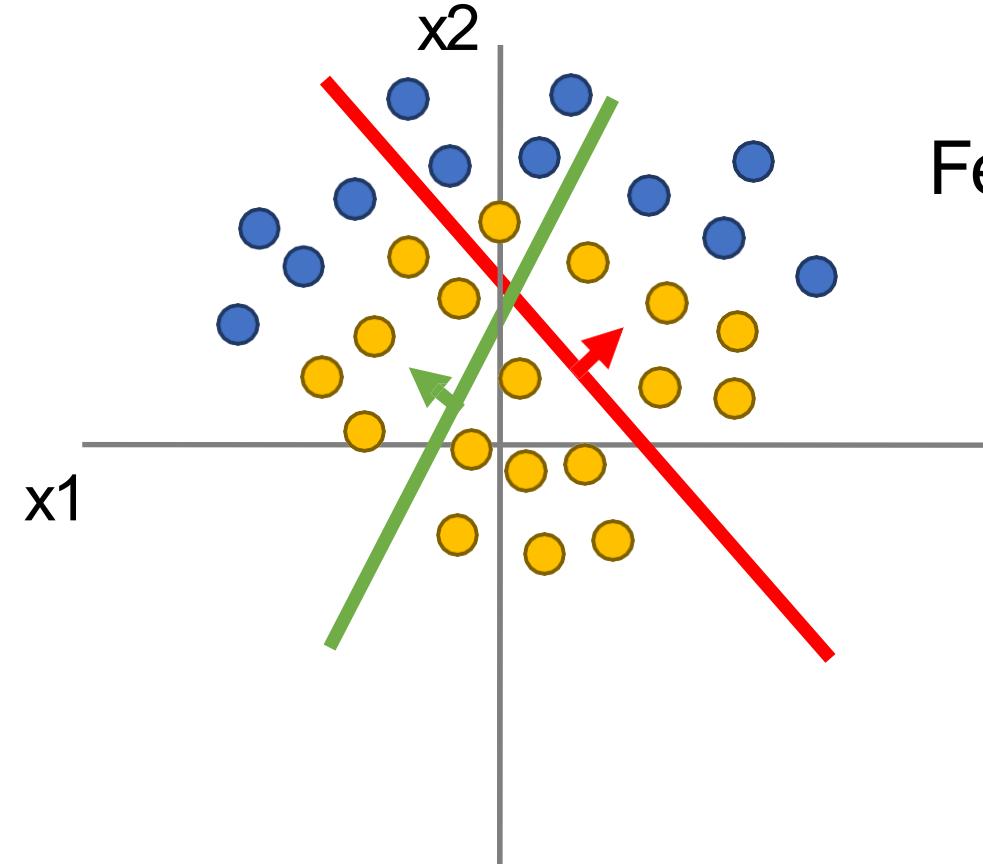
The power of non-linear activation function

Points not linearly
separable in original space



Consider a linear transform: $h = Wx$
Where x, h are both 2-dimensional

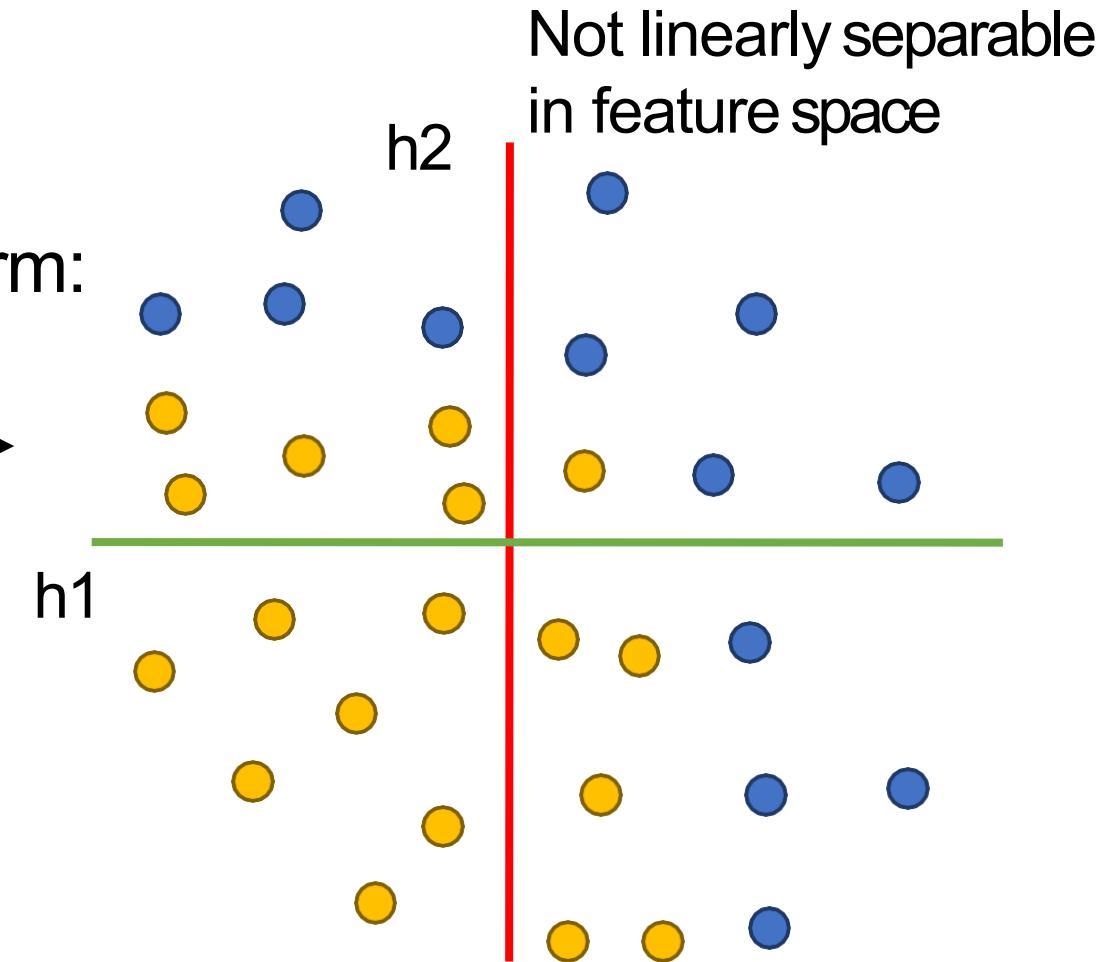
Points not linearly
separable in original space



Feature transform:
 $h = Wx$

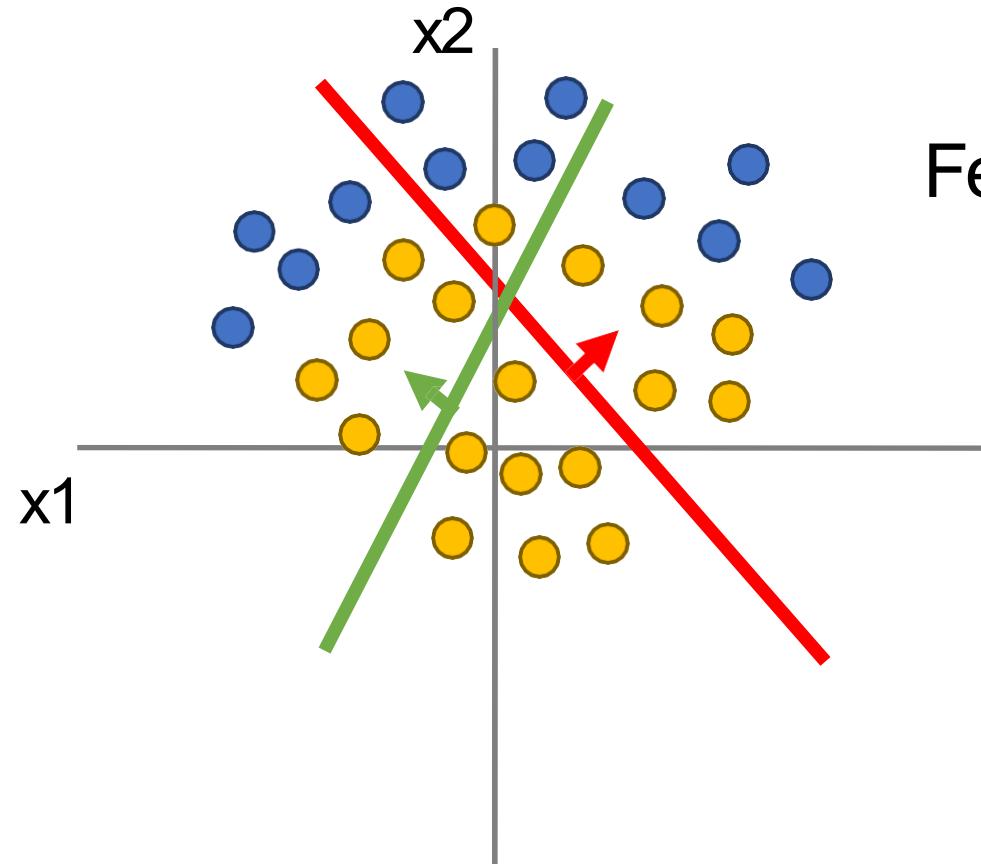


Consider a linear transform: $h = Wx$
Where x, h are both 2-dimensional



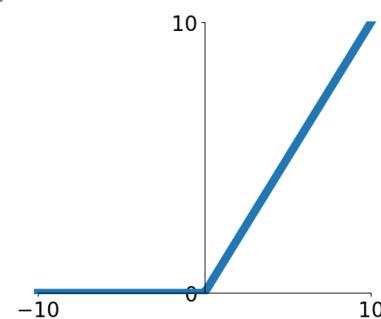
Not linearly separable
in feature space

Points not linearly
separable in original space

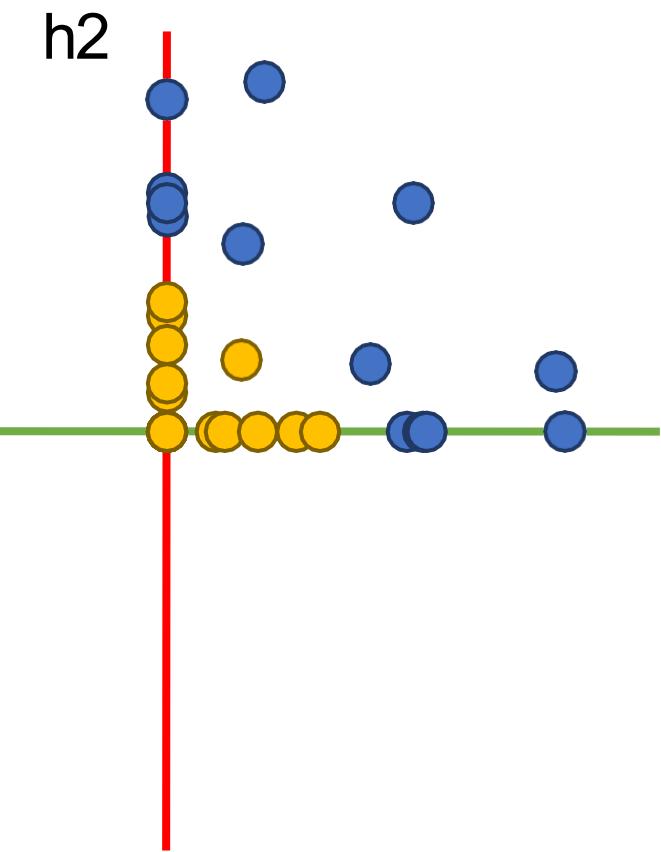


Feature transform:

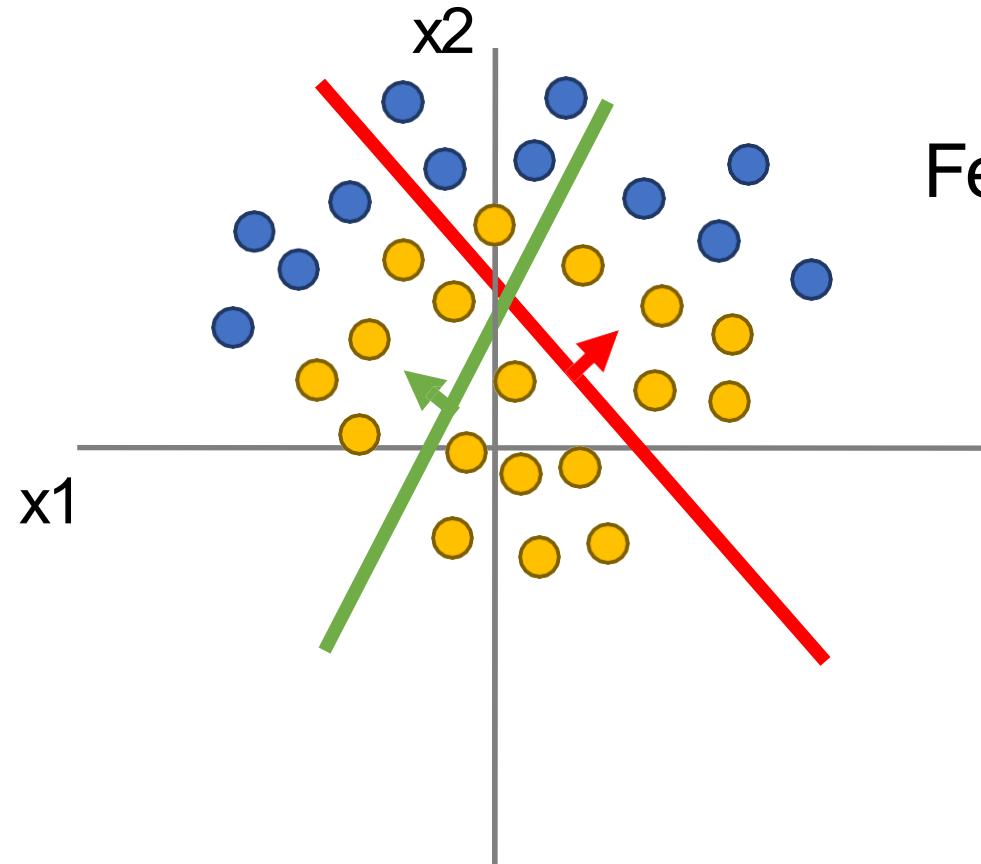
$$h = \text{ReLU}(Wx)$$



Consider a neural net hidden layer:
 $h = \text{ReLU}(Wx) = \max(0, Wx)$
Where x, h are both 2-dimensional

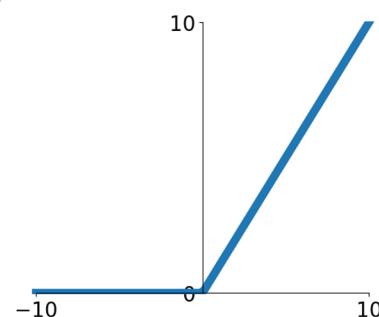


Points not linearly
separable in original space



Feature transform:

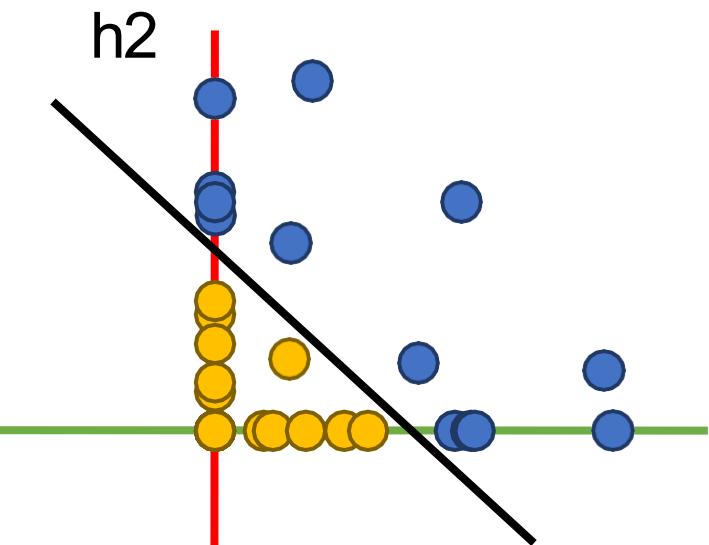
$$h = \text{ReLU}(Wx)$$



Consider a neural net hidden layer:

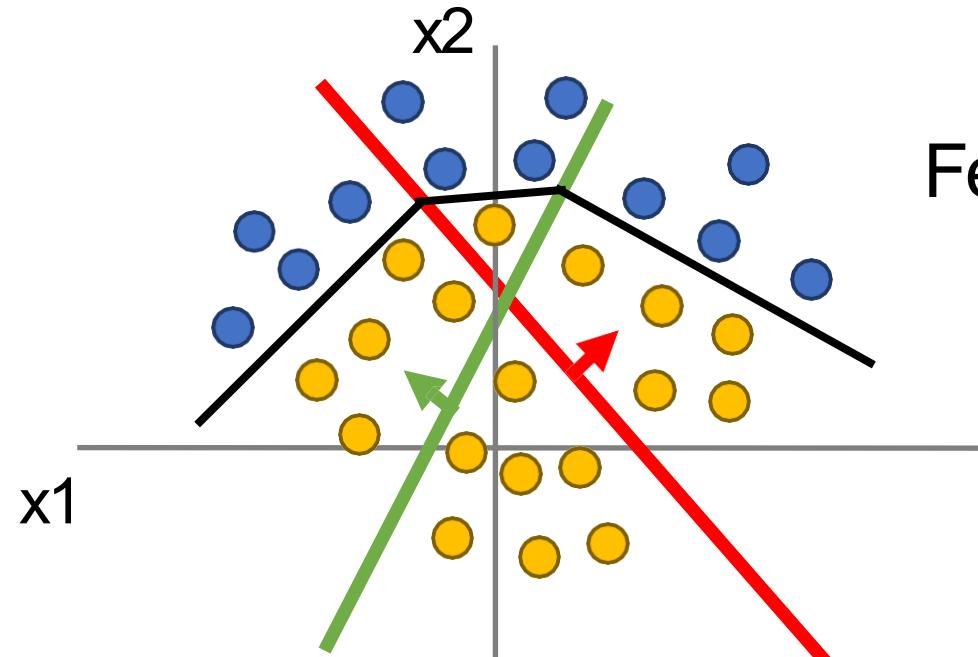
$$h = \text{ReLU}(Wx) = \max(0, Wx)$$

Where x, h are both 2-dimensional



Points are linearly
separable in features space!

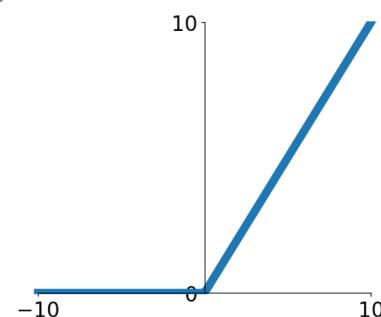
Points not linearly
separable in original space



Linear classifier in feature
space gives nonlinear
classifier in original space

Feature transform:

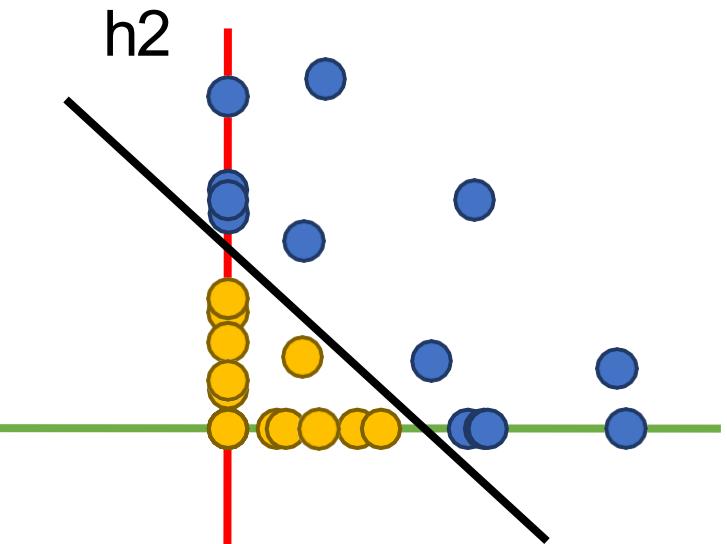
$$h = \text{ReLU}(Wx)$$



Consider a neural net hidden layer:

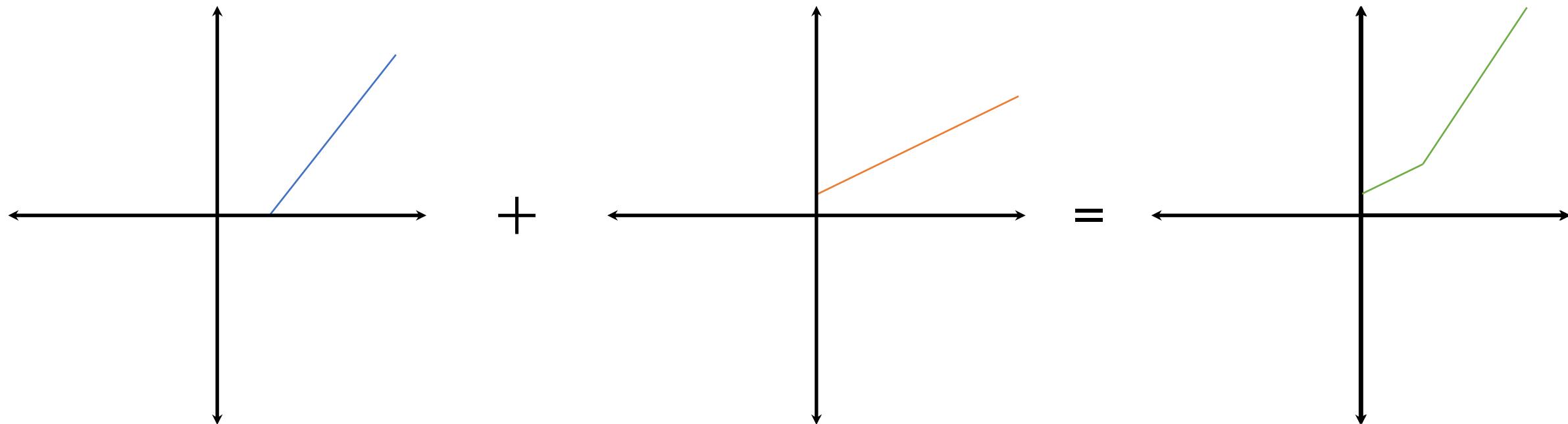
$$h = \text{ReLU}(Wx) = \max(0, Wx)$$

Where x, h are both 2-dimensional

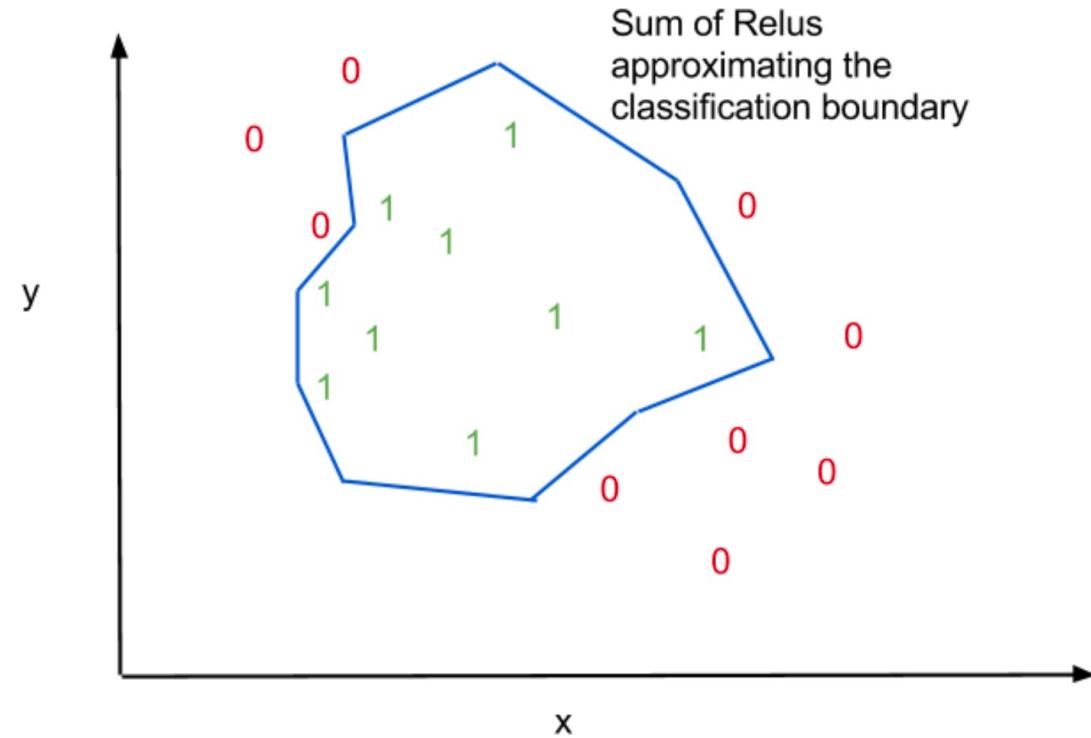
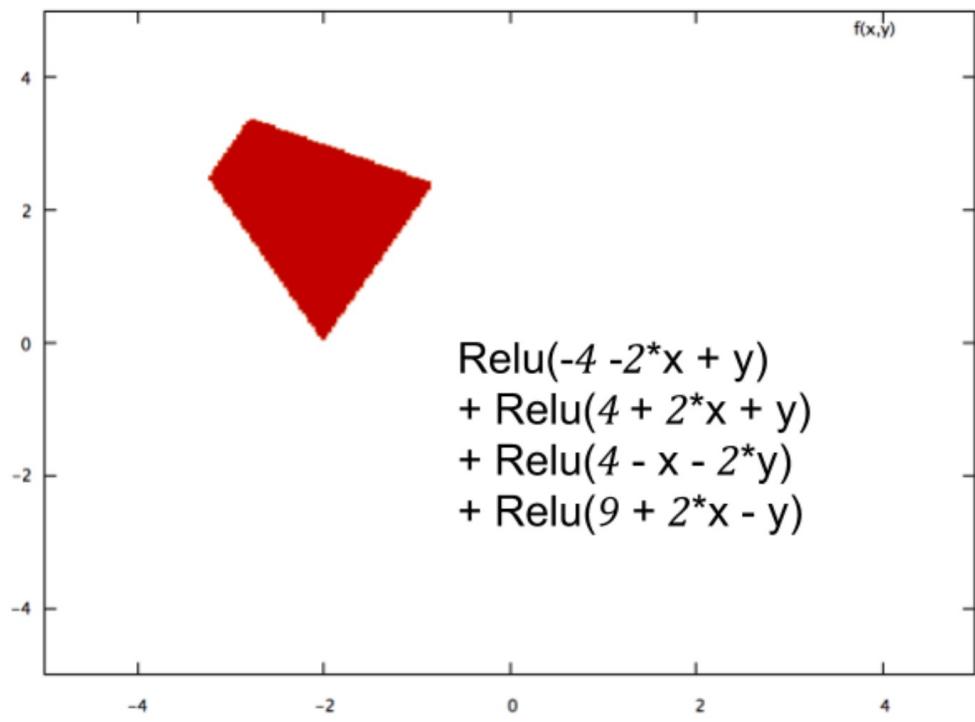


Points are linearly
separable in features space!

ReLU: Rectified Linear Units



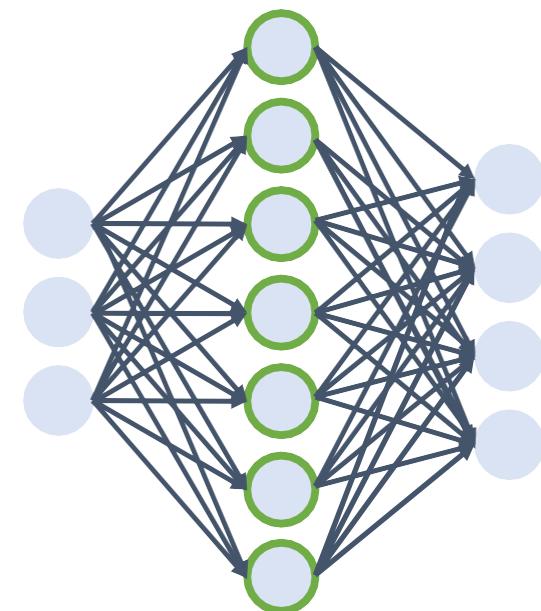
[ReLU](#), half-wave rectification, etc.



Universal Approximation

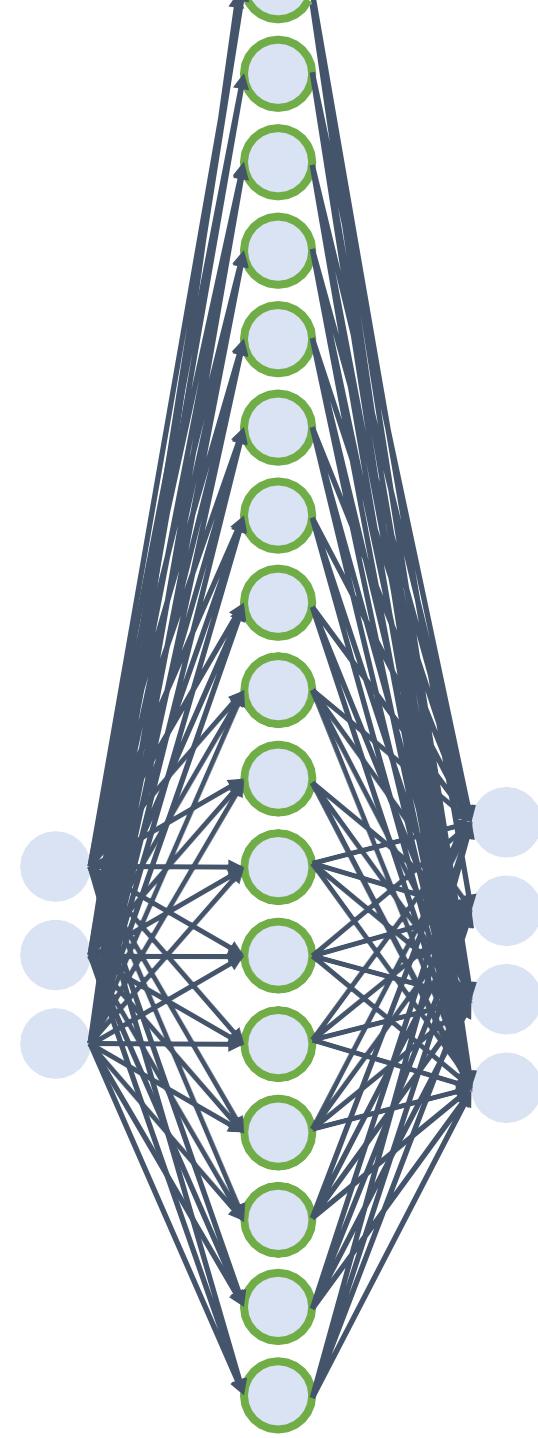
- A simple FC-Non-Linearity-FC network can approximate any function on a finite support, under:
 - mild assumptions
 - mild restrictions on non-linearity
- First proven for Sigmoid by Cybenko 1989
 - Later extended for ReLU and other non-linearities

$$\hat{y} = w_2^T f(w_1^T x + b_1) + b_2$$



Universal Approximation

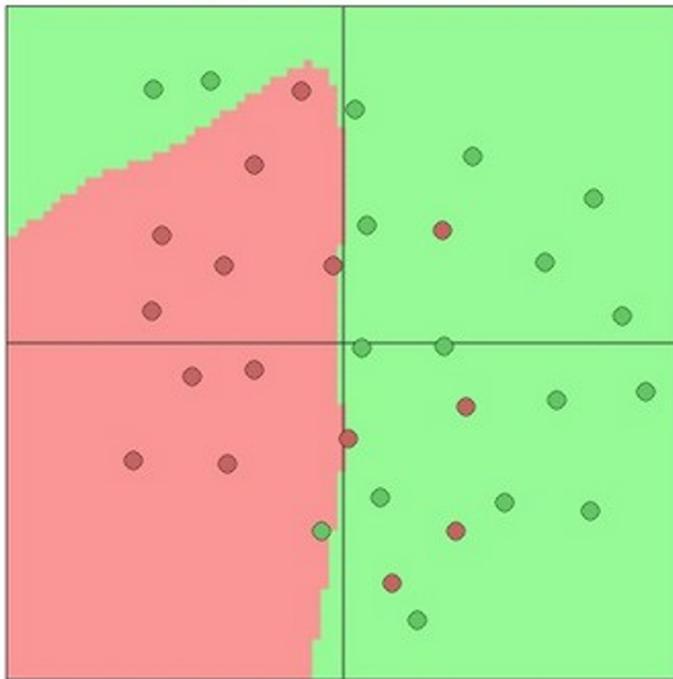
- A simple FC-Non-Linearity-FC network can approximate any function on a finite support, under:
 - mild assumptions
 - mild restrictions on non-linearity
- Requires extremely large number of neurons in the FC layer
 - Hard to train
 - Computational infeasible
 - Theoretically possible



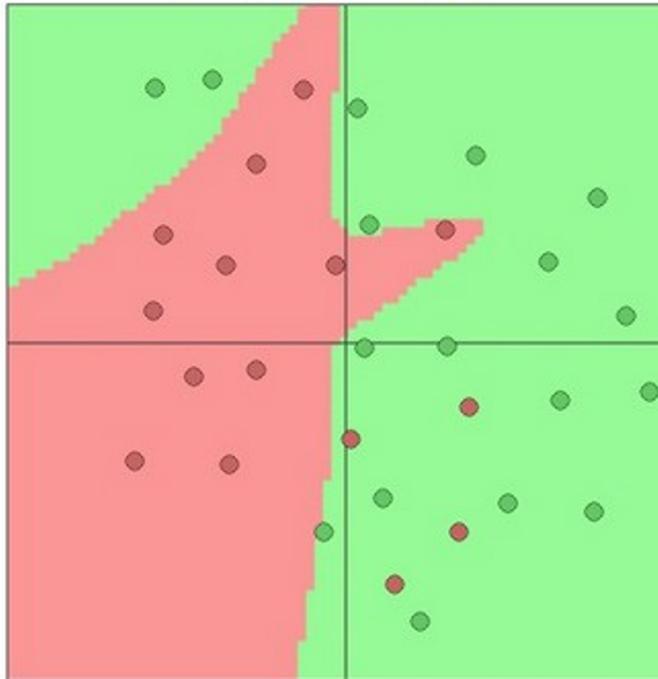
Setting the number of layers and their sizes

Setting the number of layers and their sizes

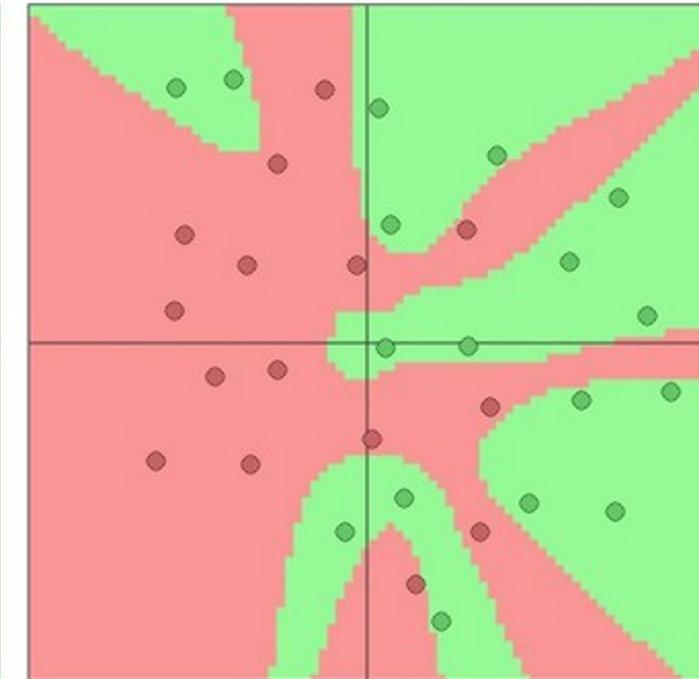
3 hidden units



6 hidden units



20 hidden units



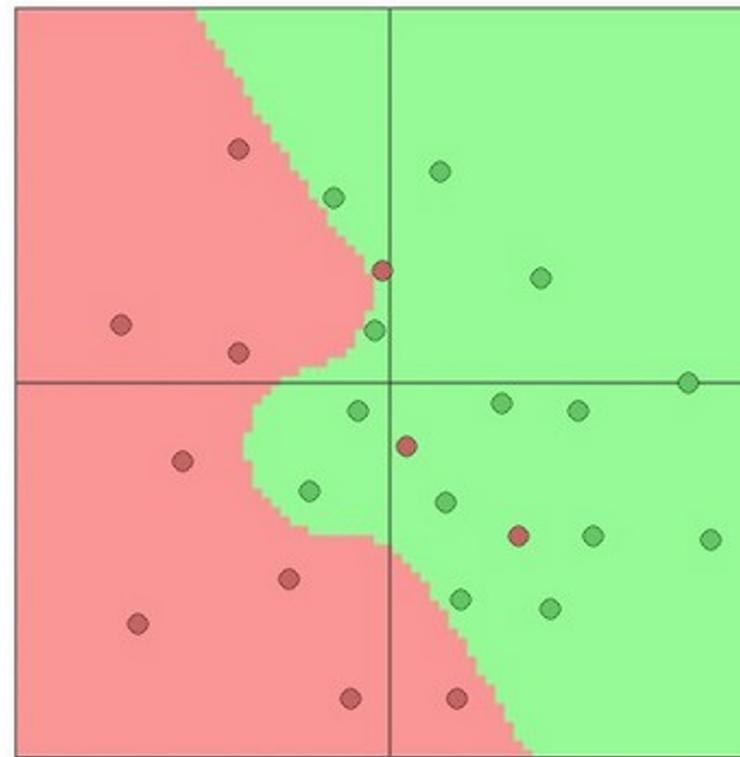
More hidden units = more capacity

Don't regularize with size; instead use stronger L2

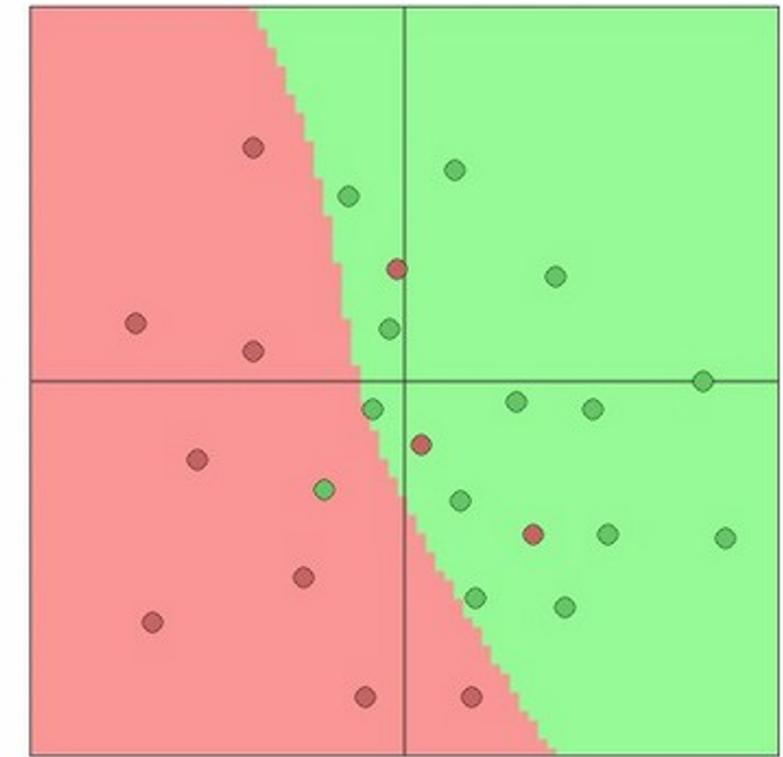
$\lambda = 0.001$



$\lambda = 0.01$



$\lambda = 0.1$



Any Question ?