

SPATIAL  
ECOLOGY

# Neural Nets

Antonio Fonseca

# Agenda

## 1) Perceptron

- Quick recap
- Hands-on tutorial
- Intro to gradient descent and optimizers

## 2) Feedforward Neural Networks

- The limitations of Perceptrons
- Multi-layer Perceptron
- Training: the forward and back-propagation
- Debugging tips

# Linear Regression Optimization

- Add an offset  $w_0$ :  $f(\mathbf{x}; \mathbf{w}) = \mathbf{w}^T \mathbf{x} + w_0$ ,  $\mathcal{D} = \{(\mathbf{x}_1, y_1), \dots (\mathbf{x}_n, y_n)\}$

$$\mathbf{w}^* = \arg \min_{\mathbf{w}} \sum_{i=1}^n (\mathbf{w}^T \mathbf{x}_i + w_0 - y_i)^2$$

$$= \arg \min_{\mathbf{w}} L(\mathbf{w}; \mathcal{D})$$

- Set  $\frac{\partial L(\mathbf{w}; \mathcal{D})}{\partial w_i} = 0$  for each  $i$

# Mean squared error loss

$$\mathbf{w}^* = \arg \min_{\mathbf{w}} \sum_{i=1}^n (\mathbf{w}^T \mathbf{x}_i + w_0 - y_i)^2$$

Rewrite:

$$\begin{aligned} (X\mathbf{w} - \mathbf{y})^T (X\mathbf{w} - \mathbf{y}) &= (\mathbf{w}^T X^T - \mathbf{y}^T)(X\mathbf{w} - \mathbf{y}) \\ &= \mathbf{w}^T X^T X \mathbf{w} - \mathbf{w}^T X^T \mathbf{y} - \mathbf{y}^T X \mathbf{w} + \mathbf{y}^T \mathbf{y} \\ &= \mathbf{w}^T X^T X \mathbf{w} - 2\mathbf{w}^T X^T \mathbf{y} + \mathbf{y}^T \mathbf{y}. \end{aligned}$$

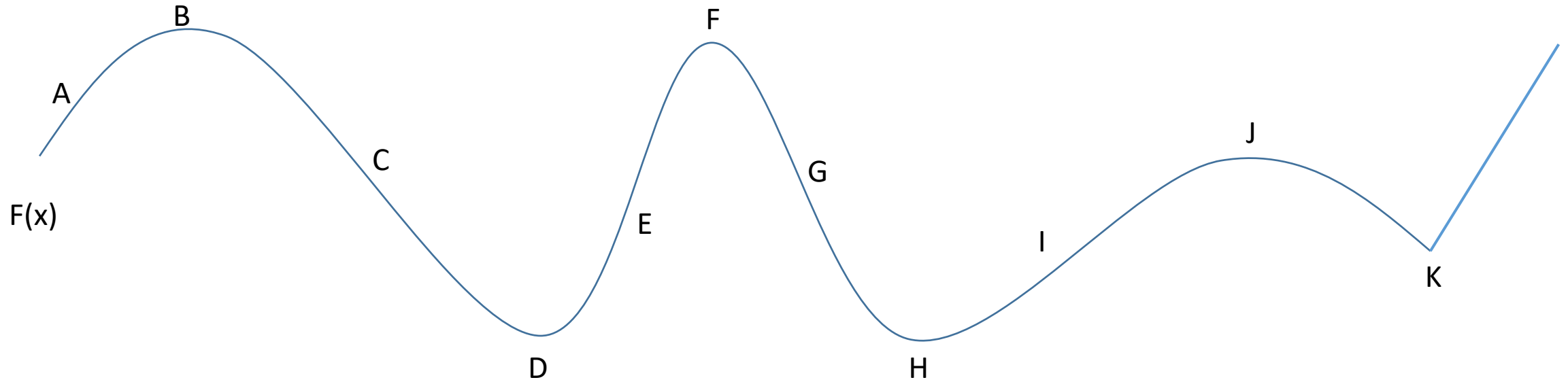
$$\frac{\partial}{\partial \mathbf{w}} \mathbf{w}^T X^T X \mathbf{w} - 2\mathbf{w}^T X^T \mathbf{y} + \mathbf{y}^T \mathbf{y} = 0$$

$$2X^T X \mathbf{w} - 2X^T \mathbf{y} = 0$$

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

$$\mathbf{w} = (X^T X)^{-1} X^T \mathbf{y}$$

# More on the derivatives



# Regularization

- Ridge regression: penalize with L2 norm

$$\mathbf{w}^* = \arg \min \sum_i L(f(\mathbf{x}_i; \mathbf{w}), y_i) + \lambda \sum_{j=1}^m w_j^2$$

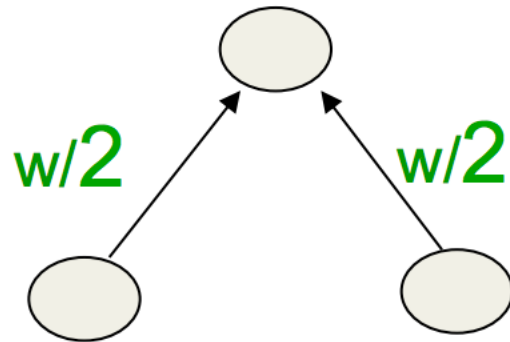
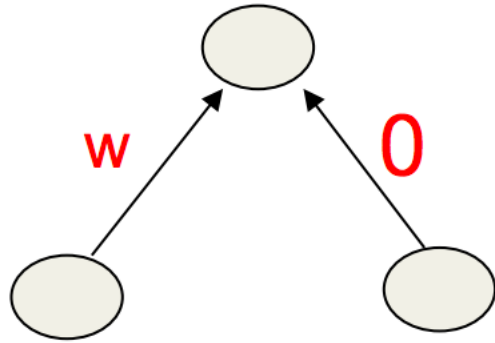
- Closed form solution exists  $\mathbf{w}^* = (\lambda I + X^T X)^{-1} X^T \mathbf{y}$

- LASSO regression: penalize with L1 norm

$$\mathbf{w}^* = \arg \min \sum_i L(f(\mathbf{x}_i; \mathbf{w}), y_i) + \lambda \sum_{j=1}^m |w_j|$$

- No closed form solution but still convex (optimal solution can be found)

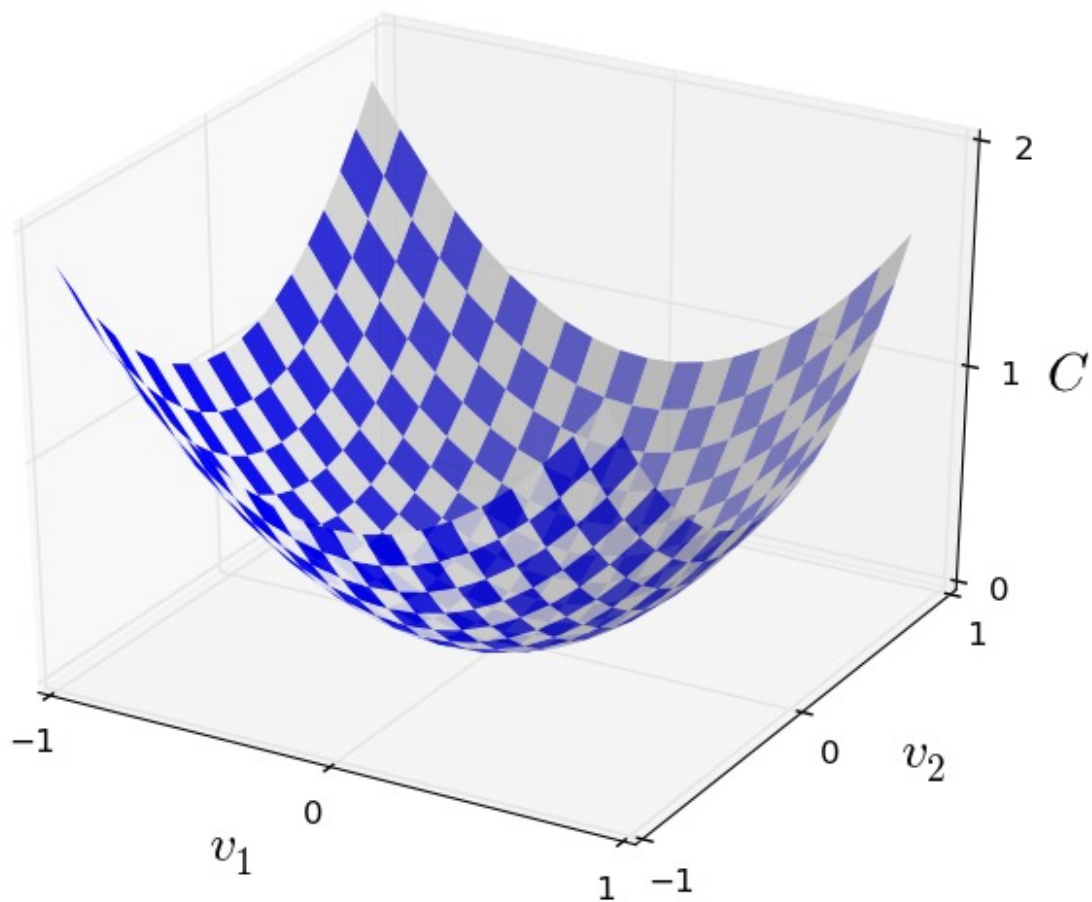
# Regularization



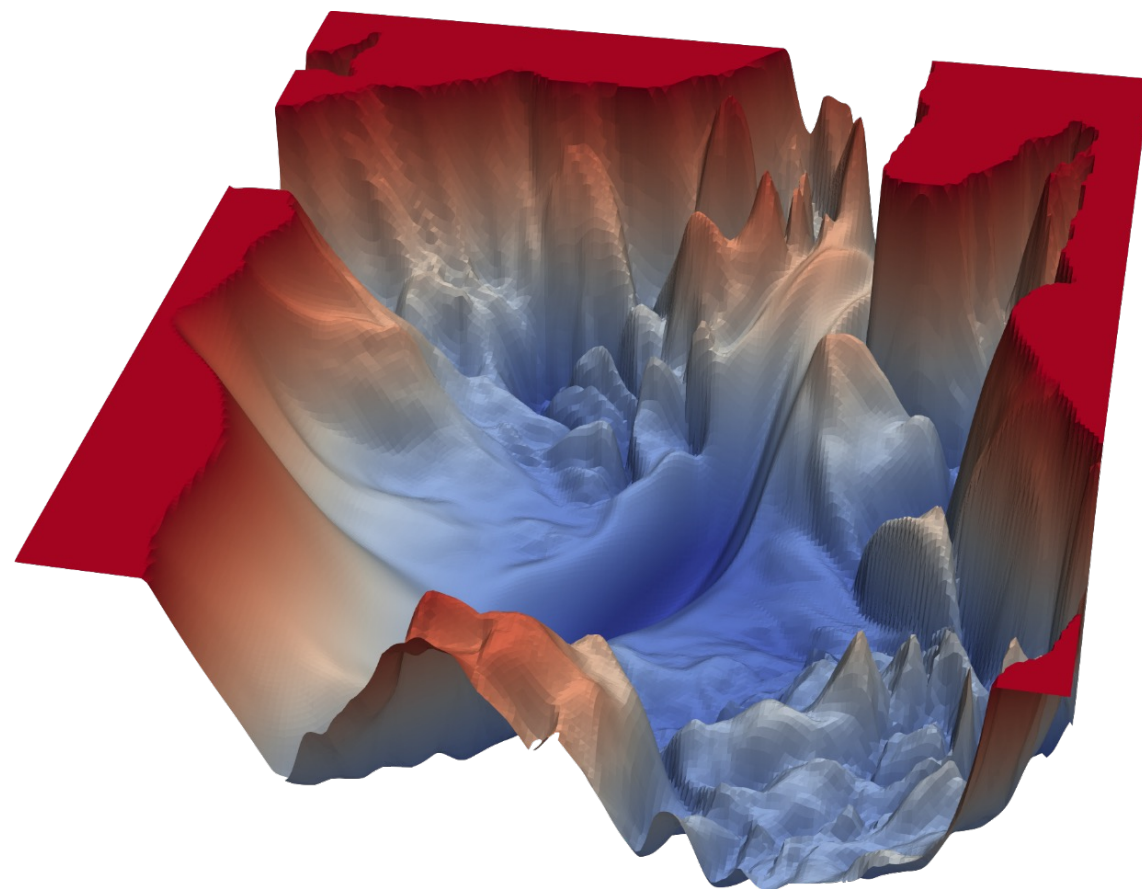
- Prefers to share smaller weights
- Makes model smoother
- More Convex



# Expectation



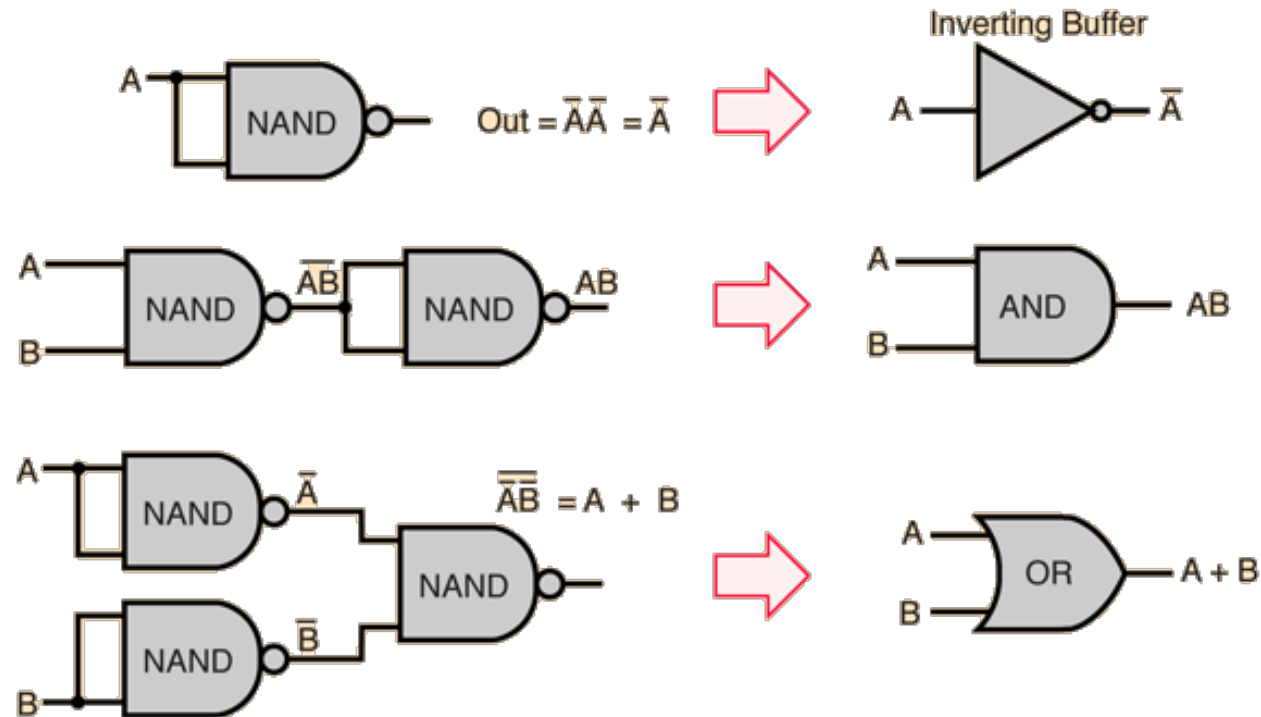
# Reality





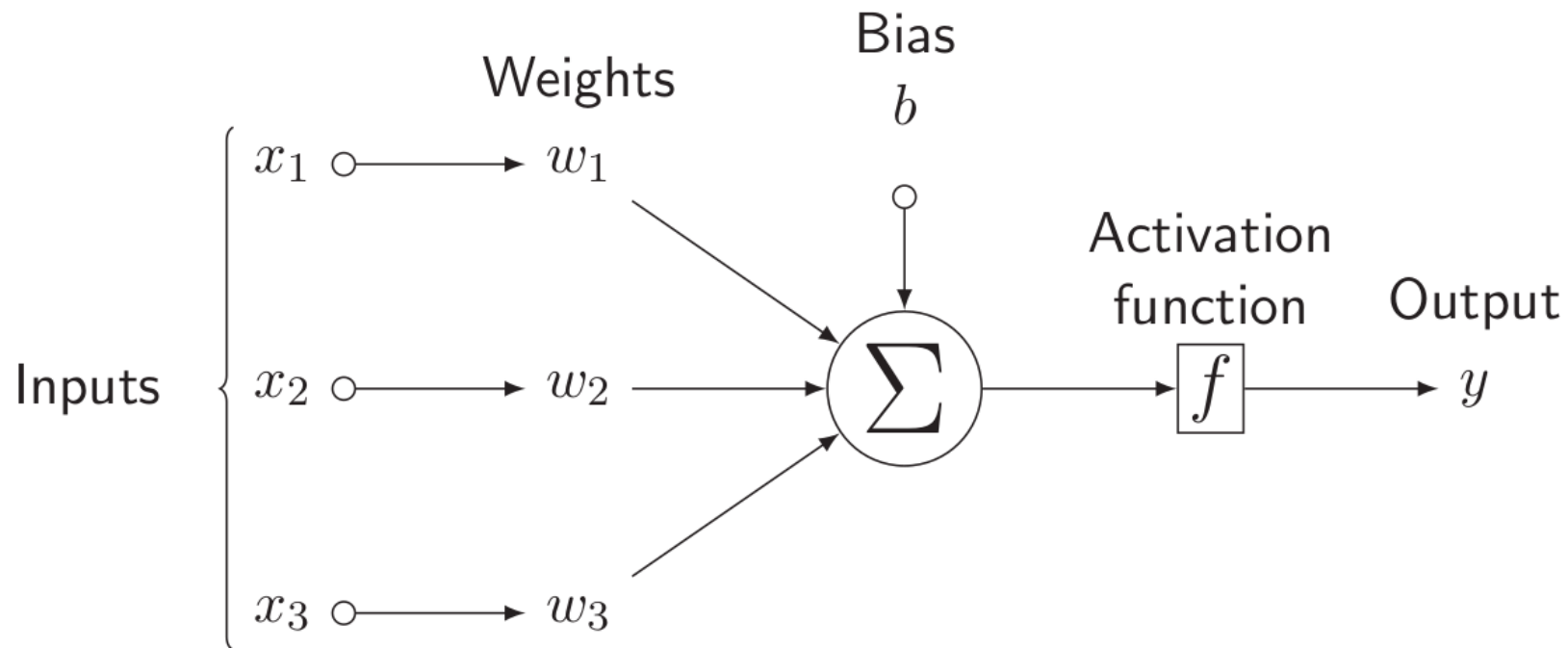
# Logic circuits with perceptrons

- NAND gates can be constructed from perceptrons
- NAND gates are universal for computation
  - Any computation can be built from NAND gates
  - Therefore, perceptrons are universal for computation



# Perceptron: Threshold Logic

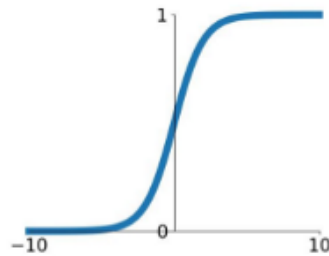
$$\mathcal{L}_{\text{perc}}(\mathbf{x}, y) = \begin{cases} 0 & \text{if } y\mathbf{w}^\top \mathbf{f}(\mathbf{x}) > 0 \\ -y\mathbf{w}^\top \mathbf{f}(\mathbf{x}) & \text{if } y\mathbf{w}^\top \mathbf{f}(\mathbf{x}) \leq 0 \end{cases}$$



# Activation functions

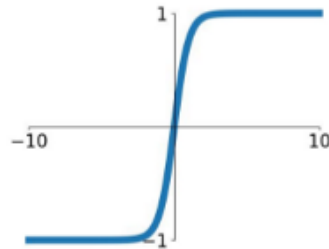
## Sigmoid

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



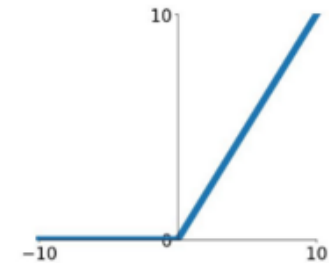
## tanh

$$\tanh(x)$$



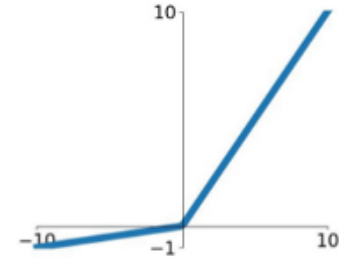
## ReLU

$$\max(0, x)$$



## Leaky ReLU

$$\max(0.1x, x)$$

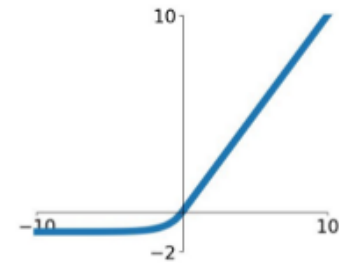


## Maxout

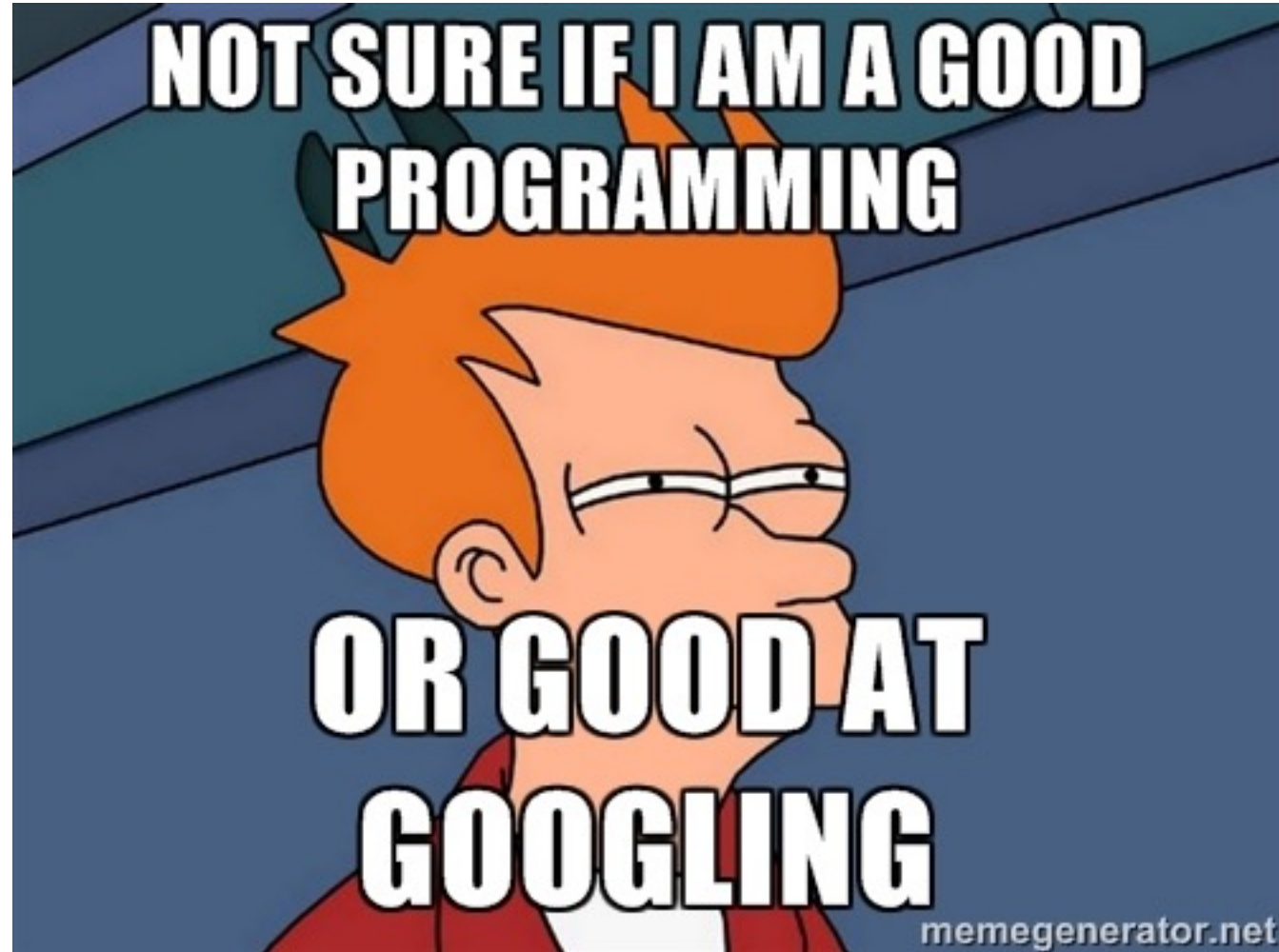
$$\max(w_1^T x + b_1, w_2^T x + b_2)$$

## ELU

$$\begin{cases} x & x \geq 0 \\ \alpha(e^x - 1) & x < 0 \end{cases}$$



Now let's get our hand dirty!



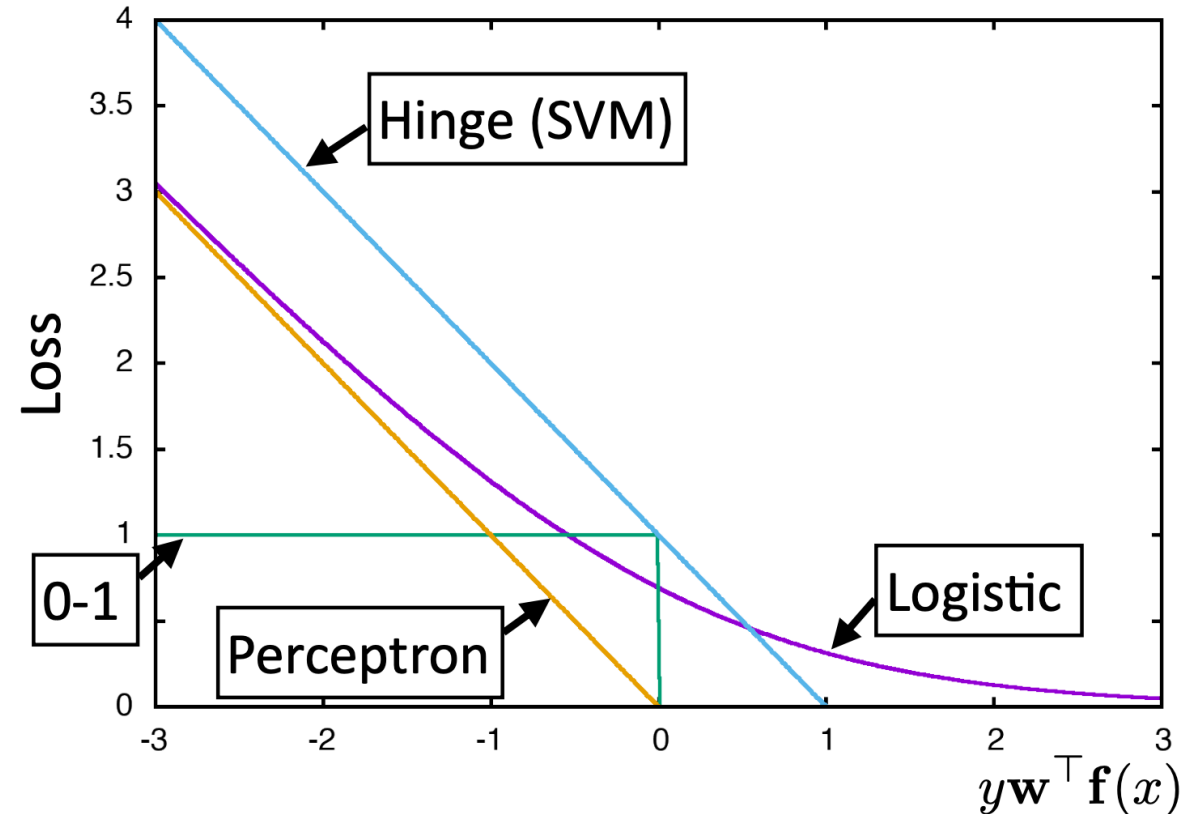
# (Putting things in perspective)

$$\mathcal{L}_{\text{lr}}(\mathbf{x}, y) = \begin{cases} -y\mathbf{w}^\top \mathbf{f}(\mathbf{x}) + \log(1 + \exp(y\mathbf{w}^\top \mathbf{f}(\mathbf{x}))) & \text{if } y = +1 \text{ (positive)} \\ \log(1 + \exp(-y\mathbf{w}^\top \mathbf{f}(\mathbf{x}))) & \text{if } y = -1 \text{ (negative)} \end{cases}$$

$$\mathcal{L}_{\text{perc}}(\mathbf{x}, y) = \begin{cases} 0 & \text{if } y\mathbf{w}^\top \mathbf{f}(\mathbf{x}) > 0 \\ -y\mathbf{w}^\top \mathbf{f}(\mathbf{x}) & \text{if } y\mathbf{w}^\top \mathbf{f}(\mathbf{x}) \leq 0 \end{cases}$$

Main differences:

- Perceptron: gradient-based optimization
- LR: probabilistic model
- Perceptron: if the data are linearly separable, perceptron is guaranteed to converge.
- LR: likelihood can never truly be maximized with a finite  $\mathbf{w}$  vector.



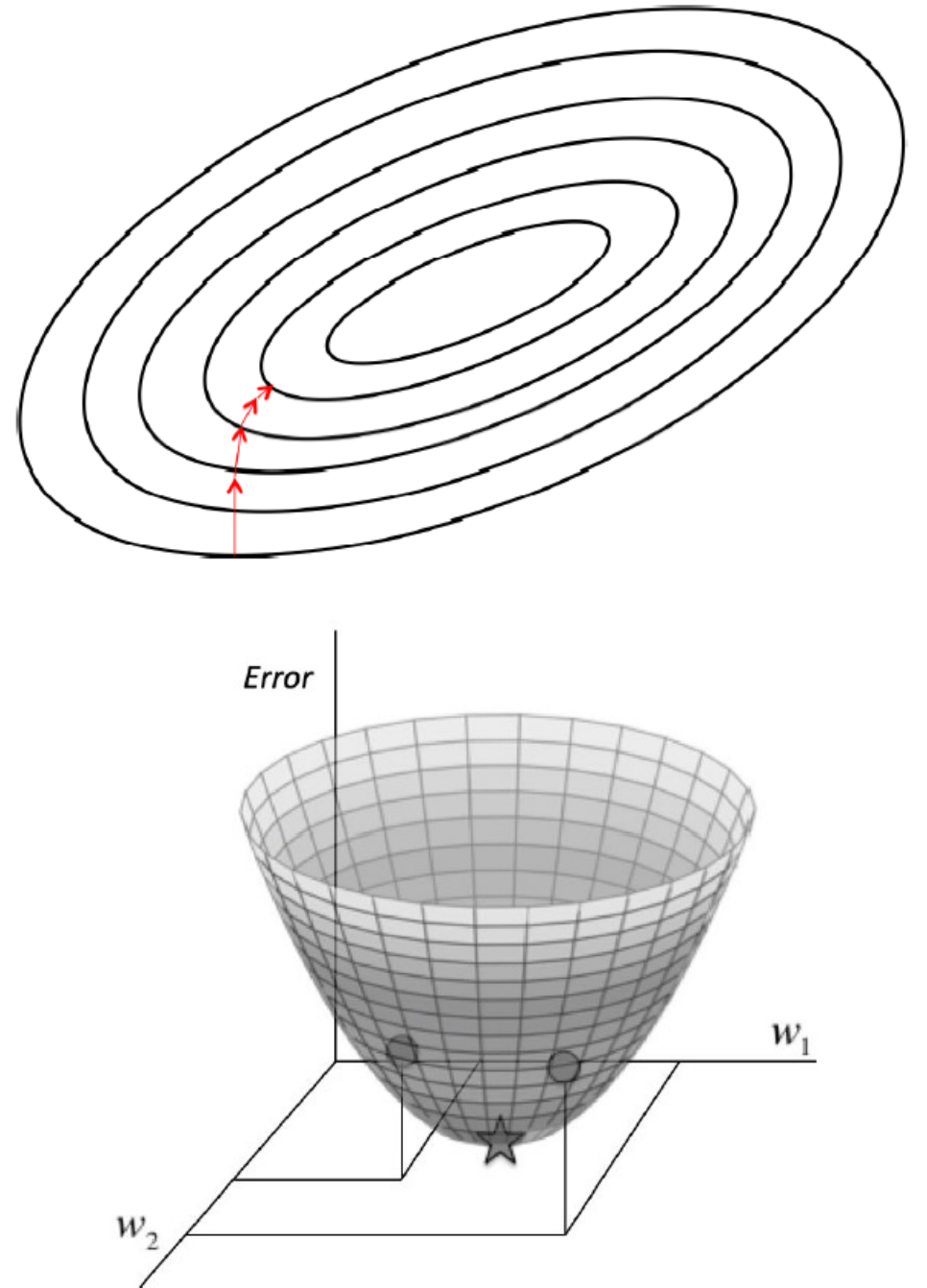
# Optimizers

## Gradient

$$\begin{aligned}\Delta w_k &= -\frac{\partial E}{\partial w_k} \\ &= -\frac{\partial}{\partial w_k} \left( \frac{1}{m} \sum_i (w^T X_i - y_i)_i^2 \right)\end{aligned}$$

$$w_{i+1} = w_i + \Delta w_k$$

## Stochastic gradient descent (**SGD**)





# Optimizers

## Hyperparameters

- Learning rate ( $\alpha$ )

$$\Delta w_k = -\alpha \frac{\partial E}{\partial w_k}$$
$$= -\alpha \frac{\partial}{\partial w_k} \left( \frac{1}{m} \sum_i (w^T X_i - y_i)_i^2 \right)$$

$$w_{i+1} = w_i + \Delta w_k$$

## Stochastic gradient descent (**SGD**)

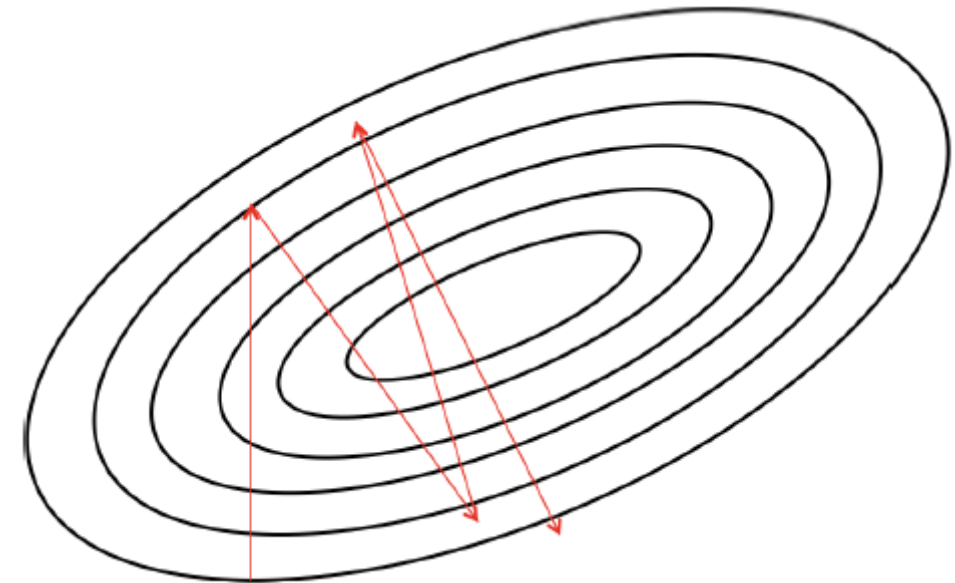
Practical test:

lr\_val = [1; 0.1; 0.01]

momentum\_val = 0

nesterov\_val = 'False'

decay\_val = 1e-6



Result of a large learning rate  $\alpha$

# Optimizers



Watch out for local minimal areas

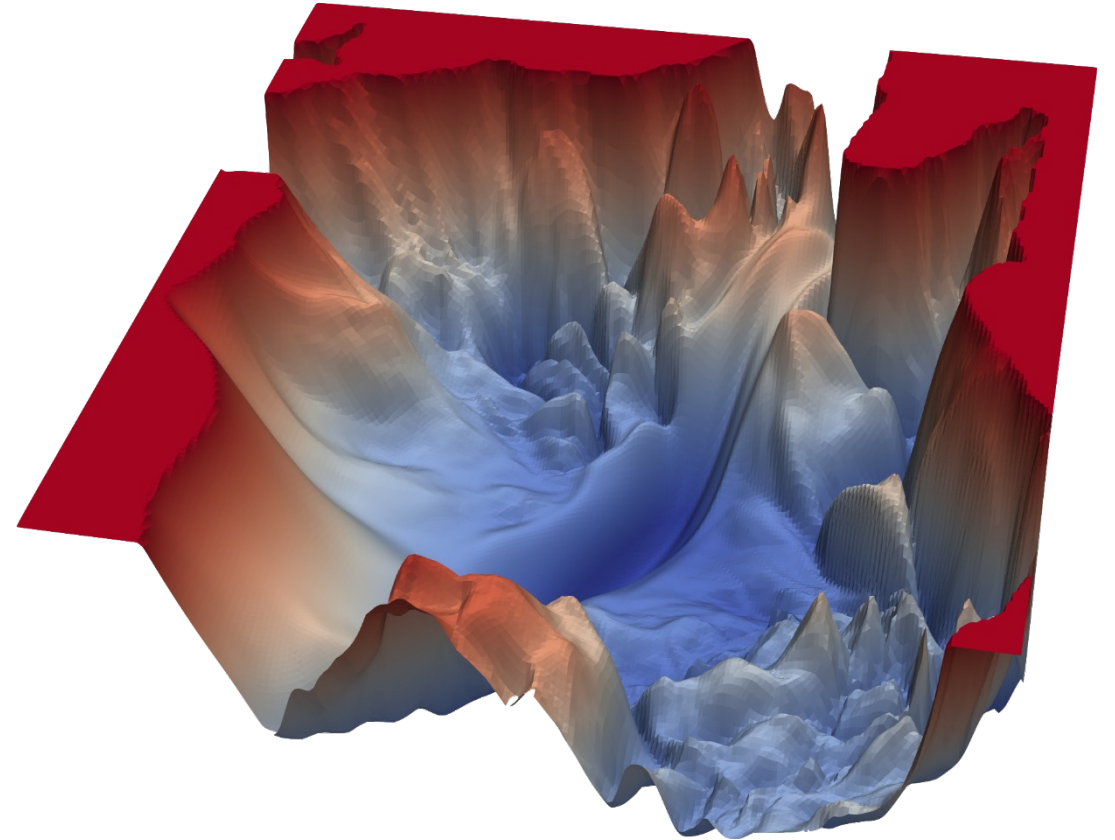
Hyperparameters

- Learning rate ( $\alpha$ )

$$\begin{aligned}\Delta w_k &= -\alpha \frac{\partial E}{\partial w_k} \\ &= -\alpha \frac{\partial}{\partial w_k} \left( \frac{1}{m} \sum_i (w^T X_i - y_i)_i^2 \right)\end{aligned}$$

$$w_{i+1} = w_i + \Delta w_k$$

Stochastic gradient descent (**SGD**)

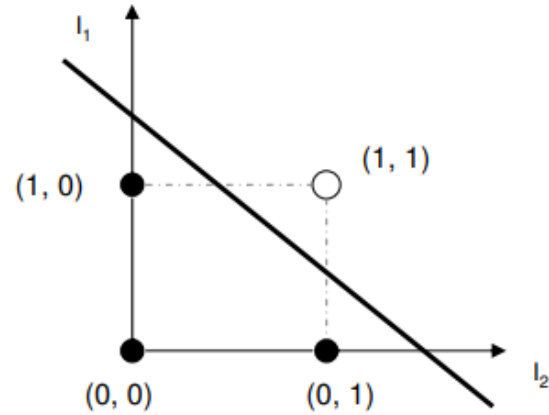


# Gradient Descent

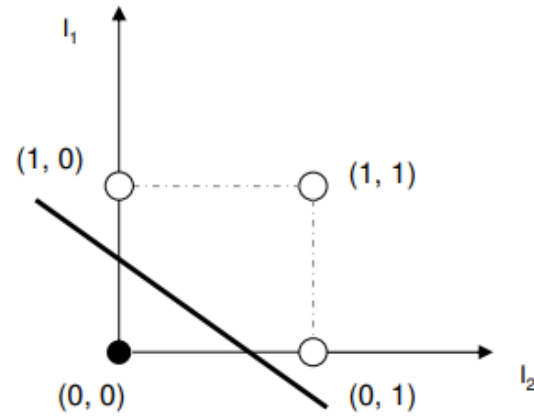
- Gradient descent refers to taking a step in the direction of the ***gradient (partial derivative)*** of a weight or bias with respect to the cost function
- Gradients are propagated backwards through the network in a process known as ***backpropagation***
- The size of the step taken in the direction of the gradient is called the ***learning rate***

# Limitations of the Perceptron

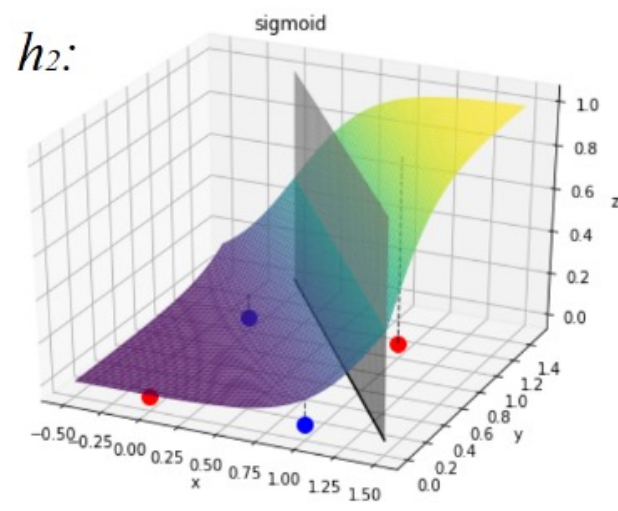
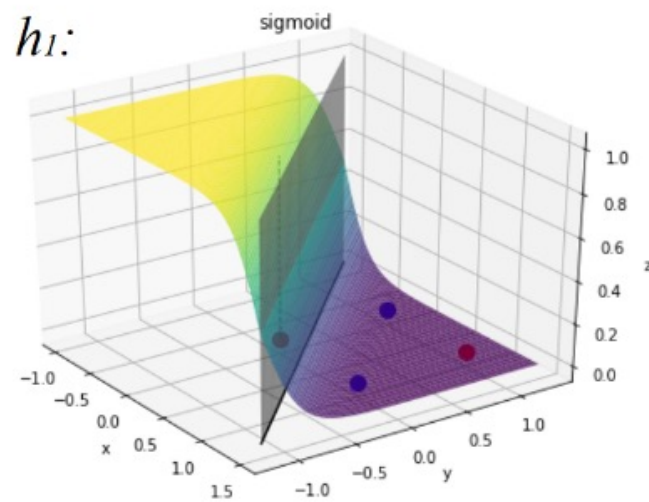
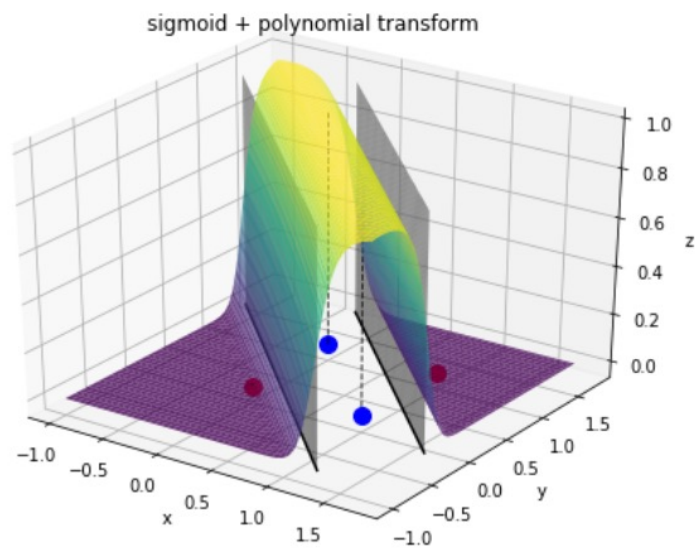
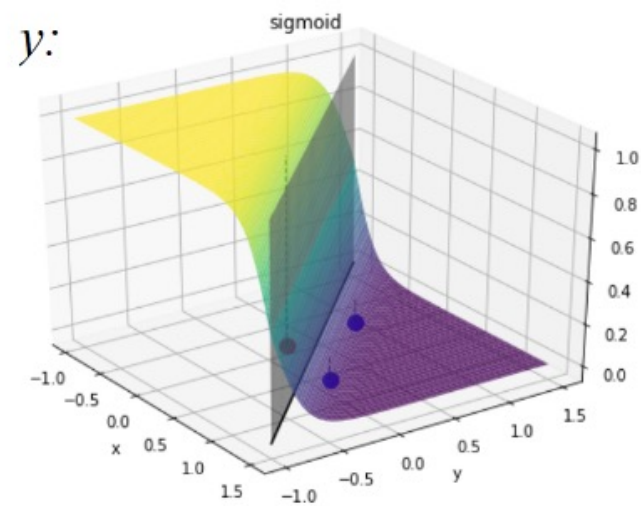
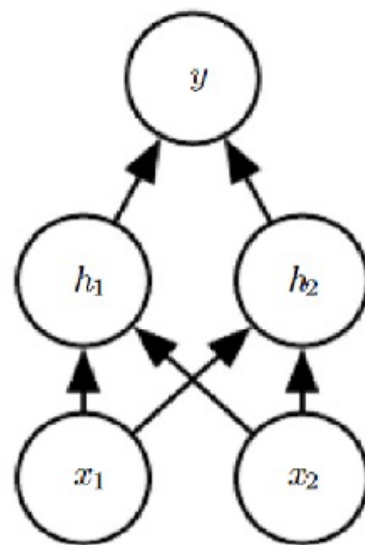
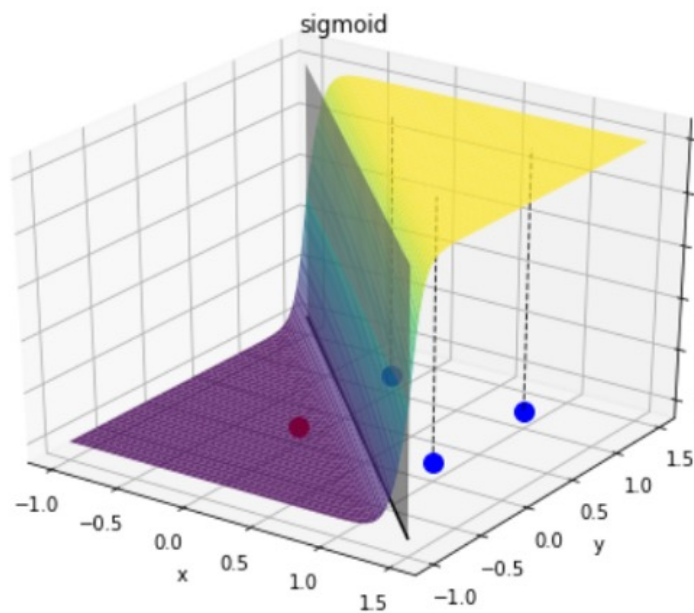
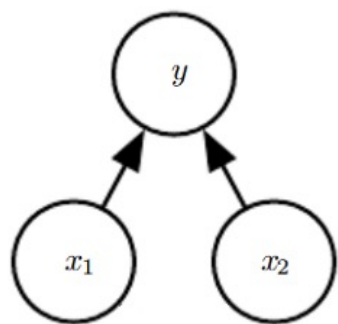
AND		
$I_1$	$I_2$	out
0	0	0
0	1	0
1	0	0
1	1	1



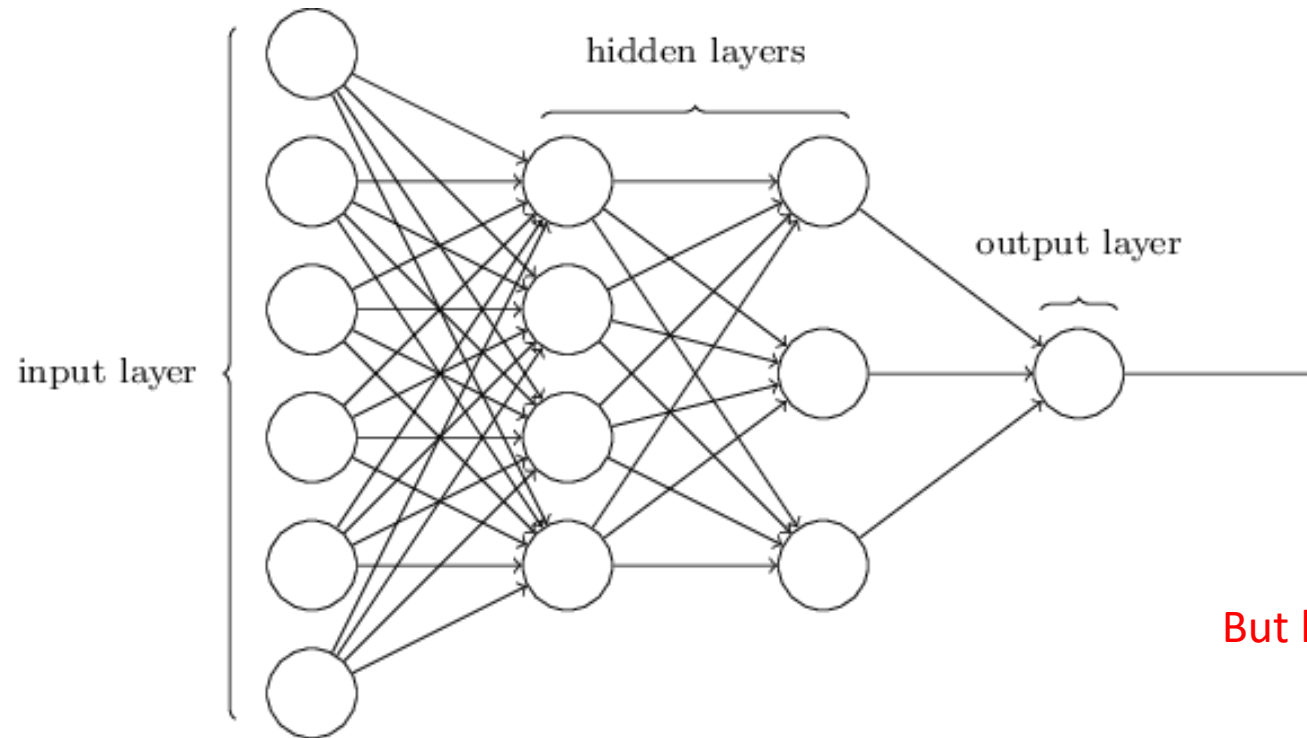
OR		
$I_1$	$I_2$	out
0	0	0
0	1	1
1	0	1
1	1	1



Perceptron



# Architecture of Neural Networks



But how do we train it?

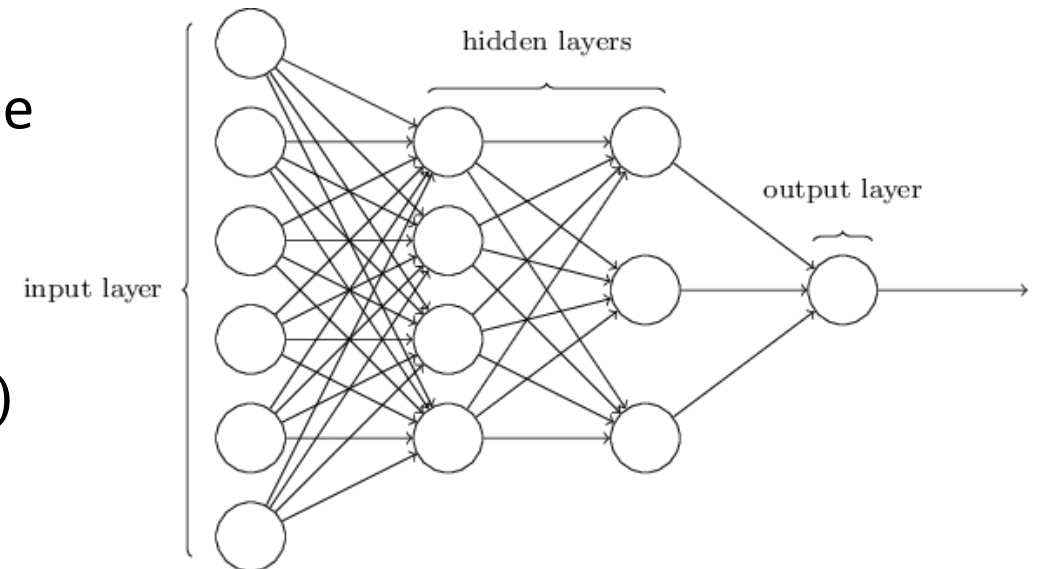
- Sometimes called multi-layer perceptron (MLP)
- Output from one layer is used as input for the next (Feedforward network)



# Forward Propagation

- Store weights and biases as matrices
- Suppose we are considering the weights from the second (hidden) layer to the third (output) layer
  - $w$  is the weight matrix with  $w_{ji}$  the weight for the connection between the  $i$ th neuron in the second layer and the  $j$ th neuron in the third layer
  - $b$  is the vector of biases in the third layer
  - $a$  is the vector of activations (output) of the 2<sup>nd</sup> layer
  - $a'$  the vector of activations (output) of the third layer

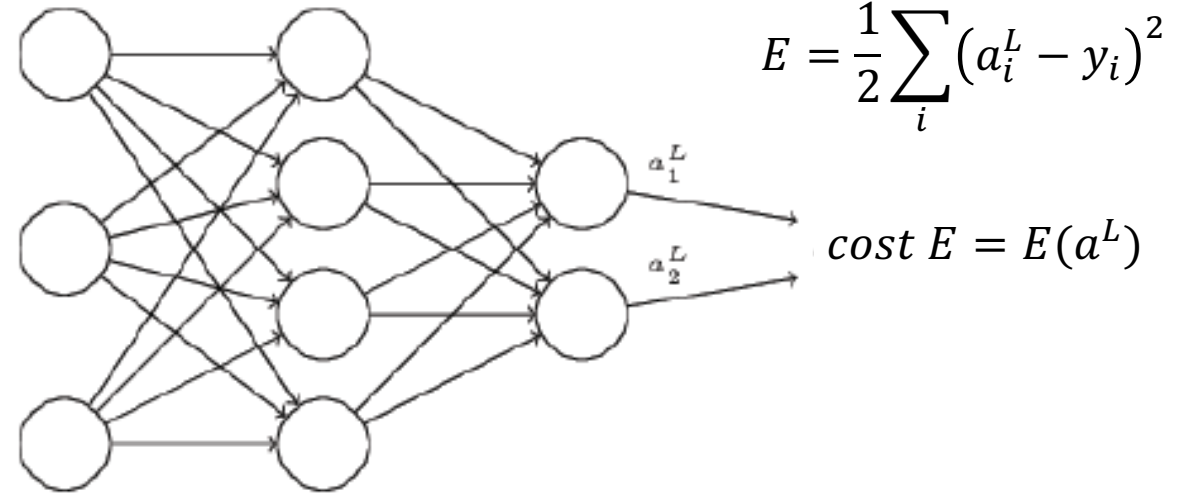
$$a' = \sigma(wa + b)$$



# Backpropagation

- 1. Input  $x$ :** Set the corresponding activation  $a^1$  for the input layer.
- 2. Feedforward:** For each  $l = 2, 3, \dots, L$  compute  $z^l = w^l a^{l-1} + b^l$  and  $a^l = \sigma(z^l)$ .
- 3. Output error  $\delta^L$ :** Compute the vector  $\delta^L = \nabla_a C \odot \sigma'(z^L)$ .
- 4. Backpropagate the error:** For each  $l = L - 1, L - 2, \dots, 2$  compute  $\delta^l = ((w^{l+1})^T \delta^{l+1}) \odot \sigma'(z^l)$ .
- 5. Output:** The gradient of the cost function is given by  $\frac{\partial C}{\partial w_{jk}^l} = a_k^{l-1} \delta_j^l$  and  $\frac{\partial C}{\partial b_j^l} = \delta_j^l$ .

$$\frac{\partial E}{\partial w_{ji}^l} = \frac{\partial E}{\partial a_i^l} \frac{\partial a_i^l}{\partial z_j^l} \frac{\partial (w_{ji}^l a_i^{l-1})}{\partial w_{ji}^l}$$



$$z_j^l = \sum_i w_{ji}^l a_i^{l-1} + b_j^l \quad a_j^l = \sigma \left( \sum_i w_{ji}^l a_i^{l-1} + b_j^l \right) = \sigma(z_j^l)$$

$$\delta_j^L \equiv \frac{\partial E}{\partial z_j^L} = \frac{\partial E}{\partial a_i^L} \frac{\partial a_i^L}{\partial z_j^L} = \frac{\partial E}{\partial a_j^L} \sigma'(z_j^L) \quad (1)$$

$$\begin{aligned} \delta_j^l &\equiv \frac{\partial E}{\partial z_j^l} = \frac{\partial E}{\partial z_i^{l+1}} \frac{\partial z_i^{l+1}}{\partial z_j^l} = \frac{\partial z_i^{l+1}}{\partial z_j^l} \delta_i^{l+1} \\ &= \frac{\partial (\sum_i w_{ij}^{l+1} a_j^l + b_i^{l+1})}{\partial z_j^l} \delta_i^{l+1} = \sum_i w_{ij}^{l+1} \delta_i^{l+1} \sigma'(z_j^l) \quad (2) \end{aligned}$$

# Back to the code (Feedforward networks)

When people want to use Machine Learning without math



# How training works

1. In each ***epoch***, randomly shuffle the training data
2. Partition the shuffled training data into ***mini-batches***
3. For each mini-batch, apply a single step of **gradient descent**
  - **Gradients** are calculated via ***backpropagation*** (the next topic)
4. Train for multiple epochs

# Debugging a neural network

- What can we do?
  - Should we change the learning rate?
  - Should we initialize differently?
  - Do we need more training data?
  - Should we change the architecture?
  - Should we run for more epochs?
  - Are the features relevant for the problem?
- Debugging is an art
  - We'll develop good heuristics for choosing good architectures and hyper parameters