

## APPLIED MACHINE LEARNING

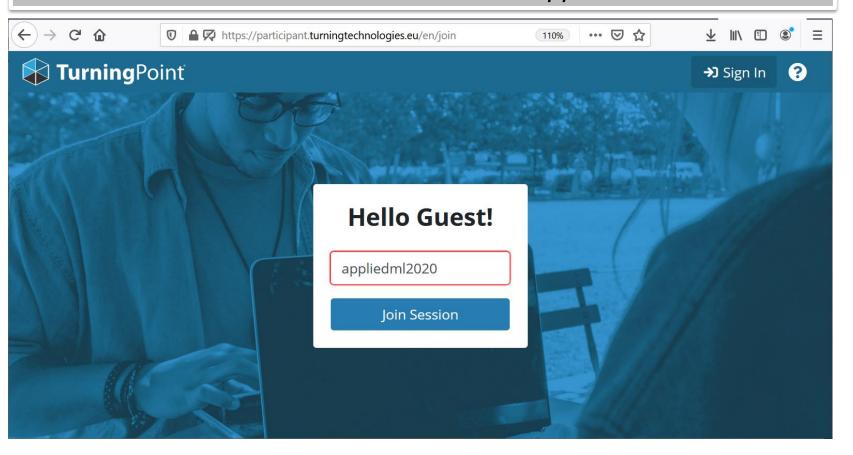
## Introduction to Neural Networks Interactive session



## Launch polling system

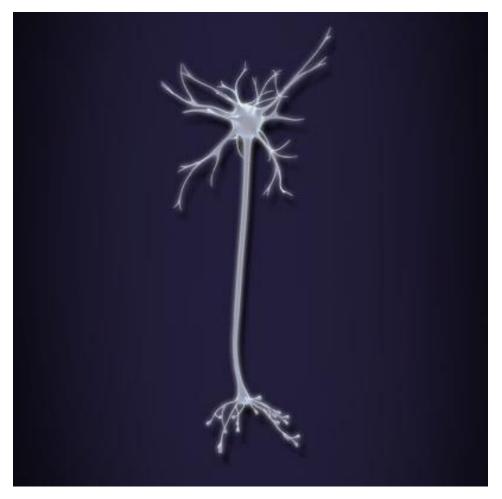
https://participant.turningtechnologies.eu/en/join

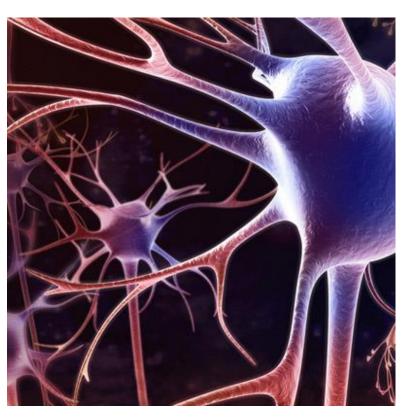
Acces as GUEST and enter the session id: appliedml2020





## The Neuron

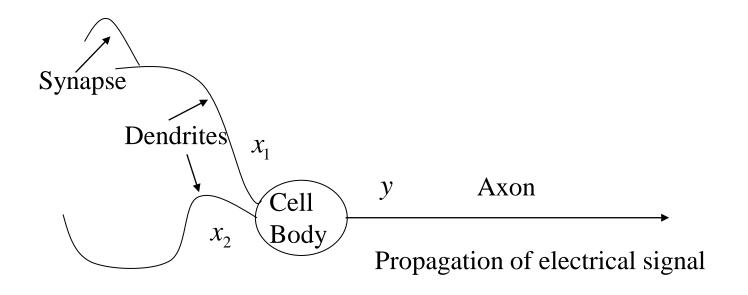




Courtesy of 3DScience.com

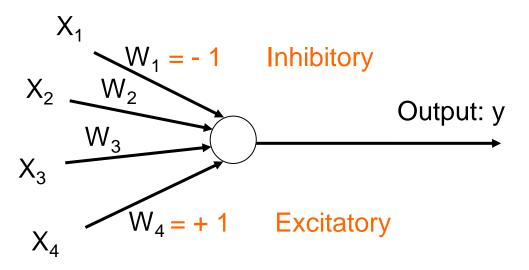


## Simple Model of Biological Neurons' Signal Processing: The Leaky – Integrator Neuron





## Simple Model – The Perceptron



X: Input vector from all other neurons

w: the strength of each synapse

y: neuron output

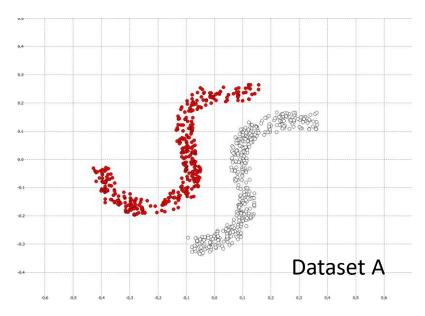
f: transfer function

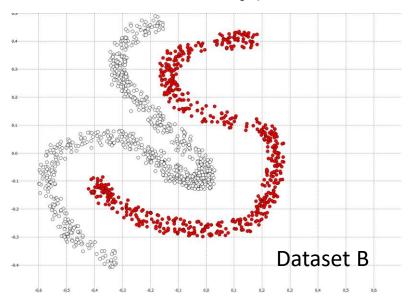
$$y = \theta \left( \sum_{i=1}^{4} w_i \cdot x_i + w_0 \right)$$

$$\theta(...) = \begin{cases} 1 \text{ if } \mathbf{w}^T \mathbf{x} + w_0 \ge 0 \\ 0 \text{ if } \mathbf{w}^T \mathbf{x} + w_0 < 0 \end{cases}$$

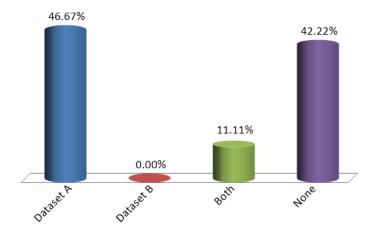


# For which dataset the perceptron model can achieve perfect classification (100% accuracy)





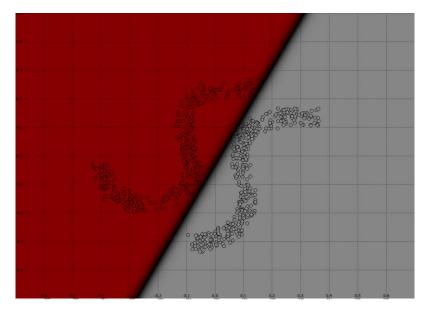
- A. Dataset A
- B. Dataset B
- C. Both
- D. None

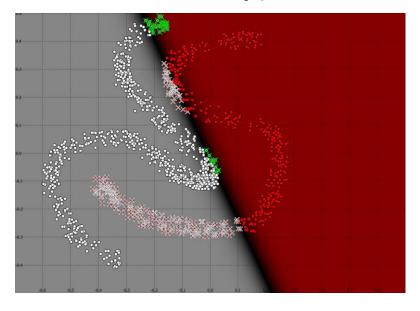






# For which dataset the perceptron model can achieve perfect classification (100% accuracy)



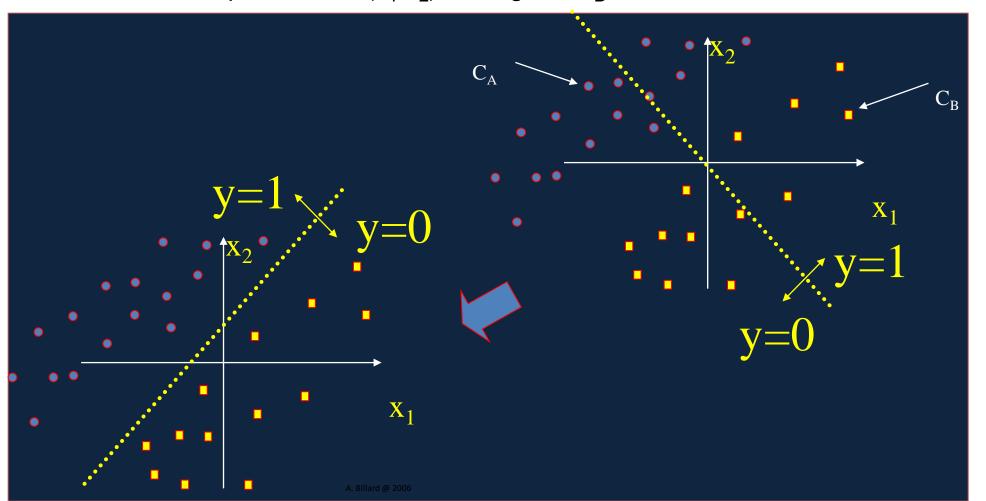


The perceptron is a simple linear classifier – identical to linear SVM



## Learning in one Neuron

Goal: we would like the neuron: to output 1 when  $(x_1,x_2)$  belongs to  $C_{A_1}$  and to output 0 when  $(x_1,x_2)$  belongs to  $C_{B_1}$ 





## Learning in one Neuron

The perceptron iteratively updates the weight vector w until the classification is correct

For each training pair  $(x_1,x_2)$ : If correctly classified:

do nothing

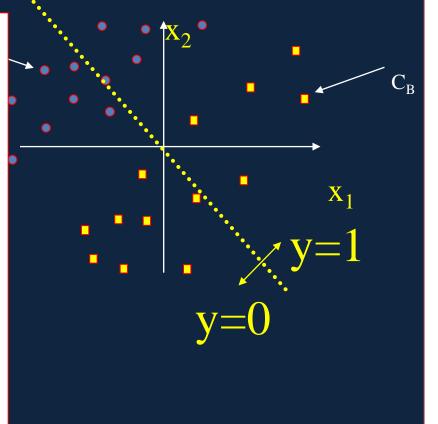
If wrongly classified:

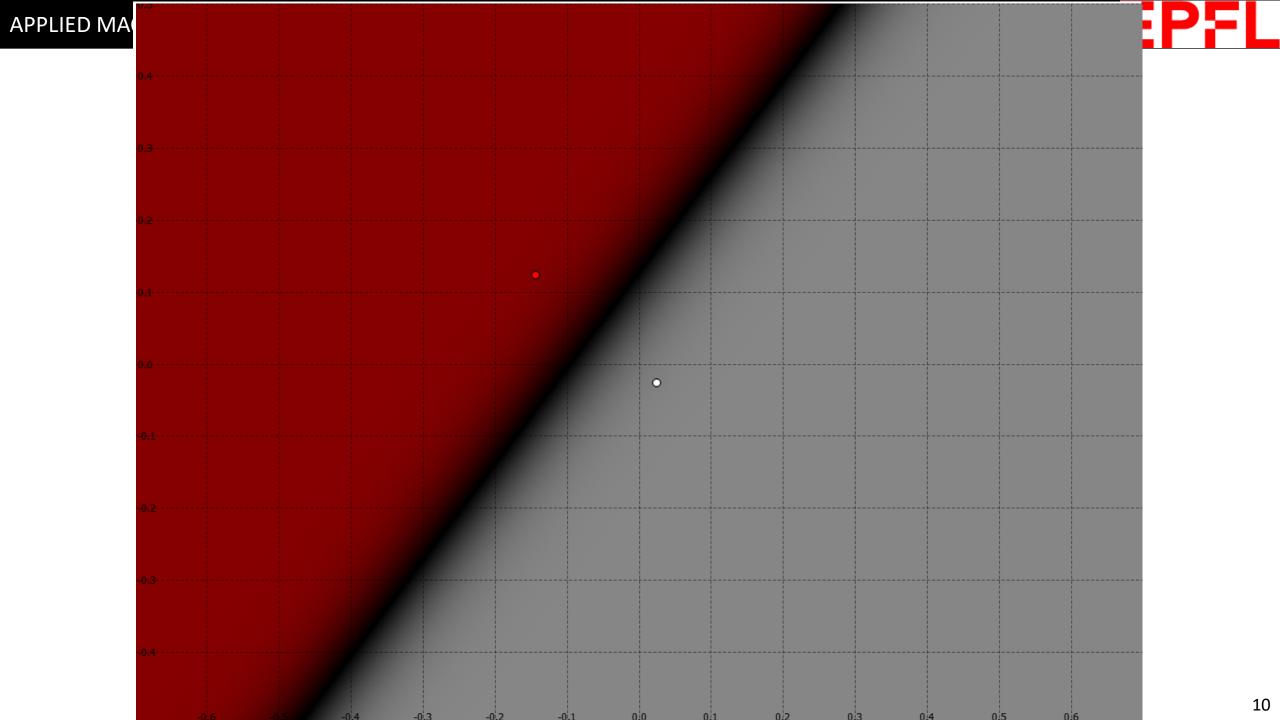
$$\vec{w}(t+1) = \vec{w}(t) - \eta \cdot \vec{x}$$

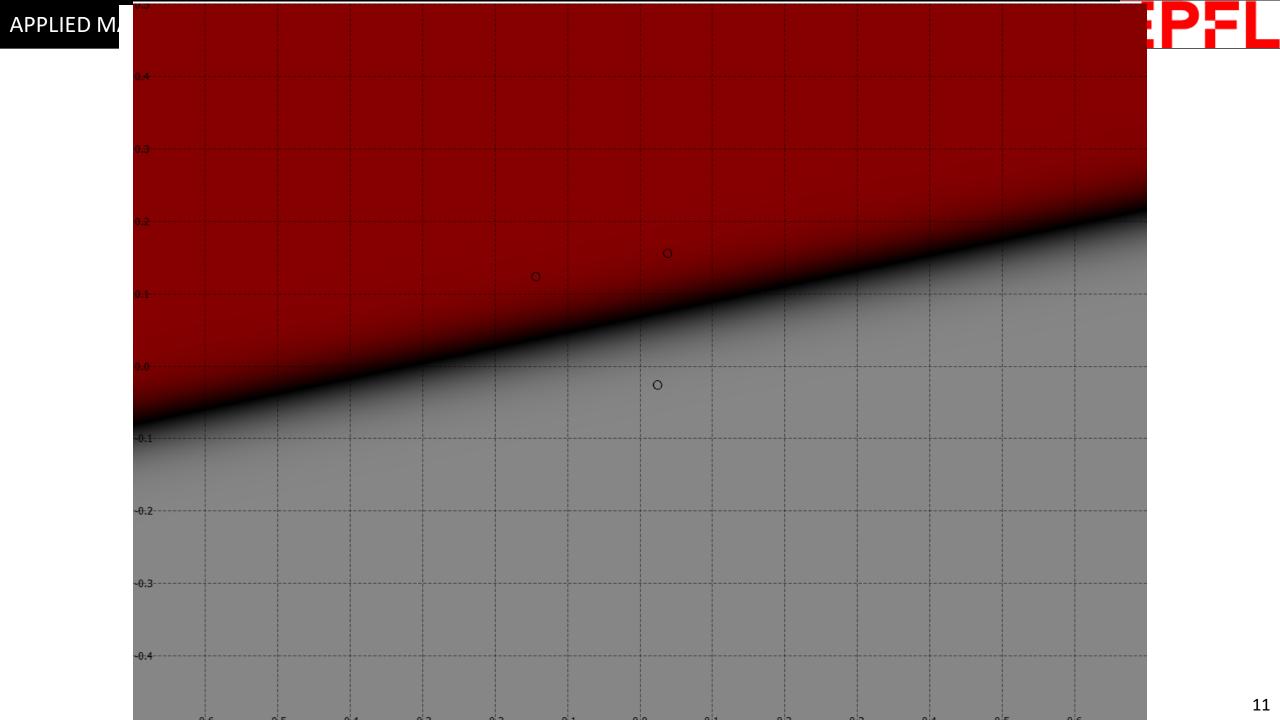
$$\text{if } \vec{w}^T \cdot \vec{x} \ge 0 \text{ and } \vec{x} \text{ belongs to } C_B$$

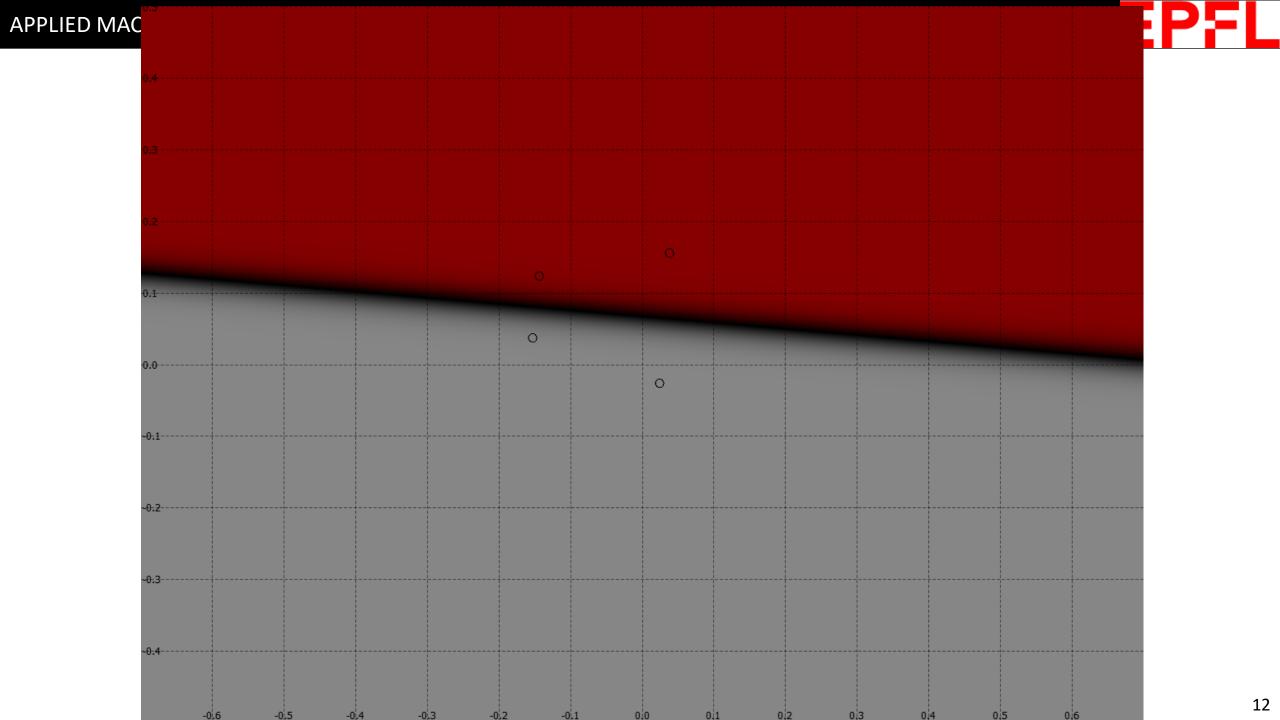
$$\vec{w}(t+1) = \vec{w}(t) + \eta \cdot \vec{x}$$

$$\text{if } \vec{w}^T \cdot \vec{x} < 0 \text{ and } \vec{x} \text{ belongs to } C_A$$



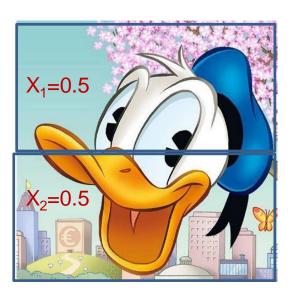


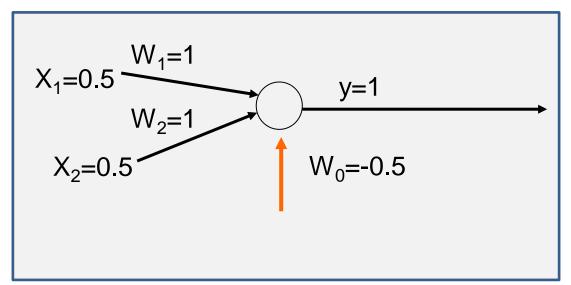






## The perceptron





$$y = \theta \left( \sum_{i=1}^{4} w_i \cdot x_i + w_0 \right)$$

**Assume** X<sub>i</sub> codes average grey scale intensity in part i of the image.

Is it Donald?

$$\theta(..) = \begin{cases} 1 & \text{Yes } \text{if } \mathbf{w}^T \mathbf{x} + w_0 \ge 0 \\ 0 & \text{No } \text{if } \mathbf{w}^T \mathbf{x} + w_0 < 0 \end{cases}$$

#### APPLIED MACHINE LEARNING







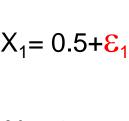




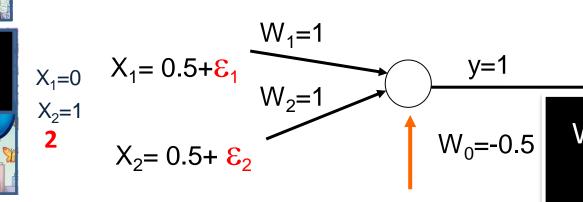




$$X_1 = 0$$
  
 $X_2 = 1$ 



$$X_2 = 0.5 + \epsilon_2$$



Which of the four types of disturbances would lead to incorrect classification?



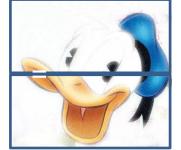
$$X_1=0.1$$
  
 $X_2=0.1$   
**3**

$$X_1 = 0.9$$
  
 $X_2 = 0.9$ 

4

Is it Donald?

$$\theta(...) = \begin{cases} 1 & \text{Yes if } \mathbf{w}^T \mathbf{x} + w_0 \ge 0 \\ 0 & \text{No if } \mathbf{w}^T \mathbf{x} + w_0 < 0 \end{cases}$$



The classifier is robust to large but local disturbances (case 2), but not to diffuse disturbances (case 3). How can we enlarge robustness?





$$X_1 = 0.4$$

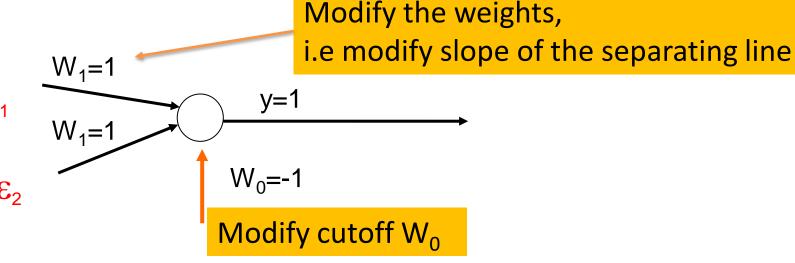
$$X_2 = 0.4$$



$$X_2 = \hat{x}$$

 $X_{1}=0$   $X_{1}=1+\varepsilon_{1}$   $X_{2}=1$   $X_{2}=1+\varepsilon_{2}$ 

$$X_2 = 1 + \frac{\epsilon_2}{\epsilon_2}$$





$$X_1 = 0.1$$

 $X_2 = 0.1$ 



$$X_2 = 0.9$$

Is it Donald?

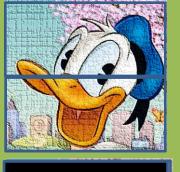
$$\theta(..) = \begin{cases} 1 & \text{Yes } \text{if } \mathbf{w}^T \mathbf{x} + w_0 \ge 0 \\ 0 & \text{No } \text{if } \mathbf{w}^T \mathbf{x} + w_0 < 0 \end{cases}$$

### Modify the encoding

4

The classifier is robust to large but local disturbances (case 2), but not to diffuse disturbances (case 3). How can we enlarge robustness?







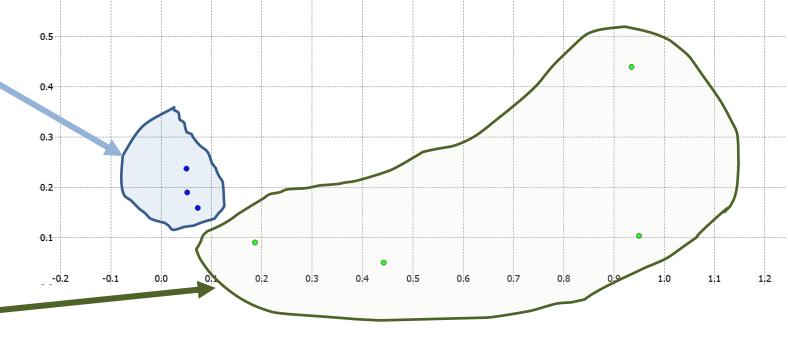




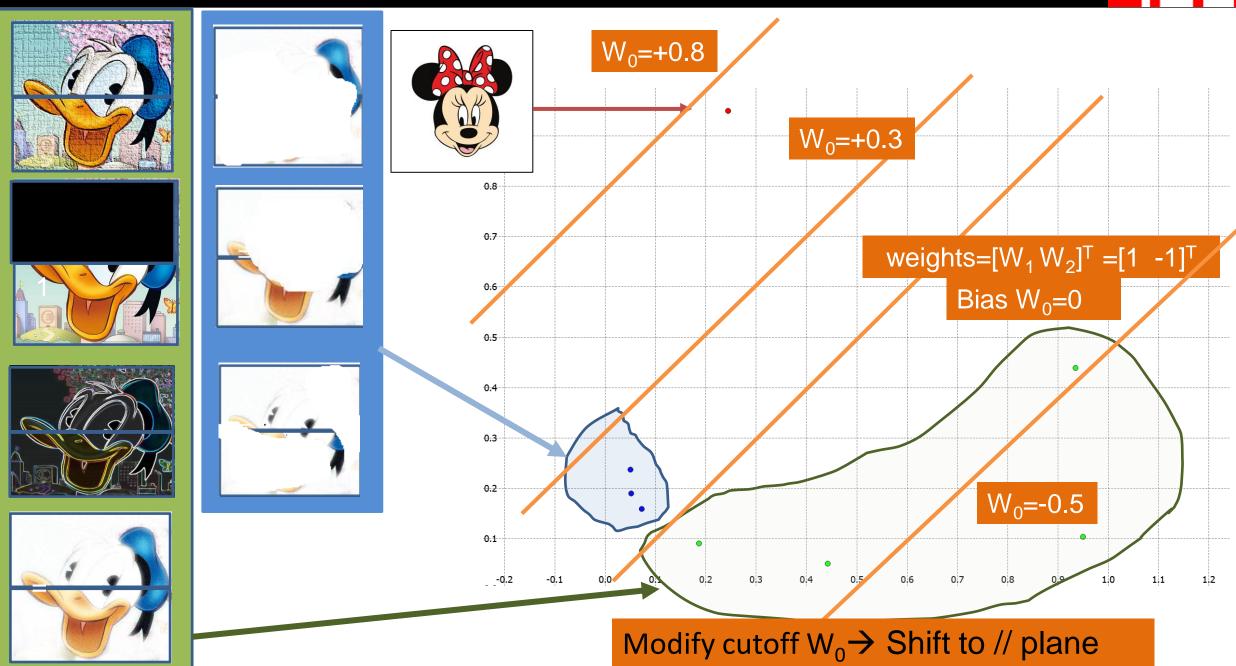




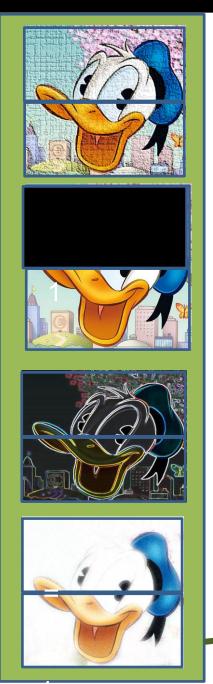
Consider now that you take real images. Perform PCA. These are the projection on the 1st two eigenvectors after PCA on real images.

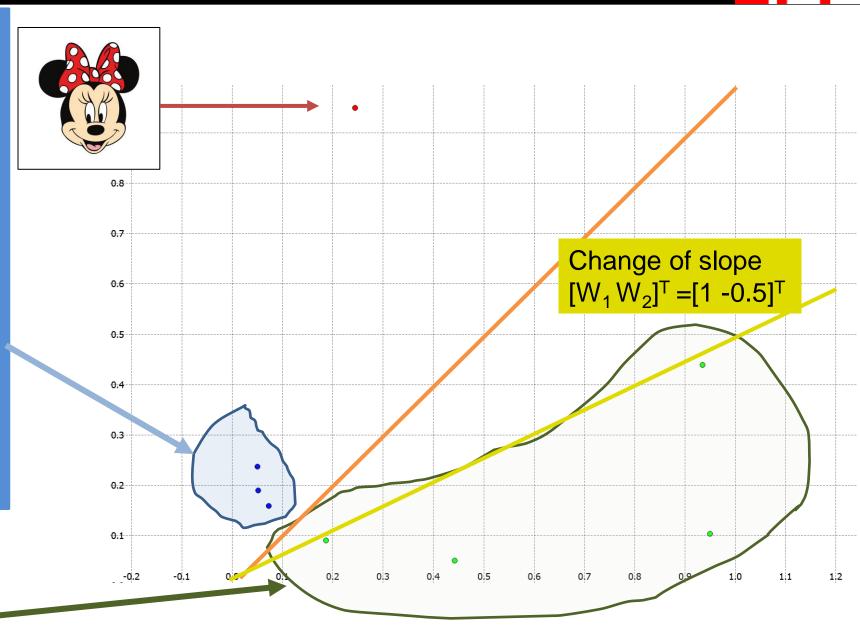








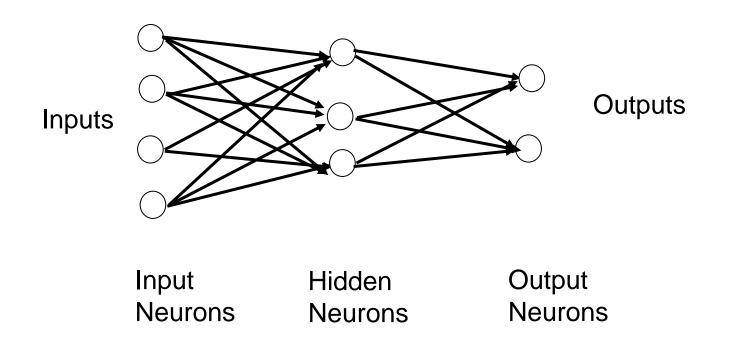






## Multi-Layer

### A two-layer Feed-Forward Neural Network



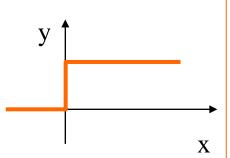


#### PERCEPTRON (1940, McCulloch & Pitts)

Threshold Function:

$$f(z) = \theta(z) = \begin{cases} 1 & \text{if } z \ge 0 \\ 0 & \text{if } z < 0 \end{cases} \quad \Longrightarrow \quad y = \theta \left( \sum_{i} w_{i} \cdot x_{i} \right) \quad \blacksquare$$

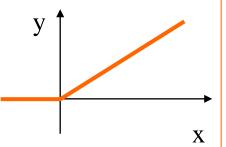
$$y = \theta \left( \sum_{i} w_{i} \cdot x_{i} \right)$$



RELU (Rectified Linear Unit) (Hahnloser & Seung 2011)

Threshold Function:

$$f(z) = \theta(z) = \max(0, z) = \begin{cases} z & \text{if } z \ge 0\\ 0 & \text{if } z < 0 \end{cases}$$



**Sigmoid** 

The steepest the sigmoid slope, the closer it is to perceptron

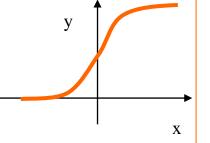
The flatter the sigmoid slope, the closer it is to RELU

Positive form:

$$f(z) = g(z) = \frac{1}{1 + e^{-D \cdot z}}$$

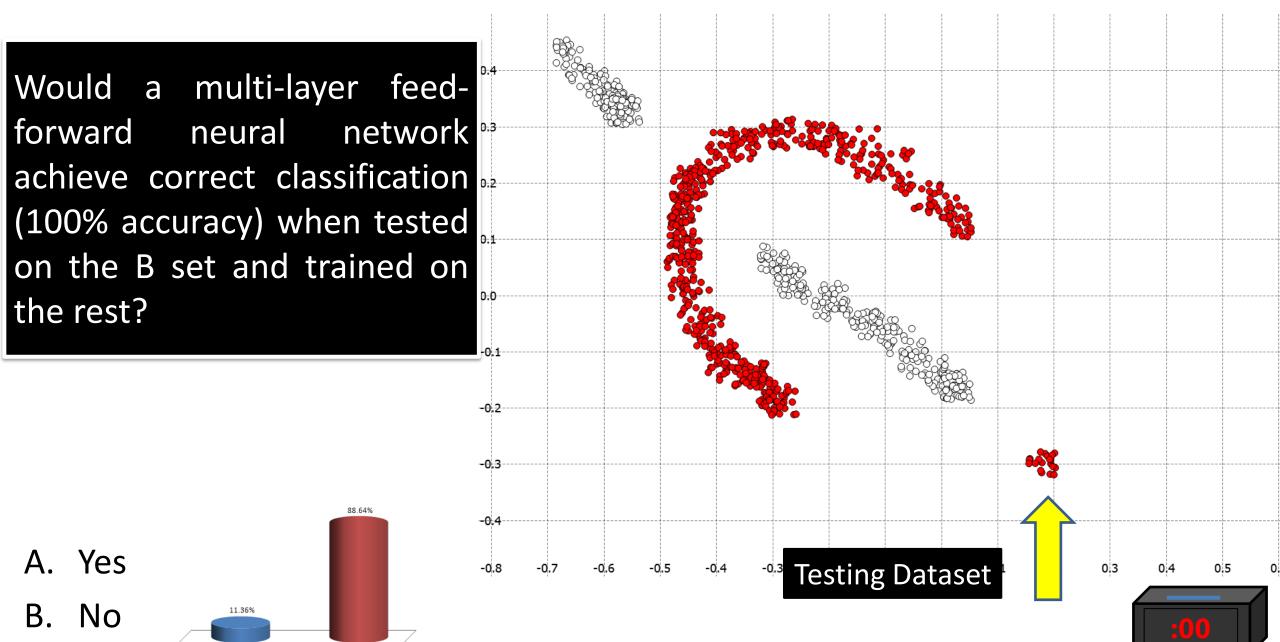


$$f(z) = g(z) = \frac{1}{1 + e^{-D \cdot z}} \qquad \Longrightarrow \qquad y = \frac{1}{1 + e^{-D \cdot \sum_{i} w_{i} \cdot x_{i}}}$$



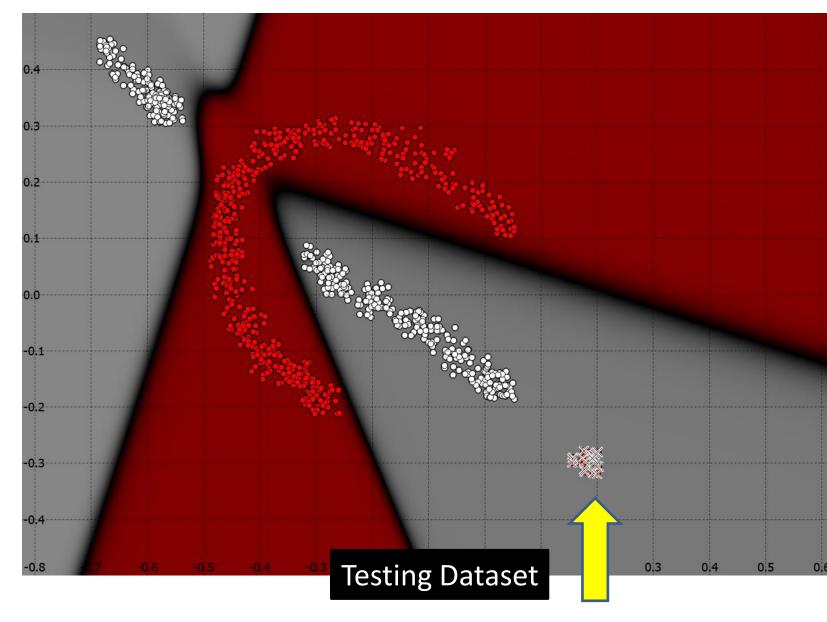
D is the slope of the function





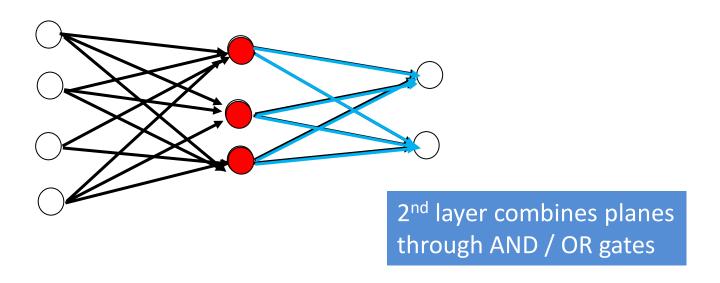


No, as the NN will learn piece-wise decomposition of the space.
Set of linear functions with ReLu or quasi linear with sigmoid activation fct.





Each hidden neuron creates a dividing plane.

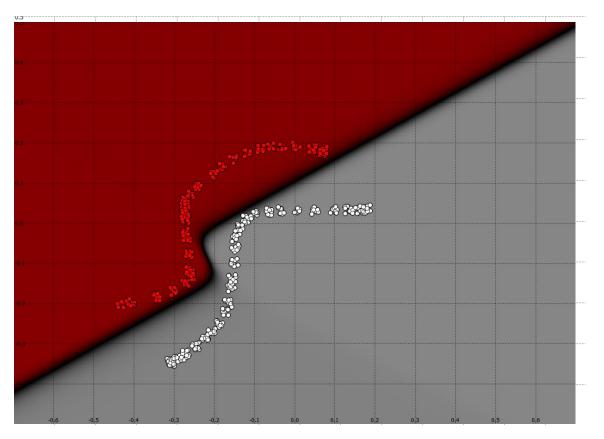


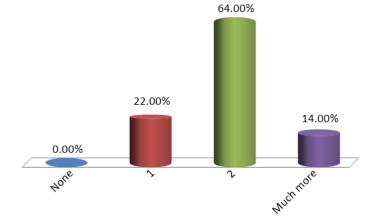


How many neurons do you need at minimum in the hidden layer to separate correctly these two classes?



- B. 1
- C. 2
- D. 3

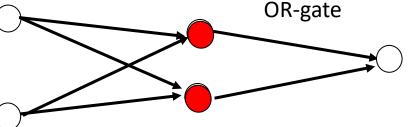












 $y>=0 \rightarrow red class$ 

 $y<0 \rightarrow$  white class

2<sup>nd</sup> cutting plane

2 neurons suffice as the boundary is composed of 2 partitions, but it depends on the slope of the activation function. 1<sup>st</sup> partition

2nd partition

Steeper sigmoid slope

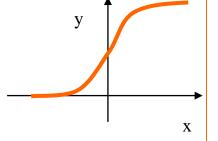
### **Sigmoid**

#### Positive form:

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

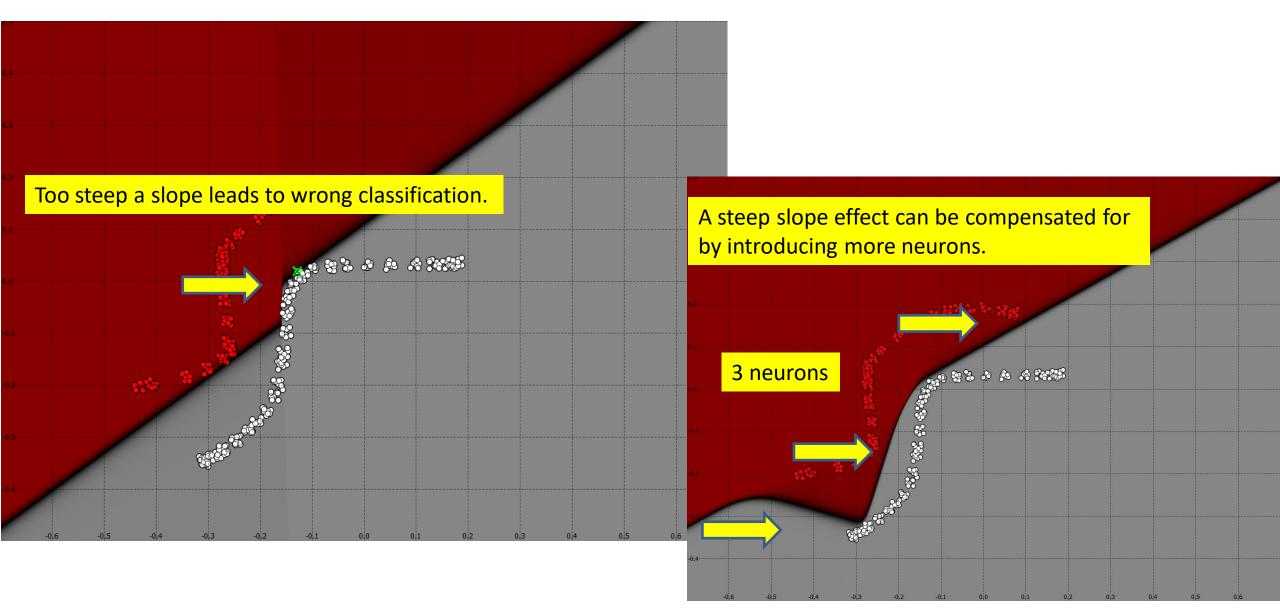


$$f(z) = g(z) = \frac{1}{1 + e^{-D \cdot z}} \qquad \Longrightarrow \qquad y = \frac{1}{1 + e^{-D \cdot \sum_{i} w_{i} \cdot x_{i}}}$$



D is the slope of the function



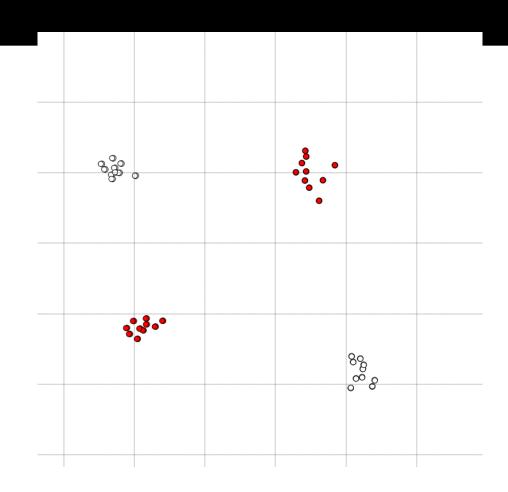


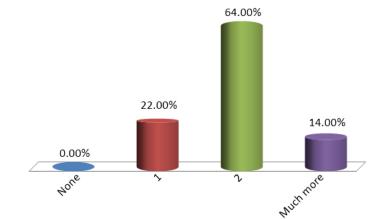
EPFL

How many neurons do you need at minimum in the hidden layer to separate correctly these two classes?



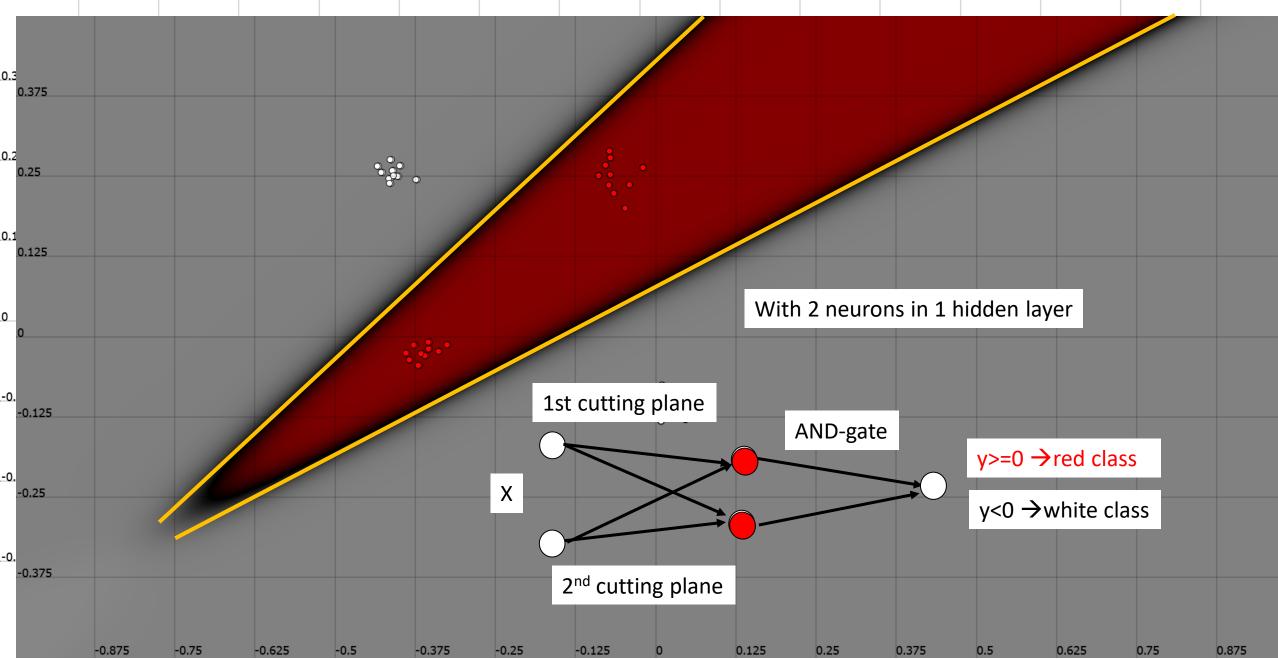
- B. 2
- C. 4
- D. 6









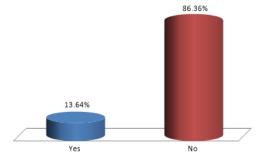




The more neurons, the more layers, the merrier?

A. Yes

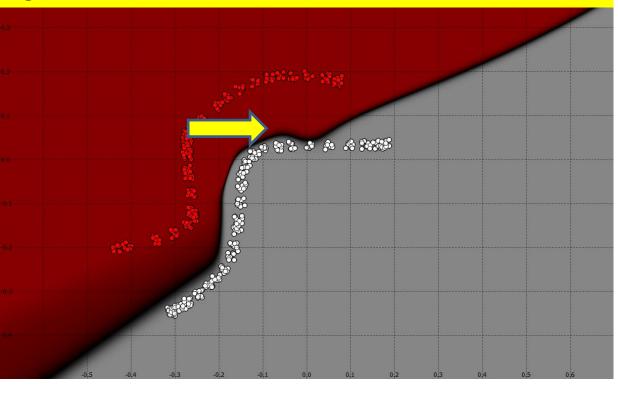
B. No

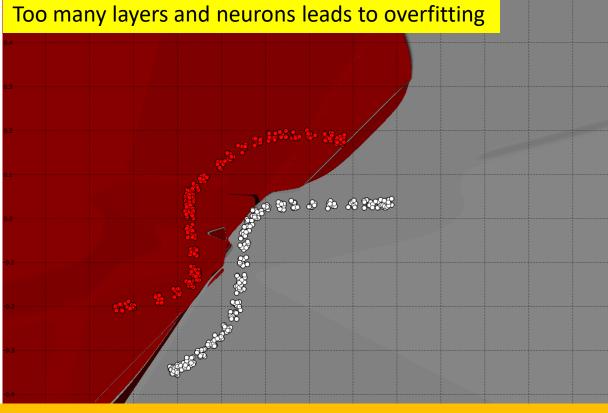






More neurons on 1-hidden layer leads to smoother and tighter curve

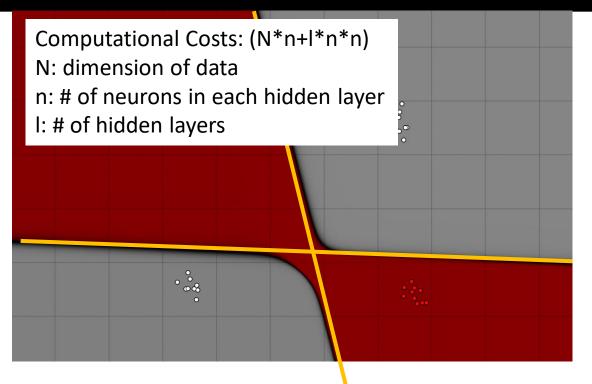




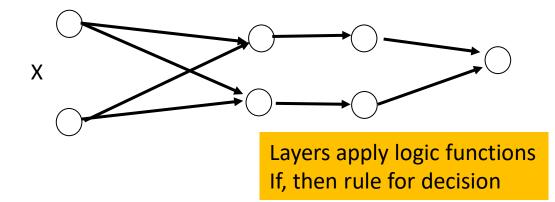
Here 5 layers, 10 neurons per layer → overfitting
With 500 datapoints total / too many parameters to estimate

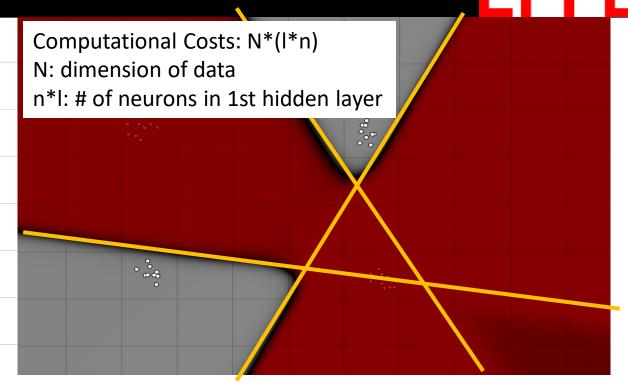
Early stopping with dropout and weight constraint can mitigate this effect.

#### APPLIED MACHINE LEARNING

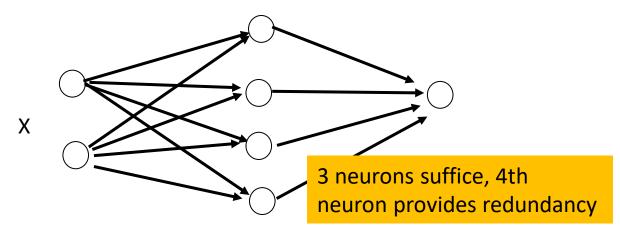


With 2 neurons in each of 2 hidden layers

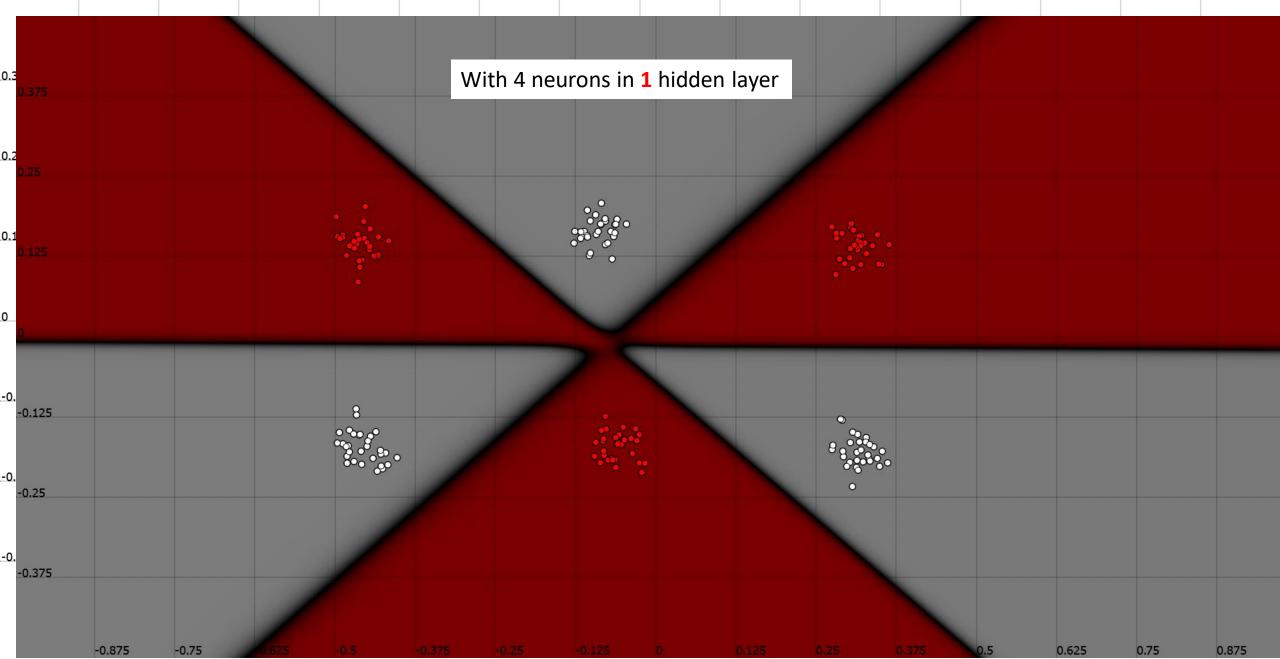




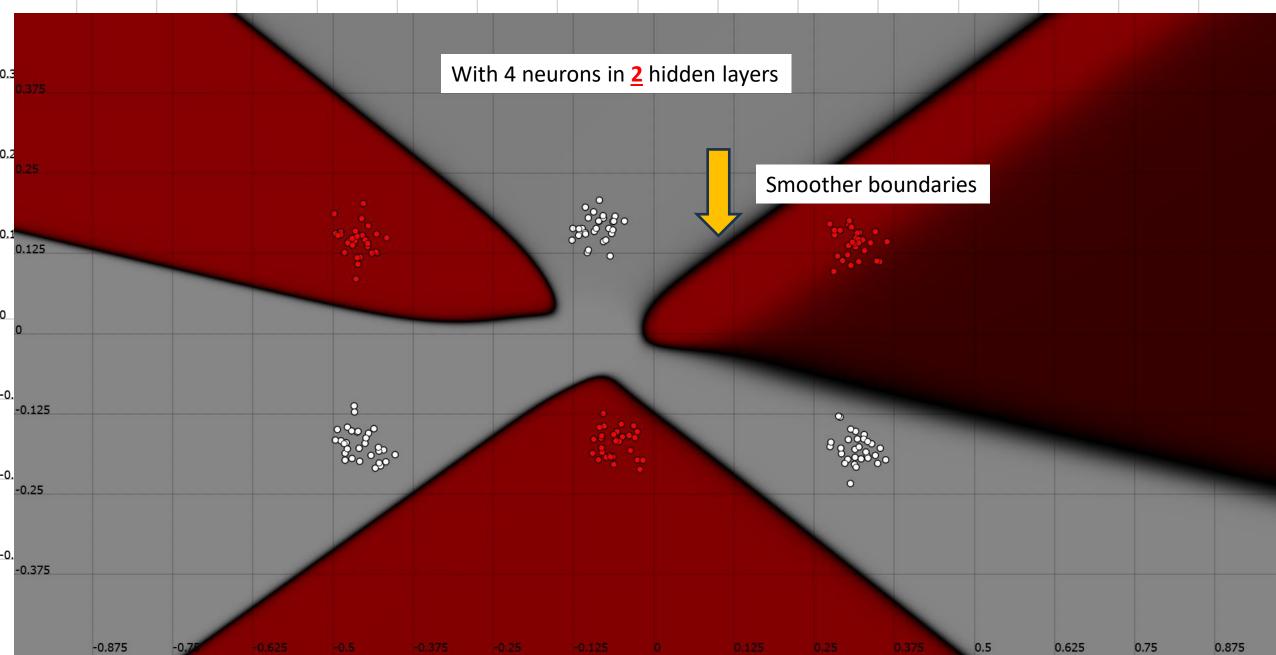
With 4 neurons in hidden layer













#### THE BACKPROPAGATION ALGORITHM

For each input-output pair, do:

- 1) Calculate the output y<sub>i</sub> of each output unit j of the net
- 2) Calculate  $\frac{\partial E}{\partial w_{ii}}$  for all weights on arcs to output nodes
- 3) Calculate  $\frac{\partial E}{\partial s_i}$  for each of these nodes
- 4) Calculate  $\frac{\partial E}{\partial s_i}$  for the last hidden layer
- 5) ... and so on, propagating  $\frac{\partial E}{\partial s_i}$  backward towards the first layer
- 6) finally, compute all the weight changes:

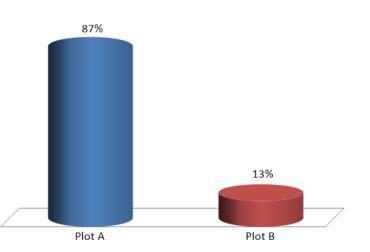
$$\Delta w_{ij} = -\eta \cdot \frac{\partial E}{\partial w_{ij}}$$

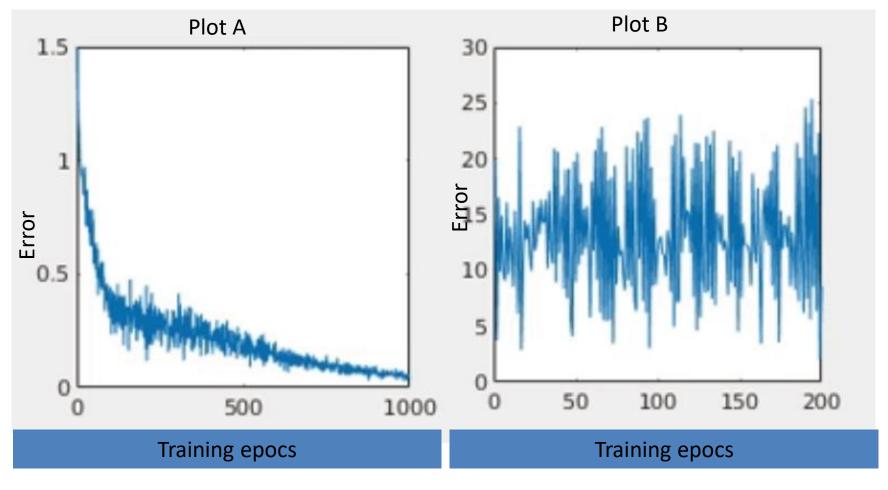


Which of the following error-plot corresponds to a NN trained with a lower learning rate?



B. Plot B

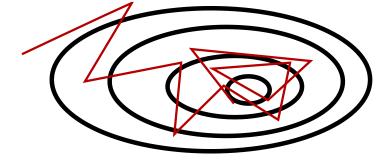






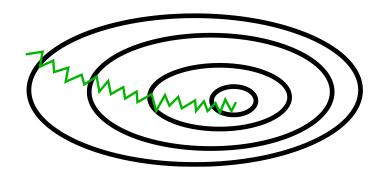
## Effects of the learning rate

large learning rate



Might not converge

small learning rate

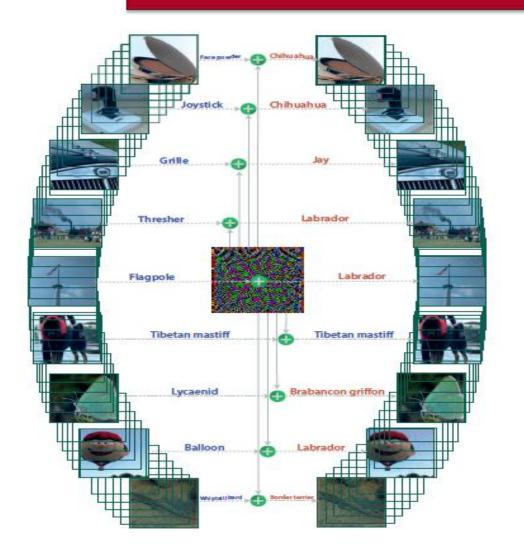


Slower learning

Adaptive learning rate (learning rate decay)



### How to Fool Neural Networks





Class 1



Class 0

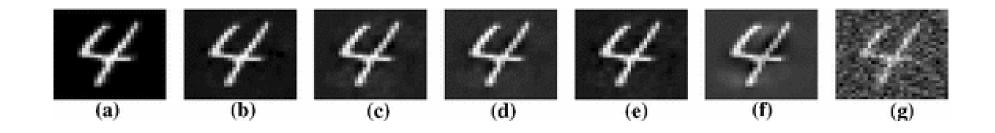
Add a small "image-agnostic" perturbation vx' = x + v

that leads to missclassification.

Algorithm that can learn from few training datapoints which "minimal" perturbation  $(\min \|\nu\|)$  that leads to missclassification with some probability p.

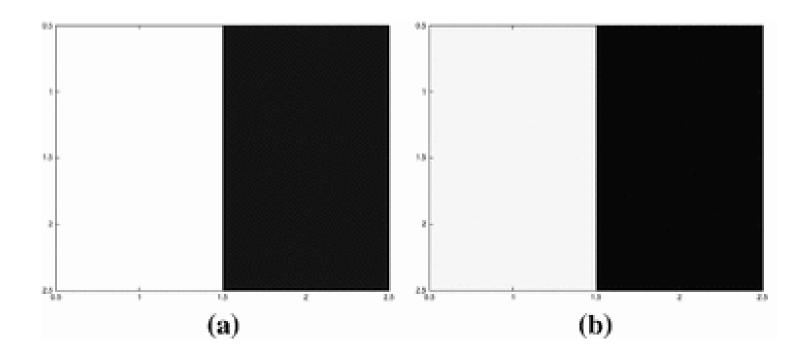


### **How to Fool Neural Networks**



Original image **a** and minimally perturbed images (**b**)–(**f**) that switch the estimated label of linear (**b**), quadratic (**c**), cubic (**d**), RBF(1) (**e**), RBF(0.1) (**f**) classifiers.





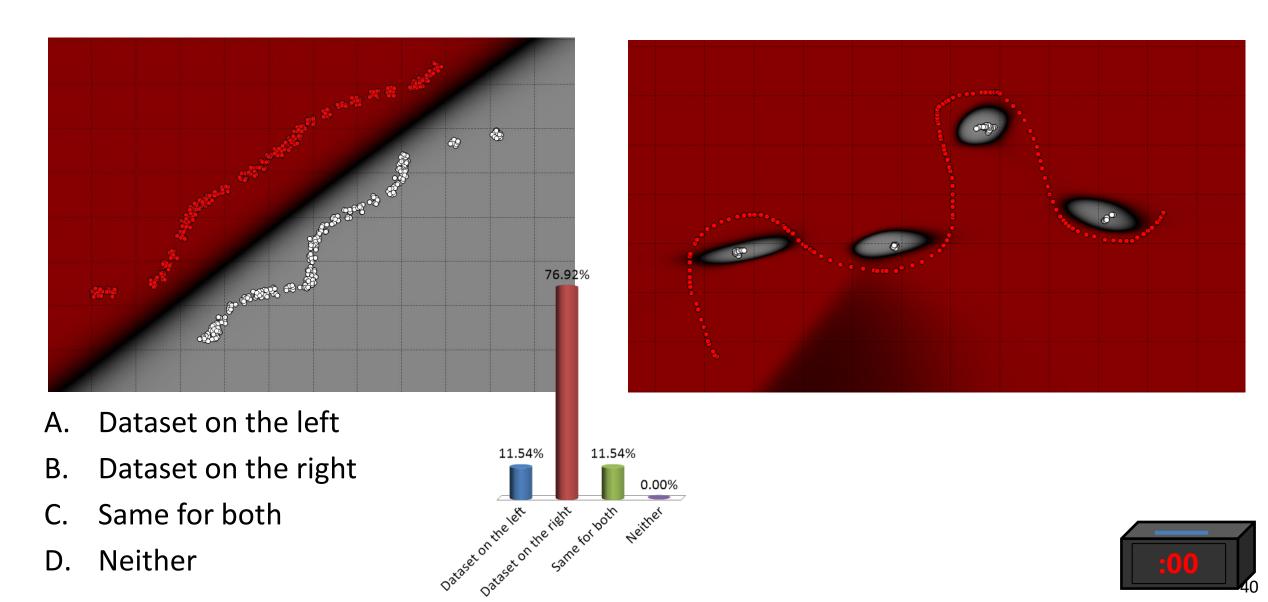
A small perturbation, quasi invisible, is applied to pattern (a) and makes it look like pattern (b). Pattern (b) changes label with a linear classifier.

Can you explain why?

The linear classifier has learned an AND-gate for all white (positive pixels). A small perturbation that reduces the grey scale value of the white pixels leads to incorrect prediction.

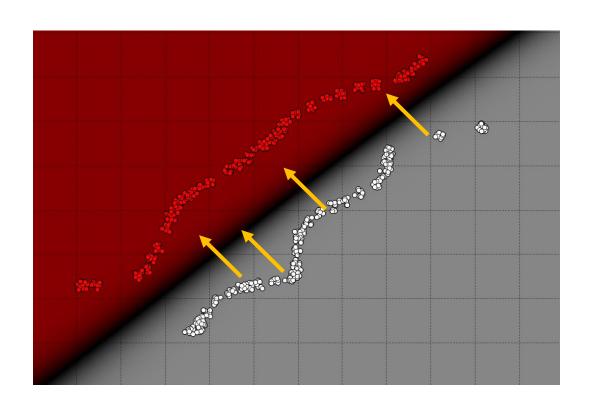


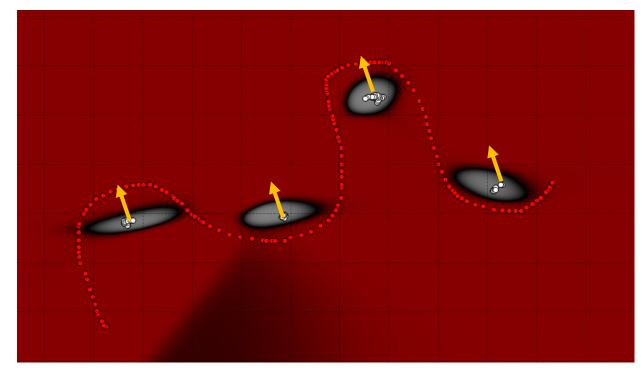
## Which of the two datasets would be easiest to fool?





### Which of the two datasets would be easiest to fool?



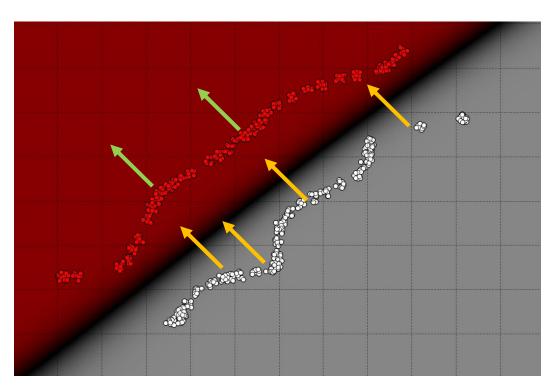


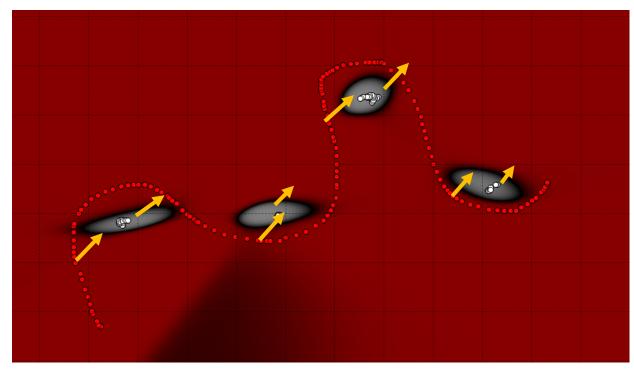
#### It is as easy if one wants to fool one class for the other.

In the linear classification case (left), one just needs to find the correct perturbation vector. In nonlinear classification case (right), it becomes easy when the two classes are tightly mingled and we have quasi overfitting across one class.



## Which of the two datasets would be easiest to fool?





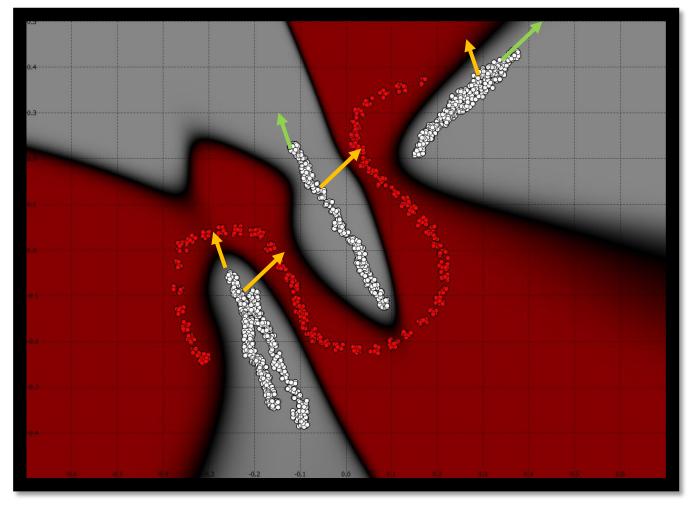
It is not as easy if one wants to fool both classes with a single vector.

In the linear classification case (left), one vector will be insufficient.

In the nonlinear classification case (right), one vector may be sufficient to lead to incorrect classifications of elements of both classes.

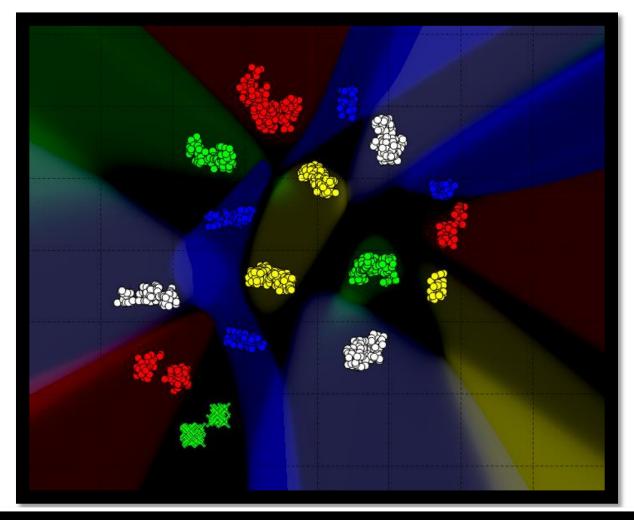
The more mingled the classes, the easier it is to fool the NN. It is hence crucial to identify features (subdimensions, subpatterns) distinct for each class and based the classification on these features.





If the curvature of the boundary is very complex, it becomes more difficult to find a single vector that works for all the elements of the classes.





The more classes, the tighter the fit, the easier to generate a perturbation that will lead to misclassification with non-negligible probability on all classes.

How to construct boundaries robust to such perturbations?



1: Build a pre-processing layer to detect « perturbations ».

Akhtar, Naveed, Jian Liu, and Ajmal Mian. "Defense against universal adversarial perturbations." *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 2018

2: Train using adversarial approach – make the network more robust to this through adversarial network

Ali Shafahi, Mahyar Najibi, Zheng Xu, John P. Dickerson, Larry S. Davis, and Tom Goldstein. Universal adversarial training. ArXiv, abs/1811.11304, 2020.



## Which of the three classification techniques are most sensitive to fooling?

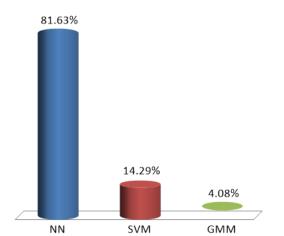
NN

Class label determined by  $y \sim \sum w_{ij} (x + v)$ , with Re Lu

B. SVM

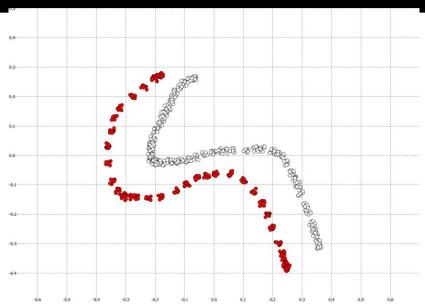
**GMM** 

Class label determined by  $y \sim \left(\sum_{i} w_{ij} \cdot k(x + v)\right)$  k : RBF or Gauss Fct

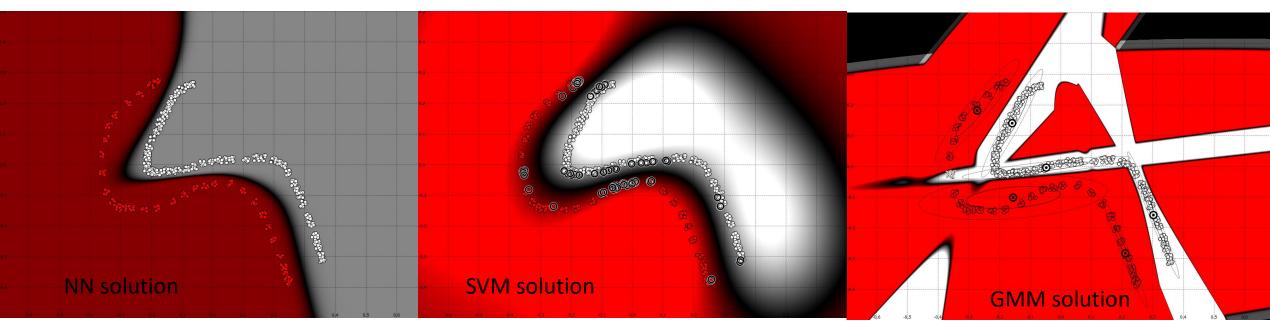








All techniques are as sensitive, as all decisions functions depend on a measure of the distance to the boundary.



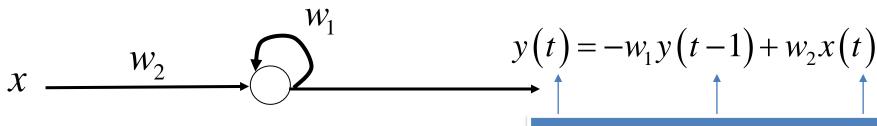


#### Recurrent Neural Networks

Most biological network have recurrent connections.

This change of direction in the flow of information is interesting, as it can allow:

- To keep a memory of the activation of the neuron
- To propagate the information across output neurons



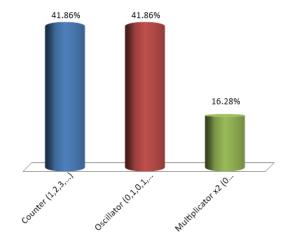
Differentiated response over time

Continuous model: 
$$\dot{y} = -\frac{1}{\tau}y + \int xdt$$



# What type of function can we embed in a single recurrent neuron?

- A. Counter (1,2,3,...)
- B. Oscillator (0,1,0,1, ...)
- C. Multiplicator x2 (0; 2)



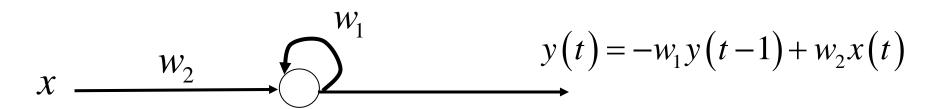


Binary input 1/0



## What type of function can we embed in a single recurrent neuron?

- A. Counter (1,2,3,...)  $w_1 = -1$  and  $w_2 = 1$
- B. Oscillator (0,1,0,1, ...)  $w_1 = w_2 = 1$
- C. Multiplicator x2 (0; 2) Not possible



Binary input 1/0