Montreal, July 8–12

# DLMI2024
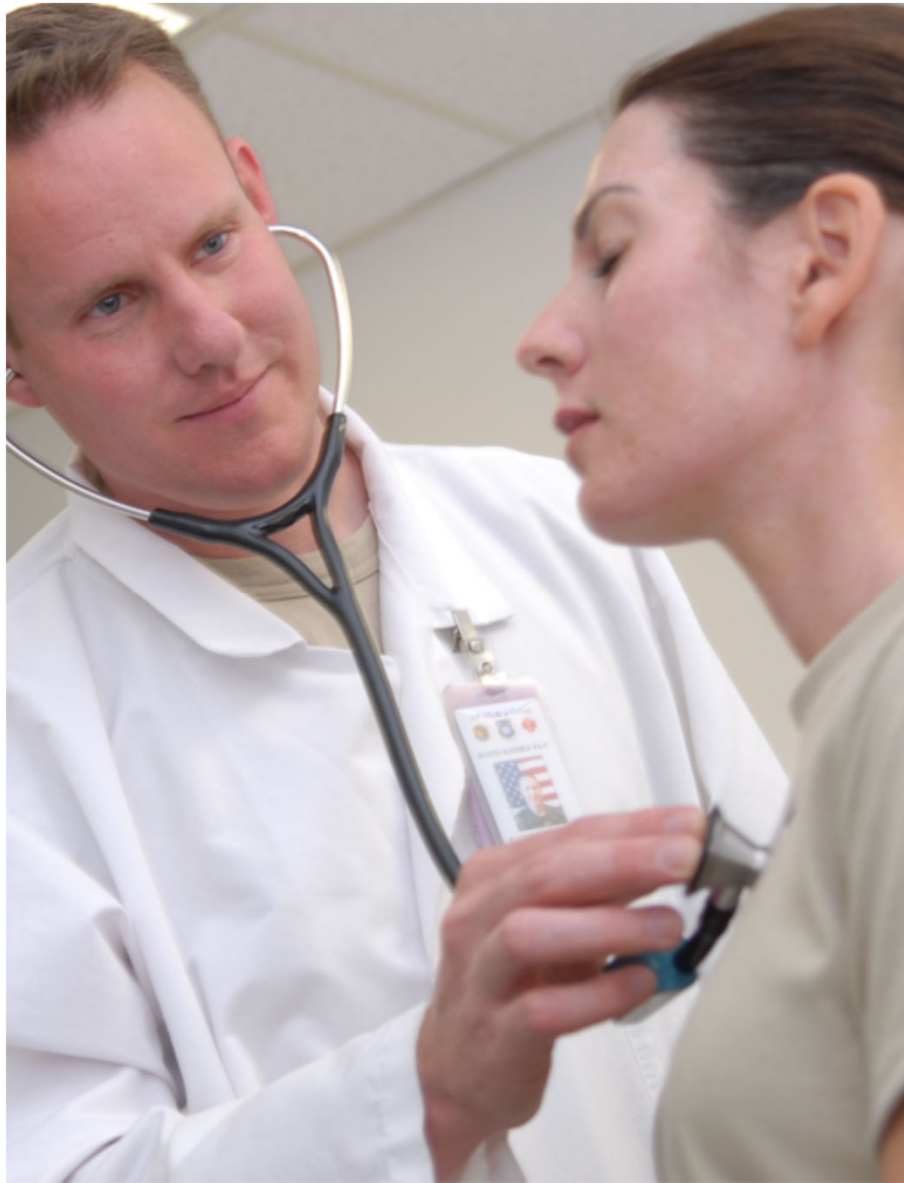
# Basics of deep learning part 2

By

Pierre-Marc Jodoin
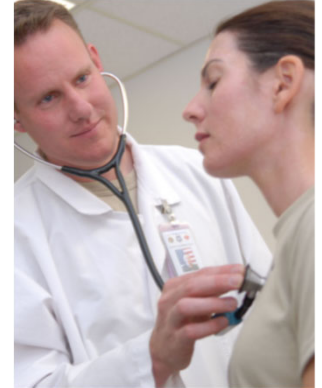
UNIVERSITÉ DE
SHERBROOKE

# Lets start with a simple example
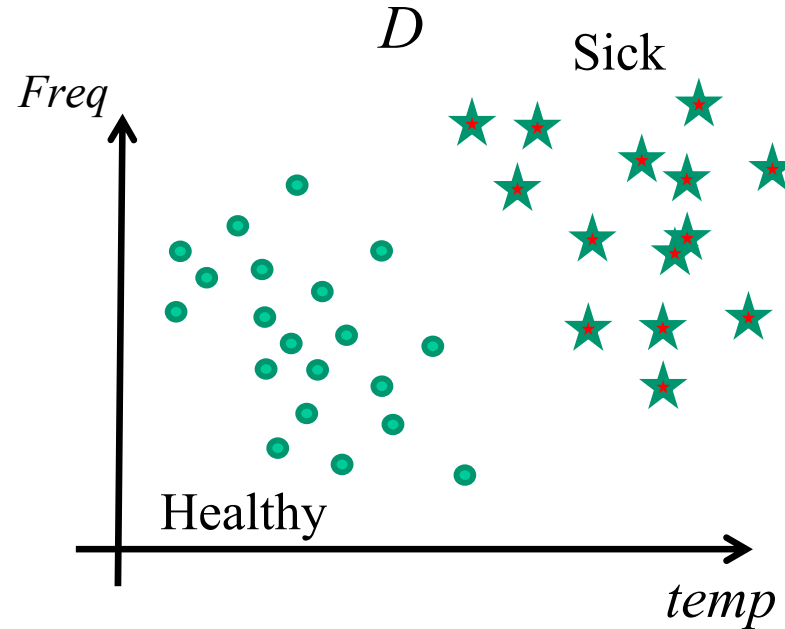


From Wikimedia Commons
the free media repository

# Lets start with a simple example

$D$

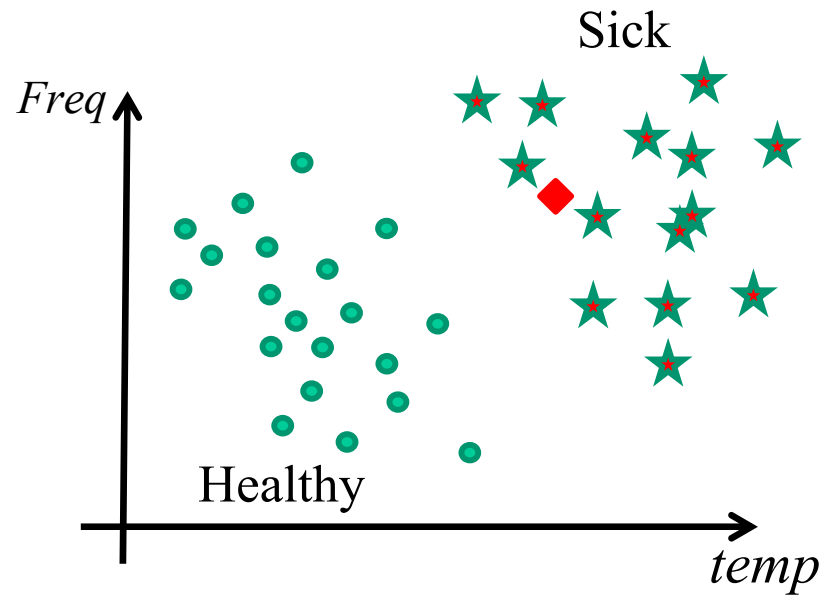| | ( temp, freq) | diagnostic |
|---|---|---|
| Patient 1 | (37.5, 72) | Healthy |
| Patient 2 | (39.1, 103) | Sick |
| Patient 3 | (38.3, 100) | Sick |
| | (…) | … |
| Patient N | (36.7, 88) | Healthy |

$\vec{x}$       $t$

$D$

Sick

*Freq*

Healthy

*temp*

# Lets start with a simple example

A new patient comes to the hospital
**How can we determine if he is sick or not?**

Sick

*Freq*

Healthy

*temp*

# Solution

**Divide the feature space in 2 regions : <span style="color:red">sick</span> and <span style="color:red">healthy</span>**

# Solution

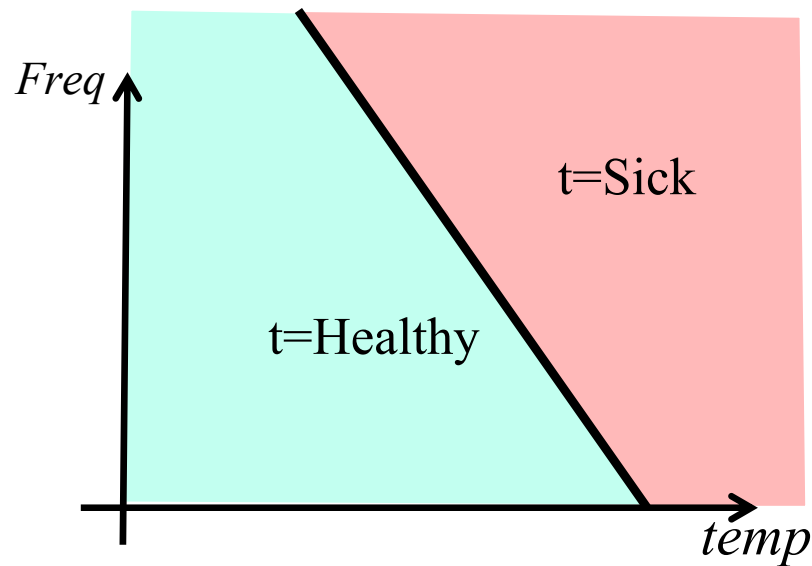**Divide the feature space in 2 regions : <span style="color:red">sick</span> and <span style="color:red">healthy</span>**

# More formally

$$y(\vec{x}) = \begin{cases} Healthy \text{ if } \vec{x} \text{ is in the green region} \\ Sick \text{ otherwise} \end{cases}$$



Freq

t=Sick

t=Healthy

temp
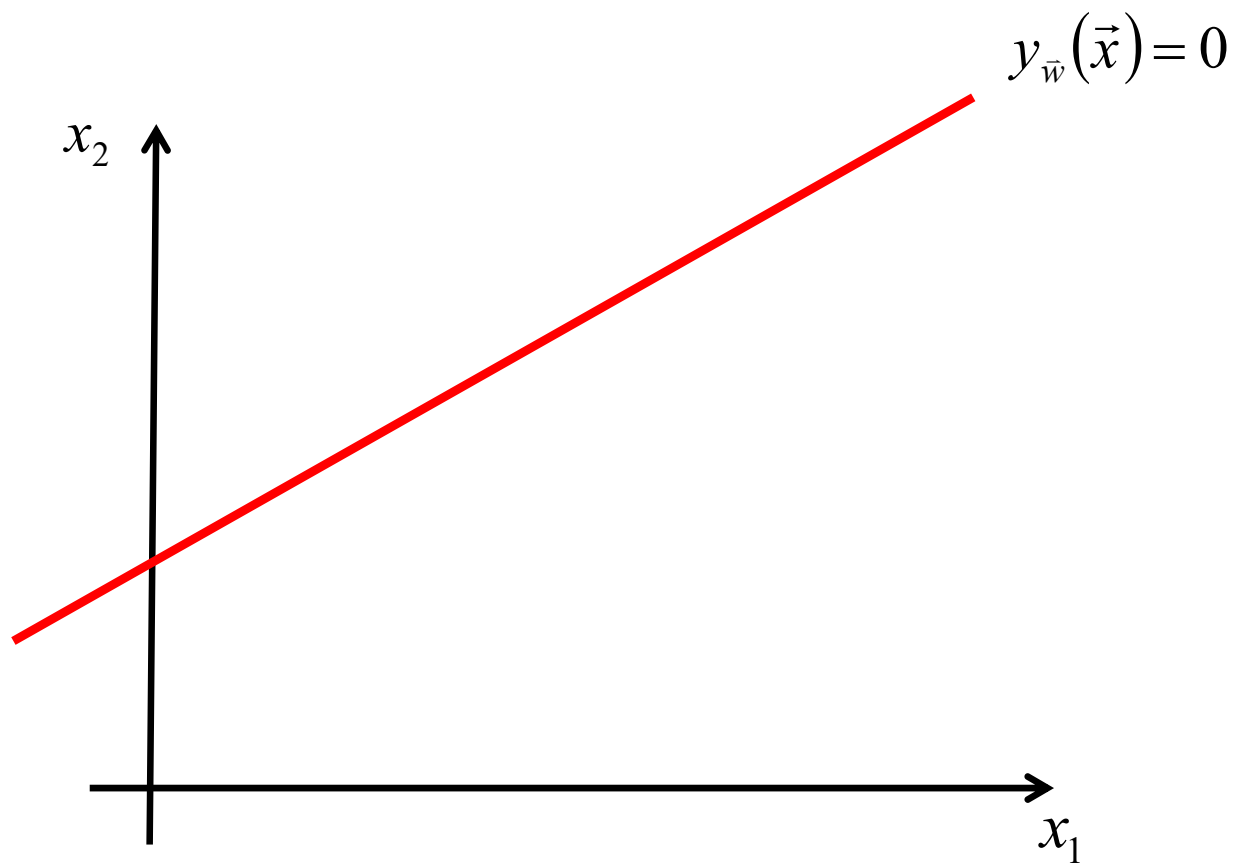
How to split the feature space?

# Linear function

$$y_{\vec{w}}(\vec{x}) = w_1 x_1 + w_2 x_2 + w_0$$

$y_{\vec{w}}(\vec{x}) = 0$

# Classification function

$$y_{\vec{w}}(\vec{x}) = w_1 x_1 + w_2 x_2 + w_0$$

$y_{\vec{w}}(\vec{x}) = 0$

$x_2$
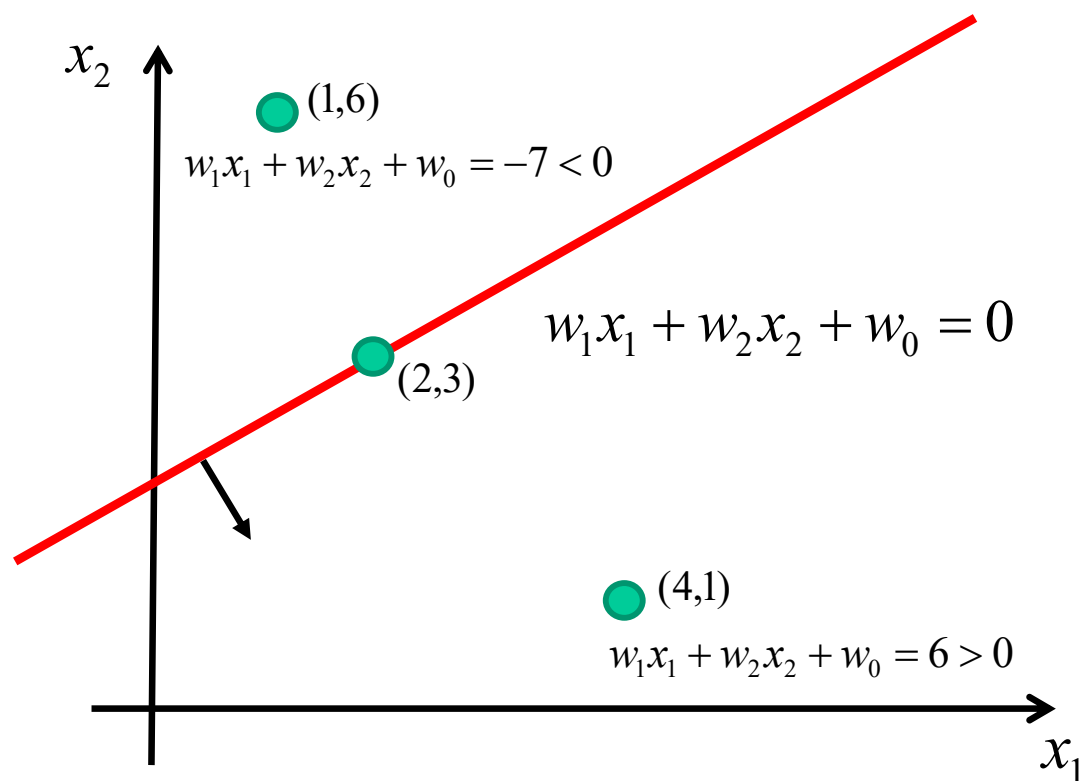
$y_{\vec{w}}(\vec{x}) < 0$

$y_{\vec{w}}(\vec{x}) > 0$

$\vec{N} = (w_1, w_2)$

$x_1$

# Classification function

$$y_{\vec{w}}(\vec{x}) = w_1 x_1 + w_2 x_2 + w_0$$

$w_1 = 1.0$

$w_2 = -2.0$

$w_0 = 4.0$

(1,6)

$w_1 x_1 + w_2 x_2 + w_0 = -7 < 0$

$w_1 x_1 + w_2 x_2 + w_0 = 0$

(2,3)

(4,1)

$w_1 x_1 + w_2 x_2 + w_0 = 6 > 0$

$x_2$

$x_1$

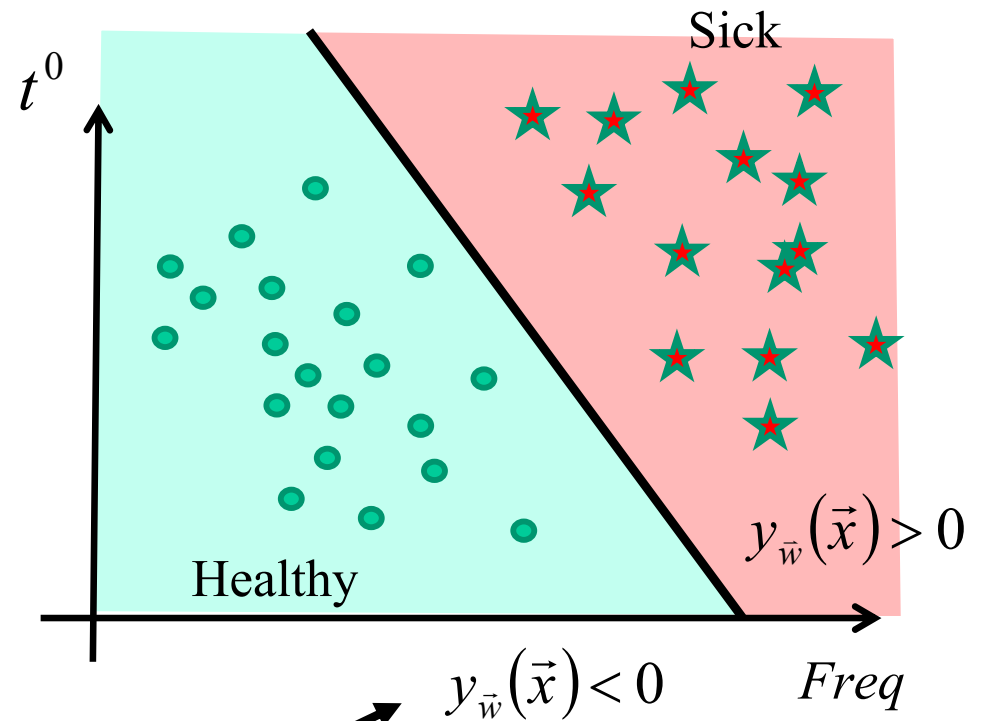linear classification = dot product with bias included

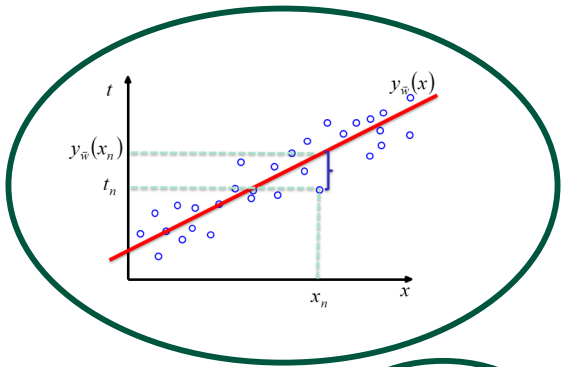$$y_{\vec{w}}(\vec{x}) = \vec{w}^T \vec{x}$$

# Learning

With the **training dataset** $D$

**the GOAL is to**

find the parameters $(w_0, w_1, w_2)$ that would best separate the two classes.

$t^0$

Sick

$y_{\vec{w}}(\vec{x}) > 0$

Healthy

$y_{\vec{w}}(\vec{x}) < 0$

*Freq*

$$\vec{w} = \arg\min_{\vec{w}} \sum_{n=1}^{N} \left(\vec{w}^{\mathrm{T}}\vec{x}_n - t_n\right)^2$$

$$\vec{w} = \left(X^{\mathrm{T}}X\right)^{-1}X^{\mathrm{T}}T$$

YES! If data follows a **Gaussian** distribution

$$\vec{w} = \arg\min_{\vec{w}} \sum_{n=1}^{N} (\vec{w}^{\mathrm{T}} \vec{x}_n - t_n)^2$$

$$\vec{w} = (X^{\mathrm{T}} X)^{-1} X^{\mathrm{T}} T$$

Otherwise, we need another solution

18

# Loss function



$$L(y_{\vec{w}}(\vec{x}), D) \approx 0$$

$$L(y_{\vec{w}}(\vec{x}), D) > 0$$

$$L(y_{\vec{w}}(\vec{x}), D) >> 0$$

**Good!**  **Medium**  **BAD!**

# So far…

1. Training dataset: $D$
2. Classification function (a line in 2D) : $y_{\vec{w}}(\vec{x}) = w_1 x_1 + w_2 x_2 + w_0$
3. Loss function: $L(y_{\vec{w}}(\vec{x}), D)$

4. Training : find $(w_0, w_1, w_2)$ that minimize $L(y_{\vec{w}}(\vec{x}), D)$

# Before deep neural nets were … **linear models**

# In this session…



Perceptron
Logistic regression
Multi-layer perceptron

# Perceptron



Rosenblatt, Frank (1958), **The Perceptron: A Probabilistic Model for Information Storage and Organization in the Brain**, Psychological Review, v65, No. 6, pp. 386–408

# Perceptron
## (2D and 2 classes)



$$y_{\vec{w}}(\vec{x}) = w_0 + w_1 x_1 + w_2 x_1$$

$$= \vec{w}^{\mathrm{T}} \vec{x}$$

24

# Perceptron
(2D and 2 classes)

$$y_{\vec{w}}(\vec{x}) = sign(\vec{w}^{\mathrm{T}}\vec{x})$$

Activation function

$y_{\vec{w}}(\vec{x}) > 0$

$y_{\vec{w}}(\vec{x}) = 0$

$y_{\vec{w}}(\vec{x}) < 0$

$\vec{w}$

$x_1$

$x_2$

$1$

$w_0$

$w_1$

$w_2$

$\vec{w}^{\mathrm{T}}\vec{x}$

$sign$

$y_{\vec{w}}(\vec{x}) \in \{-1, +1\}$

# Perceptron

## (2D and 2 classes)



$$y_{\vec{w}}(\vec{x}) = sign(\vec{w}^\mathrm{T}\vec{x}) \in \{-1,+1\}$$

**Neuron**

Dot product + activation function

# So far…

1. Training dataset: $D$
2. Classification function (a line in 2D): $y_{\vec{w}}(\vec{x}) = w_1 x_1 + w_2 x_2 + w_0$
3. Loss function: $L\left(y_{\vec{w}}(\vec{x}), D\right)$

# So far…

1. Training dataset: $D$

2. Classification function (a line in 2D) :



3. Loss function: $L(y_{\vec{w}}(\vec{x}), D)$

4. Training : find $(w_0, w_1, w_2)$ that minimize $L(y_{\vec{w}}(\vec{x}), D)$

# Linear classifiers have **limits**

# Non-linearly separable training data



Linear classifier = **<u>large error rate</u>**

# Non-linearly separable training data

Three classical solutions

1. Acquire more observations
2. Use a non-linear classifier
3. Transform the data

# Non-linearly separable training data



Three classical solutions

1. **More observations**
2. Use a non-linear classifier
3. Transform the data

# Acquire more data



$D$

| | ( temp, freq) | diagnostic |
|---|---|---|
| Patient 1 | (37.5, 72) | healthy |
| Patient 2 | (39.1, 103) | sick |
| Patient 3 | (38.3, 100) | sick |
| | (…) | … |
| Patient N | (36.7, 88) | healthy |

$\vec{x}$        $t$

$D$

| | ( temp, freq, headache) | Diagnostic |
|---|---|---|
| Patient 1 | (37.5, 72, 2) | healthy |
| Patient 2 | (39.1, 103, 8) | sick |
| Patient 3 | (38.3, 100, 6) | sick |
| | (…) | … |
| Patient N | (36.7, 88, 0) | healthy |

$\vec{x}$        $t$

# Non-linearly separable training data



$$y_{\vec{w}}(\vec{x}) = w_1 x_1 + w_2 x_2 + w_0$$

(line)

$$y_{\vec{w}}(\vec{x}) = w_1 x_1 + w_2 x_2 + w_3 x_3 + w_0$$

(plane)

# Perceptron
(2D and 2 classes)



$$y_{\vec{w}}(\vec{x}) = sign(\vec{w}^{\mathrm{T}}\vec{x}) \in \{-1, +1\}$$

$$y_{\vec{w}}(\vec{x}) = w_1 x_1 + w_2 x_2 + w_0$$

(line)

# Perceptron
## (3D and 2 classes)

$y_{\vec{w}}(\vec{x}) = +1$

$y_{\vec{w}}(\vec{x}) = -1$

1 $\quad w_0$

$x_1 \quad w_1$

$x_2 \quad w_2$

$x_3 \quad w_3$

$\vec{w}^{\mathrm{T}}\vec{x} \;\; sg$

$y_{\vec{w}}(\vec{x}) = sign(\vec{w}^{\mathrm{T}}\vec{x}) \in \{-1,+1\}$

$$y_{\vec{w}}(\vec{x}) = w_1 x_1 + w_2 x_2 + w_3 x_3 + w_0$$

(plane)

36

# Perceptron
## (K-D and 2 classes)



$$y_{\vec{w}}(\vec{x}) = sign\left(\vec{w}^{\mathrm{T}}\vec{x}\right) \in \{-1, +1\}$$

$$y_{\vec{w}}(\vec{x}) = w_1 x_1 + w_2 x_2 + w_3 x_3 + \ldots + w_K x_K + w_0$$
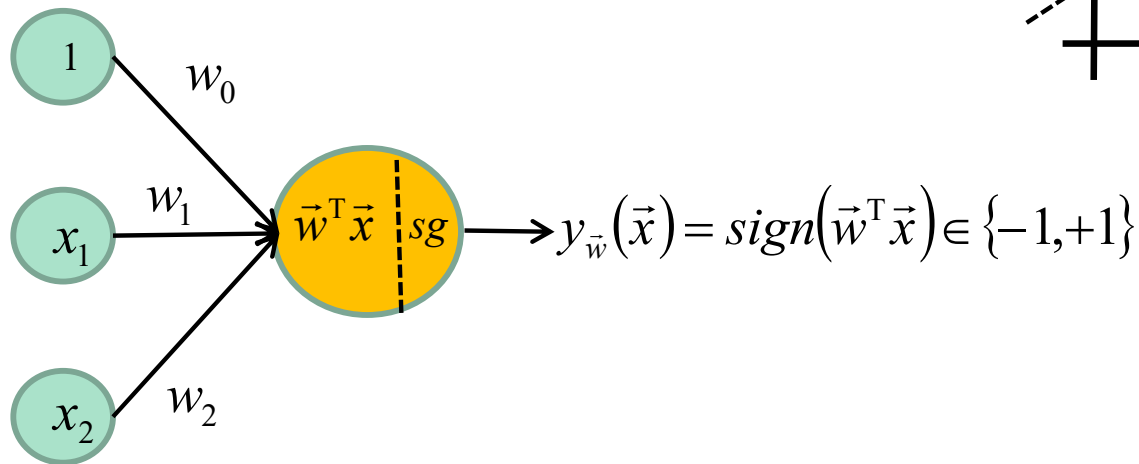
(hyperplane)

# How do we train a Perceptron?

# Loss function



$$L\left(y_{\vec{w}}(\vec{x}), D\right) \approx 0 \qquad L\left(y_{\vec{w}}(\vec{x}), D\right) > 0 \qquad L\left(y_{\vec{w}}(\vec{x}), D\right) >> 0$$

$L\!\left(y_{\vec{w}}(\vec{x}), D\right)$

BEST SOLUTION

$w_1$

$w_2$

# Perceptron

**Question**: how to find the best solution? $\nabla L\left(y_{\vec{w}}(\vec{x}), D\right) = 0$

Random initialization

```
[w1,w2]=np.random.randn(2)
```

# Gradient descent

**Question**: how to find the best solution?   $\nabla L\left(y_{\vec{w}}(\vec{x}), D\right) = 0$

$$\vec{w}^{[k+1]} = \vec{w}^{[k]} - \eta \nabla L\left(y_{\vec{w}^{[k]}}(\vec{x}), D\right)$$

→ Gradient of the loss function

→ *Learning rate*

# Optimisation

$$\vec{\mathrm{w}}^{[k+1]} = \vec{\mathrm{w}}^{[k]} - \eta^{[k]}\nabla L$$

→ Gradient of the loss function

→ learning rate

**Stochastic gradient descent (SGD)**

Init $\vec{\mathrm{w}}$

k=0

DO k=k+1

    FOR n = 1 to N

$$\vec{w} = \vec{w} - \eta^{[k]}\nabla L(\vec{x}_n)$$

UNTIL every data is well classified or k== MAX_ITER

# Perceptron Criterion (loss)

**Observation**

A **wrongly classified sample** is when

$$\vec{w}^{\mathrm{T}}\vec{x}_n > 0 \text{ and } t_n = -1$$

or

$$\vec{w}^{\mathrm{T}}\vec{x}_n < 0 \text{ and } t_n = +1.$$

Consequently $-\vec{w}^{\mathrm{T}}\vec{x}_n t_n$ is **ALWAYS positive for wrongly classified samples**



46

# Perceptron gradient descent

$$L\left(y_{\vec{w}}(\vec{x}), D\right) = \sum_{\vec{x}_n \in V} - \vec{w}^{\mathrm{T}} \vec{x}_n t_n$$

$$\nabla L\left(y_{\vec{w}}(\vec{x}), D\right) = \sum_{\vec{x}_n \in V} - \vec{x}_n t_n$$

**Stochastic gradient descent (SGD)**

```
Init  w⃗
k=0
DO k=k+1
      FOR n = 1 to N
```
$$\qquad \text{IF} \quad \vec{w}^{\mathrm{T}} \vec{x}_n t_n < 0 \quad \text{THEN} \quad /* \text{ wrongly classified } */$$
$$\qquad \vec{w} = \vec{w} + \eta t_n \vec{x}_n$$

UNTIL every data is well classified OR k=k_MAX

NOTE : **learning rate** $\eta$ :

- **Too low** => slow convergence
- **Too large** => might not converge (even diverge)
- Can **decrease** at each iteration  (e.g. $\eta^{[k]} = cst / k$)

47

# So far…

1. Training dataset: $D$

2. Linear classification function: $y_{\vec{w}}(\vec{x}) = w_1 x_1 + w_2 x_2 + \ldots + w_M x_M + w_0$

3. Loss function: $L\left(y_{\vec{w}}(\vec{x}), D\right) = \sum_{\vec{x}_n \in V} -\vec{w}^{\mathrm{T}} \vec{x}_n t_n$

# So far…

1. Training dataset: $D$

2. Linear classification function:

3. Loss function: $L\left(y_{\vec{w}}(\vec{x}), D\right) = \sum_{\vec{x}_n \in V} - \vec{w}^{\mathrm{T}} \vec{x}_n t_n$



$y_{\vec{w}}(\vec{x}) = sign\left(\vec{w}^{\mathrm{T}} \vec{x}\right) \in \{-1, +1\}$

4. Training : find $\vec{w}$ that minimizes $L\left(y_{\vec{w}}(\vec{x}), D\right)$

$$\nabla L\left(y_{\vec{w}}(\vec{x}), D\right) = 0$$

# Multiclass Perceptron

(2D and 3 classes)



$$y_{\vec{w},0}(\vec{x}) = \vec{w}_0^T \vec{x} = w_{0,0} + w_{0,1}x_1 + w_{0,2}x_2$$

$$y_{\vec{w},1}(\vec{x}) = \vec{w}_1^T \vec{x} = w_{1,0} + w_{1,1}x_1 + w_{1,2}x_2$$

$$y_{\vec{w},2}(\vec{x}) = \vec{w}_2^T \vec{x} = w_{2,0} + w_{2,1}x_1 + w_{2,2}x_2$$

51

# Multiclass Perceptron
(2D and 3 classes)

$\vec{w}_0^T \vec{x}$

$\vec{w}_1^T \vec{x}$

$\vec{w}_2^T \vec{x}$

$\arg\max_i y_{W,i}(\vec{x})$

$w_{0,1}$
$w_{1,1}$
$w_{2,1}$
$w_{0,2}$
$w_{1,2}$
$w_{2,2}$
$w_{0,0}$
$w_{1,0}$
$w_{2,0}$

$y_{\vec{w},2}(\vec{x})$ is max

$y_{\vec{w},1}(\vec{x})$ is max

$y_{\vec{w},0}(\vec{x})$ is max

$$y_W(\vec{x}) = W\vec{x}$$

$$y_W(\vec{x}) = \begin{bmatrix} w_{0,0} & w_{0,1} & w_{0,2} \\ w_{1,0} & w_{1,1} & w_{1,2} \\ w_{2,0} & w_{2,1} & w_{2,2} \end{bmatrix} \begin{bmatrix} 1 \\ x_1 \\ x_2 \end{bmatrix}$$

# Multiclass Perceptron
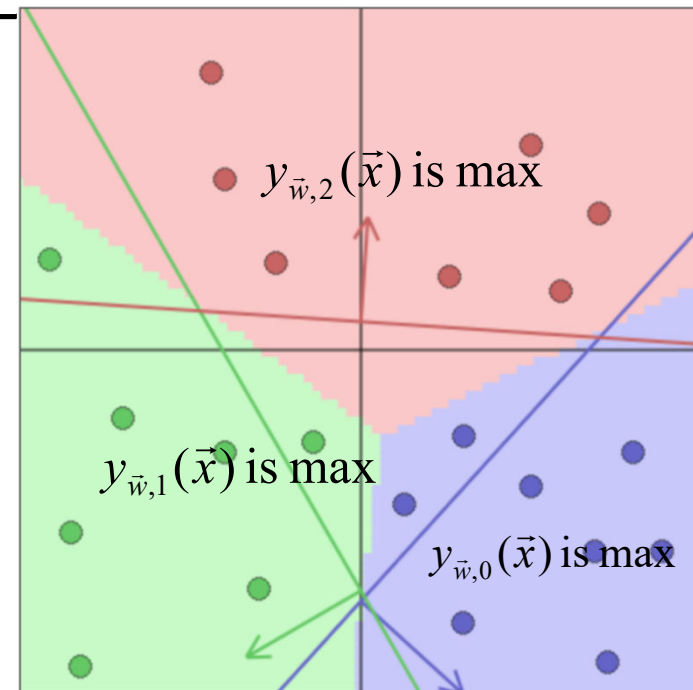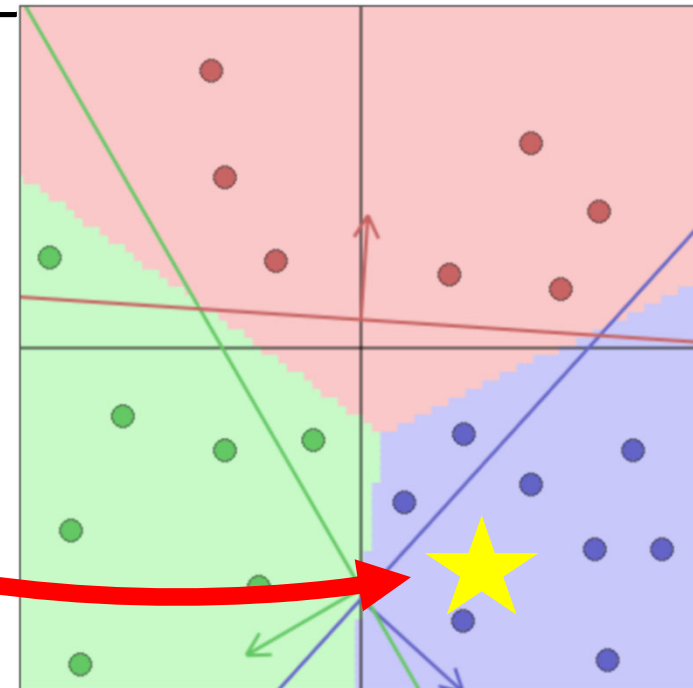
(2D and 3 classes)

Example

$(1.1, -2.0)$

$$y_W(\vec{x}) = \begin{bmatrix} -2 & -3.6 & 0.5 \\ -4 & 2.4 & 4.1 \\ -6 & 4 & -4.9 \end{bmatrix} \begin{bmatrix} 1 \\ 1.1 \\ -2 \end{bmatrix} = \begin{bmatrix} -6.9 \\ -9.6 \\ 8.2 \end{bmatrix} \begin{array}{l} \text{Class 0} \\ \text{Class 1} \\ \text{Class 2} \end{array}$$

# Multiclass Perceptron

Loss function

$$L\big(y_W(\vec{x}), D\big) = \sum_{\vec{x}_n \in V} \big(\vec{w}_j^T \vec{x}_n - \vec{w}_{t_n}^T \vec{x}_n\big)$$

Sum over all wrongly
classified samples

Score of the wrong class

Score of the true class

$$\nabla L\big(y_{\vec{w}}(\vec{x}), D\big) = \sum_{\vec{x}_n \in V} \vec{x}_n$$

# Multiclass Perceptron

Stochastic gradient descent (SGD)

init **W**
k=0, i=0
DO k=k+1
      FOR n = 1 to N
            $j = \arg\max \mathrm{W}^T \vec{x}_n$
           IF $j \neq t_i$ THEN /* wrongly classified sample */

$$\vec{w}_j = \vec{w}_j - \eta \vec{x}_n$$

$$\vec{w}_{t_n} = \vec{w}_{t_n} + \eta \vec{x}_n$$

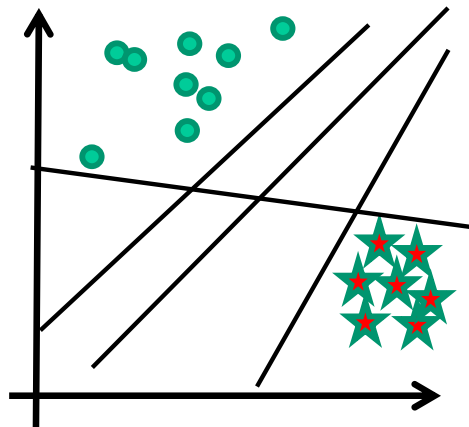UNTIL every data is well classified or k > K_MAX.

# Perceptron

**Advantages**:
- Very simple
- Does **NOT** assume the data follows a **Gaussian distribution**.
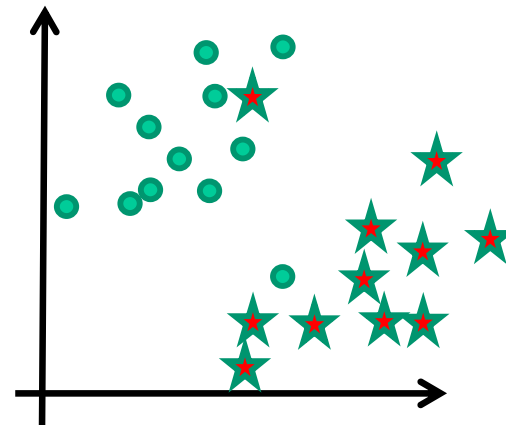- If data is **linearly separable**, convergence is **guaranteed**.

**Limitations**:
- Zero gradient for many solutions => several "perfect solutions"
- Data must be **linearly separable**



Many "optimal" solutions

Will never converge

Two famous ways of improving the Perceptron

1.  New **activation function** + new **Loss**

     → **Logistic regression**
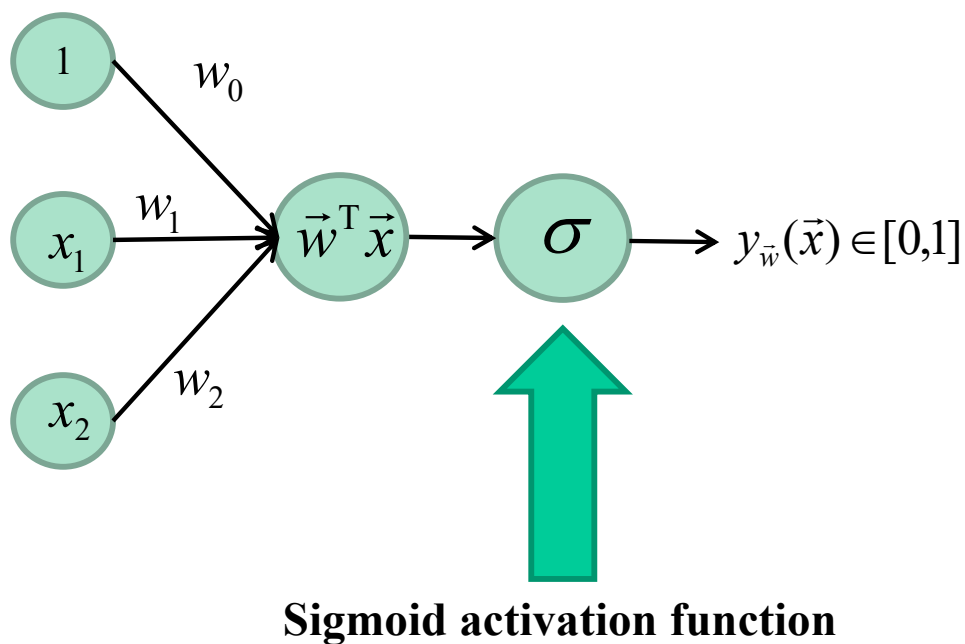
1.  **New network** architecture

     → **Multilayer Perceptron / CNN**

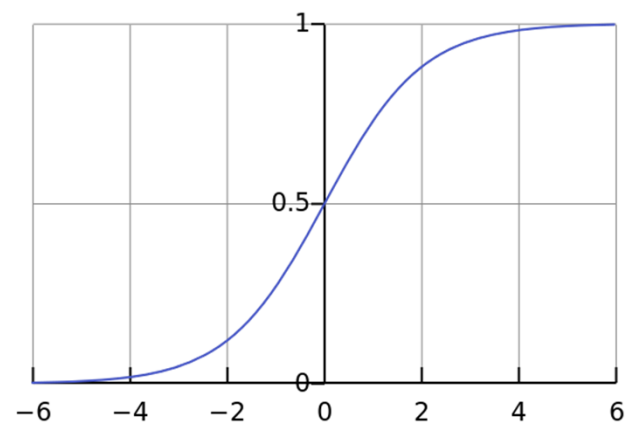# Logistic regression

# Logistic regression

(2D, 2 classes)

New activation function: **sigmoid**



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

$\vec{w}^{\mathrm{T}}\vec{x}$    $\sigma$    $y_{\vec{w}}(\vec{x}) \in [0,1]$

**Sigmoid activation function**
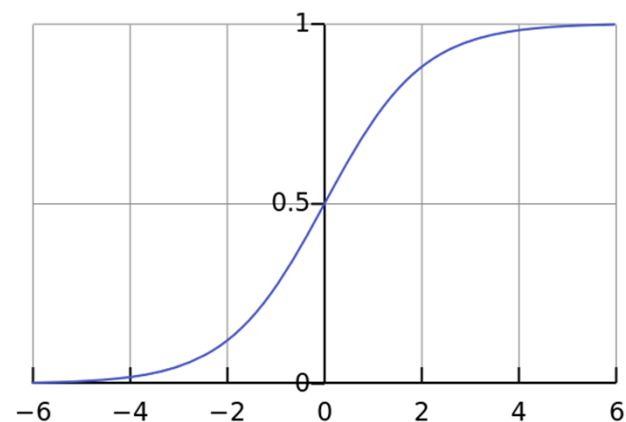
# Logistic regression

(2D, 2 classes)

New activation function: **sigmoid**



**Neuron**
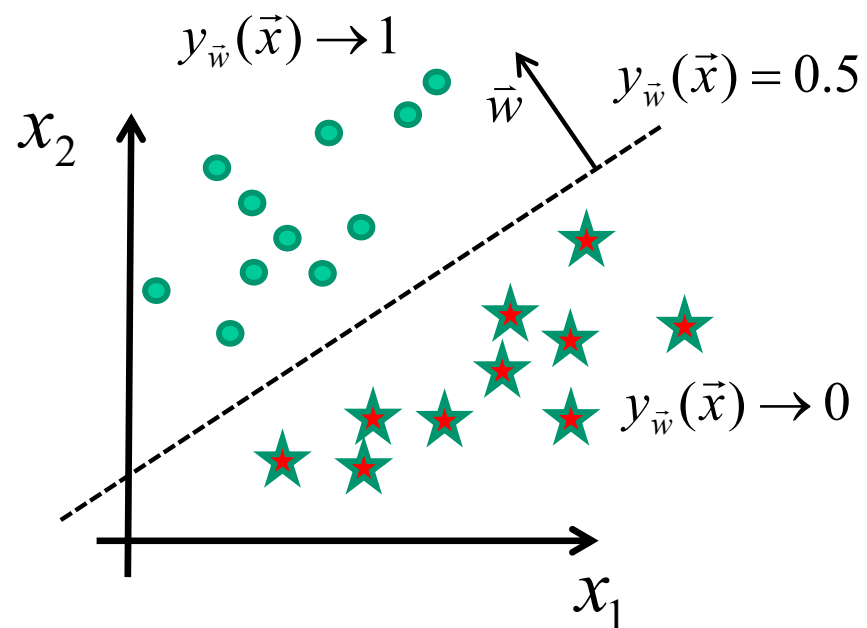
$$\sigma(t) = \frac{1}{1 + e^{-t}}$$

$$y_{\vec{w}}(\vec{x}) = \sigma(\vec{w}^T \vec{x}) = \frac{1}{1 + e^{-\vec{w}^T \vec{x}}}$$

# Logistic regression

(2D, 2 classes)

New activation function: **sigmoid**



$$y_{\vec{w}}(\vec{x}) = \sigma\left(\vec{w}^T \vec{x}\right)$$

# Logistic regression

(2D, 2 classes)

Example

$$\vec{x}_n = (0.4, -1.0), \vec{w} = [2.0, -3.6, 0.5]$$



$$\vec{w}^T \vec{x}_n = 2 - 3.6 * 0.4 - 0.5 = -1.94$$

$$y_{\vec{w}}(\vec{x}) = \sigma(-1.94) = \frac{1}{1 + e^{1.94}} = 0.125$$

Since 0.125 is lower than 0.5, $\vec{x}_n$ is **behind** the plane.

# Logistic regression

(K-D, 2 classes)

Like the Perceptron the logistic regression accomodates for K-D vectors



$$y_{\vec{w}}(\vec{x}) = \sigma\left(\vec{w}^T \vec{x}\right)$$

$$y_{\vec{w}}(\vec{x}) = \sigma\left(\vec{w}^T \vec{x}\right)$$

The neuron computes $\vec{w}^T \vec{x}$ then applies $\sigma$, with inputs $1, x_1, x_2, \ldots, x_K$ weighted by $w_0, w_1, w_2, \ldots, w_K$.

## What is the loss function?

With a sigmoid, we can **simulate a conditional probability** of $c_1$ GIVEN $\vec{x}$



$$y_{\vec{w}}(\vec{x}) = \sigma\left(\vec{w}^T \vec{x}\right) \approx P\left(c_1 \mid \vec{x}\right)$$

With a sigmoid, we can **simulate a conditional probability** of $c_1$ GIVEN $\vec{x}$



$$y_{\vec{w}}(\vec{x}) = \sigma(\vec{w}^T \vec{x}) \approx P(c_1 \mid \vec{x})$$

$$\Rightarrow P(c_0 \mid \vec{w}, \vec{x}) = 1 - y_{\vec{w}}(\vec{x})$$

Cost function is **–ln of the prediction**

$$L(y_{\vec{w}}(\vec{x}), D) = -\sum_{n=1}^{N} t_n \ln\left(y_{\vec{w}}(\vec{x}_n)\right) + (1 - t_n) \ln\left(1 - y_{\vec{w}}(\vec{x}_n)\right)$$

*2 Class Cross entropy*

We can also show that

$$\nabla_{\vec{w}} L(y_{\vec{w}}(\vec{x}), D) = \sum_{n=1}^{N} \left(y_{\vec{w}}(\vec{x}_n) - t_n\right) \vec{x}_n$$
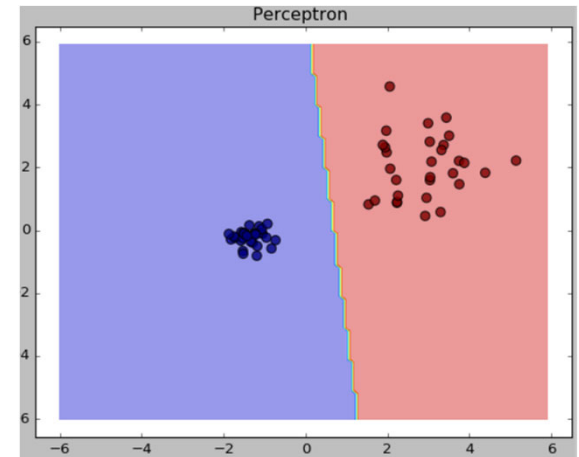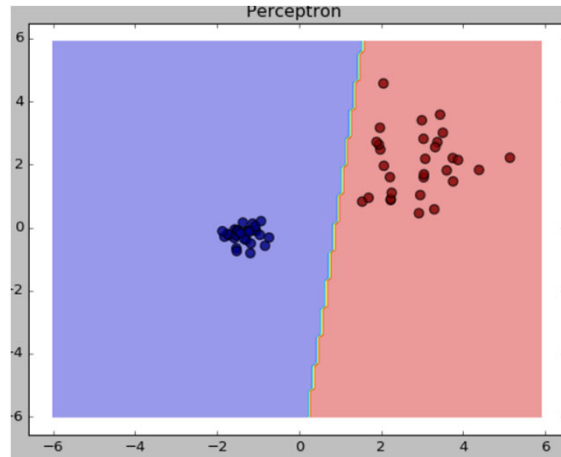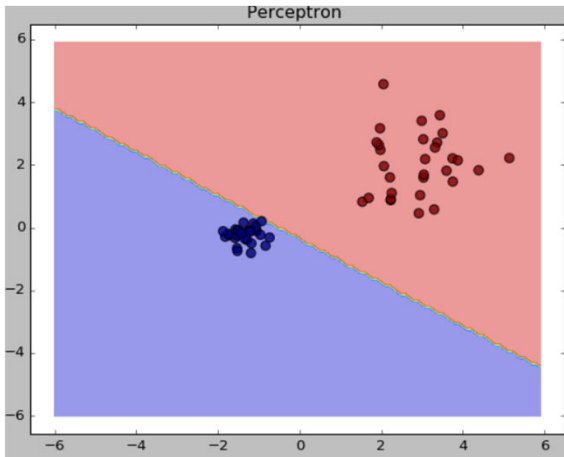
As opposed to the Perceptron the gradient does not depend on the wrongly classified samples

# Logistic Network

Advantages:

- **More stable than the Perceptron**
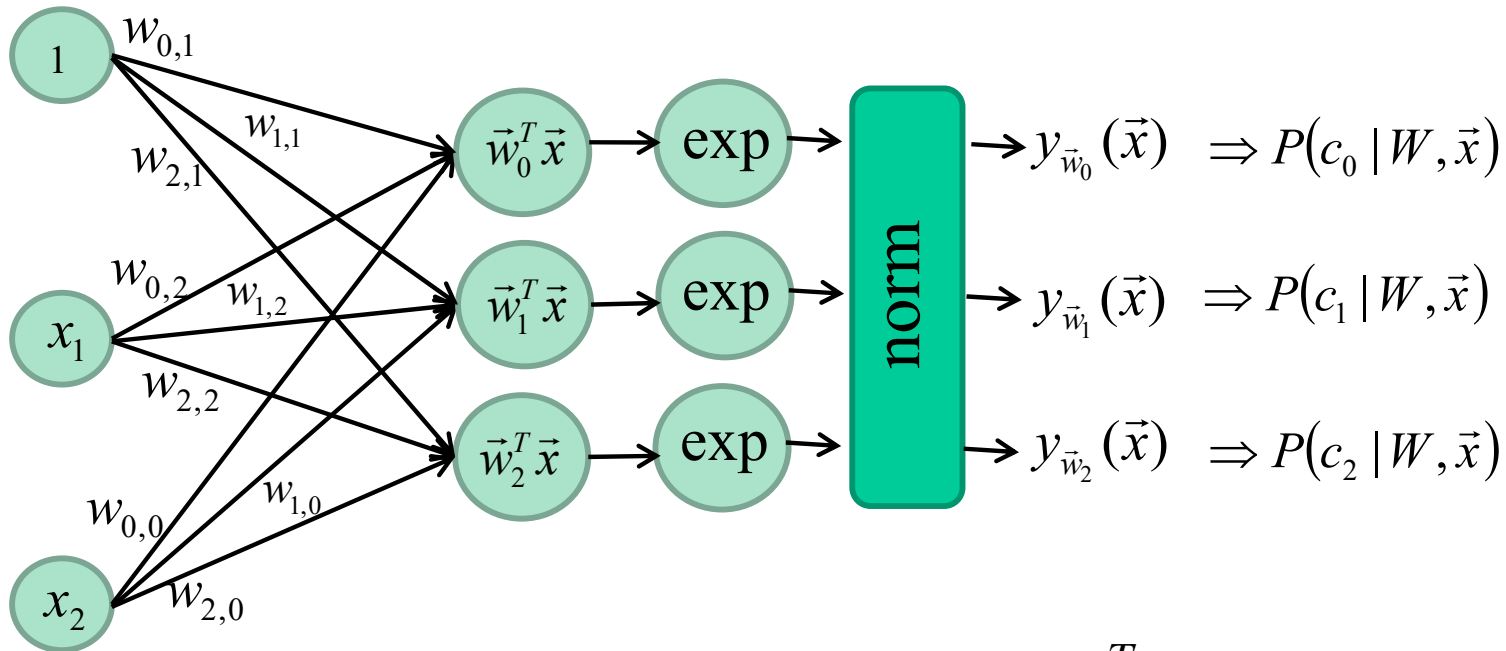- More effective when the data is **non separable**
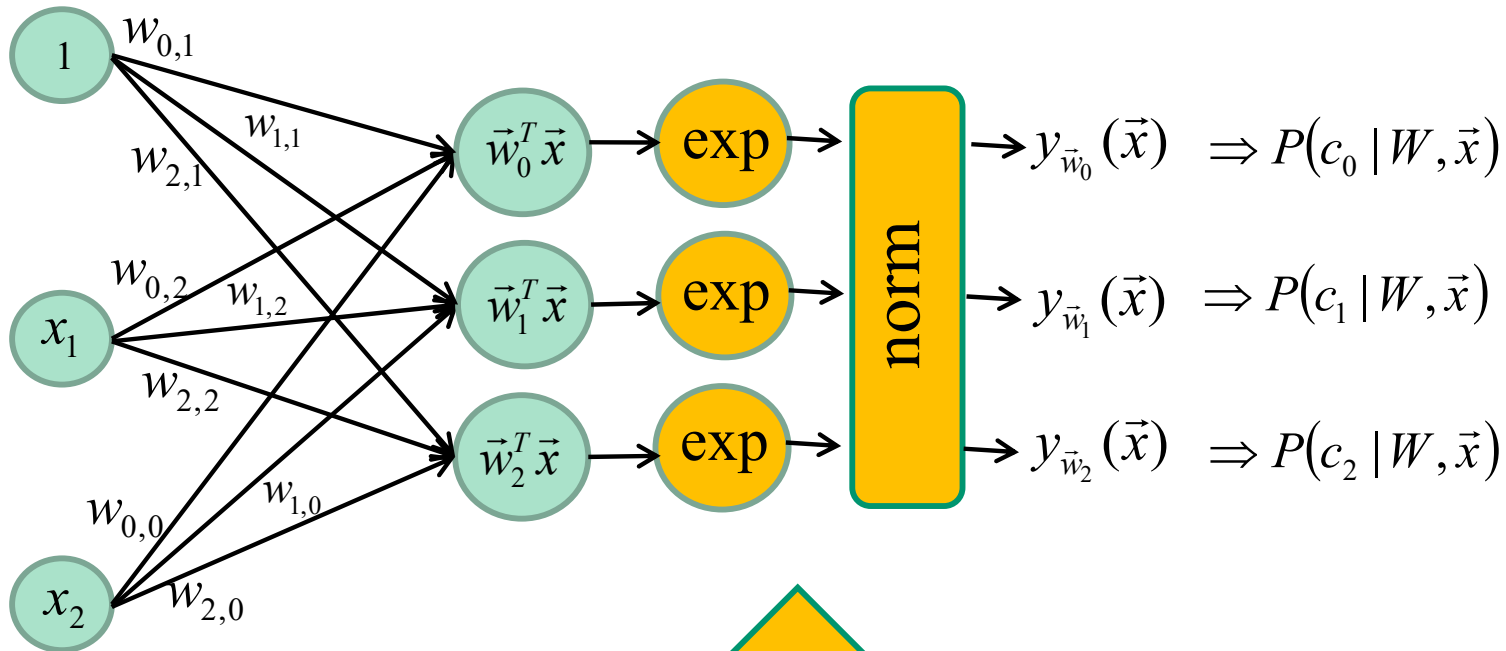


69

# And for K>2 classes?

New activation function : **Softmax**



$$y_{\vec{w}_i}(\vec{x}) = \frac{e^{\vec{w}_i^T \vec{x}}}{\sum_c e^{\vec{w}_c^T \vec{x}}}$$

# And for K>2 classes?

New activation function : **Softmax**



**Softmax**

$$y_{\vec{w}_i}(\vec{x}) = \frac{e^{\vec{w}_i^T \vec{x}}}{\sum_c e^{\vec{w}_c^T \vec{x}}}$$

# And for K>2 classes?



Cifar10

$$\text{'airplane'} \Rightarrow t = \begin{bmatrix} 1000000000 \end{bmatrix}$$

$$\text{'automobile'} \Rightarrow t = \begin{bmatrix} 0100000000 \end{bmatrix}$$

$$\text{'bird'} \Rightarrow t = \begin{bmatrix} 0010000000 \end{bmatrix}$$

$$\text{'cat'} \Rightarrow t = \begin{bmatrix} 0001000000 \end{bmatrix}$$

$$\text{'deer'} \Rightarrow t = \begin{bmatrix} 0000100000 \end{bmatrix}$$

$$\text{'dog'} \Rightarrow t = \begin{bmatrix} 0000010000 \end{bmatrix}$$

$$\text{'frog'} \Rightarrow t = \begin{bmatrix} 0000001000 \end{bmatrix}$$

$$\text{'horse'} \Rightarrow t = \begin{bmatrix} 0000000100 \end{bmatrix}$$

$$\text{'ship'} \Rightarrow t = \begin{bmatrix} 0000000010 \end{bmatrix}$$

$$\text{'truck'} \Rightarrow t = \begin{bmatrix} 0000000001 \end{bmatrix}$$

Class labels : **__one-hot vectors__**

# K>2 classes

*Cross entropy Loss*

$$L\left(y_W(\vec{x}), D\right) = -\sum_{n=1}^{N}\sum_{k=1}^{K} t_{kn} \ln y_{W_k}(\vec{x}_n)$$

# Regularization

Different weights may give the same score

$$\vec{x} = (1.0, 1.0, 1.0)$$

$$\vec{w}_1^T = [1, 0, 0]$$

$$\vec{w}_2^T = [1/3, 1/3, 1/3]$$

Which weights are the best?

Solution: Maximum a posteriori

$$\vec{w}_1^T \vec{x} = \vec{w}_2^T \vec{x} = 1$$

# Maximum *a posteriori*

Regularization

Constant

$$\arg\min_W = L(y_{\vec{w}}(\vec{x}), D) + \lambda R(W)$$

**Loss function**

**Regularization**

In general L1 or L2 $\quad R(\theta) = \|W\|_1 \text{ ou } \|W\|_2$

# Wow! Loooots of information!

# Lets recap…

# Neural networks

2 classes



Sign activation

Perceptron

Sigmoid activation

Logistic regression

# K-Class Neural networks



Perceptron

Logistic regression

Softmax activation

# Loss functions

2 classes

$$L(y_{\vec{w}}(\vec{x}), D) = \sum_{\vec{x}_n \in V} -t_n \vec{w}^\mathrm{T} \vec{x}_n \qquad \text{where V is the set of wrongly classified samples}$$

$$L(y_{\vec{w}}(\vec{x}), D) = -\sum_{n=1}^{N} t_n \ln(y_{\vec{w}}(\vec{x}_n)) + (1 - t_n)\ln(1 - y_{\vec{w}}(\vec{x}_n)) \qquad \text{Cross entropy loss}$$

# Loss functions

K classes

$$L(y_{\vec{w}}(\vec{x}), D) = \sum_{\vec{x}_n \in V} \left( \vec{w}_j^T \vec{x}_n - \vec{w}_{t_n}^T \vec{x}_n \right) \quad \text{where V } is \text{ the set of wrongly classified samples}$$

$$L(y_{\vec{w}}(\vec{x}), D) = -\sum_{n=1}^{N} \sum_{k=1}^{K} t_{kn} \ln y_{W_k}(\vec{x}_n) \qquad Cross\ entropy \text{ loss with a Softmax}$$

$$L\left(y_{\vec{w}}(\vec{x}), D\right) = \sum_{n=1}^{N} l\left(y_W(\vec{x}_n), t_n\right) + \lambda R(W)$$

**Constant**

**Loss function**

**Regularization**

$$R(W) = \|W\|_1 \ \ or \ \ \|W\|_2$$

81

# Now, lets go
# **DEEPER**

# Non-linearly separable training data

Three classical solutions

1. Acquire more data
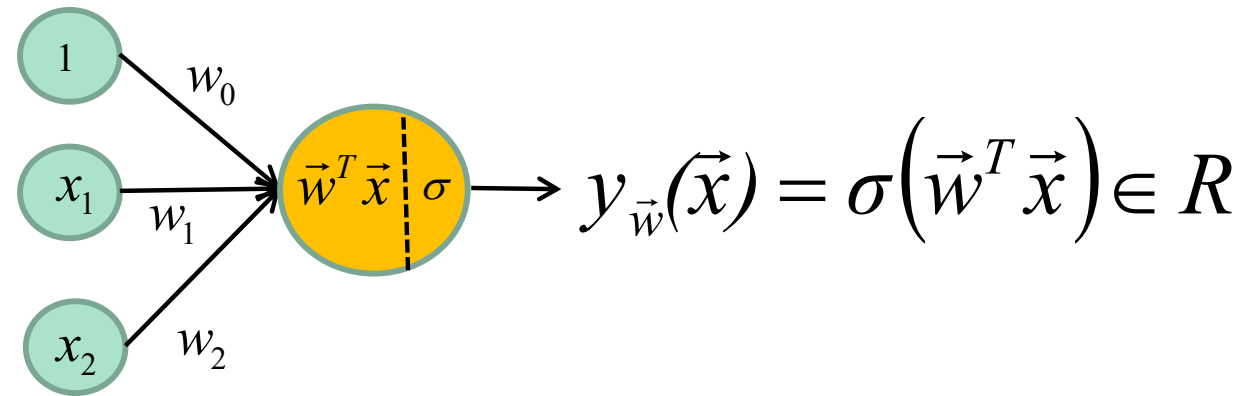2. Use a non-linear classifier
3. Transform the data

# Non-linearly separable training data

Three classical solutions

1. Acquire more data
2. Use a non-linear classifier
3. **Transform the data**

# 2D, 2Classes, Linear logistic regression
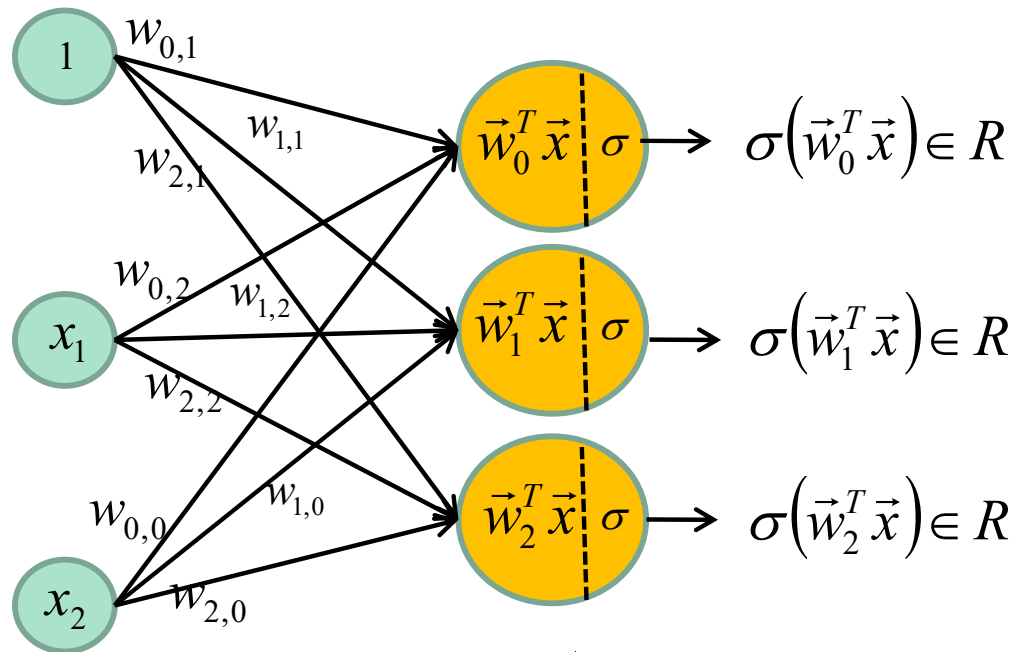


$$y_{\vec{w}}(\vec{x}) = \sigma\left(\vec{w}^T \vec{x}\right) \in R$$
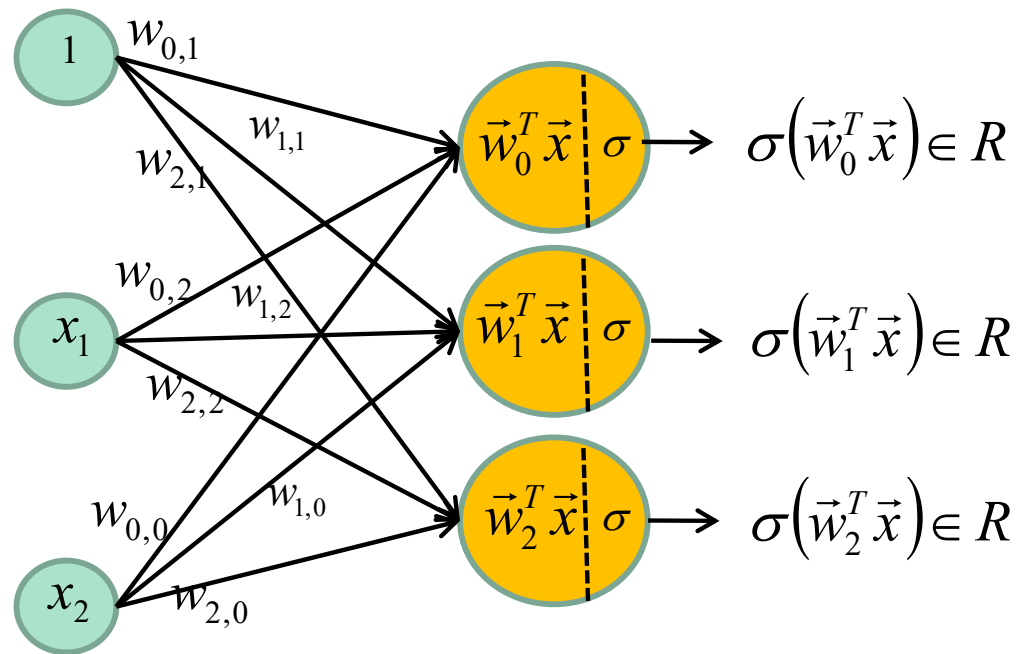
**Input layer
(3 "neurons")**

**Output layer
(1 neuron with sigmoid)**
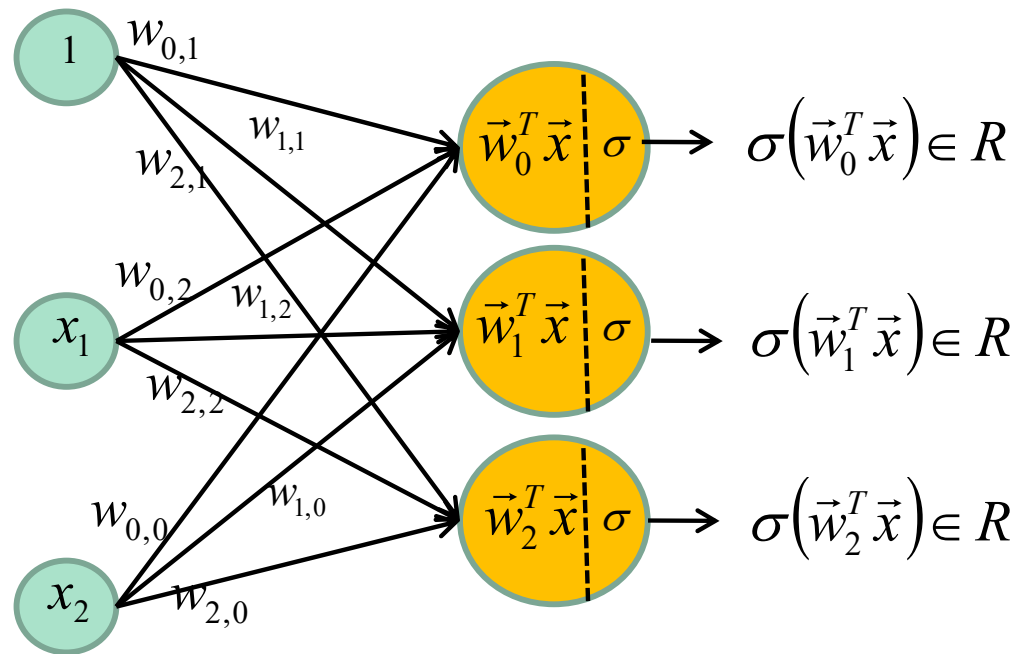
# Let's add 3 neurons



**Input layer (3 "neurons")**

**First layer (3 neurons)**

**NOTE:** The output of the first layer is a vector of **3 real** values

$$\sigma\left(\begin{bmatrix} w_{0,0} & w_{0,1} & w_{0,2} \\ w_{1,0} & w_{1,1} & w_{1,2} \\ w_{2,0} & w_{2,1} & w_{2,2} \end{bmatrix}\begin{bmatrix} 1 \\ x_1 \\ x_2 \end{bmatrix}\right) \in R^3$$
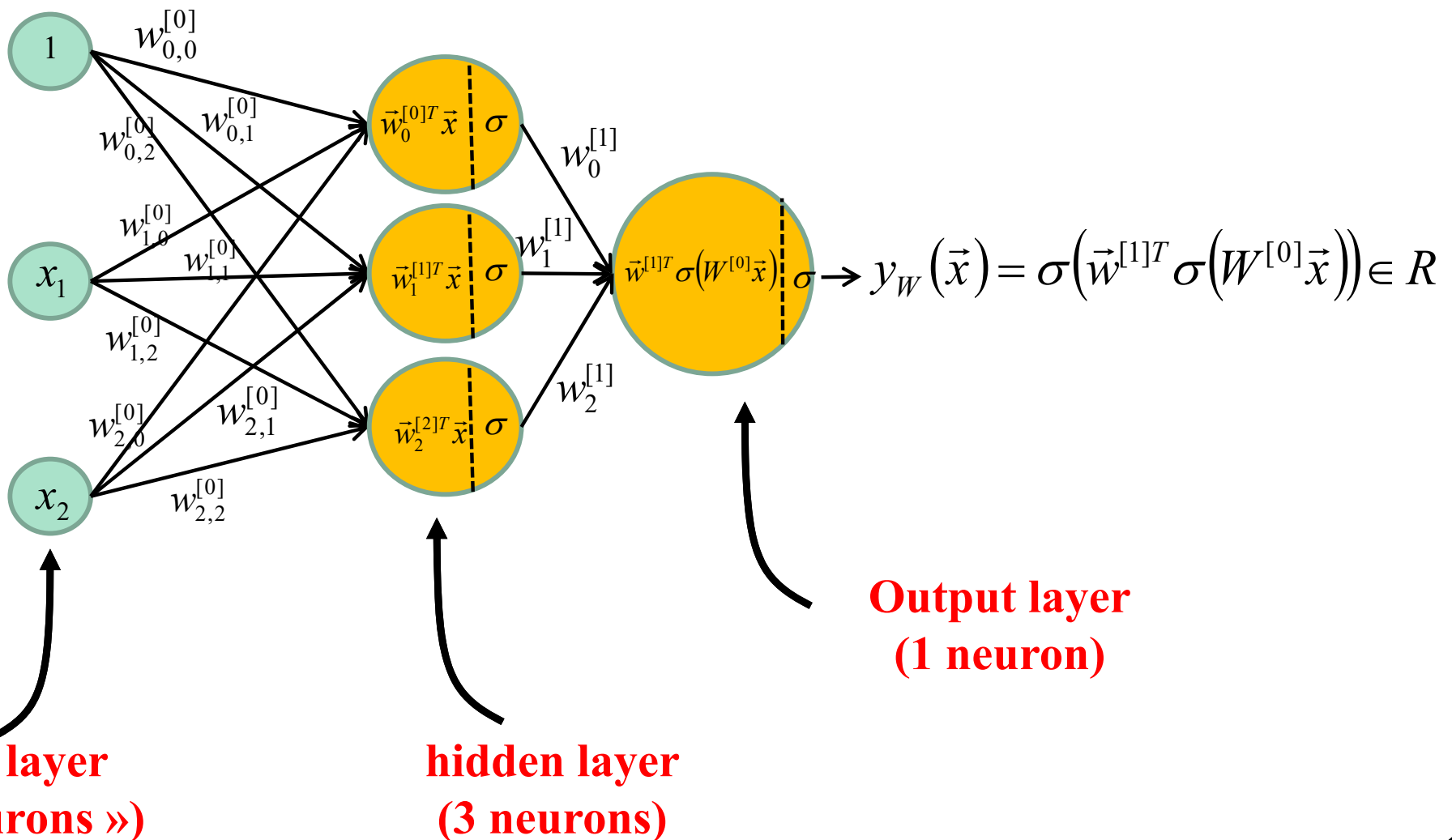
**NOTE:** The output of the first layer is a vector of **3 real** values
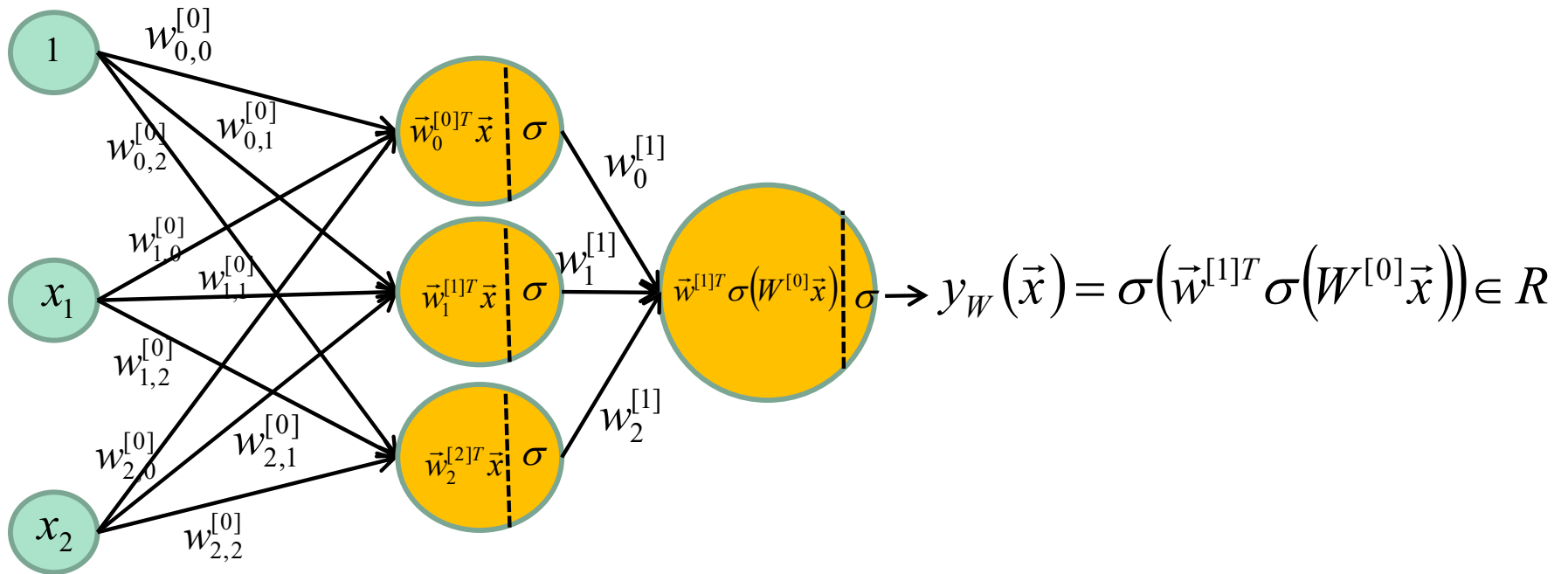
$$\sigma\left(W^{[0]}\vec{x}\right)$$

# 2-D, 2-Class, 1 hidden layer

If we want a **2-class Classification** via a **logistic regression** (a **cross entropy loss**) we must add an **output neuron.**



$$y_W(\vec{x}) = \sigma\left(\vec{w}^{[1]T}\sigma\left(W^{[0]}\vec{x}\right)\right) \in R$$

**Output layer
(1 neuron)**

**Input layer
(3 « neurons »)**

**hidden layer
(3 neurons)**

# 2-D, 2-Class, 1 hidden layer



$$y_W(\vec{x}) = \sigma\left(\vec{w}^{[1]T}\sigma\left(W^{[0]}\vec{x}\right)\right) \in R$$

Visual simplification

$$y_W(\vec{x}) = \sigma\left(\vec{w}^{[1]T}\sigma\left(W^{[0]}\vec{x}\right)\right)$$

# 2-D, 2-Class, 1 hidden layer



Input layer  Hidden layer  Output layer

$W^{[0]}$  $\vec{w}^{[1]}$

$$y_W(\vec{x}) = \sigma\left(\vec{w}^{[1]T}\sigma\left(W^{[0]}\vec{x}\right)\right)$$

**bias** neuron.

This network contains a total of **13 parameters**

3x3                     1x4

Input layer        Hidden layer
Hidden layer      Output layer

# 2-D, 2-Class, 1 hidden layer

Input layer

Hidden layer

Output layer



$$y_W(\vec{x}) = \sigma\left(\vec{w}^{[1]T} \sigma\left(W^{[0]}\vec{x}\right)\right)$$
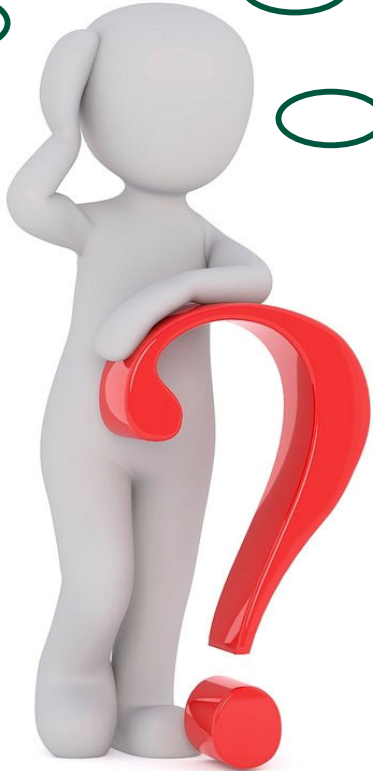
Increasing the number of neurons = increasing the **capacity of the model**

This network has 5x3+1x6=**21 parameters**

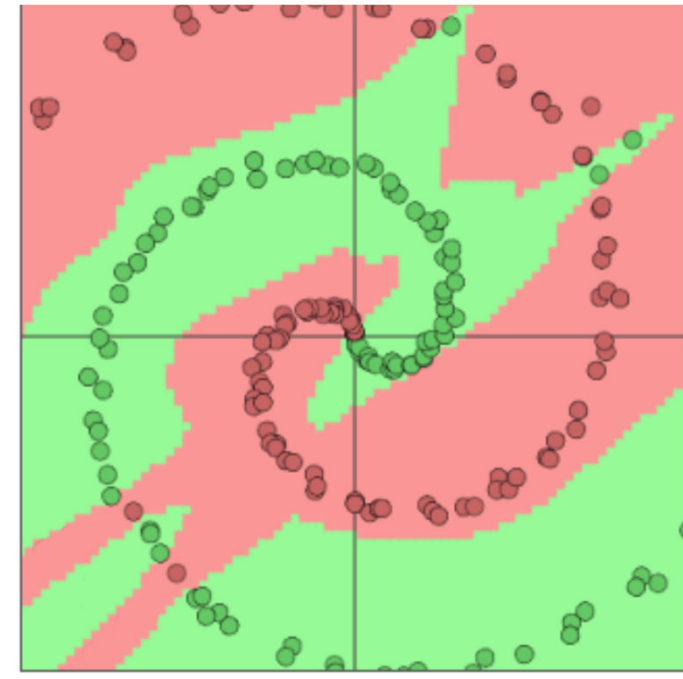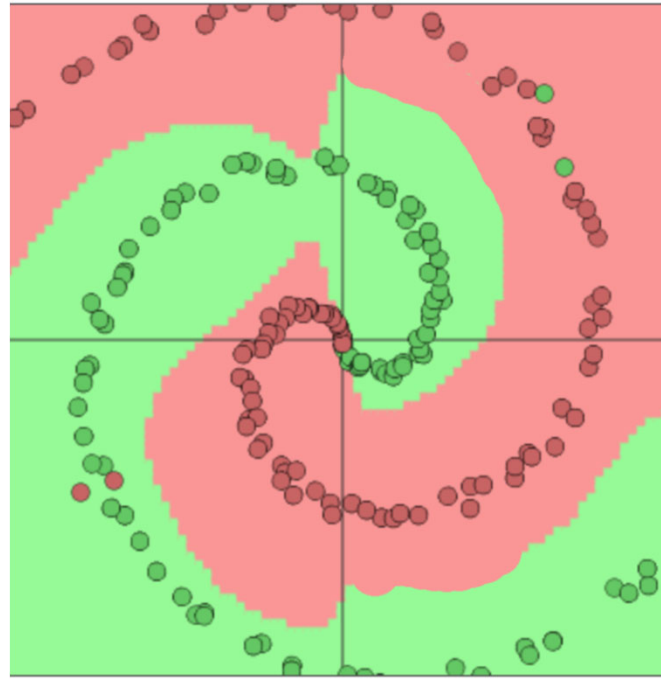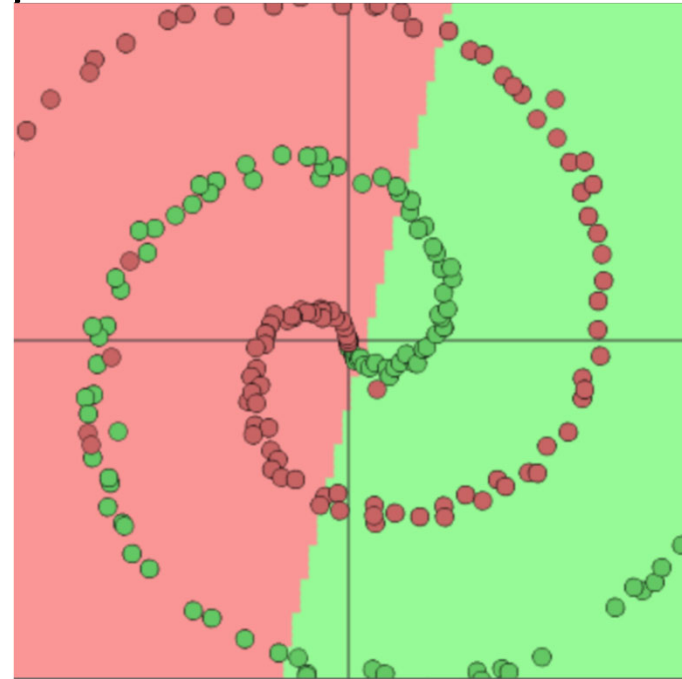$\vec{w}^{\mathrm{T}}\vec{x}$

Line

93

# Nb neurons VS Capacity

No hidden neuron

12 hidden neurons

60 hidden neurons



Linear classification
**Underfitting**
(low capacity)

Non linear classification
**Good result**
(good capacity)

Non linear classification
**Over fitting**
(too large capacity)

http://cs.stanford.edu/people/karpathy/convnetjs/demo/classify2d.html

# kD, 2Classes, 1 hidden layer



Input layer | Hidden layer | Output layer

$$y_W(\vec{x}) = \sigma\left(\vec{w}^{[1]T}\sigma\left(W^{[0]}\vec{x}\right)\right)$$

Increasing the dimensionality of the data = **more columns in** $W^{[0]}$

This network has 5x(k+1)+1x6 **parameters**

# kD, 2Classes, <u>**2 hidden layers**</u>

Input layer    Hidden Layer 1    Hidden Layer 2    Output layer

$$W^{[0]} \in R^{5 \times k+1}$$

$$W^{[1]} \in R^{3 \times 6}$$

$$\vec{w}^{[2]} \in R^4$$

$$W^{[0]} \qquad W^{[1]} \qquad \vec{w}^{[2]T}$$

$$y_W(\vec{x}) = \sigma\left(\vec{w}^{[2]T} \sigma\left(W^{[1]} \sigma\left(W^{[0]} \vec{x}\right)\right)\right)$$

Adding an hidden layer = Adding a matrix multiplication

This network has 5x(k+1)+6x3 + 1x4 **parameters**

# kD, 2 Classes, 4 hidden layer network

Input layer   Hidden Layer 1   Hidden Layer 2   Hidden Layer 3   Hidden Layer 4   Output layer

$W^{[0]}$

$W^{[1]}$

$W^{[2]}$

$W^{[3]}$

$\vec{w}^{[4]}$

$x_1$

$x_2$

$x_3$

$(...)$

$x_k$

$1$

$$W^{[0]} \in R^{5 \times k + 1}$$

$$W^{[1]} \in R^{3 \times 6}$$

$$W^{[2]} \in R^{4 \times 4}$$

$$W^{[3]} \in R^{7 \times 5}$$

$$\vec{w}^{[4]} \in R^{8}$$

$$y_W(\vec{x}) = \sigma\left(\vec{w}^{[4]T} \sigma\left(W^{[3]} \sigma\left(W^{[2]} \sigma\left(W^{[1]} \sigma\left(W^{[0]} \vec{x}\right)\right)\right)\right)\right)$$

This network has 5x(k+1)+6x3 + 4x4+7x5+1x8 **parameters**

97

# kD, 2 Classes, 4 hidden layer network

Input layer    Hidden Layer 1    Hidden Layer 2    Hidden Layer 3    Hidden Layer 4    Output layer

$W^{[0]}$

$W^{[1]}$

$W^{[2]}$

$W^{[3]}$

$\vec{w}^{[4]}$

$x_1$

$x_2$

$x_3$

$(...)$

$x_k$

$1$

$W^{[0]} \in R^{5 \times k+1}$

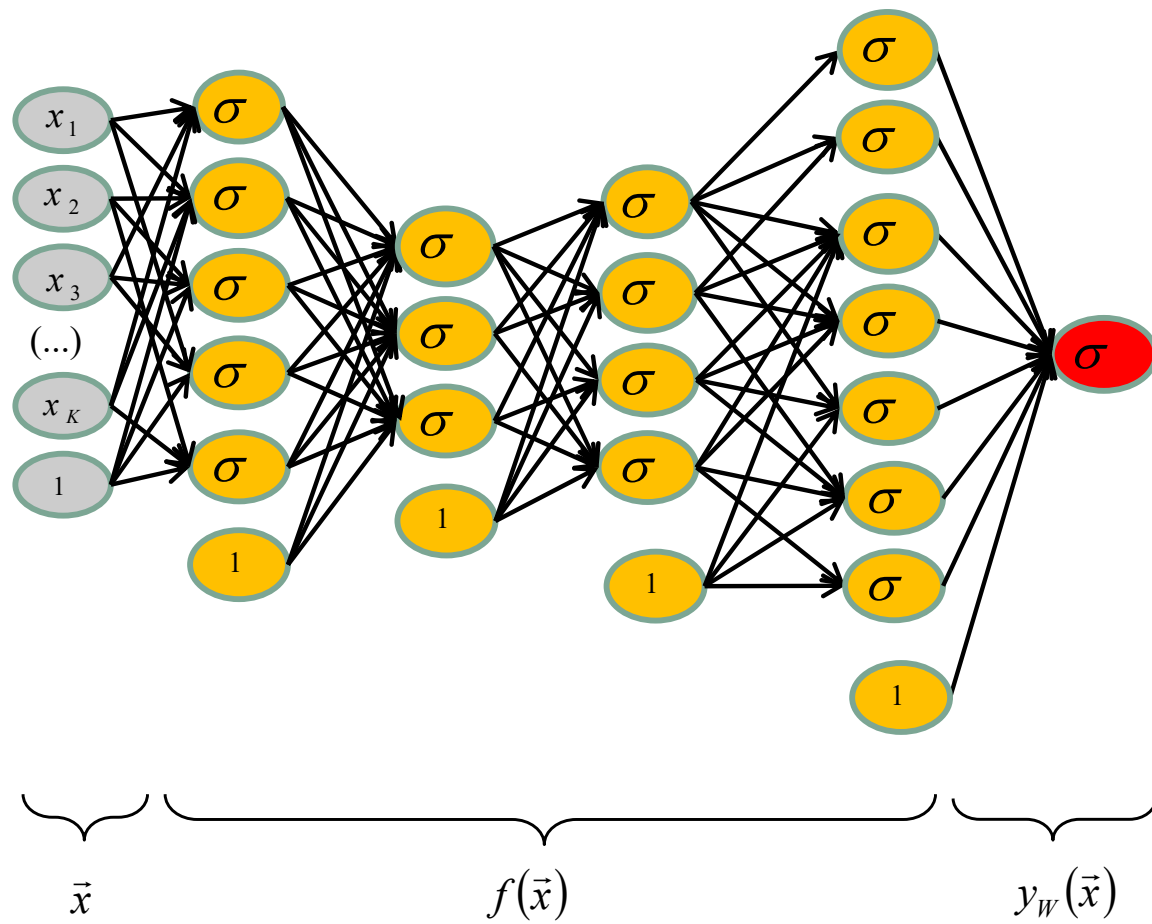$W^{[1]} \in R^{3 \times 6}$

$W^{[2]} \in R^{4 \times 4}$
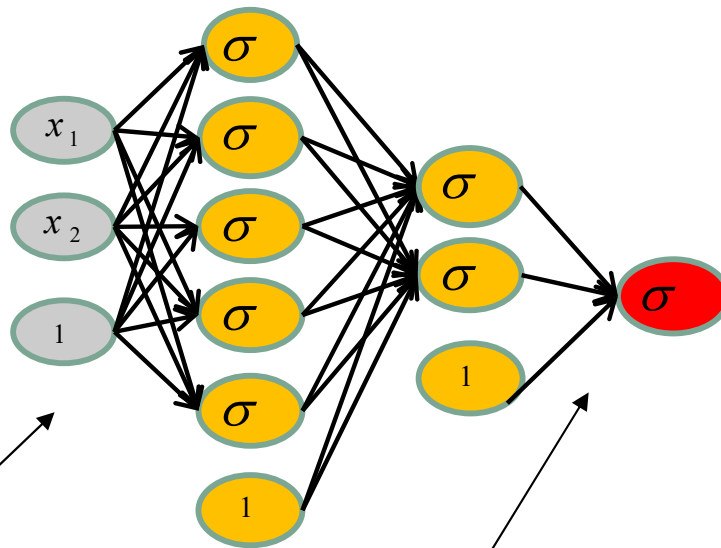
$W^{[3]} \in R^{7 \times 5}$

$\vec{w}^{[4]} \in R^{8}$

$$y_W(\vec{x}) = \sigma\left(\vec{w}^{[4]T}\sigma\left(W^{[3]}\sigma\left(W^{[2]}\sigma\left(W^{[1]}\sigma\left(W^{[0]}\vec{x}\right)\right)\right)\right)\right)$$

NOTE : More hidden layers = **deeper** network = **more capacity**.

# Multilayer Perceptron

# Example



Input data

Output of the last layer

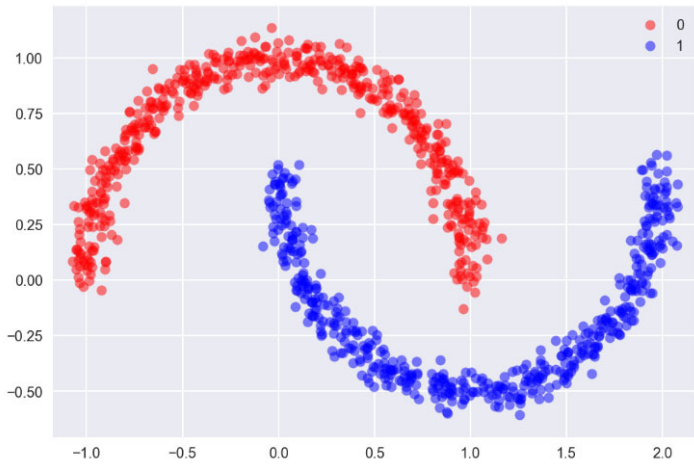Output of the network
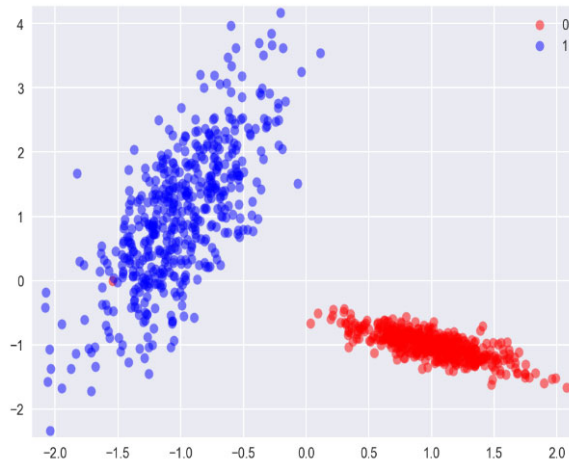
$\vec{x}$

$y_W(\vec{x})$

# Example



$\vec{x}$
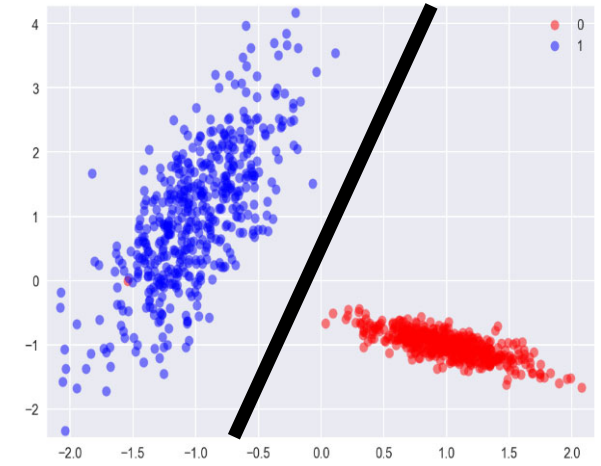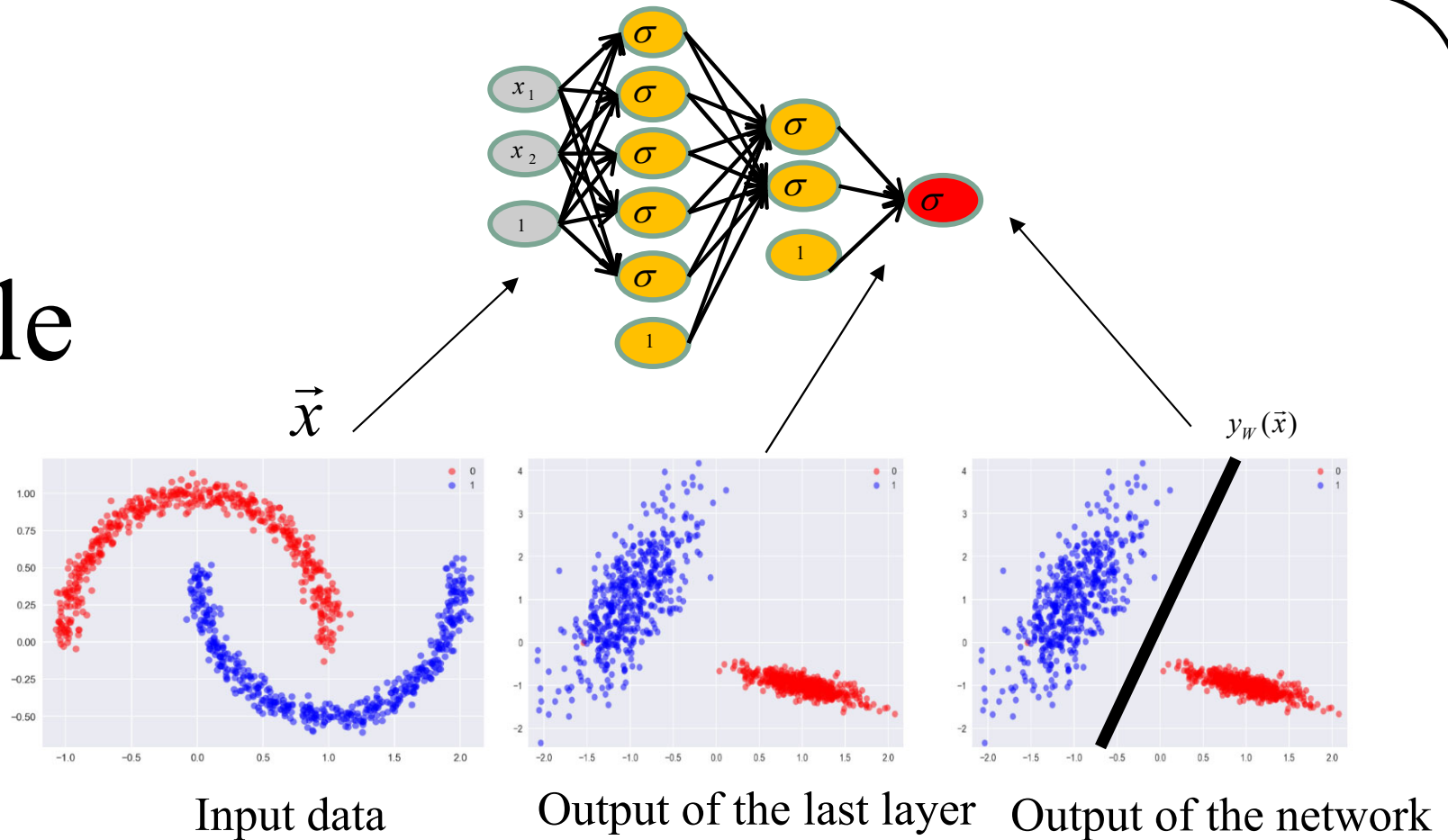
$y_W(\vec{x})$

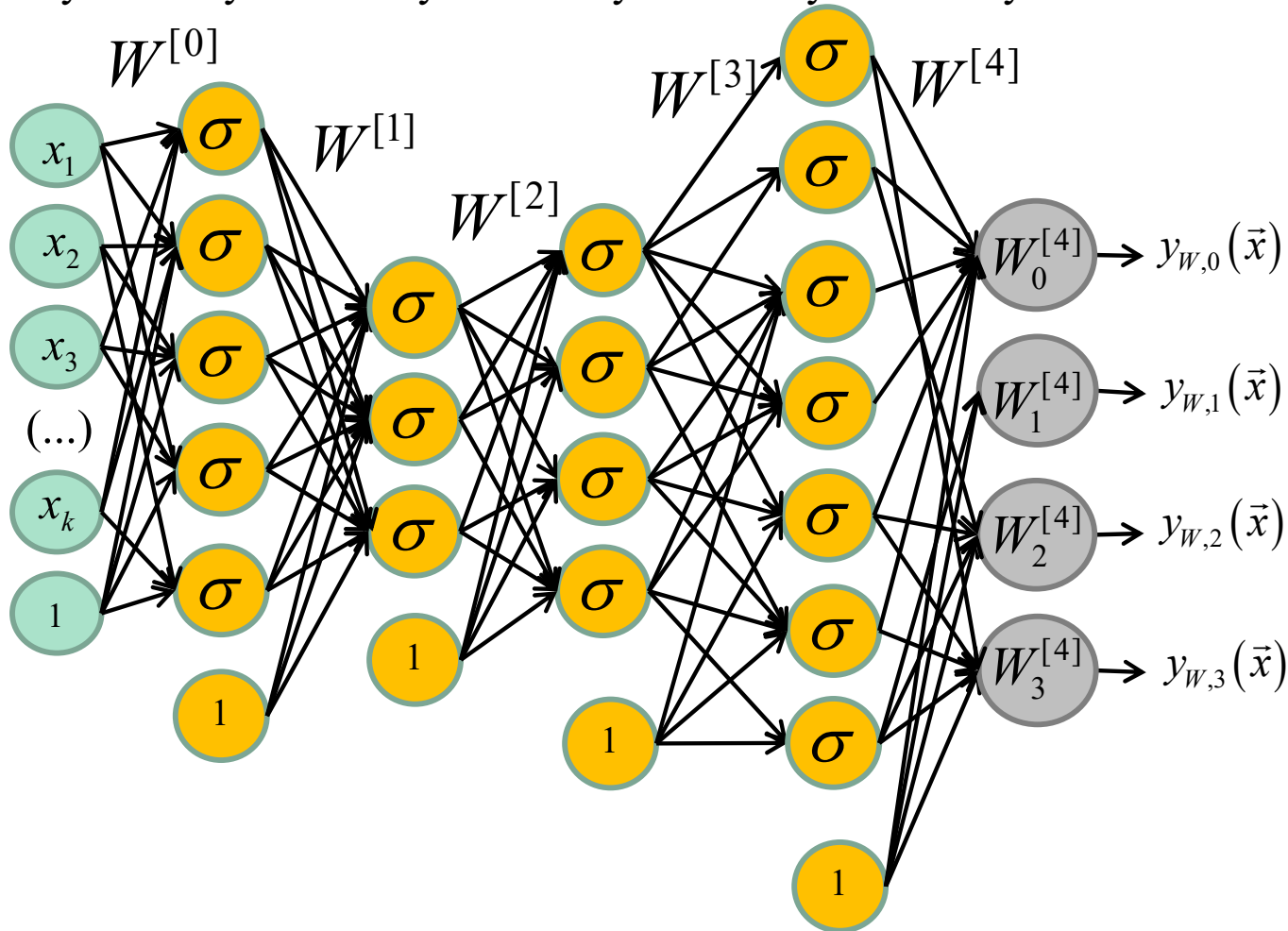Input data          Output of the last layer          Output of the network

A classification neural network is a **linear classifier** with a bunch of neurons that act as a **basis function**.

A **K-Class** neural network has **K output** neurons.

# kD, **4 Classes**, 4 hidden layer network

Input layer

Hidden Layer 1

Hidden Layer 2

Hidden Layer 3

Hidden Layer 4

Output layer

$W^{[0]}$

$W^{[1]}$

$W^{[2]}$

$W^{[3]}$

$W^{[4]}$

$x_1$

$x_2$

$x_3$

$(...)$

$x_k$

$1$

$\sigma$ ... (hidden units)

$1$

$W_0^{[4]} \rightarrow y_{W,0}(\vec{x})$

$W_1^{[4]} \rightarrow y_{W,1}(\vec{x})$

$W_2^{[4]} \rightarrow y_{W,2}(\vec{x})$

$W_3^{[4]} \rightarrow y_{W,3}(\vec{x})$

$$W^{[0]} \in R^{5 \times k+1}$$
$$W^{[1]} \in R^{3 \times 6}$$
$$W^{[2]} \in R^{4 \times 4}$$
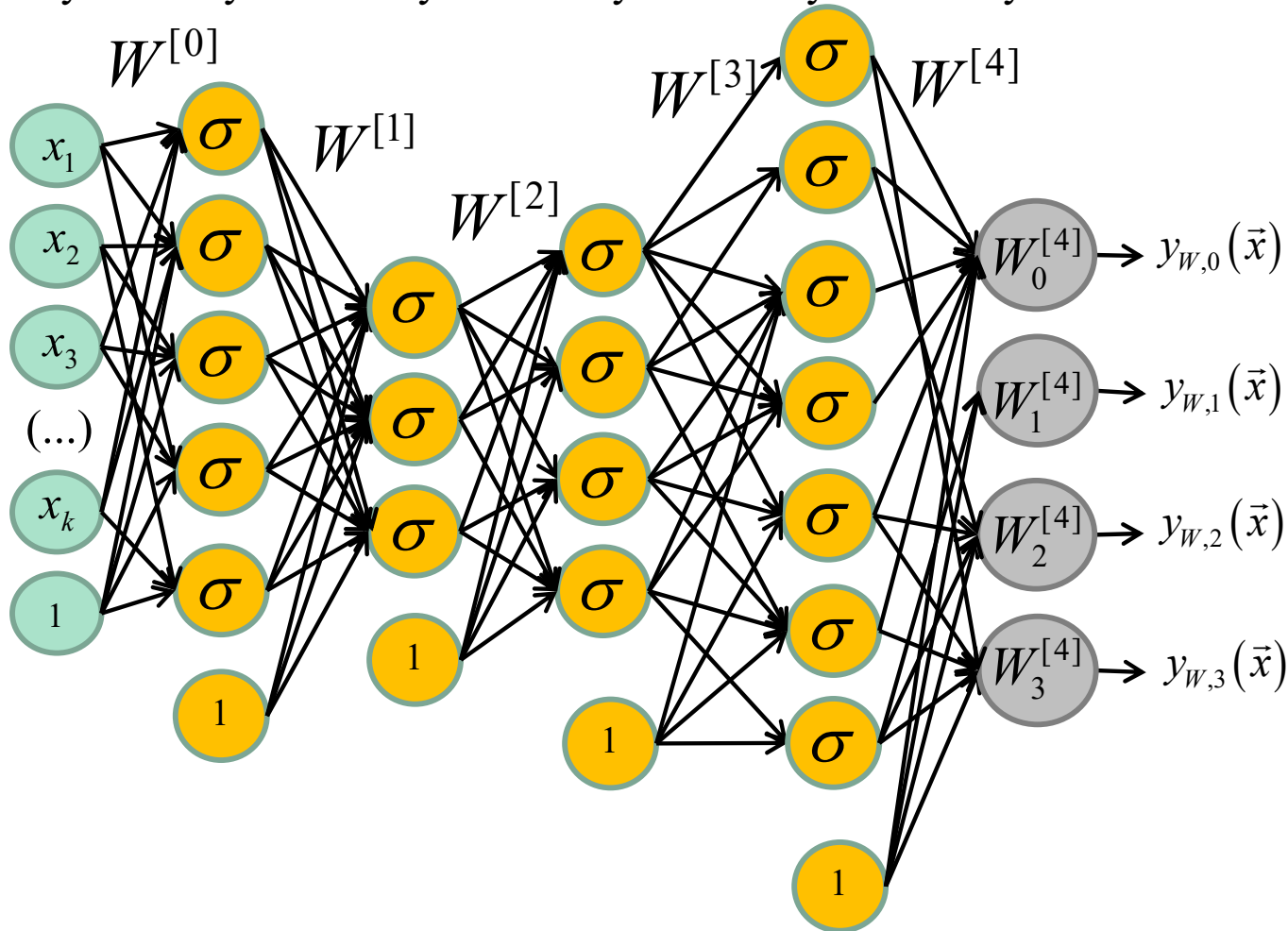$$W^{[3]} \in R^{7 \times 5}$$
$$W^{[4]} \in R^{8 \times 4}$$

$$y_W(\vec{x}) = W^{[4]} \sigma \left( W^{[3]} \sigma \left( W^{[2]} \sigma \left( W^{[1]} \sigma \left( W^{[0]} \vec{x} \right) \right) \right) \right)$$

# kD, **4 Classes**, 4 hidden layer network

Input layer    Hidden Layer 1    Hidden Layer 2    Hidden Layer 3    Hidden Layer 4    Output layer

*Perceptron loss*



$$y_W\left(\vec{x}\right) = W^{[4]}\sigma\left(W^{[3]}\sigma\left(W^{[2]}\sigma\left(W^{[1]}\sigma\left(W^{[0]}\vec{x}\right)\right)\right)\right)$$

# kD, **4 Classes**, 4 hidden layer network

Input layer  Hidden Layer 1  Hidden Layer 2  Hidden Layer 3  Hidden Layer 4  Output layer
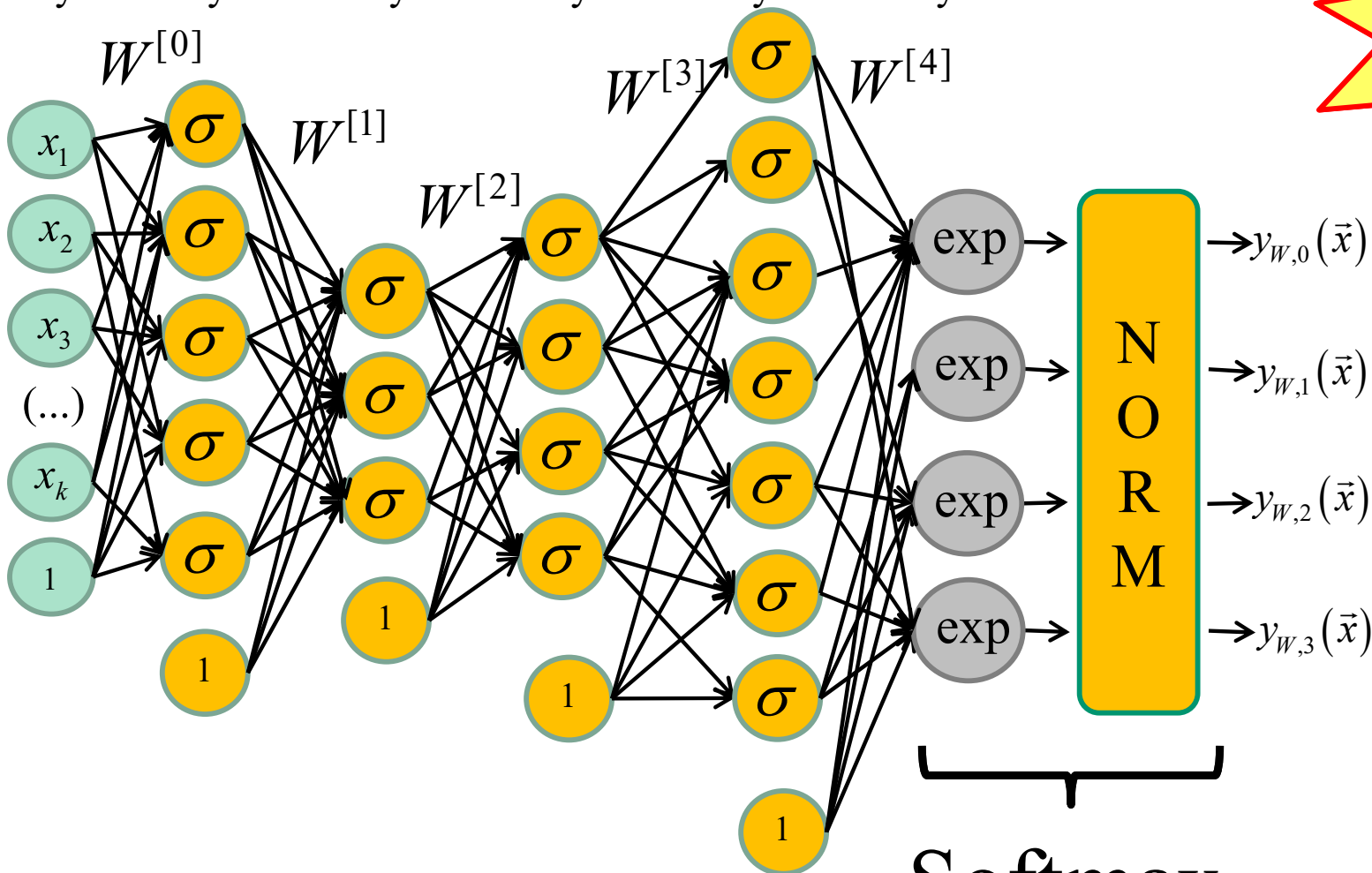
**Cross entropy**

$W^{[0]}$

$W^{[1]}$

$W^{[2]}$

$W^{[3]}$

$W^{[4]}$

$x_1$
$x_2$
$x_3$
$(...)$
$x_k$
$1$

$\sigma$

exp

N O R M

$\rightarrow y_{W,0}(\vec{x})$
$\rightarrow y_{W,1}(\vec{x})$
$\rightarrow y_{W,2}(\vec{x})$
$\rightarrow y_{W,3}(\vec{x})$

## Softmax

$$y_W(\vec{x}) = \text{softmax}\left(W^{[4]}\sigma\left(W^{[3]}\sigma\left(W^{[2]}\sigma\left(W^{[1]}\sigma\left(W^{[0]}\vec{x}\right)\right)\right)\right)\right)$$

# In conclusion

- Linear classifiers
  - Perceptron
  - Logistic regression
- 2-Class vs K-Class neural nets
- Loss function
  - Hinge Loss
  - Cross-entropy loss
- Gradient descent
- Multi-layer perceptron.

# Evaluation metrics

# How to evaluate a model?

**True negative (10)**

Hypertension

Healthy

# False positive (5)

Hypertension

Healthy

# False negative (3)

# Confusion matrix

Hypertension



Healthy

Ground truth

|  | Positive | Negative |
|---|---|---|
| **Positive** | TP = 11 | FP=5 |
| **Negative** | FN = 3 | TN=10 |

Model prediction

TP  +  FN = 14 = TOTAL # positive
TN  + FP = 15 = TOTAL # negative

TP + FP = 16 = TOTAL # of patients classified +1
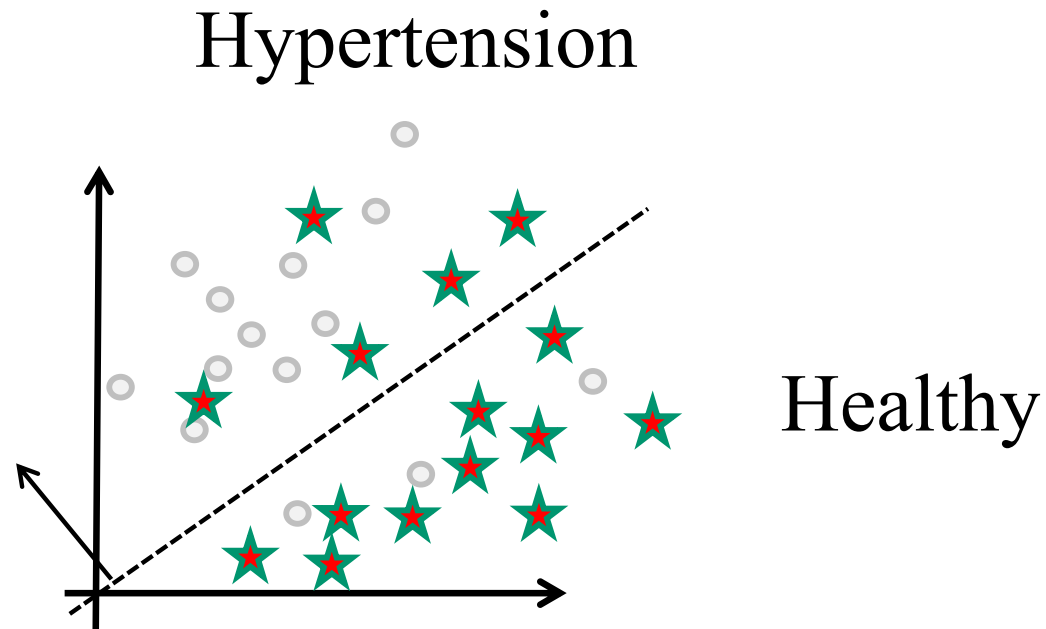FN + TN = 13= TOTAL # of patients classified -1

# False positive rate

| TP = 11 | FP=5 |
|---------|------|
| FN = 3 | TN=10 |



Hypertension

Healthy

$$FPR = FP/(FP+TN) = 5/15$$

# Specificity
**(true negative rate)**

Hypertension

Healthy

$$Sp = TN/(FP+TN)=11/15$$

116

# False negative rate

| TP = 11 | FP=5 |
|---------|------|
| FN = 3 | TN=10 |

Hypertension

Healthy

$$FNR = FN/(FN+TP) = 3/14$$

# **Recall**
**(True positive rate)**

| TP = 11 | FP=5 |
|---------|------|
| FN = 3 | TN=10 |

Hypertension

Healthy

$$\mathbf{Re} = TP/(FN+TP)=1-FNR=11/14$$

# Precision

| TP = 11 | FP=5 |
|---------|------|
| FN = 3  | TN=10 |

Hypertension

Healthy

$$\textbf{Pr} = TP/(TP+FP) = 11/16$$

# False Discovery Rate

Hypertension

Healthy

$$\mathbf{FDR} = FP/(TP+FP) = 5/16$$

# Accuracy

Hypertension

Healthy

Rate of good classification = (TP+TN)/(FP+FN+TP+TN) = 21/29

# In short

Ground truth

|  | Positive | Negative |
|---|---|---|
| **Positive** | TP = 11 | FP=5 |
| **Negative** | FN=3 | TN=10 |

Model prediction

TN + FP = 15 = TOTAL # negative
TP + FN = 14 = TOTAL # positive

TP + FP = 16 = TOTAL # of patients classified +1
FN + TN = 13= TOTAL # of patients classified -1

False positive rate = FP/(FP+TN) = 5/15
False negative rate = FN/(FN+TP) = 3/14

Specificity (**Sp**) = TN/(FP+TN)=1-FPR=10/15
Recall (**Re**) = TP/(TP+TN)=11/14
Precision (**Pr**) = TP/(TP+FP) = 11/16

Accuracy = (TP+TN)/(FP+FN+TP+TN) = 21/29
**F-measure** = 2* (Re*Pr)/(Pr+Re)=0.73

# Different thresholds, different results



FPR:low
FNR:high

FPR:med
FNR:med

FPR:high
FNR:low

# ROC curves

Compute Recall and FPR for **different thresholds**

# ROC curves

Good way for comparing methods

**Which method is best?**

**Ultimate goal**

Recall (1-FNR)

False positive rate

# ROC curves

**Example :** 10 motion detection methods

# Precision recall curve

**Sampe spirit that the ROC curve**

**Which one is best?**



Precision (y-axis)

Recall (x-axis)

**Bull's eye**
Pr=Re=1

# Segmentation metrics

Sørensen–Dice coefficient

Ground truth

|  | Positive | Negative |
|---|---|---|
| **Positive** | TP = 11 | FP=5 |
| **Negative** | FN=3 | TN=10 |

Model prediction

$$\text{Dice} = 2TP/(TP+FP+FN)$$
$$= 2*11/(2*11+5+10)$$
$$= 0.59$$

# Segmentation metrics

Sørensen–Dice coefficient

Ground truth

|  | Positive | Negative |
|---|---|---|
| **Positive** | TP = 11 | FP=5 |
| **Negative** | FN=3 | **TN=10,0000** |

Model prediction

**Useful when TN is very large**

Dice = 2TP/(TP+FP+FN)
  = 2*11/(2*11+5+10)
  = 0.59

131

# Right ventricle segmentation

Sørensen–Dice coefficient



Prediction           Ground truth

# Right ventricle segmentation

Sørensen–Dice coefficient



133

# Right ventricle segmentation

Sørensen–Dice coefficient



$$Dice = \frac{2\left|GT \cap Pred\right|}{\left|GT\right| + \left|Pred\right|}$$

134

# Right ventricle segmentation

Limit of the Dice coefficient



Dice=0.95

Dice=0.95

# Right ventricle segmentation

Hausdorff distance

dist = 0.02 mm

dist = 1mm

dist = 7mm

# Right ventricle segmentation

Hausdorff distance

dist = 0.02 mm

dist = 1mm

$$Hd = \max_{a \in GT} \left[ \min_{b \in pred} \left[ dist(a,b) \right] \right]$$

dist = 7mm

# Right ventricle segmentation

Hausdorff distance

dist = 0.02 mm

dist = 1mm

**dist = 7mm**

$Hd = 7\,mm$

# Right ventricle segmentation

Hausdorff distance



Hd=7mm

Hd=1.6mm

# Extra slides

# Better understand

*Cross entropy* vs *Hinge loss*

# Cross entropy vs Hinge loss

Different *loss* = different **network output**

- Hinge loss : output = <u>matrix-vector</u>
- Cross entrpty: sortie = <u>softmax</u>



$$y_{W_i}\left(\vec{x}_n\right) = \vec{W}_i^{\mathrm{T}}\vec{x}$$

$$y_{W_i}\left(\vec{x}_n\right) = \frac{e^{\vec{W}_i^{\mathrm{T}}\vec{x}}}{\sum_j e^{\vec{W}_j^{\mathrm{T}}\vec{x}}}$$

143

$$\vec{x} = \begin{bmatrix} -15 \\ 22 \\ -44 \\ 56 \end{bmatrix}, t = 2$$

$$\begin{bmatrix} 0.0 & 0.01 & -0.05 & 0.1 & 0.05 \\ 0.2 & 0.7 & 0.2 & 0.05 & 0.16 \\ -0.3 & 0.0 & -0.45 & -0.2 & 0.03 \end{bmatrix}$$

$$W$$

$$\begin{bmatrix} 1 \\ -15 \\ 22 \\ -44 \\ 56 \end{bmatrix}$$

$$\vec{x}$$

**Hinge loss**

*Score*

$$\begin{bmatrix} -2.85 \\ 0.86 \\ 0.28 \end{bmatrix}$$

$$\max(0, -2.85 - 0.28 + 1) +$$
$$\max(0, 0.86 - 0.28 + 1)$$
$$=$$
$$\max(0, -2.13) + \max(0, 1.58)$$
$$=$$

**1.58**

**Cross entropy**

*Score*

$$\begin{bmatrix} -2.85 \\ 0.86 \\ 0.28 \end{bmatrix} \xrightarrow{\exp} \begin{bmatrix} 0.06 \\ 2.36 \\ 1.32 \end{bmatrix} \xrightarrow{\text{norm}} \begin{bmatrix} 0.02 \\ 0.63 \\ 0.35 \end{bmatrix}$$

$$-\ln(0.35)$$
$$=$$
**0.452**

**(Softmax)**

$$\vec{x} = \begin{bmatrix} -15 \\ 22 \\ -44 \\ 56 \end{bmatrix}, t = 2$$

$$\begin{bmatrix} 0.0 & 0.01 & -0.05 & 0.1 & 0.05 \\ 0.2 & 0.7 & 0.2 & 0.05 & 0.16 \\ -0.3 & 0.0 & -0.45 & -0.2 & 0.03 \end{bmatrix} \begin{bmatrix} 1 \\ -15 \\ 22 \\ -44 \\ 56 \end{bmatrix}$$

**Hinge loss**

*Score*

$$\begin{bmatrix} -2.85 \\ 0.86 \\ 0.28 \end{bmatrix}$$

$$\max(0, -2.85 - 0.28 + 1) +$$
$$\max(0, 0.86 - 0.28 + 1)$$
$$=$$
$$\max(0, -2.13) + \max(0, 1.58)$$
$$=$$

***1.58***

**Cross entropy**

*Score*

$$\begin{bmatrix} -2.85 \\ 0.86 \\ 0.28 \end{bmatrix} \xrightarrow{\exp} \begin{bmatrix} 0.06 \\ 2.36 \\ 1.32 \end{bmatrix} \xrightarrow{\text{norm}} \begin{bmatrix} 0.02 \\ 0.63 \\ 0.35 \end{bmatrix}$$

$$-\ln(0.35)$$
$$=$$
***0.452***

**(Softmax)**

**Q1**: What happens if we increase the score of class 0(-2.85)?

$$\vec{x} = \begin{bmatrix} -15 \\ 22 \\ -44 \\ 56 \end{bmatrix}, t = 2$$

$$\begin{bmatrix} 0.0 & 0.01 & -0.05 & 0.1 & 0.05 \\ 0.2 & 0.7 & 0.2 & 0.05 & 0.16 \\ -0.3 & 0.0 & -0.45 & -0.2 & 0.03 \end{bmatrix} \begin{bmatrix} 1 \\ -15 \\ 22 \\ -44 \\ 56 \end{bmatrix}$$

**Hinge loss**

*Score*

$$\begin{bmatrix} -2.85 \\ 0.86 \\ 0.28 \end{bmatrix}$$

$$\max(0, -2.85 - 0.28 + 1) +$$
$$\max(0, 0.86 - 0.28 + 1)$$
$$=$$
$$\max(0, -2.13) + \max(0, 1.58)$$
$$=$$

***1.58***

**Cross entropy**

*Score*

$$\begin{bmatrix} -2.85 \\ 0.86 \\ 0.28 \end{bmatrix} \xrightarrow{\exp} \begin{bmatrix} 0.06 \\ 2.36 \\ 1.32 \end{bmatrix} \xrightarrow{\text{norm}} \begin{bmatrix} 0.02 \\ 0.63 \\ 0.35 \end{bmatrix}$$

$$-\ln(0.35)$$
$$=$$
***0.452***

**(Softmax)**

**Q2**: What happens if we increase the score of class 1(0.86)?

$$\vec{x} = \begin{bmatrix} -15 \\ 22 \\ -44 \\ 56 \end{bmatrix}, t = 2$$

$$\begin{bmatrix} 0.0 & 0.01 & -0.05 & 0.1 & 0.05 \\ 0.2 & 0.7 & 0.2 & 0.05 & 0.16 \\ -0.3 & 0.0 & -0.45 & -0.2 & 0.03 \end{bmatrix} \begin{bmatrix} 1 \\ -15 \\ 22 \\ -44 \\ 56 \end{bmatrix}$$

**Q3**: what is the MIN/MAX values of those two losses?

*Hinge loss*

*Score*

$$\begin{bmatrix} -2.85 \\ 0.86 \\ 0.28 \end{bmatrix}$$

$$\max(0, -2.85 - 0.28 + 1) +$$
$$\max(0, 0.86 - 0.28 + 1)$$
$$=$$
$$\max(0, -2.13) + \max(0, 1.58)$$
$$=$$
***1.58***

**Cross entropy**

*Score*

$$\begin{bmatrix} -2.85 \\ 0.86 \\ 0.28 \end{bmatrix} \xrightarrow{\text{exp}} \begin{bmatrix} 0.06 \\ 2.36 \\ 1.32 \end{bmatrix} \xrightarrow{\text{norm}} \begin{bmatrix} 0.02 \\ 0.63 \\ 0.35 \end{bmatrix}$$

$$-\ln(0.35)$$
$$=$$
***0.452***

**(Softmax)**

$$\vec{x} = \begin{bmatrix} -15 \\ 22 \\ -44 \\ 56 \end{bmatrix}, t = 2$$

**Hinge loss**

*Score*

$$\begin{bmatrix} -2.85 \\ 0.86 \\ 0.28 \end{bmatrix}$$

$$\max(0, -2.85 - 0.28 + 1) +$$
$$\max(0, 0.86 - 0.28 + 1)$$
$$=$$
$$\max(0, -2.13) + \max(0, 1.58)$$
$$=$$
*1.58*

$$\begin{bmatrix} 0.0 & 0.01 & -0.05 & 0.1 & 0.05 \\ 0.2 & 0.7 & 0.2 & 0.05 & 0.16 \\ -0.3 & 0.0 & -0.45 & -0.2 & 0.03 \end{bmatrix} \begin{bmatrix} 1 \\ -15 \\ 22 \\ -44 \\ 56 \end{bmatrix}$$

**Cross entropy**

*Score*

$$\begin{bmatrix} -2.85 \\ 0.86 \\ 0.28 \end{bmatrix} \xrightarrow{\exp} \begin{bmatrix} 0.06 \\ 2.36 \\ 1.32 \end{bmatrix} \xrightarrow{\text{norm}} \begin{bmatrix} 0.02 \\ 0.63 \\ 0.35 \end{bmatrix}$$

$$-\ln(0.35)$$
$$=$$
*0.452*

**(Softmax)**

**Q4**: what would happed to the total loss if we were to add an L2 regularization?

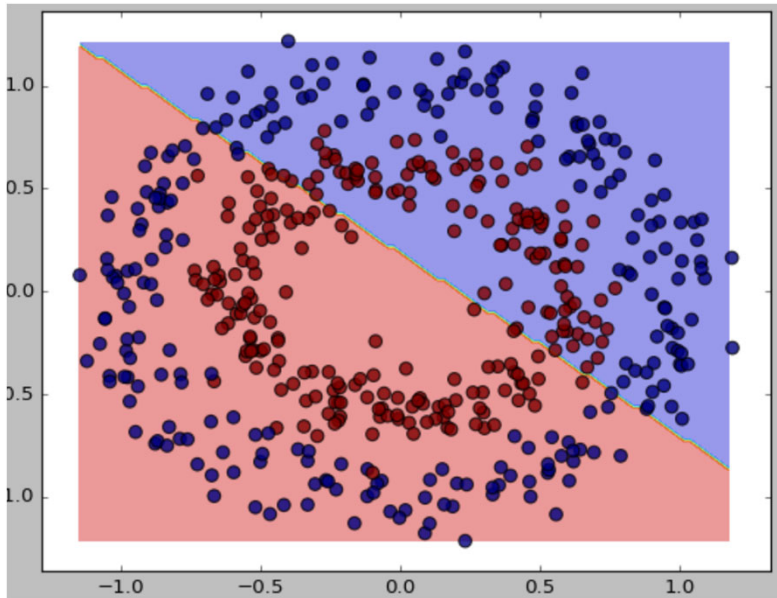Increase the number
of neurons

Increase the number
of layers
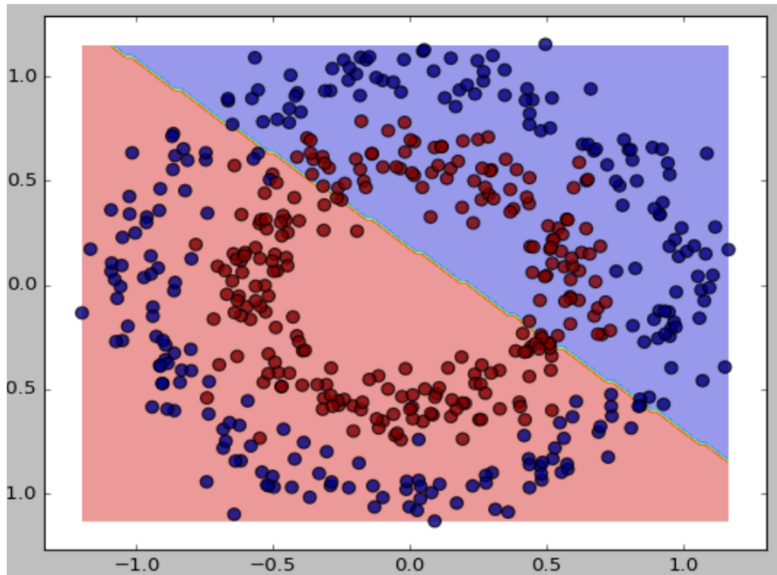
**Increase the
capacity of the network**

Increasing the capacity of the network
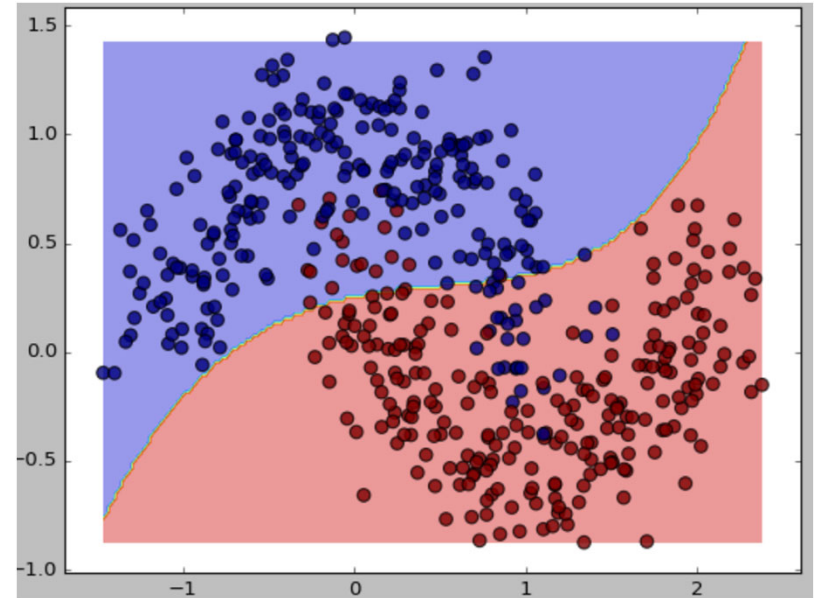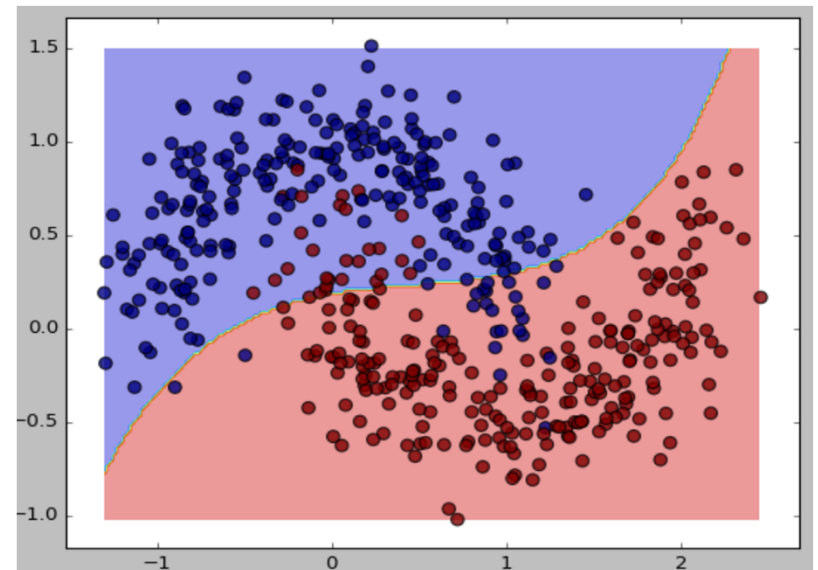can lead to **over-fitting**

# Under-fitting

# Could do better...



Precision on the training set = 52.2%

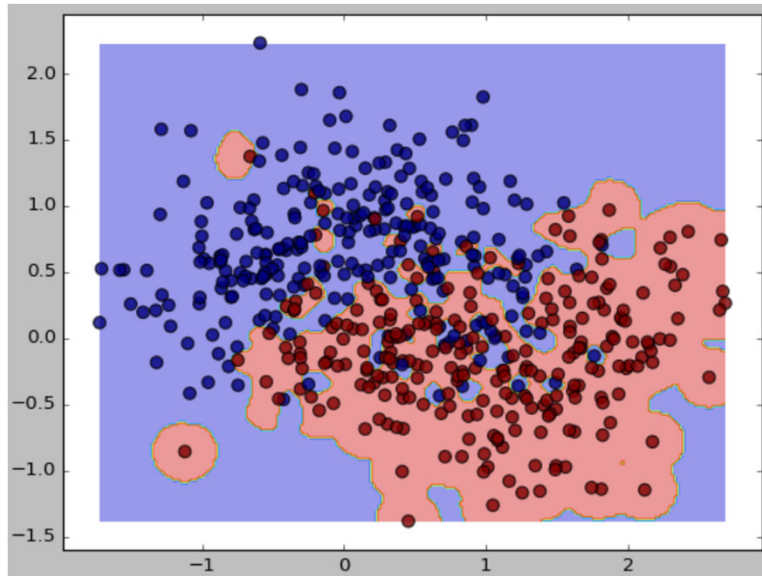Precision on the training set = 82%
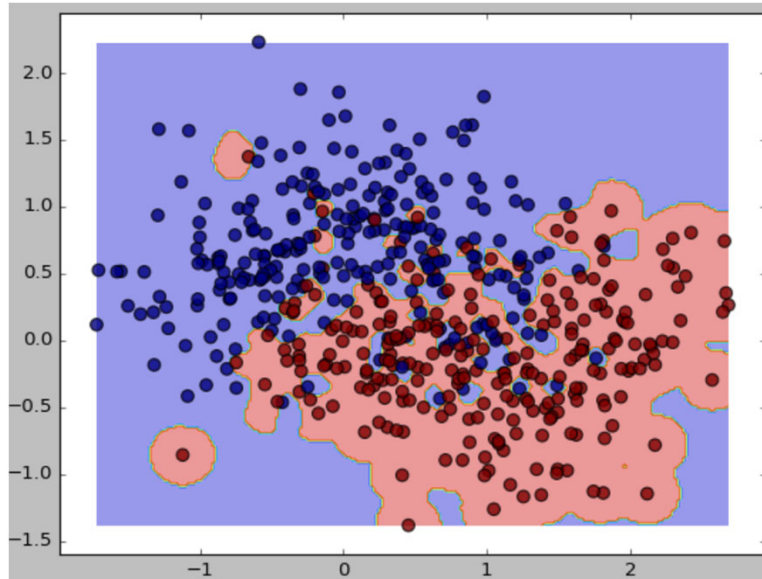
Precision on the test set = 51.2%

Precision on the test set = 80%
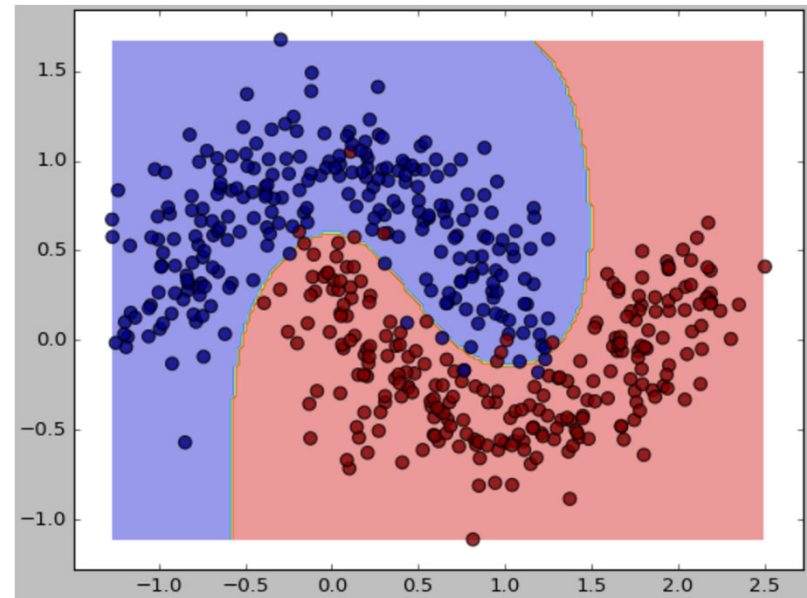
**Overfitting**

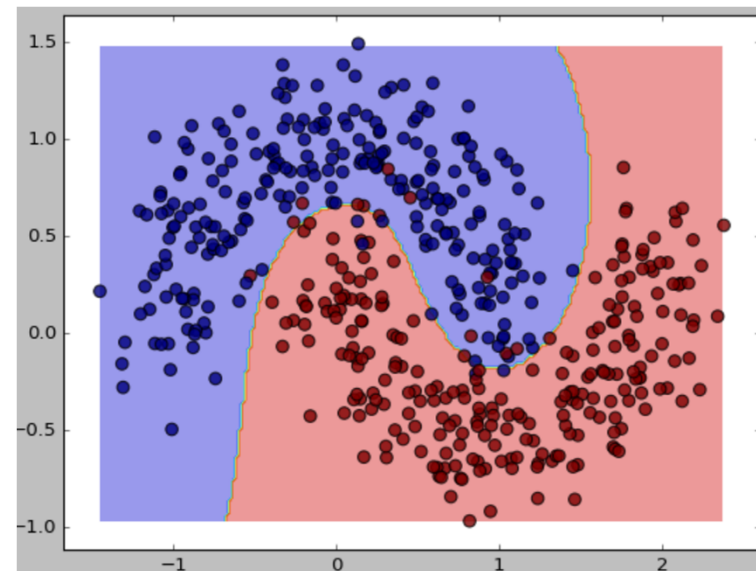Precision on the training set = 99.6%

Precision on the test setb= 78%
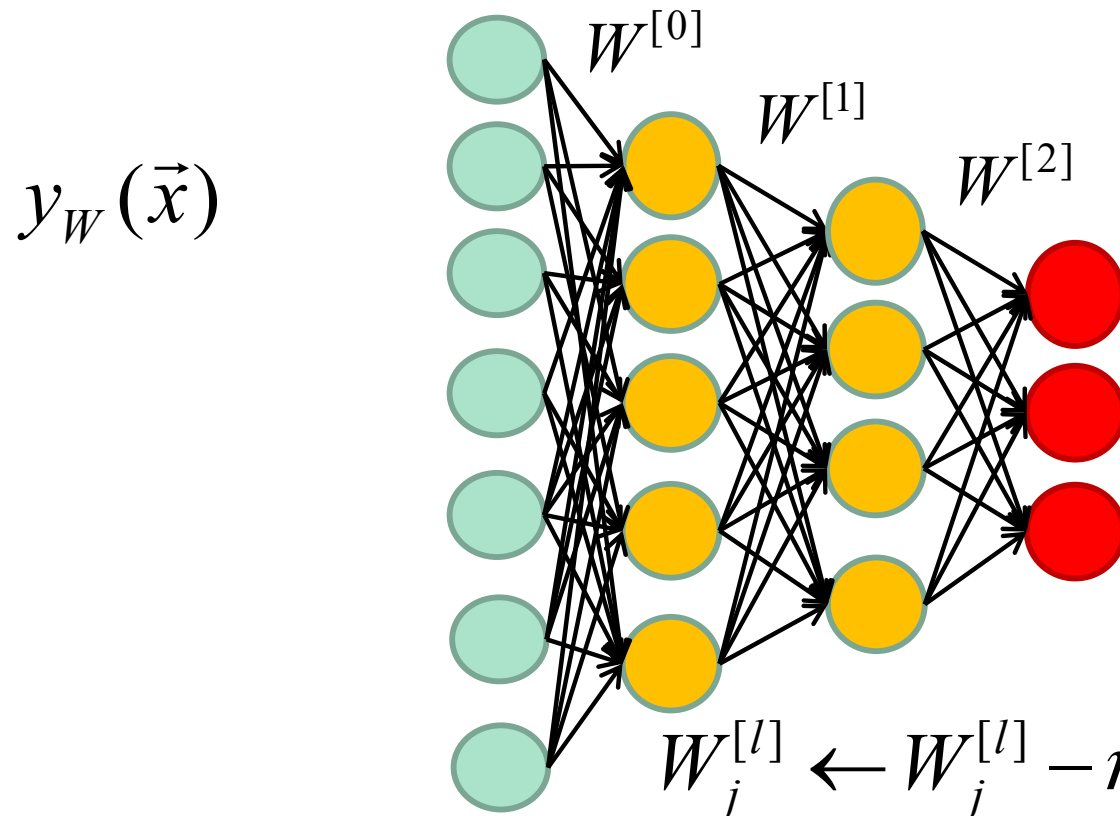
**SUPER !!!**

Precision on the training set =97.8%

Precision on the test set = 96.2%

# kD, 2 Classes, 4 hidden layer network

Input layer    Hidden Layer 1    Hidden Layer 2    Hidden Layer 3    Hidden Layer 4    Output layer

$$W^{[0]} \in R^{5 \times k+1}$$

$$W^{[1]} \in R^{3 \times 6}$$

$$W^{[2]} \in R^{4 \times 4}$$

$$W^{[3]} \in R^{7 \times 5}$$

$$\vec{w}^{[4]} \in R^{8}$$

$$y_W(\vec{x})$$

$$W^{[0]}$$

$$W^{[1]}$$

$$W^{[2]}$$

$$L\left(y_{\vec{w}}(\vec{x}), D\right)$$

$$W_j^{[l]} \leftarrow W_j^{[l]} - \eta \frac{\partial\left(L\left(y_W(\vec{x})\right)\right)}{\partial W_{j}^{[l]}(\vec{x}), D)}$$