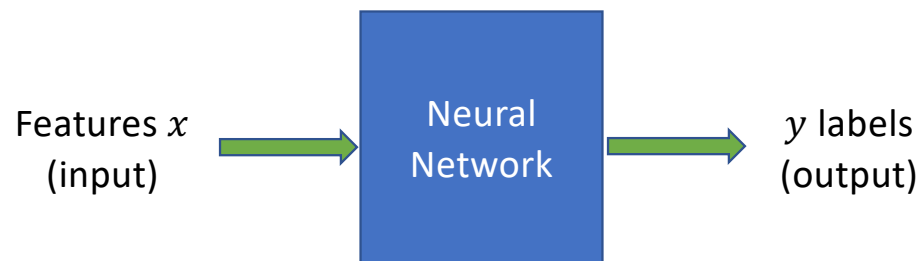


## Chapter 2

# Regression and classification

Neural networks and deep learning

# Regression and classification



Primarily, neural network are used to *predict* output **labels** from input **features**.

Prediction tasks can be classified into two categories:

**Regression:** the labels are continuous (age, income, height, etc.)

**Classification:** the labels are discrete (sex, digits, type of flowers, etc.),

Training finds network weights and biases that are optimal for **prediction** of labels from features.

# Linear neuron

Synaptic input  $u$  to a neuron is given by

$$u = \mathbf{w}^T \mathbf{x} + b$$

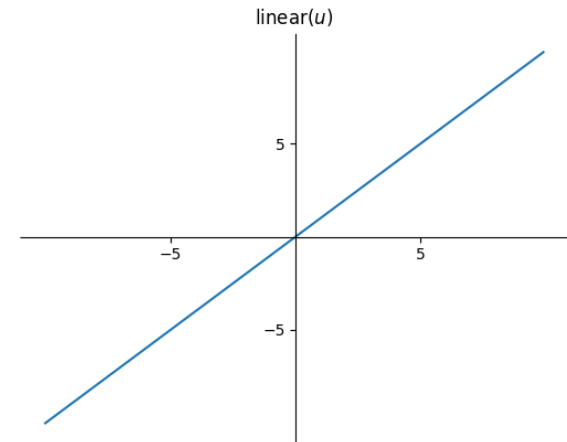
A linear neuron has a *linear activation function*. That is,

$$y = f(u) = u$$

A linear neuron with weights  $\mathbf{w} = (w_1 \ w_2 \ \cdots \ w_n)^T$  and bias  $b$  has an output:

$$y = \mathbf{w}^T \mathbf{x} + b$$

where input  $\mathbf{x} = (x_1 \ x_2 \ \cdots \ x_n)^T \in \mathbf{R}^n$  and output  $y \in \mathbf{R}$ .



# Linear neuron performs linear regression

Representing a dependent (output) variable as a linear combination of independent (input) variables is known as **linear regression**.

The output of a linear neuron can be written as

$$y = w_1x_1 + w_2x_2 \cdots + w_nx_n + b$$

where  $x_1, x_2, \cdots x_n$  are the inputs. That is, a linear neuron performs linear regression and the weights and biases (that is,  $b$  and  $w_1, \dots, w_n$ ) act as regression coefficients. The above function forms a **hyperplane** in Euclidean space  $\mathbf{R}^n$ .

Given a training examples  $\{(x_p, d_p)\}_{p=1}^P$  where input  $x_p \in \mathbf{R}^n$  and target  $d_p \in \mathbf{R}$ , training a linear neuron finds a linear mapping  $\phi: \mathbf{R}^n \rightarrow \mathbf{R}$  given by:

$$y = \mathbf{w}^T \mathbf{x} + b$$

# Stochastic gradient descent (SGD) for linear neuron

The **cost function**  $J$  for regression is usually given as the *square error* (s.e.) between neuron outputs and targets.

Given a training pattern  $(\mathbf{x}, d)$ ,  $\frac{1}{2}$  square error cost  $J$  is defined as

$$J = \frac{1}{2} (d - y)^2$$

where  $y$  is neuron output for input pattern  $\mathbf{x}$  and  $d$  is the target label.

$$y = \mathbf{w}^T \mathbf{x} + b$$

The  $\frac{1}{2}$  in the cost function is introduced to simplify learning equations and does not affect the optimal values of the parameters (weights and bias).

# SGD for linear neuron



$$J = \frac{1}{2} (d - y)^2$$
$$y = u = \mathbf{w}^T \mathbf{x} + b$$

$$\frac{\partial J}{\partial u} = \frac{\partial J}{\partial y} = -(d - y) \quad (\text{A})$$

$$\nabla_{\mathbf{w}} J = \frac{\partial J}{\partial \mathbf{w}} = \frac{\partial J}{\partial u} \frac{\partial u}{\partial \mathbf{w}} \quad (\text{B})$$

# SGD for linear neuron

$$u = \mathbf{w}^T \mathbf{x} + b = w_1 x_1 + w_2 x_2 \cdots + w_n x_n + b$$

$$\frac{\partial u}{\partial \mathbf{w}} = \begin{pmatrix} \frac{\partial u}{\partial w_1} \\ \frac{\partial u}{\partial w_2} \\ \vdots \\ \frac{\partial u}{\partial w_n} \end{pmatrix} = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix} = \mathbf{x} \quad (\text{C})$$

$$\frac{\partial u}{\partial \mathbf{w}} = \mathbf{x}$$

Substituting (A) and (C) in (B),

$$\nabla_{\mathbf{w}} J = -(d - y) \mathbf{x} \quad (\text{D})$$

Similarly, since  $\frac{\partial u}{\partial b} = 1$ ,

$$\nabla_b J = \frac{\partial J}{\partial u} \frac{\partial u}{\partial b} = -(d - y) \quad (\text{E})$$

# SGD for linear neuron

Gradient learning equations:

$$\mathbf{w} \leftarrow \mathbf{w} - \alpha \nabla_{\mathbf{w}} J$$

$$b \leftarrow b - \alpha \nabla_b J$$

By substituting from (D) and (E), SGD equations for a linear neuron are given by

$$\mathbf{w} \leftarrow \mathbf{w} + \alpha (d - y) \mathbf{x}$$

$$b \leftarrow b + \alpha (d - y)$$



# SGD algorithm for linear neuron

Given a training dataset  $\{(\mathbf{x}_p, d_p)\}_{p=1}^P$

Set learning parameter  $\alpha$

Initialize  $\mathbf{w}$  and  $b$

Repeat until convergence:

For every training pattern  $(\mathbf{x}_p, d_p)$ :

$$y_p = \mathbf{w}^T \mathbf{x}_p + b$$

$$\mathbf{w} \leftarrow \mathbf{w} + \alpha(d_p - y_p)\mathbf{x}_p$$

$$b \leftarrow b + \alpha(d_p - y_p)$$

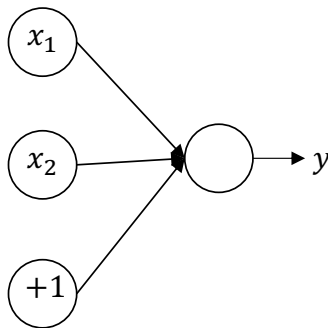
## Example 1: SGD on a linear neuron

Train a linear neuron to perform the following mapping, using stochastic gradient descent (SGD) learning:

$\mathbf{x}^T = (x_1, x_2)$	$d$
(0.54, -0.96)	1.33
(0.27, 0.50)	0.45
(0.00, -0.55)	0.56
(-0.60, 0.52)	-1.66
(-0.66, -0.82)	-1.07
(0.37, 0.91)	0.30

Use a learning factor  $\alpha = 0.01$ .

# Example 1



Let's initialize weights randomly and biases to zeros

$$\mathbf{w} = \begin{pmatrix} 0.92 \\ 0.71 \end{pmatrix} \text{ and } b = 0.0$$

$$\alpha = 0.01$$

Need to shuffle the patterns before every epoch.

# Example 1: epoch 1

## Epoch 1 begins

Shuffle the patterns ....

First pattern  $p = 1$  is applied:

$$\mathbf{x}_p = \begin{pmatrix} 0.54 \\ -0.96 \end{pmatrix} \text{ and } d_p = 1.33$$

$$y_p = \mathbf{w}^T \mathbf{x}_p + b = (0.92 \quad 0.71) \begin{pmatrix} 0.54 \\ -0.96 \end{pmatrix} + 0.0 = -0.19$$

$$\text{s.e.} = (d_p - y_p)^2 = 2.292$$

$$\mathbf{w} = \mathbf{w} + \alpha(d_p - y_p)\mathbf{x}_p = \begin{pmatrix} 0.92 \\ 0.71 \end{pmatrix} + 0.01 \times (1.33 + 0.19) \begin{pmatrix} 0.54 \\ -0.96 \end{pmatrix} = \begin{pmatrix} 0.93 \\ 0.70 \end{pmatrix}$$

$$b = b + \alpha(d_p - y_p) = 0 + 0.01 \times (1.33 + 0.19) = 0.02$$

## Example 1: epoch 1

Second pattern  $p = 2$  is applied:

$$\mathbf{x}_p = \begin{pmatrix} -0.66 \\ -0.82 \end{pmatrix} \text{ and } d_p = -1.07$$

$$y_p = \mathbf{w}^T \mathbf{x}_p + b = (0.93 \quad 0.70) \begin{pmatrix} -0.66 \\ -0.82 \end{pmatrix} + 0.02 = -0.19$$

$$\text{s.e.} = (d_p - y_p)^2 = 0.01$$

$$\mathbf{w} = \mathbf{w} + \alpha(d_p - y_p)\mathbf{x}_p = \begin{pmatrix} 0.93 \\ 0.70 \end{pmatrix} + 0.01 \times (-1.07 + 0.19) \begin{pmatrix} -0.66 \\ -0.82 \end{pmatrix} = \begin{pmatrix} 0.93 \\ 0.70 \end{pmatrix}$$

$$b = b + \alpha(d_p - y_p) = 0.02 + 0.01 \times (-1.07 + 0.19) = 0.02$$

Iterations continues for patterns  $p = 3, \dots 6$ .

**the second epoch starts .....**

Shuffle the patterns

Apply patterns  $p = 1, 2, \dots 6$

**Training epochs continue until convergence.**

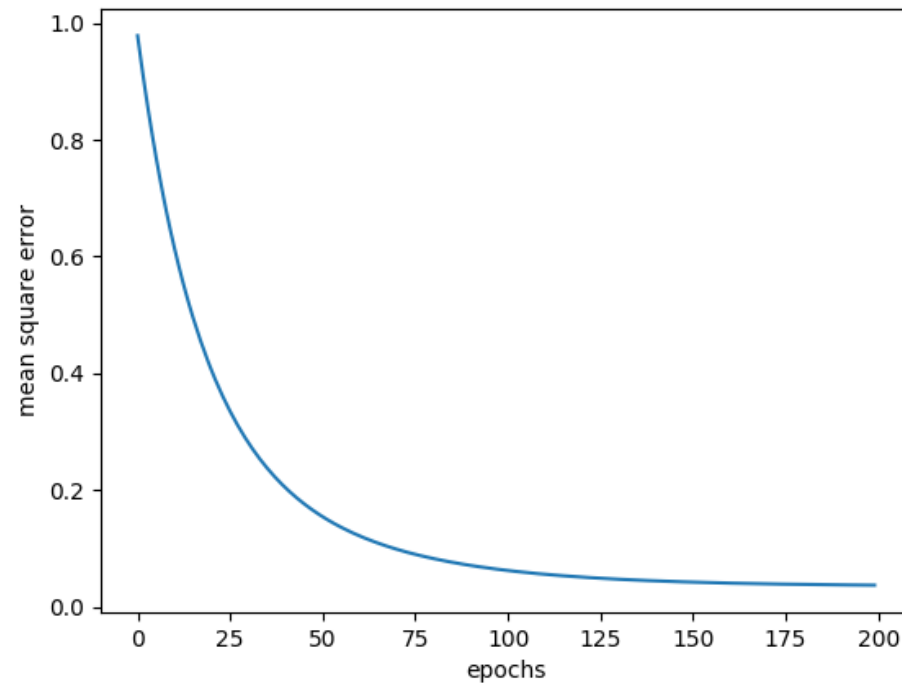
## Example 1: epoch 1

$x_p$	$y_p$	s.e.	$w$	$b$
$x_1 = \begin{pmatrix} 0.54 \\ -0.96 \end{pmatrix}$	-0.19	2.29	$\begin{pmatrix} 0.93 \\ 0.70 \end{pmatrix}$	0.02
$x_2 = \begin{pmatrix} -0.66 \\ -0.82 \end{pmatrix}$	-1.17	0.01	$\begin{pmatrix} 0.93 \\ 0.70 \end{pmatrix}$	0.03
$x_3 = \begin{pmatrix} 0.00 \\ -0.55 \end{pmatrix}$	-0.37	0.87	$\begin{pmatrix} 0.93 \\ 0.69 \end{pmatrix}$	0.03
$x_4 = \begin{pmatrix} 0.27 \\ 0.50 \end{pmatrix}$	0.62	0.03	$\begin{pmatrix} 0.92 \\ 0.69 \end{pmatrix}$	0.02
$x_5 = \begin{pmatrix} -0.60 \\ 0.52 \end{pmatrix}$	-0.17	2.21	$\begin{pmatrix} 0.93 \\ 0.69 \end{pmatrix}$	0.01
$x_6 = \begin{pmatrix} 0.37 \\ 0.91 \end{pmatrix}$	0.98	0.45	$\begin{pmatrix} 0.93 \\ 0.68 \end{pmatrix}$	0.00

## Example 1: epoch 200

$x_p$	$y_p$	s.e.	$w$	$b$
$x_1 = \begin{pmatrix} 0.54 \\ -0.96 \end{pmatrix}$	1.49	0.03	$\begin{pmatrix} 2.00 \\ -0.44 \end{pmatrix}$	-0.01
$x_2 = \begin{pmatrix} 0.00 \\ -0.55 \end{pmatrix}$	0.22	0.12	$\begin{pmatrix} 2.00 \\ -0.44 \end{pmatrix}$	-0.01
$x_3 = \begin{pmatrix} -0.66 \\ -0.82 \end{pmatrix}$	-0.98	0.01	$\begin{pmatrix} 2.00 \\ -0.44 \end{pmatrix}$	-0.01
$x_4 = \begin{pmatrix} 0.37 \\ 0.91 \end{pmatrix}$	0.33	0.00	$\begin{pmatrix} 2.00 \\ -0.44 \end{pmatrix}$	-0.01
$x_5 = \begin{pmatrix} 0.27 \\ 0.50 \end{pmatrix}$	0.30	0.02	$\begin{pmatrix} 2.00 \\ -0.44 \end{pmatrix}$	-0.01
$x_6 = \begin{pmatrix} -0.60 \\ 0.52 \end{pmatrix}$	-1.45	0.04	$\begin{pmatrix} 2.00 \\ -0.44 \end{pmatrix}$	-0.01

# Example 1



$$\text{m.s.e.} = \frac{1}{6} \sum_{p=1}^6 (d_p - y_p)^2$$



# Example 1

At convergence:

$$\mathbf{w} = \begin{pmatrix} 2.00 \\ -0.44 \end{pmatrix}$$
$$b = -0.013$$

Mean square error = 0.037

The regression equation:

$$y = \mathbf{x}^T \mathbf{w} + b = (x_1 \quad x_2) \begin{pmatrix} 2.00 \\ -0.44 \end{pmatrix} - 0.013$$

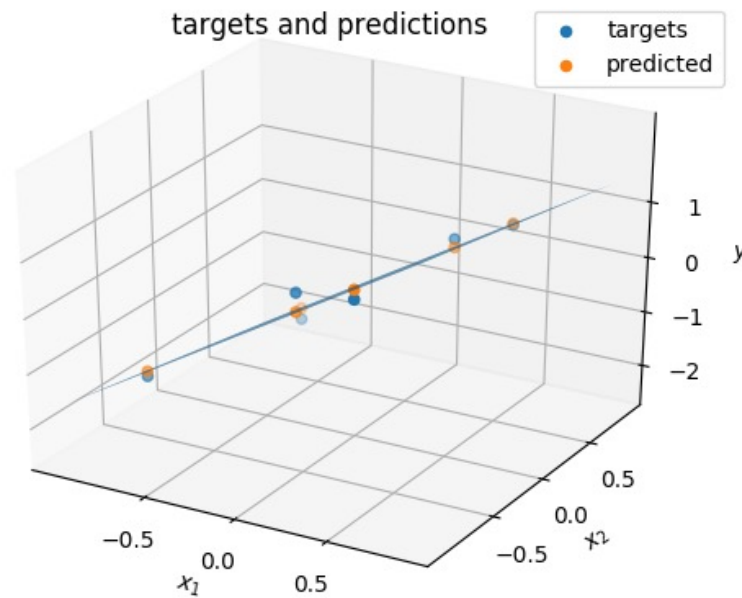
The mapping learnt by the linear neuron:

$$y = 2.00x_1 - 0.44x_2 - 0.013$$

# Example 1

<i>inputs</i> $x = (x_1, x_2)$	<i>predictions</i> $y = 2.00x_1 - 0.44x_2 - 0.01$	<i>targets</i> $d$
(0.54, -0.96)	1.49	1.33
(0.27, 0.50)	0.30	0.45
(0.00, -0.55)	0.22	0.56
(-0.60, 0.52)	-1.45	-1.66
(-0.66, -0.82)	-0.98	-1.07
(0.37, 0.91)	0.33	0.30

# Example 1



The mapping portrays a hyperplane in the 3-dimensional space:

$$y = 2.00x_1 - 0.44x_2 - 0.013$$

# Gradient descent (GD) for linear neuron

Given a training dataset  $\{(\mathbf{x}_p, d_p)\}_{p=1}^P$ , cost function  $J$  is given by the sum of square errors (s.s.e.):

$$J = \frac{1}{2} \sum_{p=1}^P (d_p - y_p)^2$$

where  $y_p$  is the neuron output for input pattern  $\mathbf{x}_p$ .

$$J = \sum_{p=1}^P J_p \quad (\text{F})$$

where  $J_p = \frac{1}{2} (d_p - y_p)^2$  is the square error for the  $p$ th pattern.

# GD for linear neuron

From (F):

$$\begin{aligned}\nabla_{\mathbf{w}} J &= \sum_{p=1}^P \nabla_{\mathbf{w}} J_p \\ &= - \sum_{p=1}^P (d_p - y_p) \mathbf{x}_p && \text{from (D)} \\ &= -((d_1 - y_1)\mathbf{x}_1 + (d_2 - y_2)\mathbf{x}_2 + \cdots + (d_P - y_P)\mathbf{x}_P) \\ &= -(\mathbf{x}_1 \quad \mathbf{x}_2 \quad \cdots \quad \mathbf{x}_P) \begin{pmatrix} (d_1 - y_1) \\ (d_2 - y_2) \\ \vdots \\ (d_P - y_P) \end{pmatrix} \\ &= -\mathbf{X}^T (\mathbf{d} - \mathbf{y})\end{aligned}$$

(G)

where  $\mathbf{X} = \begin{pmatrix} \mathbf{x}_1^T \\ \mathbf{x}_2^T \\ \vdots \\ \mathbf{x}_P^T \end{pmatrix}$  is the data matrix,  $\mathbf{d} = \begin{pmatrix} d_1 \\ d_2 \\ \vdots \\ d_P \end{pmatrix}$  is the target vector, and  $\mathbf{y} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_P \end{pmatrix}$  is the output vector.

## GD for linear neuron

Similarly,  $\nabla_b J$  can be obtained by considering inputs of +1 and substituting a vector of +1 in (G):

$$\nabla_b J = -\mathbf{1}_P^T (\mathbf{d} - \mathbf{y}) \quad (\text{H})$$

where  $\mathbf{1}_P = \begin{pmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{pmatrix}$  has  $P$  elements of 1.

The output vector  $\mathbf{y}$  for the batch of  $P$  patterns is given by

$$\mathbf{y} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_P \end{pmatrix} = \begin{pmatrix} u_1 \\ u_2 \\ \vdots \\ u_P \end{pmatrix} = \begin{pmatrix} \mathbf{x}_1^T \mathbf{w} + b \\ \mathbf{x}_2^T \mathbf{w} + b \\ \vdots \\ \mathbf{x}_P^T \mathbf{w} + b \end{pmatrix} = \begin{pmatrix} \mathbf{x}_1^T \\ \mathbf{x}_2^T \\ \vdots \\ \mathbf{x}_P^T \end{pmatrix} \mathbf{w} + b \begin{pmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{pmatrix} = \mathbf{X} \mathbf{w} + b \mathbf{1}_P$$

# GD for linear neuron

Substituting (G) and (H) in gradient descent equations:

$$\begin{aligned}\mathbf{w} &\leftarrow \mathbf{w} - \alpha \nabla_{\mathbf{w}} J \\ b &\leftarrow b - \alpha \nabla_b J\end{aligned}$$

We get GD learning equations for the linear neuron as

$$\begin{aligned}\mathbf{w} &\leftarrow \mathbf{w} + \alpha \mathbf{X}^T (\mathbf{d} - \mathbf{y}) \\ b &\leftarrow b + \alpha \mathbf{1}_p^T (\mathbf{d} - \mathbf{y})\end{aligned}$$

And  $\alpha$  is the learning factor.

Where:

$$\mathbf{y} = \mathbf{X}\mathbf{w} + b \mathbf{1}_p$$

# GD for linear neuron

Given a training dataset  $(X, d)$

Set learning parameter  $\alpha$

Initialize  $\mathbf{w}$  and  $b$

Repeat until convergence:

$$\mathbf{y} = X\mathbf{w} + b\mathbf{1}_p$$

$$\mathbf{w} \leftarrow \mathbf{w} + \alpha X^T(d - \mathbf{y})$$

$$b \leftarrow b + \alpha \mathbf{1}_p^T(d - \mathbf{y})$$



## GD and SGD for a linear neuron

GD	SGD
$(X, d)$	$(x_p, d_p)$
$J = \frac{1}{2} \sum_{p=1}^P (d_p - y_p)^2$	$J = \frac{1}{2} (d_p - y_p)^2$
$y = u = Xw + b\mathbf{1}_P$	$y_p = u_p = x_p^T w + b$
$w \leftarrow w + \alpha X^T (d - y)$	$w \leftarrow w + \alpha (d_p - y_p) x_p$
$b \leftarrow b + \alpha \mathbf{1}_P^T (d - y)$	$b \leftarrow b + \alpha (d_p - y_p)$

# Perceptron

Perceptron is a neuron having a **sigmoid** activation function and has an output

.

$$y = f(u)$$

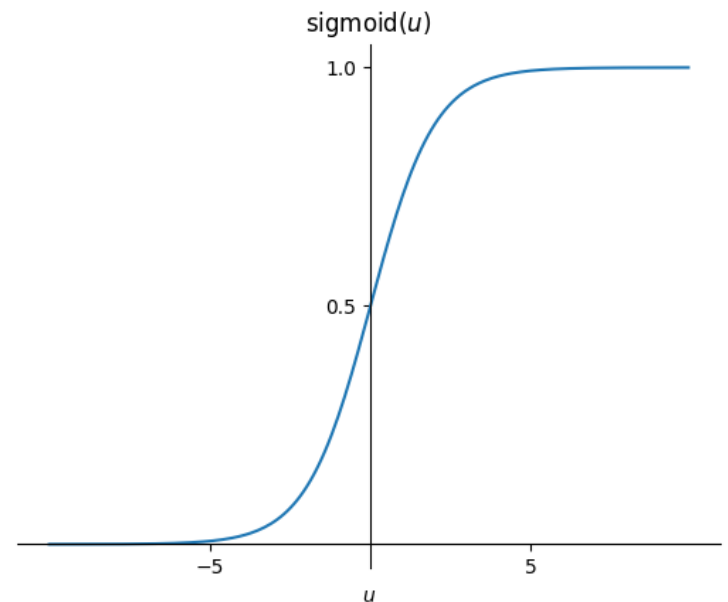
Where:

$$f(u) = \frac{1}{1 + e^{-u}} = \text{sigmoid}(u)$$

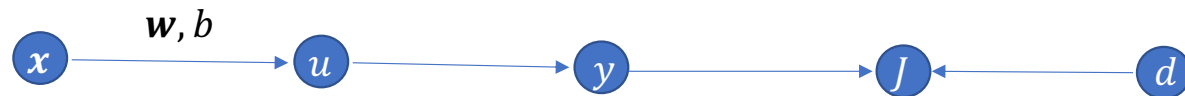
And  $u = \mathbf{w}^T \mathbf{x} + b$

The square error is used as cost function for learning.

Perceptron performs a *non-linear regression* of inputs.



# SGD for perceptron



Cost function  $J$  is given by

$$J = \frac{1}{2} (d - y)^2$$

where  $y = f(u)$  and  $u = \mathbf{w}^T \mathbf{x} + b$

The gradient with respect to the synaptic input:

$$\frac{\partial J}{\partial u} = \frac{\partial J}{\partial y} \frac{\partial y}{\partial u} = -(d - y) f'(u)$$

From (C),  $\frac{\partial u}{\partial \mathbf{w}} = \mathbf{x}$  and  $\frac{\partial u}{\partial b} = 1$ .

$$\nabla_{\mathbf{w}} J = \frac{\partial J}{\partial u} \frac{\partial u}{\partial \mathbf{w}} = -(d - y) f'(u) \mathbf{x} \quad (I)$$

$$\nabla_b J = \frac{\partial J}{\partial u} \frac{\partial u}{\partial b} = -(d - y) f'(u) \quad (J)$$

# SGD for perceptron

Gradient learning equations:

$$\mathbf{w} \leftarrow \mathbf{w} - \alpha \nabla_{\mathbf{w}} J$$

$$b \leftarrow b - \alpha \nabla_b J$$

Substituting gradients from (I) and (J), SGD equations for a perceptron are given by

$$\mathbf{w} \leftarrow \mathbf{w} + \alpha(d - y)f'(u)\mathbf{x}$$

$$b \leftarrow b + \alpha(d - y)f'(u)$$

# SGD algorithm for perceptron

Given a training dataset  $\{(\mathbf{x}_p, d_p)\}_{p=1}^P$

Set learning parameter  $\alpha$

Initialize  $\mathbf{w}$  and  $b$

Repeat until convergence:

For every training pattern  $(\mathbf{x}_p, d_p)$ :

$$u_p = \mathbf{w}^T \mathbf{x}_p + b$$

$$y_p = f(u_p) = \frac{1}{1 + e^{-u_p}}$$

$$\mathbf{w} \leftarrow \mathbf{w} + \alpha(d_p - y_p)f'(u_p)\mathbf{x}_p$$

$$b \leftarrow b + \alpha(d_p - y_p)f'(u_p)$$

## GD for perceptron

Given a training dataset  $\{(\mathbf{x}_p, d_p)\}_{p=1}^P$ , cost function  $J$  is given by the sum of square errors (s.s.e) over all the patterns:

$$J = \frac{1}{2} \sum_{p=1}^P (d_p - y_p)^2 = \sum_{p=1}^P J_p \quad (\text{F})$$

where  $J_p = \frac{1}{2} (d_p - y_p)^2$  is the square error for the  $p$ th pattern.

# GD for perceptron

From (F):

$$\begin{aligned}\nabla_{\mathbf{w}} J &= \sum_{p=1}^P \nabla_{\mathbf{w}} J_p \\ &= - \sum_{p=1}^P (d_p - y_p) f'(u_p) \mathbf{x}_p && \text{From (J)} \\ &= -((d_1 - y_1) f'(u_1) \mathbf{x}_1 + (d_2 - y_2) f'(u_2) \mathbf{x}_2 + \cdots + (d_P - y_P) f'(u_P) \mathbf{x}_P) \\ &= -(\mathbf{x}_1 \quad \mathbf{x}_2 \quad \cdots \quad \mathbf{x}_P) \begin{pmatrix} (d_1 - y_1) f'(u_1) \\ (d_2 - y_2) f'(u_2) \\ \vdots \\ (d_P - y_P) f'(u_P) \end{pmatrix} \\ &= -\mathbf{X}^T (\mathbf{d} - \mathbf{y}) \cdot f'(\mathbf{u})\end{aligned}$$

(K)

$$\text{where } \mathbf{X} = \begin{pmatrix} \mathbf{x}_1^T \\ \mathbf{x}_2^T \\ \vdots \\ \mathbf{x}_P^T \end{pmatrix}, \mathbf{d} = \begin{pmatrix} d_1 \\ d_2 \\ \vdots \\ d_P \end{pmatrix}, \mathbf{y} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_P \end{pmatrix}, \text{ and } f'(\mathbf{u}) = \begin{pmatrix} f'(u_1) \\ f'(u_2) \\ \vdots \\ f'(u_P) \end{pmatrix}$$

## GD for perceptron

Substituting  $\mathbf{X}^T$  by  $\mathbf{1}_P^T$  in (K), we get

$$\nabla_b J = -\mathbf{1}_P^T (\mathbf{d} - \mathbf{y}) \cdot f'(\mathbf{u}) \quad (\text{L})$$

where  $\mathbf{1}_P = (1 \ 1 \ \dots \ 1)^T$ .

The gradient descent learning is given by

$$\begin{aligned} \mathbf{w} &\leftarrow \mathbf{w} - \alpha \nabla_{\mathbf{w}} J \\ b &\leftarrow b - \alpha \nabla_b J \end{aligned}$$

Substituting (K) and (L), we get the learning equations:

$$\begin{aligned} \mathbf{w} &\leftarrow \mathbf{w} + \alpha \mathbf{X}^T (\mathbf{d} - \mathbf{y}) \cdot f'(\mathbf{u}) \\ b &\leftarrow b + \alpha \mathbf{1}_P^T (\mathbf{d} - \mathbf{y}) \cdot f'(\mathbf{u}) \end{aligned}$$

Note that  $\cdot$  is the element-wise product.



# GD for perceptron

Given a training dataset( $X, \mathbf{d}$ )

Set learning parameter  $\alpha$

Initialize  $\mathbf{w}$  and  $b$

Repeat until convergence:

$$\mathbf{u} = X\mathbf{w} + b\mathbf{1}_P$$

$$y = f(\mathbf{u}) = \frac{1}{1+e^{-u}}$$

$$\mathbf{w} \leftarrow \mathbf{w} + \alpha X^T(\mathbf{d} - y) \cdot f'(\mathbf{u})$$

$$b \leftarrow b + \alpha \mathbf{1}_P^T(\mathbf{d} - y) \cdot f'(\mathbf{u})$$

# Gradient descent for perceptron

GD	SGD
$(X, \mathbf{d})$	$(\mathbf{x}_p, d_p)$
$J = \frac{1}{2} \sum_{p=1}^P (d_p - y_p)^2$	$J = \frac{1}{2} (d_p - y_p)^2$
$\mathbf{u} = X\mathbf{w} + b\mathbf{1}_P$	$u_p = \mathbf{x}_p^T \mathbf{w} + b$
$\mathbf{y} = f(\mathbf{u})$	$y_p = f(u_p)$
$\mathbf{w} = \mathbf{w} + \alpha X^T (\mathbf{d} - \mathbf{y}) \cdot f'(\mathbf{u})$	$\mathbf{w} = \mathbf{w} + \alpha (d_p - y_p) f'(u_p) \mathbf{x}_p$
$b = b + \alpha \mathbf{1}_P^T (\mathbf{d} - \mathbf{y}) \cdot f'(\mathbf{u})$	$b = b + \alpha (d_p - y_p) f'(u_p)$

# Derivatives of Sigmoid

The activation function of the (continuous) perceptron is *sigmoid function* (i.e., unipolar sigmoidal function with  $a = 1.0$  and  $b = 1.0$ ) :

$$y = f(u) = \frac{1}{1 + e^{-u}}$$

The derivative is given by

$$f'(u) = \frac{-1}{(1+e^{-u})^2} \frac{\partial(e^{-u})}{\partial u} = \frac{e^{-u}}{(1+e^{-u})^2} = \frac{1}{1+e^{-u}} - \frac{1}{(1+e^{-u})^2} = y(1-y)$$

For *Tanh* function (bipolar sigmoid):

$$y = f(u) = \frac{e^{+u} - e^{-u}}{e^{+u} + e^{-u}}$$
$$f'(u) = \frac{(e^{+u} + e^{-u})(e^{+u} + e^{-u}) - (e^{+u} - e^{-u})(e^{+u} - e^{-u})}{(e^{+u} + e^{-u})^2} = \left(1 - \left(\frac{e^{+u} - e^{-u}}{e^{+u} + e^{-u}}\right)^2\right) = 1 - y^2$$

## Example 2

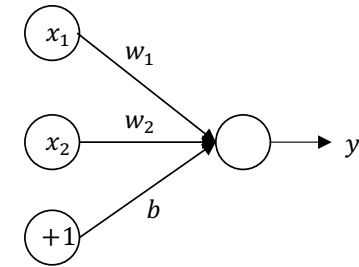
Design a perceptron to learn the following mapping by using gradient descent (GD):

$x = (x_1, x_2)$	$d$
(0.77, 0.02)	2.91
(0.63, 0.75)	0.55
(0.50, 0.22)	1.28
(0.20, 0.76)	-0.74
(0.17, 0.09)	0.88
(0.69, 0.95)	0.30
(0.00, 0.51)	-0.28

Use learning factor  $\alpha = 0.01$ .

## Example 2

$$\mathbf{X} = \begin{pmatrix} 0.77 & 0.02 \\ 0.63 & 0.75 \\ 0.50 & 0.22 \\ 0.20 & 0.76 \\ 0.17 & 0.09 \\ 0.69 & 0.95 \\ 0.00 & 0.51 \end{pmatrix} \text{ and } \mathbf{d} = \begin{pmatrix} 2.91 \\ 0.55 \\ 1.28 \\ -0.74 \\ 0.88 \\ 0.30 \\ -0.28 \end{pmatrix}$$



Initially,  $\mathbf{w} = \begin{pmatrix} 0.81 \\ 0.61 \end{pmatrix}$  and  $b = 0.0$

$\alpha = 0.01$

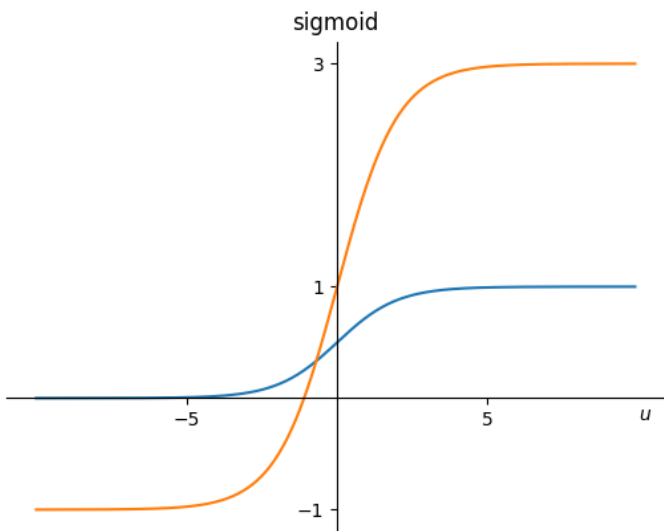
Output  $y \in [-0.74, 2.91] \subset [-1.0, 3.0]$

Note that the sigmoidal should have an amplitude = 4 and shifted downwards by 1.0.

So, the activation function should be

$$y = f(u) = \frac{4.0}{1 + e^{-u}} - 1.0$$

$$f'(u) = \frac{4e^{-u}}{(1+e^{-u})^2} = (y+1) \frac{e^{-u}}{(1+e^{-u})} = (y+1) \left(1 - \frac{1}{1+e^{-u}}\right) = \frac{1}{4} (y+1)(3-y)$$



## Example 2

Epoch 1 begins ...

$$\mathbf{u} = \mathbf{X}\mathbf{w} + b\mathbf{1}_p = \begin{pmatrix} 0.77 & 0.02 \\ 0.63 & 0.75 \\ 0.50 & 0.22 \\ 0.20 & 0.76 \\ 0.17 & 0.09 \\ 0.69 & 0.95 \\ 0.00 & 0.51 \end{pmatrix} \begin{pmatrix} 0.81 \\ 0.61 \end{pmatrix} + 0.0 \begin{pmatrix} 1.0 \\ 1.0 \\ 1.0 \\ 1.0 \\ 1.0 \\ 1.0 \\ 1.0 \end{pmatrix} = \begin{pmatrix} 0.64 \\ 0.97 \\ 0.54 \\ 0.63 \\ 0.19 \\ 1.14 \\ 0.32 \end{pmatrix}$$

$$\mathbf{y} = f(\mathbf{u}) = \frac{4.0}{1 + e^{-u}} - 1.0 = \begin{pmatrix} 1.61 \\ 1.90 \\ 1.53 \\ 1.61 \\ 1.19 \\ 2.03 \\ 1.31 \end{pmatrix}$$

$$m.s.e. = \frac{1}{7} \sum_{p=1}^7 (d_p - y_p)^2 = \frac{1}{7} ((2.91 - 1.61)^2 + (0.55 - 1.90)^2 + \dots) = 2.11$$

## Example 2

$$f'(\mathbf{u}) = \frac{1}{4}(\mathbf{y} + \mathbf{1}) \cdot (\mathbf{3} - \mathbf{y}) = \frac{1}{4} \left( \begin{pmatrix} 1.61 \\ 1.90 \\ 1.53 \\ 1.61 \\ 1.19 \\ 2.03 \\ 1.31 \end{pmatrix} + \begin{pmatrix} 1.0 \\ 1.0 \\ 1.0 \\ 1.0 \\ 1.0 \\ 1.0 \\ 1.0 \end{pmatrix} \right) \cdot \left( \begin{pmatrix} 3.0 \\ 3.0 \\ 3.0 \\ 3.0 \\ 3.0 \\ 3.0 \\ 3.0 \end{pmatrix} - \begin{pmatrix} 1.61 \\ 1.90 \\ 1.53 \\ 1.61 \\ 1.19 \\ 2.03 \\ 1.31 \end{pmatrix} \right) = \begin{pmatrix} 0.90 \\ 0.80 \\ 0.93 \\ 0.91 \\ 0.99 \\ 0.73 \\ 0.98 \end{pmatrix}$$

## Example 2

$$\mathbf{w} = \mathbf{w} + \alpha \mathbf{X}^T (\mathbf{d} - \mathbf{y}) \cdot f'(\mathbf{u})$$

$$= \begin{pmatrix} 0.81 \\ 0.61 \end{pmatrix} + 0.01 \begin{pmatrix} 0.77 & 0.63 & 0.50 & 0.20 & 0.17 & 0.69 & 0.00 \\ 0.02 & 0.75 & 0.22 & 0.76 & 0.09 & 0.95 & 0.51 \end{pmatrix} \left( \begin{pmatrix} 2.91 \\ 0.55 \\ 1.28 \\ -0.74 \\ 0.88 \\ 0.30 \\ -0.28 \end{pmatrix} - \begin{pmatrix} 1.61 \\ 1.90 \\ 1.53 \\ 1.61 \\ 1.19 \\ 2.03 \\ 1.31 \end{pmatrix} \right) \cdot \begin{pmatrix} 0.90 \\ 0.80 \\ 0.93 \\ 0.91 \\ 0.99 \\ 0.73 \\ 0.98 \end{pmatrix} = \begin{pmatrix} 0.80 \\ 0.57 \end{pmatrix}$$

$$b = b + \alpha \mathbf{1}_p^T (\mathbf{d} - \mathbf{y}) \cdot f'(\mathbf{u})$$

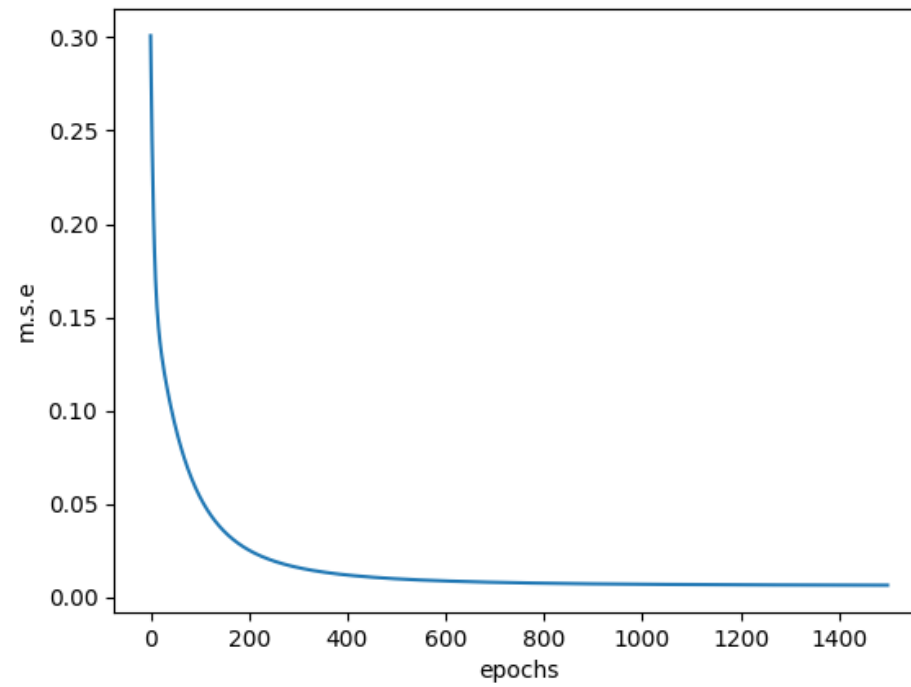
$$= 0.0 + 0.01 \times (1.0 \quad 1.0 \quad 1.0 \quad 1.0 \quad 1.0 \quad 1.0 \quad 1.0) \left( \begin{pmatrix} 2.91 \\ 0.55 \\ 1.28 \\ -0.74 \\ 0.88 \\ 0.30 \\ -0.28 \end{pmatrix} - \begin{pmatrix} 1.61 \\ 1.90 \\ 1.53 \\ 1.61 \\ 1.19 \\ 2.03 \\ 1.31 \end{pmatrix} \right) \cdot \begin{pmatrix} 0.90 \\ 0.80 \\ 0.93 \\ 0.91 \\ 0.99 \\ 0.73 \\ 0.98 \end{pmatrix} = -0.05$$



## Example 2

<i>iter</i>	<i>u</i>	<i>y</i>	<i>f'(u)</i>	<i>mse</i>	<i>w</i>	<i>b</i>
1	$\begin{pmatrix} 0.64 \\ 0.97 \\ 0.54 \\ 0.63 \\ 0.19 \\ 1.14 \\ 0.32 \end{pmatrix}$	$\begin{pmatrix} 1.61 \\ 1.90 \\ 1.53 \\ 1.61 \\ 1.19 \\ 2.03 \\ 1.31 \end{pmatrix}$	$\begin{pmatrix} 0.90 \\ 0.80 \\ 0.93 \\ 0.91 \\ 0.99 \\ 0.73 \\ 0.98 \end{pmatrix}$	2.11	$\begin{pmatrix} 0.80 \\ 0.57 \end{pmatrix}$	-0.05
2	$\begin{pmatrix} 0.57 \\ 0.88 \\ 0.47 \\ 0.54 \\ 1.04 \\ 1.95 \\ 0.24 \end{pmatrix}$	$\begin{pmatrix} 1.56 \\ 1.83 \\ 1.46 \\ 1.52 \\ 1.13 \\ 1.95 \\ 1.24 \end{pmatrix}$	$\begin{pmatrix} 0.92 \\ 0.83 \\ 0.95 \\ 0.93 \\ 1.00 \\ 0.77 \\ 0.99 \end{pmatrix}$	1.96	$\begin{pmatrix} 0.79 \\ 0.52 \end{pmatrix}$	-0.11
1500	$\begin{pmatrix} 2.05 \\ -0.45 \\ 0.56 \\ -1.94 \\ -0.16 \\ -0.85 \\ -1.89 \end{pmatrix}$	$\begin{pmatrix} 2.54 \\ 0.56 \\ 1.55 \\ -0.50 \\ 0.84 \\ 0.20 \\ -0.48 \end{pmatrix}$	$\begin{pmatrix} 0.40 \\ 0.95 \\ 0.92 \\ 0.44 \\ 0.99 \\ 0.84 \\ 0.45 \end{pmatrix}$	0.046	$\begin{pmatrix} 3.35 \\ -2.80 \end{pmatrix}$	-0.47

## Example 2



## Example 2

At convergence:

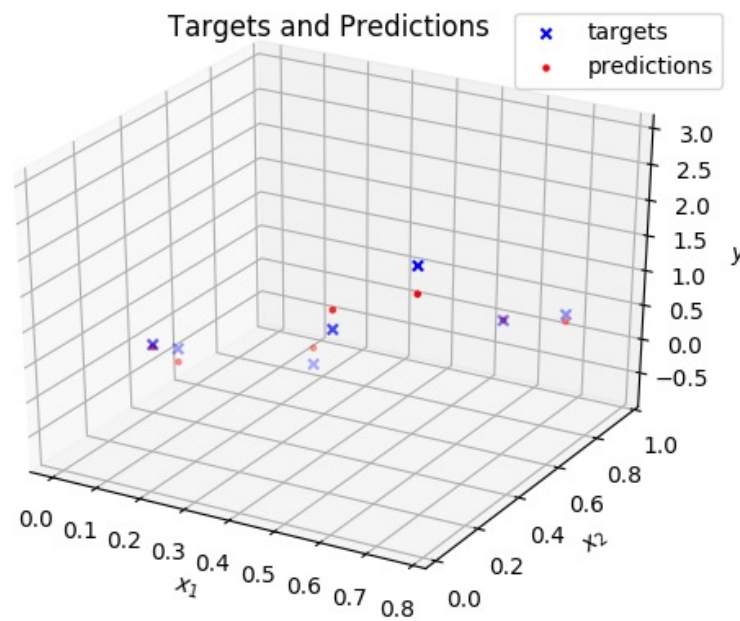
$$\mathbf{w} = \begin{pmatrix} 3.35 \\ -2.80 \end{pmatrix}$$
$$b = -0.47$$

Mean square error = 0.01

$$u = \mathbf{x}^T \mathbf{w} + b = (x_1 \quad x_2) \begin{pmatrix} 3.35 \\ -2.80 \end{pmatrix} - 0.47 = 3.35x_1 - 2.8x_2 - 0.47$$

$$y = \frac{4.0}{1 + e^{-u}} - 1.0$$
$$y = \frac{4.0}{1 + e^{-3.35x_1 + 2.8x_2 + 0.47}} - 1.0$$

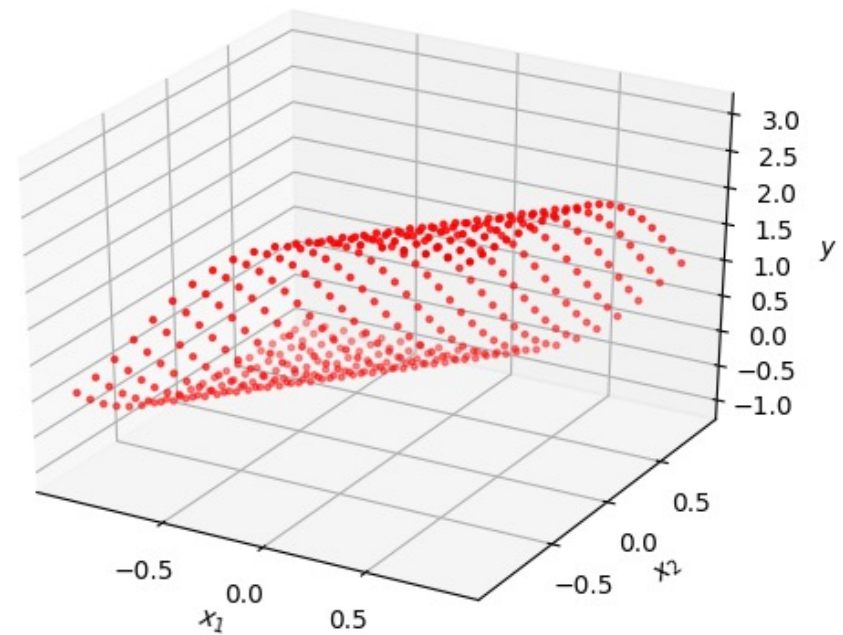
## Example 2



## Example 2

Non-linear function learnt by the perceptron

$$y = \frac{4.0}{1 + e^{-3.35x_1 + 2.8x_2 + 0.47}} - 1.0$$



# Classification Example

Classification is to identify or distinguish classes or groups of objects.

*Example:* To identify *ballet dancers* from *rugby players*.

Two distinctive *features* that can aid in classification:

- *weight*
- *height*

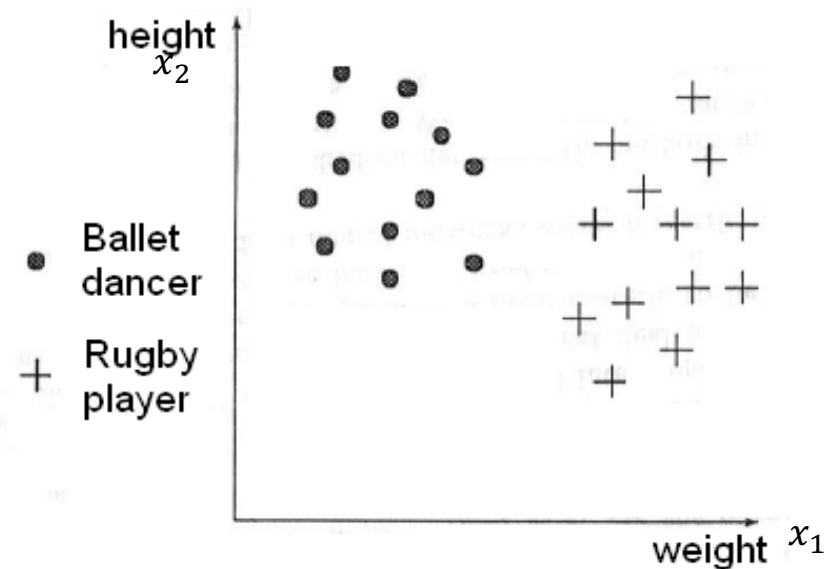


Figure: A 2-dimensional feature space

Let  $x_1$  denote weight and  $x_2$  denote height. Every individual is represented as a point  $\mathbf{x} = (x_1, x_2)$  in the feature feature space.

## Classification Example

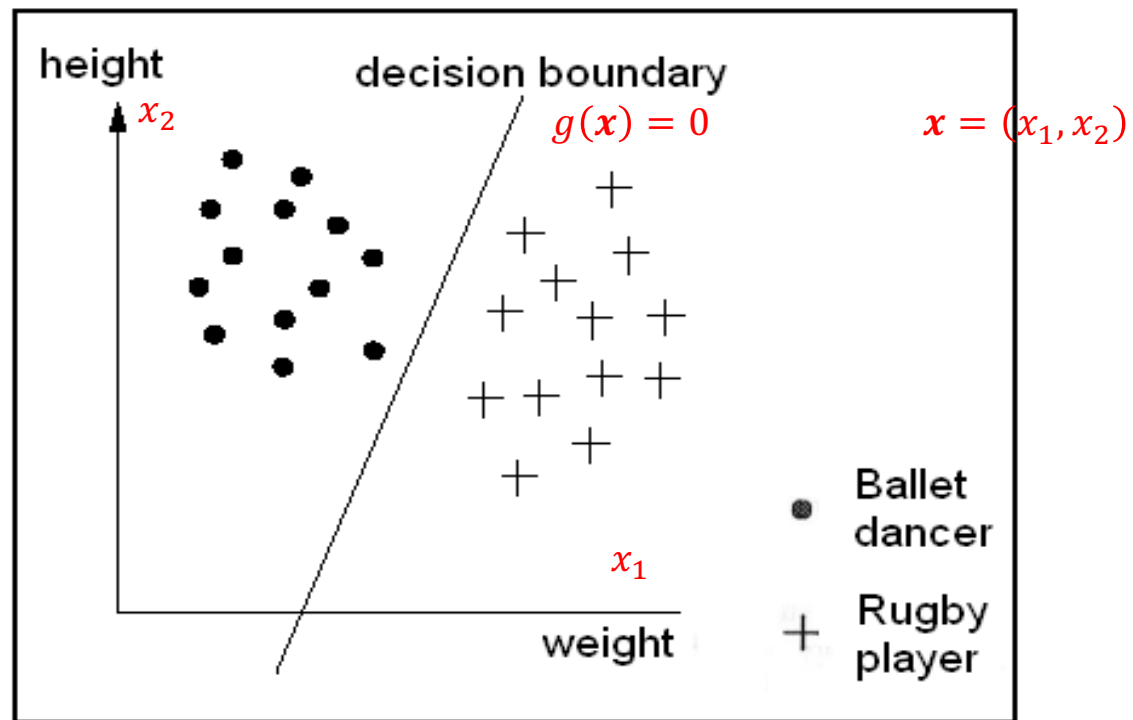


Figure: A linear classification decision boundary

# Decision boundary

The **decision boundary** of the classifier,  $g(\mathbf{x}) = 0$  where function  $g(\mathbf{x})$  is referred to as the **discriminant function**.

A classifier finds a decision boundary separating the two classes in the feature space. On one side of the decision boundary, discriminant function is positive and on other side, discriminant function is negative.

Therefore, the following class definition may be employed:

If  $g(\mathbf{x}) > 0 \Rightarrow$  Ballet dancer

If  $g(\mathbf{x}) \leq 0 \Rightarrow$  Rugby player



# Linear Classifier

If the two classes can be separated by a straight line, the classification is said to be **linearly separable**. For linear separable classes, one can design a **linear classifier**.

A linear classifier implements discriminant function or a decision boundary that is represented by a straight line (hyper plane) in the multidimensional **feature space**. Generally, the feature space is multidimensional. In the multidimensional space, a straight line or **hyperplane** is indicated by a linear sum of coordinates.

Given an input (features),  $\mathbf{x} = (x_1 \ x_2 \ \cdots \ x_n)^T$ . A linear description function is given by

$$g(\mathbf{x}) = w_0 + w_1x_1 + w_2x_2 + \cdots + w_nx_n$$

where  $\mathbf{w} = (w_1 \ w_2 \ \cdots \ w_n)^T$  are the coefficient/weights and  $w_0$  is the constant term.

# Discrete perceptron as a linear two-class classifier

The linear discriminant function can be implemented by the synaptic input to a neuron

$$g(\mathbf{x}) = u = \mathbf{w}^T \mathbf{x} + b$$

And with a threshold activation function  $1(u)$ :

$$u = g(\mathbf{x}) > 0 \rightarrow y = 1 \rightarrow \text{class1}$$

$$u = g(\mathbf{x}) \leq 0 \rightarrow y = 0 \rightarrow \text{class2}$$

That is, two-class linear classifier (or a dichotomizer) can be implemented with an artificial neuron with a threshold (unit step) activation function (**discrete perceptron**).

The output, 0 or 1, of the binary neuron represents the **label** of the class.

# Classification error

In classification, the error is expressed as total mismatches between the target and output labels.

$$\text{Classification error} = \sum_{p=1}^P 1(d_p \neq y_p)$$

# Logistic regression neuron

A **logistic regression neuron** performs a binary classification of inputs. That is, it classifies inputs into two classes with labels '0' and '1'.

The activation function of the logistic regression neuron is given by the sigmoid function:

$$f(u) = \frac{1}{1 + e^{-u}}$$

where  $u = \mathbf{w}^T \mathbf{x} + b$  is the synaptic input to the neuron.

The activation of a logistic regression neuron gives the probability of the neuron output belonging to class '1'.

$$P(y = 1|\mathbf{x}) = f(u)$$

Then

$$P(y = 0|\mathbf{x}) = 1 - P(y = 1|\mathbf{x}) = 1 - f(u)$$

# Logistic Regression Neuron

A logistic regression neuron receives an input  $\mathbf{x} \in \mathbf{R}^n$  and produces a class label  $y \in \{0, 1\}$  as the output.

$$f(u) = \frac{1}{1 + e^{-u}}$$

When  $u = 0$ ,  $f(u) = P(y = 1|\mathbf{x}) = P(y = 0|\mathbf{x}) = 0.5$ .  
That is,  $y = 1$  if  $f(u) > 0.5$ , else  $y = 0$ .

The output  $y$  of the neuron is given by:

$$y = 1(f(u) > 0.5) = 1(u > 0)$$

Note that for logistic neuron, the output and activation are different. It finds a linear boundary  $u = 0$  separating the two classes.

# SGD for logistic regression neuron

Given a training pattern  $(\mathbf{x}, d)$  where  $\mathbf{x} \in \mathbf{R}^n$  and  $d \in \{0,1\}$ .

The cost function for classification is given by the **cross-entropy**:

$$J = -d\log(f(u)) - (1 - d)\log(1 - f(u))$$

where  $u = \mathbf{w}^T \mathbf{x} + b$  and  $f(u) = \frac{1}{1+e^{-u}}$ .

The cost function  $J$  is minimized using the gradient descent procedure.

$$J = \begin{cases} -\log(f(u)) & \text{if } d = 1 \\ -\log(1 - f(u)) & \text{if } d = 0 \end{cases}$$

# SGD for logistic regression neuron

$$J = -d \log(f(u)) - (1 - d) \log(1 - f(u))$$

where  $u = \mathbf{w}^T \mathbf{x} + b$  and  $f(u) = \frac{1}{1 + e^{-u}}$ .

Gradient with respect to  $u$  :

$$\begin{aligned} \frac{\partial J}{\partial u} &= -\frac{\partial}{\partial f(u)} \left( d \log(f(u)) + (1 - d) \log(1 - f(u)) \right) \frac{\partial f(u)}{\partial u} \\ &= -\left( \frac{d}{f(u)} - \frac{(1 - d)}{1 - f(u)} \right) f'(u) \end{aligned}$$

Substituting  $f'(u) = f(u)(1 - f(u))$  for sigmoid activation function,

$$\frac{\partial J}{\partial u} = -\frac{d - f(u)}{f(u)(1 - f(u))} f(u)(1 - f(u)) = -(d - f(u))$$

## SGD for logistic regression neuron

Substituting  $\frac{\partial J}{\partial u}, \frac{\partial u}{\partial \mathbf{w}} = \mathbf{x}$ , and  $\frac{\partial u}{\partial b} = 1$ .

$$\nabla_{\mathbf{w}} J = \frac{\partial J}{\partial u} \frac{\partial u}{\partial \mathbf{w}} = -(d - f(u))\mathbf{x} \quad (\text{A})$$

$$\nabla_b J = \frac{\partial J}{\partial u} \frac{\partial u}{\partial b} = -(d - f(u)) \quad (\text{B})$$

Gradient learning equations:

$$\begin{aligned} \mathbf{w} &\leftarrow \mathbf{w} - \alpha \nabla_{\mathbf{w}} J \\ b &\leftarrow b - \alpha \nabla_b J \end{aligned}$$

Substituting  $\nabla_{\mathbf{w}} J$  and  $\nabla_b J$  for logistic regression neuron

$$\begin{aligned} \mathbf{w} &\leftarrow \mathbf{w} + \alpha (d - f(u))\mathbf{x} \\ b &\leftarrow b + \alpha (d - f(u)) \end{aligned}$$



# SGD for logistic regression neuron

Given a training dataset  $\{(\mathbf{x}_p, d_p)\}_{p=1}^P$

Set learning rate  $\alpha$

Initialize  $\mathbf{w}$  and  $b$

Iterate until convergence:

For every pattern  $(\mathbf{x}_p, d_p)$ :

$$u_p = \mathbf{w}^T \mathbf{x}_p + b$$

$$f(u_p) = \frac{1}{1+e^{-u_p}}$$

$$\mathbf{w} \leftarrow \mathbf{w} + \alpha (d_p - f(u_p)) \mathbf{x}_p$$

$$b \leftarrow b + \alpha (d_p - f(u_p))$$

## GD for logistic regression neuron

Given a training dataset  $\{(\mathbf{x}_p, d_p)\}_{p=1}^P$  where  $\mathbf{x}_p \in \mathbf{R}^n$  and  $d_p \in \{0,1\}$ .

The cost function for logistic regression is given by the *cross-entropy* (or *negative log-likelihood*) over all the training patterns:

$$J = - \sum_{p=1}^P d_p \log(f(u_p)) + (1 - d_p) \log(1 - f(u_p))$$

where  $u_p = \mathbf{w}^T \mathbf{x}_p + b$  and  $f(u_p) = \frac{1}{1+e^{-u_p}}$ .

The cost function  $J$  can be written as

$$J = \sum_{p=1}^P J_p$$

where  $J_p = -d_p \log(f(u_p)) - (1 - d_p) \log(1 - f(u_p))$  is cross-entropy due to  $p$  th pattern.

## GD for logistic regression neuron

$$\begin{aligned}\nabla_{\mathbf{w}} J &= \sum_{p=1}^P \nabla_{\mathbf{w}} J_p \\ &= - \sum_{p=1}^P (d_p - f(u_p)) \mathbf{x}_p && \text{From (A)} \\ &= - \left( (d_1 - f(u_1)) \mathbf{x}_1 + (d_2 - f(u_2)) \mathbf{x}_2 + \cdots + (d_P - f(u_P)) \mathbf{x}_P \right) \\ &= - (\mathbf{x}_1 \quad \mathbf{x}_2 \quad \cdots \quad \mathbf{x}_P) \begin{pmatrix} (d_1 - f(u_1)) \\ (d_2 - f(u_2)) \\ \vdots \\ (d_P - f(u_P)) \end{pmatrix} \\ &= -\mathbf{X}^T (\mathbf{d} - f(\mathbf{u}))\end{aligned}$$

$$\text{where } \mathbf{X} = \begin{pmatrix} \mathbf{x}_1^T \\ \mathbf{x}_2^T \\ \vdots \\ \mathbf{x}_P^T \end{pmatrix}, \mathbf{d} = \begin{pmatrix} d_1 \\ d_2 \\ \vdots \\ d_P \end{pmatrix}, \text{ and } f(\mathbf{u}) = \begin{pmatrix} f(u_1) \\ f(u_2) \\ \vdots \\ f(u_P) \end{pmatrix}$$

## GD for logistic regression neuron

By substituting  $\mathbf{1}_p$  for  $\mathbf{X}$  in above equation:

$$\nabla_b J = -\mathbf{1}_p^T (\mathbf{d} - f(\mathbf{u}))$$

Substituting the gradients in

$$\begin{aligned}\mathbf{w} &\leftarrow \mathbf{w} - \alpha \nabla_{\mathbf{w}} J \\ b &\leftarrow b - \alpha \nabla_b J\end{aligned}$$

the gradient descent learning for logistic regression neuron is given by

$$\begin{aligned}\mathbf{w} &\leftarrow \mathbf{w} + \alpha \mathbf{X}^T (\mathbf{d} - f(\mathbf{u})) \\ b &\leftarrow b + \alpha \mathbf{1}_p^T (\mathbf{d} - f(\mathbf{u}))\end{aligned}$$

Note that  $\mathbf{y}$  in the discrete perceptron is now replaced with  $f(\mathbf{u})$  in logistic regression learning equations.

# GD for logistic regression neuron

Given training data  $(\mathbf{X}, \mathbf{d})$

Set learning rate  $\alpha$

Initialize  $\mathbf{w}$  and  $b$

Iterate until convergence:

$$\mathbf{u} = \mathbf{X}\mathbf{w} + b\mathbf{1}_p$$

$$f(\mathbf{u}) = \frac{1}{1+e^{-u}}$$

$$\mathbf{w} \leftarrow \mathbf{w} + \alpha \mathbf{X}^T (\mathbf{d} - f(\mathbf{u}))$$

$$b \leftarrow b + \alpha \mathbf{1}_p^T (\mathbf{d} - f(\mathbf{u}))$$

# Learning for logistic regression neuron

GD	SGD
$(\mathbf{X}, \mathbf{d})$	$(\mathbf{x}_p, d_p)$
$J = - \sum_{p=1}^P d_p \log(f(u_p)) + (1 - d_p) \log(1 - f(u_p))$	$J_p = -d_p \log(f(u_p)) - (1 - d_p) \log(1 - f(u_p))$
$\mathbf{u} = \mathbf{X}\mathbf{w} + b\mathbf{1}_P$	$u_p = \mathbf{w}^T \mathbf{x}_p + b$
$f(\mathbf{u}) = \frac{1}{1 + e^{-\mathbf{u}}}$	$f(u_p) = \frac{1}{1 + e^{-u_p}}$
$\mathbf{y} = 1(f(\mathbf{u}) > 0.5)$	$y_p = 1(f(u_p) > 0.5)$
$\mathbf{w} \leftarrow \mathbf{w} + \alpha \mathbf{X}^T (\mathbf{d} - f(\mathbf{u}))$	$\mathbf{w} \leftarrow \mathbf{w} + \alpha (d_p - f(u_p)) \mathbf{x}_p$
$b \leftarrow b + \alpha \mathbf{1}_P^T (\mathbf{d} - f(\mathbf{u}))$	$b \leftarrow b + \alpha (d_p - f(u_p))$

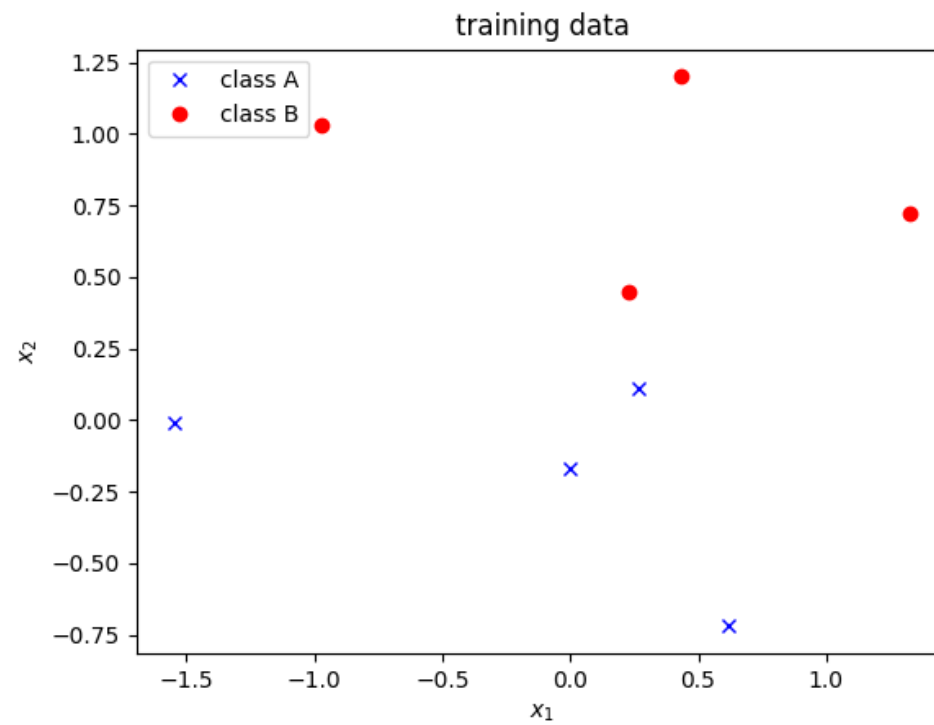
## Example 3: GD for logistic regression neuron

Train a logistic regression neuron to perform the following classification, using GD:

$(1.33 \quad 0.72) \rightarrow \text{class } B$   
 $(-1.55 \quad -0.01) \rightarrow \text{class } A$   
 $(0.62 \quad -0.72) \rightarrow \text{class } A$   
 $(0.27 \quad 0.11) \rightarrow \text{class } A$   
 $(0.0 \quad -0.17) \rightarrow \text{class } A$   
 $(0.43 \quad 1.2) \rightarrow \text{class } B$   
 $(-0.97 \quad 1.03) \rightarrow \text{class } B$   
 $(0.23 \quad 0.45) \rightarrow \text{class } B$

User a learning factor  $\alpha = 0.04$ .

## Example 3





## Example 3

Let  $y = 1$  for class A and  $y = 0$  for class B.

$$\mathbf{X} = \begin{pmatrix} 1.33 & 0.72 \\ -1.55 & -0.01 \\ 0.62 & -0.72 \\ 0.27 & 0.11 \\ 0.0 & -0.17 \\ 0.43 & 1.2 \\ -0.97 & 1.03 \\ 0.23 & 0.45 \end{pmatrix} \text{ and } \mathbf{d} = \begin{pmatrix} 0 \\ 1 \\ 1 \\ 1 \\ 1 \\ 0 \\ 0 \\ 0 \end{pmatrix}$$

Initially,  $w = \begin{pmatrix} 0.77 \\ 0.02 \end{pmatrix}$ ,  $b = 0.0$  and  $\alpha = 0.4$

## Example 3

Epoch 1:

$$\mathbf{u} = \mathbf{X}\mathbf{w} + b\mathbf{1} = \begin{pmatrix} 1.33 & 0.72 \\ -1.55 & 0.01 \\ 0.62 & -0.72 \\ 0.27 & 0.11 \\ 0.0 & -0.17 \\ 0.43 & 1.2 \\ -0.97 & 1.03 \\ 0.23 & 0.45 \end{pmatrix} \begin{pmatrix} 0.77 \\ 0.02 \end{pmatrix} + 0.0 \begin{pmatrix} 1 \\ 1 \\ 1 \\ 1 \\ 1 \\ 1 \\ 1 \\ 1 \end{pmatrix} = \begin{pmatrix} 1.04 \\ -1.2 \\ 0.46 \\ 0.21 \\ 0.00 \\ 0.36 \\ -0.73 \\ 0.19 \end{pmatrix}$$

$$f(\mathbf{u}) = \frac{1}{1 + e^{(-\mathbf{u})}} = \begin{pmatrix} 0.74 \\ 0.23 \\ 0.61 \\ 0.55 \\ 0.5 \\ 0.59 \\ 0.33 \\ 0.55 \end{pmatrix}$$

## Example 3

$$\mathbf{y} = 1(f(\mathbf{u}) > 0.5) = 1 \left( \begin{pmatrix} 0.74 \\ 0.23 \\ 0.61 \\ 0.55 \\ 0.5 \\ 0.59 \\ 0.33 \\ 0.55 \end{pmatrix} > 0.5 \right) = \begin{pmatrix} 1 \\ 0 \\ 1 \\ 1 \\ 0 \\ 1 \\ 0 \\ 1 \end{pmatrix} \quad \mathbf{d} = \begin{pmatrix} 0 \\ 1 \\ 1 \\ 1 \\ 1 \\ 0 \\ 0 \\ 0 \end{pmatrix}$$

$$\text{Classification error} = \sum_{p=1}^8 1(d_p \neq y_p) = 5$$

$$\begin{aligned} \text{Cross-entropy} &= -\sum_{p=1}^P d_p \log(f(u_p)) + (1 - d_p) \log(1 - f(u_p)) \\ &= -\log(1 - f(u_1)) - \log f(u_2) - \log f(u_3) - \dots - \log(1 - f(u_8)) \\ &= -\log(1 - 0.74) - \log 0.23 - \log 0.61 - \dots - \log(1 - 0.55) \\ &= 6.653 \end{aligned}$$

## Example 3

$$\mathbf{w} = \mathbf{w} + \alpha \mathbf{X}^T (\mathbf{d} - f(\mathbf{u}))$$

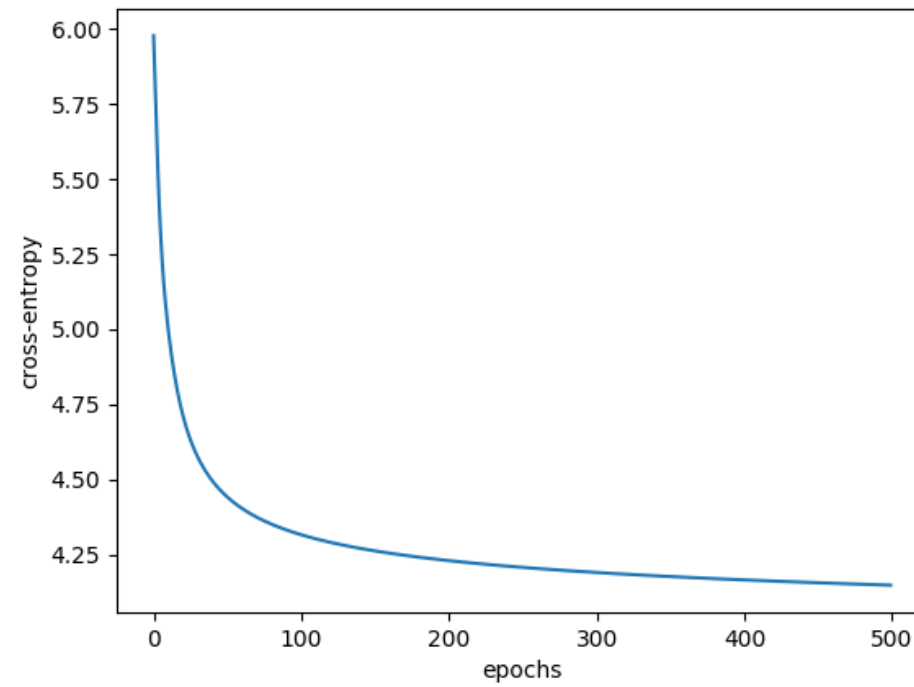
$$= \begin{pmatrix} 0.77 \\ 0.02 \end{pmatrix} + 0.04 \begin{pmatrix} 1.33 & -1.55 & 0.62 & 0.27 & 0 & 0.43 & -0.97 & 0.23 \\ 0.72 & -0.01 & -0.72 & 0.11 & -0.17 & 1.2 & 1.03 & 0.45 \end{pmatrix} \begin{pmatrix} 0 \\ 1 \\ 1 \\ 1 \\ 1 \\ 0 \\ 0 \\ 0 \end{pmatrix} - \begin{pmatrix} 0.74 \\ 0.23 \\ 0.61 \\ 0.55 \\ 0.5 \\ 0.59 \\ 0.33 \\ 0.55 \end{pmatrix}$$

$$= \begin{pmatrix} 0.69 \\ -0.2 \end{pmatrix}$$

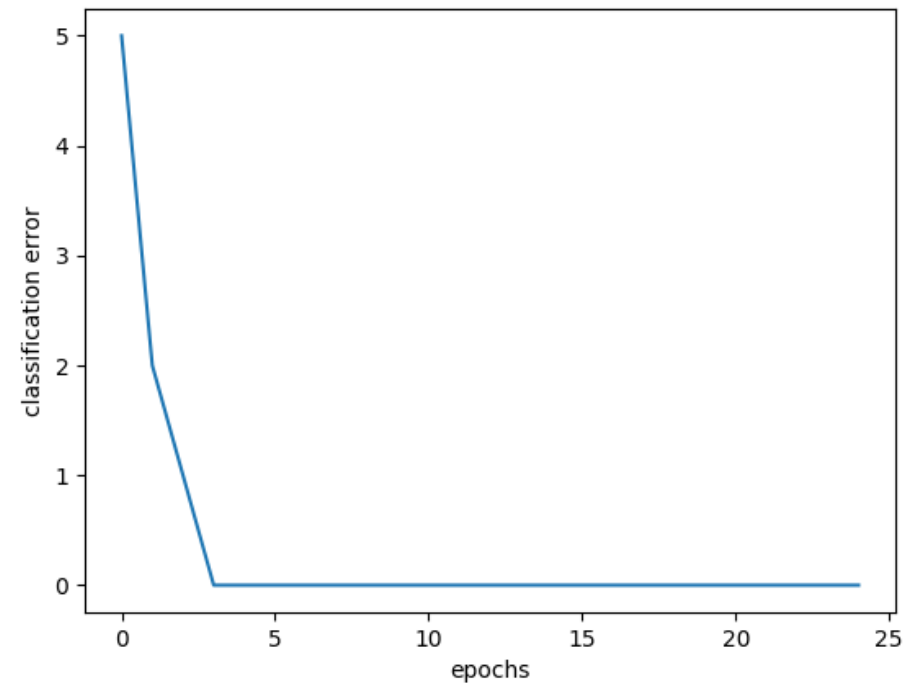
$$b = b + \alpha \mathbf{1}_p^T (\mathbf{d} - f(\mathbf{u}))$$

$$= 0.0 + 0.04 (1 \quad 1 \quad 1 \quad 1 \quad 1 \quad 1 \quad 1 \quad 1) \begin{pmatrix} 0 \\ 1 \\ 1 \\ 1 \\ 1 \\ 0 \\ 0 \\ 0 \end{pmatrix} - \begin{pmatrix} 0.74 \\ 0.23 \\ 0.61 \\ 0.55 \\ 0.5 \\ 0.59 \\ 0.33 \\ 0.55 \end{pmatrix} = -0.09$$

## Example 3



## Example 3



## Example 3

At convergence,  $\mathbf{w} = \begin{pmatrix} -1.20 \\ -15.02 \end{pmatrix}$ ,  $b = 4.47$

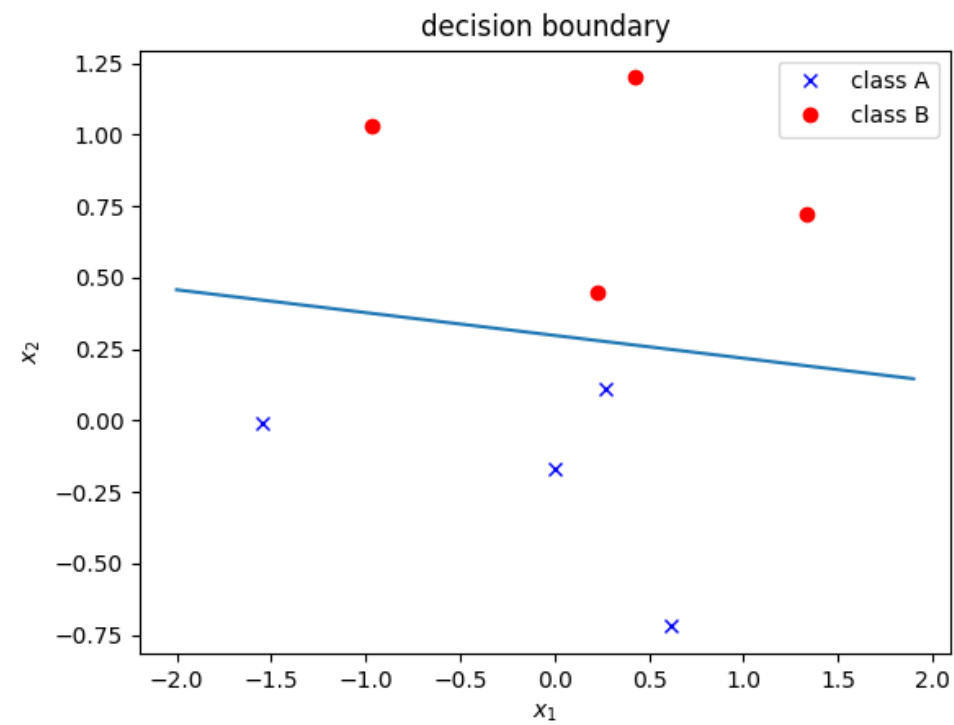
The decision boundary is given by:  $u = \mathbf{x}^T \mathbf{w} + b = 0$

$$(x_1 \ x_2)^T \begin{pmatrix} -1.20 \\ -15.02 \end{pmatrix} + 4.47 = 0$$

$$-1.20x_1 - 15.02x_2 + 4.47 = 0$$

Decision boundary:

$$-1.20x_1 - 15.02x_2 + 4.47 = 0$$



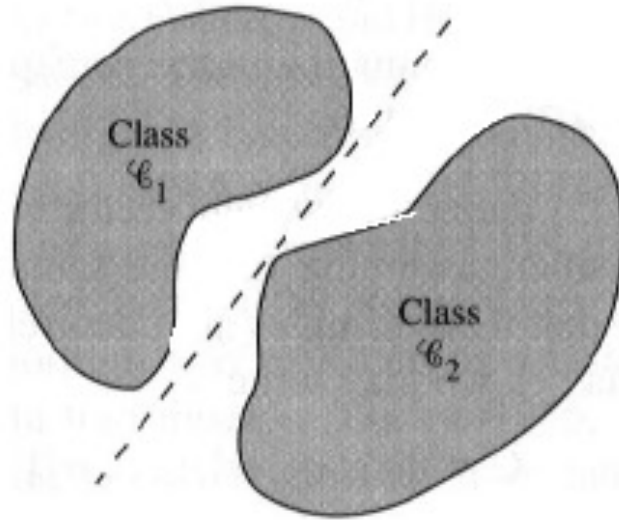


# Limitations of logistic regression neuron

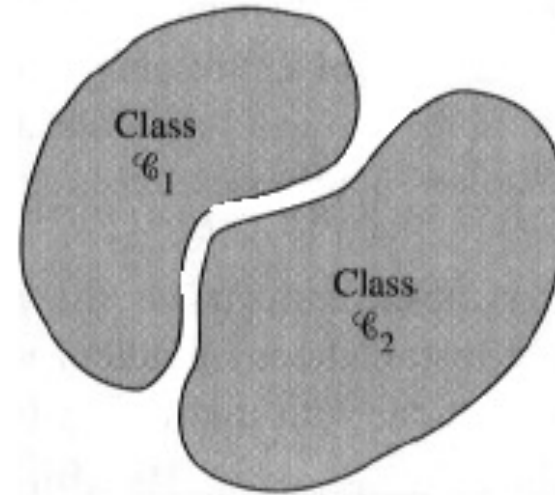
As long as the neuron is a *linear combiner* followed by a *non-linear activation function*, then regardless of the form of non-linearity used, the neuron can perform pattern classification *only* on *linearly separable* patterns.

*Linear separability* requires that the patterns to be classified must be sufficiently separated from each other to ensure that the decision boundaries are hyperplanes.

# Limitations of logistic regression neuron



(a) Linearly  
Separable Pattern



(b) Non-Linearly  
Separable Pattern

Discrete perceptron and logistic regression neuron can create only linear decision boundaries.

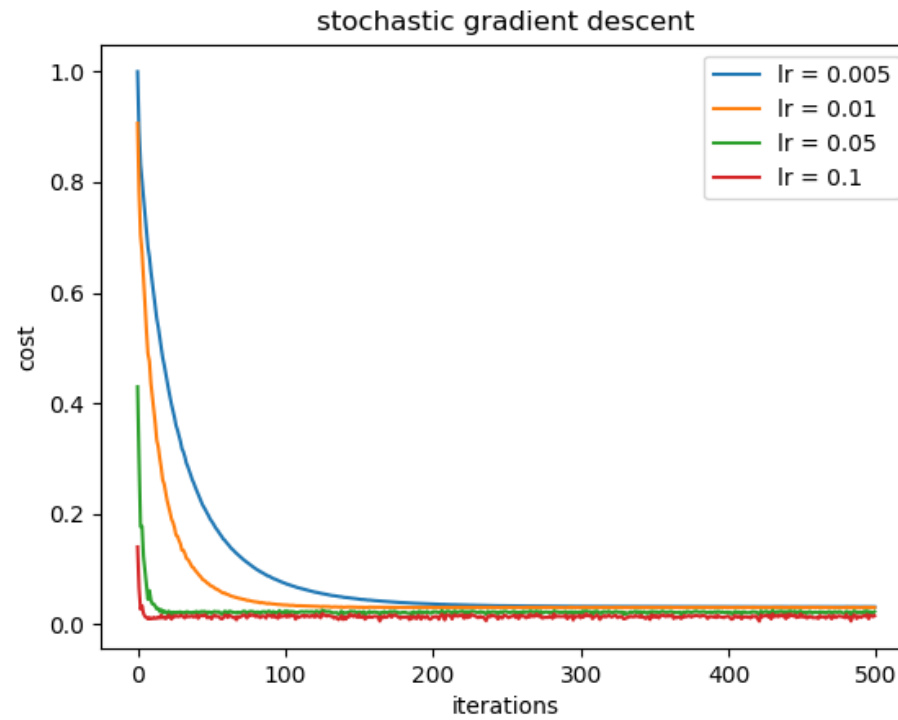
## Example 4: effects of leaning rate

Design a perceptron to learn the following mapping by using gradient descent (GD):

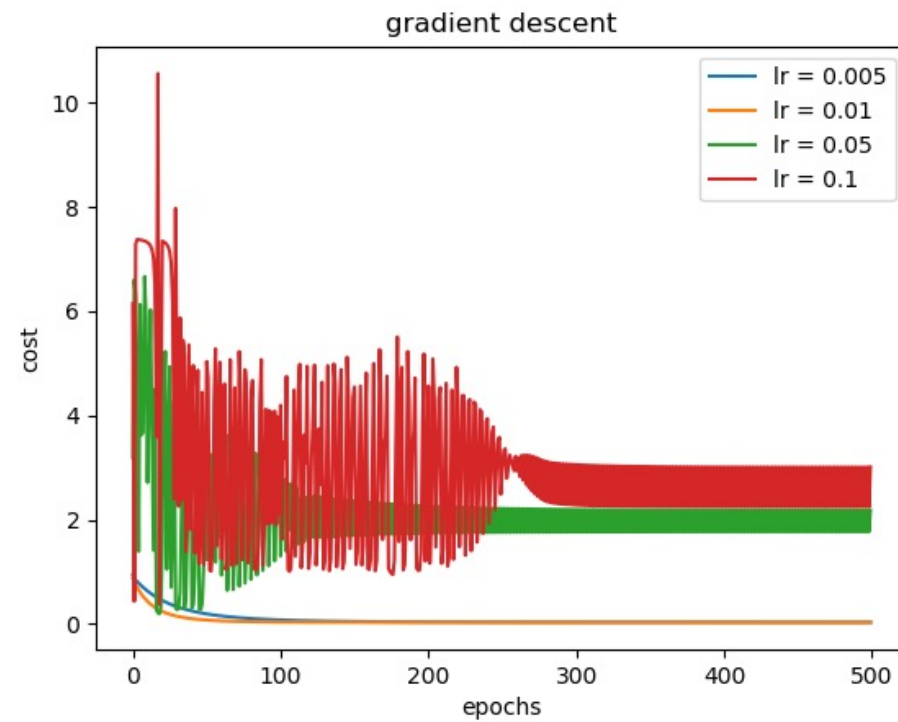
$x = (x_1, x_2)$	$d$
(0.77, 0.02)	2.91
(0.63, 0.75)	0.55
(0.50, 0.22)	1.28
(0.20, 0.76)	-0.74
(0.17, 0.09)	0.88
(0.69, 0.95)	0.30
(0.00, 0.51)	-0.28

Use learning factor  $\alpha = 0.01$ .

## Example 4: Learning rates with SGD



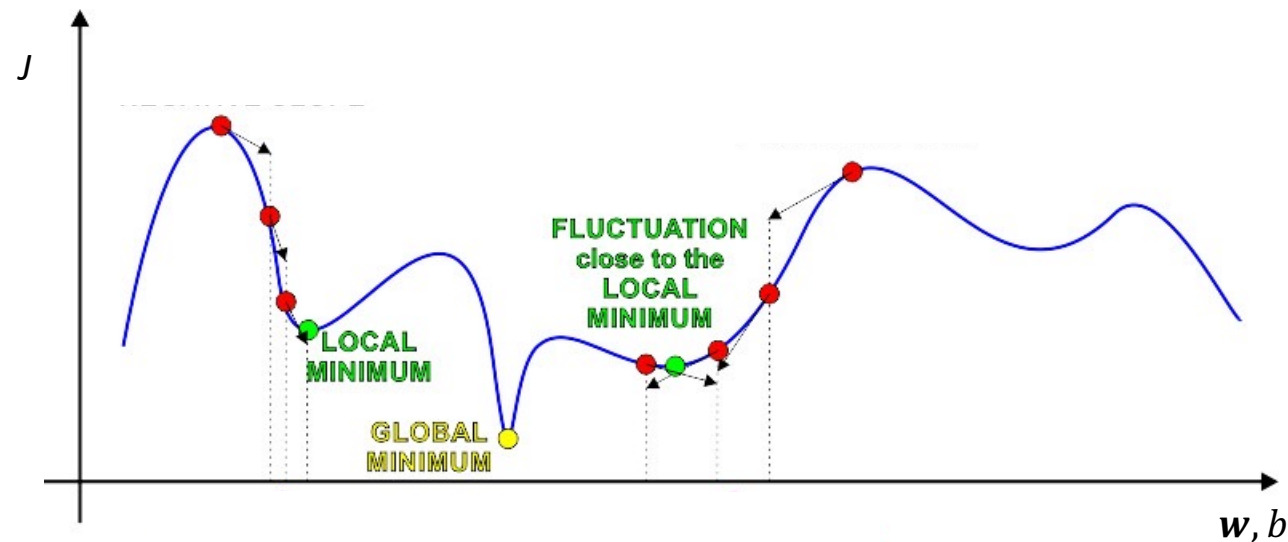
## Example 4: Learning rates with GD



# Learning rates

- At higher learning rates, convergence is faster but may not be stable.
- The *optimal learning rate* is the largest rate at which learning does not diverge.
- Generally, SGD convergence to a better solution (lower error) as it capitalizes on randomness of data. SGD takes a longer time to converge
- Usually, GD can use a higher learning rate compared to SGD; The time for one add/multiply computation per a weight update is less when patterns are trained in a (small) batch.
- In practice, *mini-batch SGD* is used. Then, the time to train a network is dependent upon
  - the learning rate
  - the batch size

# Local minima problem in gradient descent learning



Algorithm may get stuck in a local minimum of error function depending on the initial weights. Gradient descent gives a suboptimal solution and does not guarantee the optimal solution.

## Summary: types of neurons

Role	Neuron
Regression (one dimensional)	Linear neuron
	Perceptron
Classification (two classes)	Logistic regression neuron



## Summary: GD for neurons

$$\begin{aligned}
 & (X, d) \\
 & \mathbf{u} = X\mathbf{w} + b\mathbf{1}_P \\
 & \mathbf{w} = \mathbf{w} - \alpha X^T \nabla_{\mathbf{u}} J \\
 & b = b - \alpha \mathbf{1}_P^T \nabla_{\mathbf{u}} J
 \end{aligned}$$

neuron	$f(\mathbf{u}), y$	$\nabla_{\mathbf{u}} J$
Logistic regression neuron	$f(\mathbf{u}) = \frac{1}{1 + e^{-u}}$ $y = 1(f(\mathbf{u}) > 0.5)$	$-(d - f(\mathbf{u}))$
Linear neuron	$y = u$	$-(d - y)$
Perceptron	$y = f(\mathbf{u}) = \frac{1}{1 + e^{-u}}$	$-(d - y) \cdot f'(\mathbf{u})$

## Summary: SGD for neurons

$$\begin{aligned}
 &(\mathbf{x}_p, d_p) \\
 &u_p = \mathbf{w}^T \mathbf{x}_p + b \\
 &\mathbf{w} = \mathbf{w} - \alpha \nabla_{u_p} J \mathbf{x}_p \\
 &b = b - \alpha \nabla_{u_p} J
 \end{aligned}$$

neuron	$f(u_p), y_p$	$\nabla_{u_p} J$
Logistic regression neuron	$f(u_p) = \frac{1}{1 + e^{-u_p}}$ $y_p = 1(f(u_p) > 0.5)$	$-(d_p - f(u_p))$
Linear neuron	$y_p = u_p$	$-(d_p - y_p)$
Perceptron	$y_p = f(u_p) = \frac{1}{1 + e^{-u_p}}$	$-(d_p - y_p) \cdot f'(u_p)$