

Chapter 3

Classification

Neural networks and deep learning

Classification Example

Classification is to distinguish *ballet dancers* from *rugby players*.

Two distinctive *features* that can aid in classification:

- ***weight***
- ***height***

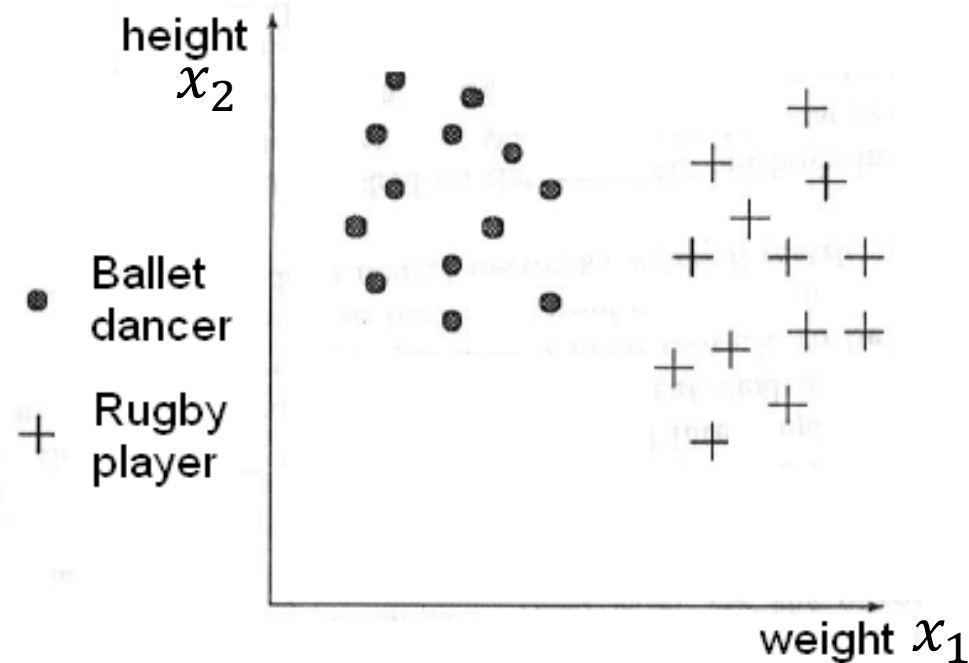


Figure: A 2-dimensional feature space

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

The figure shows the distribution of height and weight measurements on each individual in the feature space.

Classification Example

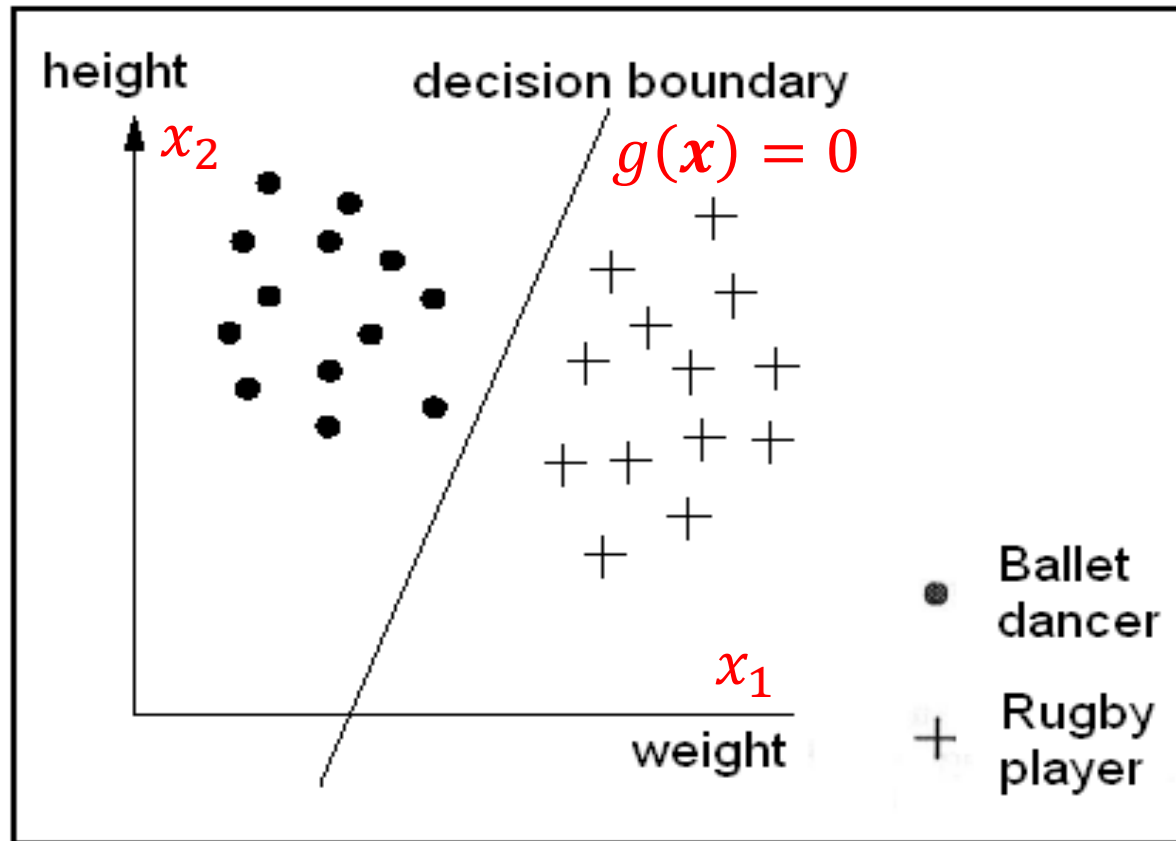


Figure: A linear classification decision boundary

Discriminant Function

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

In the example, the decision boundary separates the two classes. 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

The objective of classification is to determine a decision boundary that would separate the two clusters, i.e., rugby players and ballet dances.

Humans can easily decide some line or curve drawn between the two classes to separate them (See the previous Figure). But how can such a decision boundary or discrimination function g between them be computed ?

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**. Two class linear classifier can be designed with an artificial neuron with a threshold activation function

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 (hyper plane) is indicated by a linear sum of coordinates.

Given an input (features), $\mathbf{x} = (x_1 \ x_2 \ \cdots \ x_n)^T$. The discriminant function of a **linear classifier** is given by a linear function of features:

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

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

Linear classifier with an artificial neuron

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

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

And with a threshold activation function:

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

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

That is, two-class linear classifier can be implemented with an artificial neuron with a threshold (unit step) activation function. The output of the binary neuron represents the **label** of the class.

Linear Classifier

Consider two feature inputs x_1 and x_2 to an artificial neuron. The discriminant boundary in 2-dimensions can be written as a straight line:

$$g(\mathbf{x}) = w_1x_1 + w_2x_2 + b = 0$$

Or $x_2 = -\frac{w_1}{w_2}x_1 - \frac{b}{w_2}.$

Comparing the equation to the Cartesian system of a straight line ($y = mx + c$), the slope of line is controlled by ratio of the weight values, w_1 and w_2 and intercept is controlled by the bias term b .

Note that the constant term is needed. Otherwise, the decision boundary always goes through the origin.

With the appropriate values for the weight vector learnt, we can indeed determine the position of the decision boundary.

Weight vector is normal to the decision boundary

At the decision boundary:

$$g(\mathbf{x}) = \mathbf{x}^T \mathbf{w} + b = 0$$

Consider two points, \mathbf{x}_1 and \mathbf{x}_2 , on the decision surface.

$$\text{Then } \mathbf{x}_1^T \mathbf{w} + b = \mathbf{x}_2^T \mathbf{w} + b = 0$$

and therefore

$$(\mathbf{x}_1 - \mathbf{x}_2)^T \mathbf{w} = 0$$

That is, the weight vector \mathbf{w} is normal to the decision boundary (given by the direction of $(\mathbf{x}_1 - \mathbf{x}_2)$).

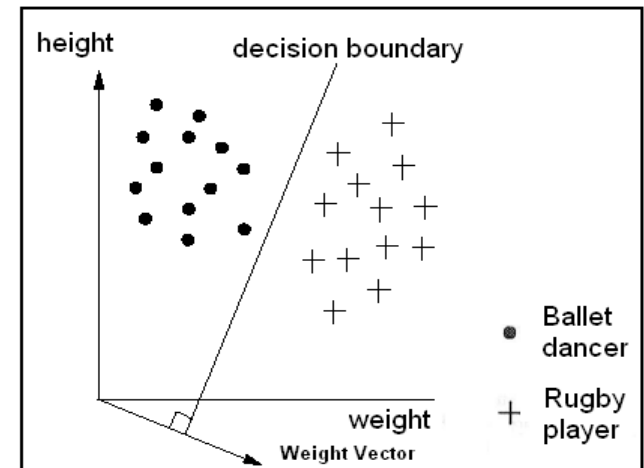
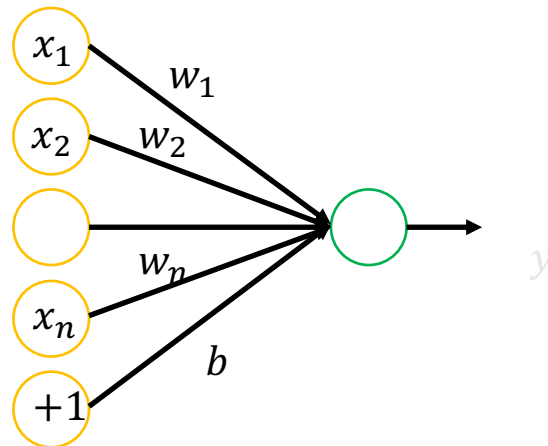


Figure: A linear classification decision boundary

Discrete Perceptron

Simple (or discrete) perceptron is a neuron that has a **threshold** or **unit-step** activation function.



Output:

$$y = f(u) = \begin{cases} 1, & u > 0 \\ 0, & u \leq 0 \end{cases}$$

where

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

Discrete perceptron classifies input patterns into two classes with a linear discriminant function. That is, discrete perceptron is capable of classifying two linearly separable classes (e.g., $y = 0$ for class 1 and $y = 1$ for class 2). It acts as a **two-class classifier** or **dichotomizer**.

Indicator function $1(\cdot)$

Indicator function $1(\cdot)$ takes value 1 when the condition given is True and value 0 when the condition is False.

$$1(x) = \begin{cases} 1, & x \text{ is True} \\ 0, & x \text{ is False} \end{cases}$$

Using the indicator function, the output of a discrete perceptron can be written as:

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

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

Discrete Perceptron Learning Algorithm

Discrete perceptron learning algorithm is a supervised scheme. It was proposed by Minsky in 1950 and its convergence can be proved. However, because of non-differentiable characteristics of the activation function, the discrete perceptron learning algorithm cannot be derived from a cost function.

Given P training pairs $\{(\mathbf{x}_p, d_p)\}_{p=1}^P$
where $\mathbf{x}_p \in \mathbf{R}^n$ is the n -dimensional input and $d_p \in \{0, 1\}$ is the binary target (desired) output of p th training pattern.

Discrete perceptron learning algorithm finds **a linear decision boundary** in the feature space.

Discrete Perceptron Learning Algorithm

The change of weights is proportional to the difference (error) between the desired output d and perceptron output y .

For training pattern (\mathbf{x}, d) :

$$\begin{aligned}u &= \mathbf{x}^T \mathbf{w} + b \\y &= 1(u > 0) \\ \delta &= d - y\end{aligned}$$

Learning:

$$\begin{aligned}\mathbf{w} &\leftarrow \mathbf{w} + \alpha \delta \mathbf{x} \\ b &\leftarrow b + \alpha \delta\end{aligned}$$

Note that $\delta = \{-1, 0, 1\}$.

The learning rate $\alpha \in (0, 1]$

When $\alpha = 1.0$, learning equations are referred to as *simple perceptron rule*.

Discrete Perceptron Learning Algorithm

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

Set the learning parameter α

Initialize \mathbf{w} and b

Repeat until convergence:

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

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

$$y_p = 1(u_p > 0)$$

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

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

Discrete Perceptron Learning

Batch learning	Stochastic learning
(X, \mathbf{d})	(\mathbf{x}_p, d_p)
$\mathbf{u} = X\mathbf{w} + b\mathbf{1}_p$	$u_p = \mathbf{x}_p^T \mathbf{w} + b$
$\mathbf{y} = 1(\mathbf{u} > 0)$	$y_p = 1(u_p > 0)$
$\mathbf{w} \leftarrow \mathbf{w} + \alpha X^T(\mathbf{d} - \mathbf{y})$	$\mathbf{w} \leftarrow \mathbf{w} + \alpha(d_p - y_p)\mathbf{x}_p$
$b \leftarrow b + \alpha \mathbf{1}_p^T(\mathbf{d} - \mathbf{y})$	$b \leftarrow b + \alpha(d_p - y_p)$

The learning equations have the same form of linear neurons.
Note that \mathbf{d} and \mathbf{y} are vectors of binary numbers.

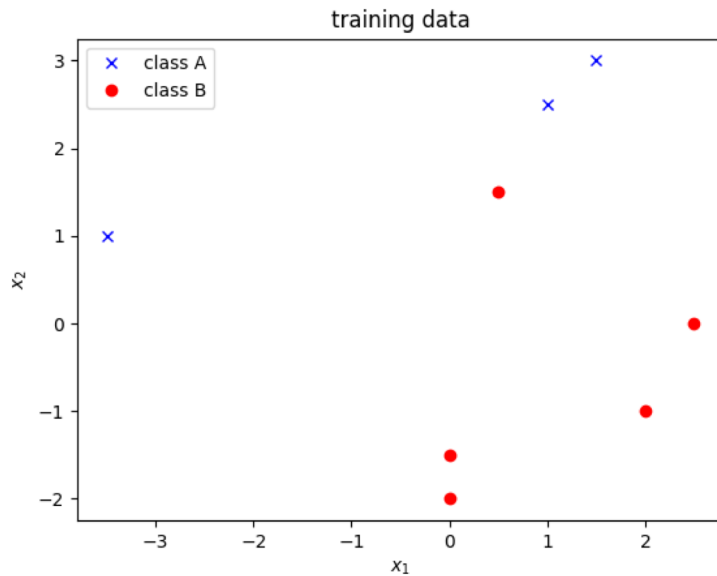
Example 1

Train a perceptron to classify the following 2-dimensional patterns, using the discrete perceptron learning algorithm:

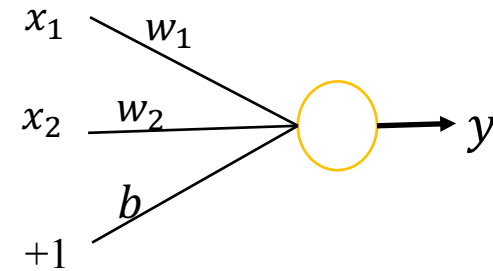
$(1.0 \ 2.5) \rightarrow \textit{class B}$
 $(2.0 \ -1.0) \rightarrow \textit{class A}$
 $(1.5 \ 3.0) \rightarrow \textit{class B}$
 $(0.0 \ -1.5) \rightarrow \textit{class A}$
 $(-3.5 \ 1.0) \rightarrow \textit{class B}$
 $(2.5 \ 0.0) \rightarrow \textit{class A}$
 $(0.5 \ 1.5) \rightarrow \textit{class A}$
 $(0.0 \ -2.0) \rightarrow \textit{class A}$

Show two iterations of SGD learning with a learning parameter $\alpha = 0.4$.

Example 1



Note that the two classes are linearly separable.



Let $y = 0$ for *class A*
 $y = 1$ for *class B*

Initially, $\mathbf{w} = \begin{pmatrix} 0.77 \\ 0.02 \end{pmatrix}$, $b = 0.0$

Learning rate $\alpha = 0.4$

Example 1

$$\mathbf{w} = \begin{pmatrix} 0.77 \\ 0.02 \end{pmatrix}, b = 0.0$$
$$y = 1(u > 0)$$

Iteration 1:

Shuffle the inputs

$$p = 1:$$

$$\mathbf{x}_1 = \begin{pmatrix} 1.5 \\ 3.0 \end{pmatrix}, d_1=1:$$

$$u_1 = \mathbf{x}_1^T \mathbf{w} + b = (1.5 \quad 3.0) \begin{pmatrix} 0.77 \\ 0.02 \end{pmatrix} + 0.0 = 1.22$$

$$y_1 = 1.0$$

$$\mathbf{w} = \mathbf{w} + \alpha(d_1 - y_1)\mathbf{x}_1 = \begin{pmatrix} 0.77 \\ 0.02 \end{pmatrix} + 0.4 \times (1 - 1) \begin{pmatrix} 1.5 \\ 3.0 \end{pmatrix} = \begin{pmatrix} 0.77 \\ 0.02 \end{pmatrix}$$

$$b = b + \alpha(d_1 - y_1) = 0.0 + 0.4 \times (1 - 1) = 0.0$$

Example 1

$p = 2$:

$$\mathbf{x}_2 = \begin{pmatrix} 0.0 \\ -2.0 \end{pmatrix}, d_2=0:$$

$$u_2 = \mathbf{x}_2^T \mathbf{w} + b = (0.0 \quad -2.0) \begin{pmatrix} 0.77 \\ 0.02 \end{pmatrix} + 0.0 = -0.04$$

$$y_2 = 0.0$$

$$\mathbf{w} = \mathbf{w} + \alpha(d_2 - y_2)\mathbf{x}_1 = \begin{pmatrix} 0.77 \\ 0.02 \end{pmatrix} + 0.4 \times (0 - 0) \begin{pmatrix} 0.0 \\ -2.0 \end{pmatrix} = \begin{pmatrix} 0.77 \\ 0.02 \end{pmatrix}$$

$$b = b + \alpha(d_1 - y_1) = 0.0 + 0.4 \times (0 - 0) = 0.0$$

Apply $\mathbf{x}_3, \mathbf{x}_4, \dots \mathbf{x}_8$, and update \mathbf{w} and b for each pattern.

Example 1

x	d	u	y	w^{new}	b^{new}
$\begin{pmatrix} 1.5 \\ 3.0 \end{pmatrix}$	1	1.22	1	$\begin{pmatrix} 0.77 \\ 0.02 \end{pmatrix}$	0.0
$\begin{pmatrix} 0.0 \\ -2.0 \end{pmatrix}$	0	-0.04	0	$\begin{pmatrix} 0.77 \\ 0.02 \end{pmatrix}$	0.0
$\begin{pmatrix} 2.5 \\ 0.0 \end{pmatrix}$	0	1.93	1	$\begin{pmatrix} -0.23 \\ 0.02 \end{pmatrix}$	-0.4
$\begin{pmatrix} 0.5 \\ 1.5 \end{pmatrix}$	0	-0.48	0	$\begin{pmatrix} 0.23 \\ 0.02 \end{pmatrix}$	-0.4
$\begin{pmatrix} -3.5 \\ 1.0 \end{pmatrix}$	1	0.43	1	$\begin{pmatrix} 0.23 \\ 0.02 \end{pmatrix}$	-0.4
$\begin{pmatrix} 0.0 \\ -1.5 \end{pmatrix}$	0	-0.43	0	$\begin{pmatrix} 0.23 \\ 0.02 \end{pmatrix}$	-0.4
$\begin{pmatrix} 0.5 \\ 1.5 \end{pmatrix}$	0	-0.88	0	$\begin{pmatrix} 0.23 \\ 0.02 \end{pmatrix}$	0.4
$\begin{pmatrix} 1.0 \\ 2.5 \end{pmatrix}$	1	-0.58	0	$\begin{pmatrix} 0.17 \\ 1.02 \end{pmatrix}$	0.0

$$\text{Classification error} = \sum_{p=1}^8 1(d_p \neq y_p) = 2$$

Example 1

Then epoch 2 begins,

Shuffle the inputs

apply $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_8$, in random order and for each input, update \mathbf{w} and b

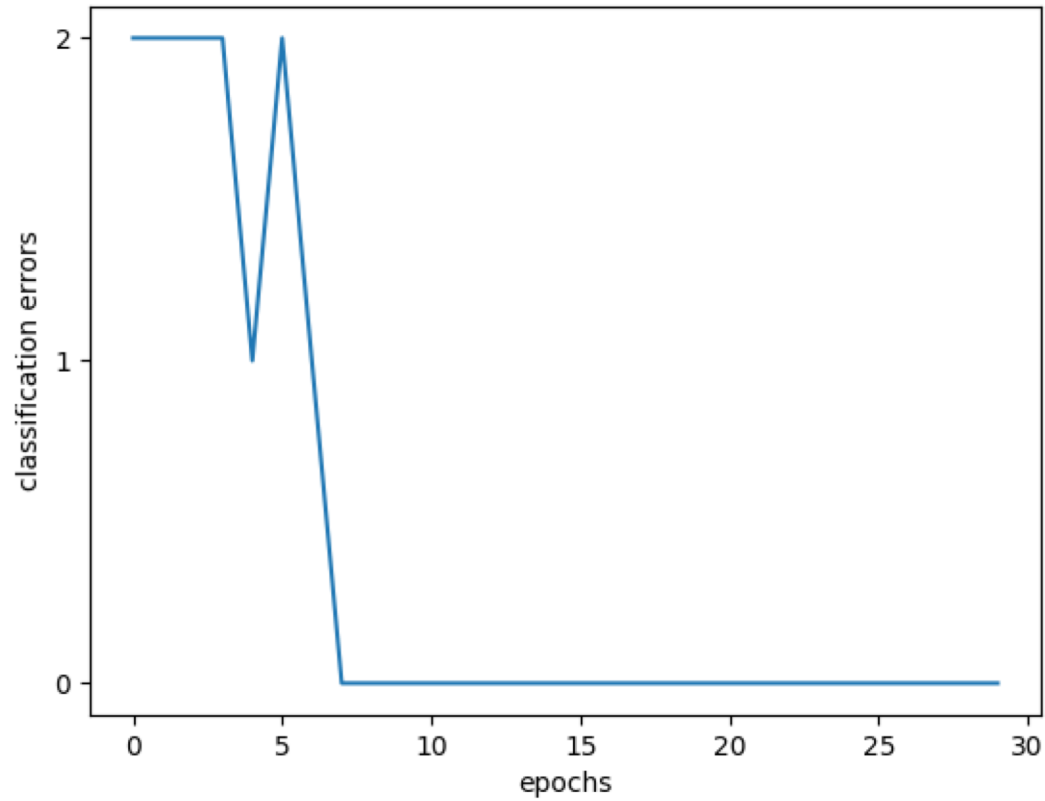
Then epoch 3 begins, apply $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_8$, in a random order and for each input, update \mathbf{w} and b

·
·
·

Iterations (epochs) continue until convergence is achieved.

Note that in each iteration, the patterns are shuffled randomly to change the presentation order.

Example 1

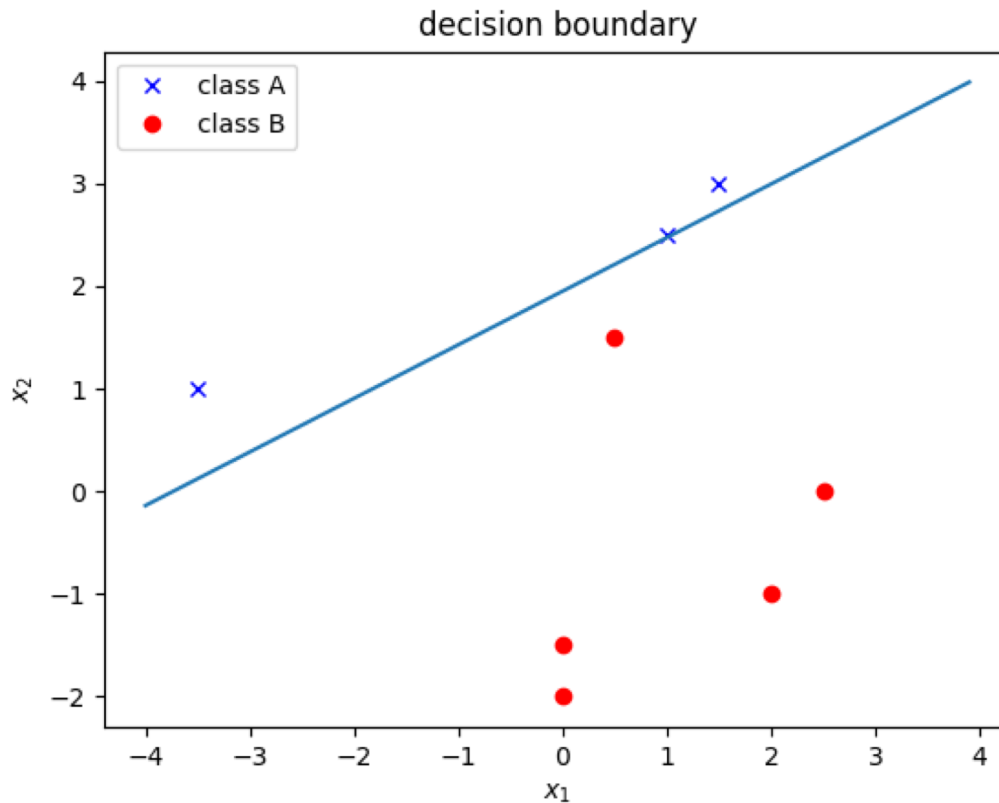


At convergence, $\mathbf{w} = \begin{pmatrix} -0.43 \\ 0.82 \end{pmatrix}$, $b = -1.6$

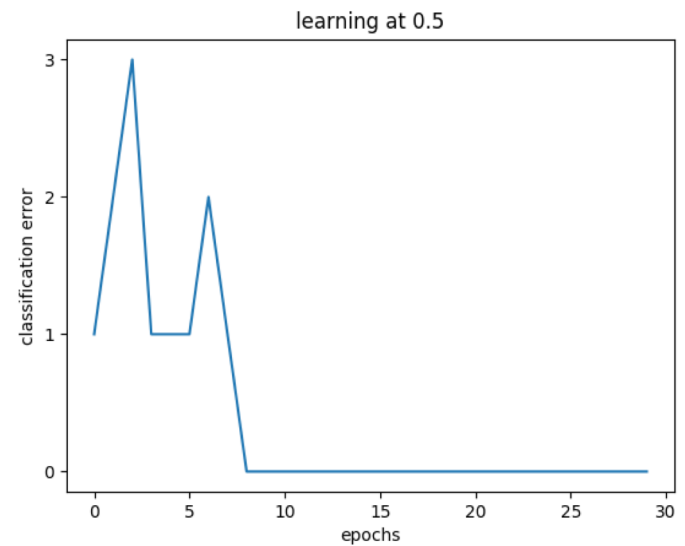
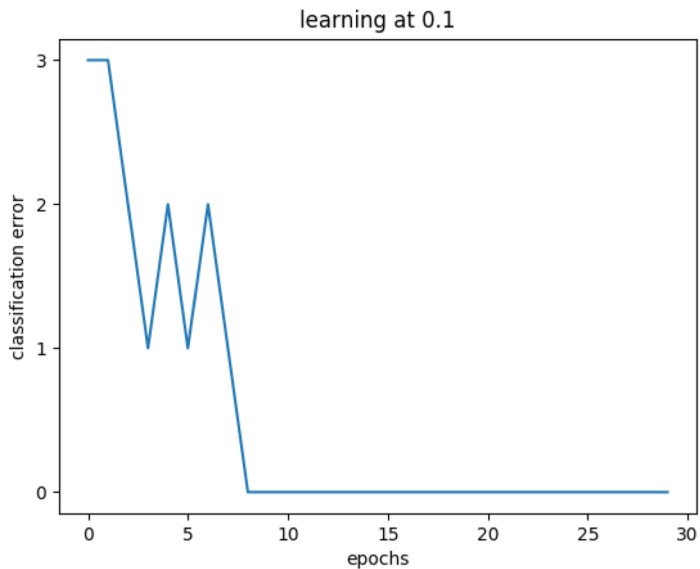
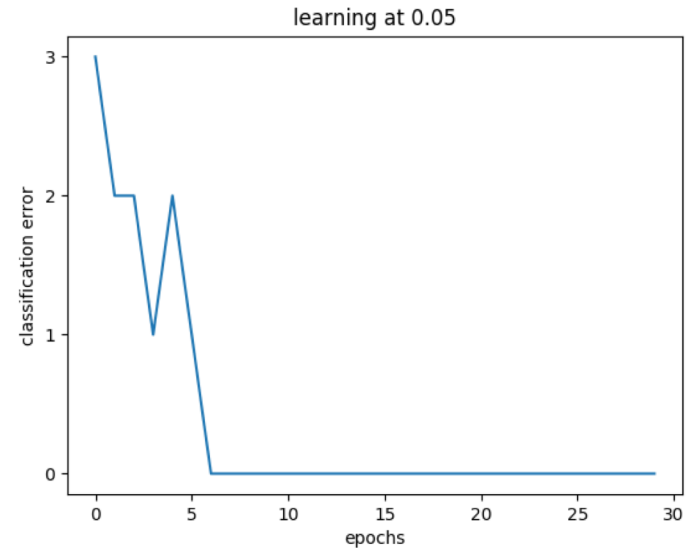
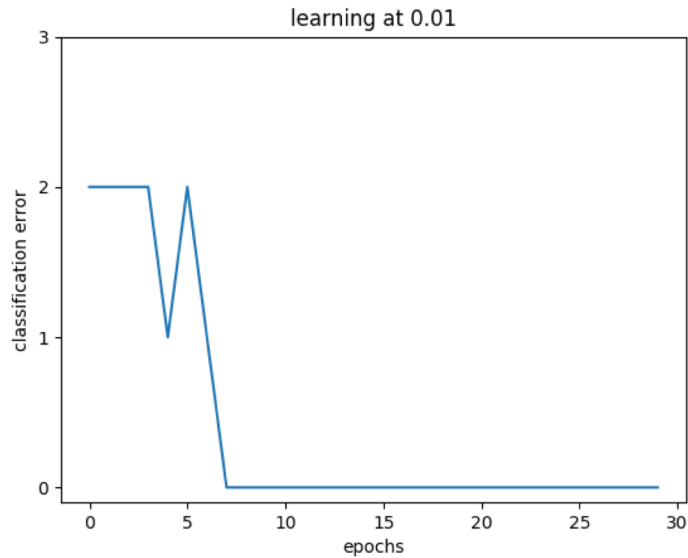
The decision boundary is given by: $\mathbf{x}^T \mathbf{w} + b = 0$

$$(x_1 \quad x_2) \begin{pmatrix} -0.43 \\ 0.82 \end{pmatrix} - 1.6 = 0$$

$$-0.43x_1 + 0.82x_2 - 1.6 = 0$$



Example 1: learning at different rates



Cross-entropy

Consider a model that produces data belonging to K classes with labels $\{1, 2, \dots, K\}$ at class probabilities p_1, p_2, \dots, p_K .

Assume n_1, n_2, \dots, n_K number of data points were observed from each class and data points are independent to one another.

The **likelihood** $p(\text{data} | \text{model})$ of data given by the model:

$$p(\text{data} | \text{model}) = p_1^{n_1} p_2^{n_2} \dots p_K^{n_K} = \prod_{k=1}^K p_k^{n_k}$$

The **negative log-likelihood** or **cross entropy** is given by:

$$-\log(p(\text{data} | \text{model})) = -\sum_{k=1}^K n_k \log p_k$$

Cross-entropy

The **negative log-likelihood** or **cross entropy** is given by:

$$-\log(p(\text{data}|\text{model})) = -\sum_{k=1}^K n_k \log p_k$$

The maximum likelihood principal states that the likelihood of data should be maximized by the model.

That is, the cross-entropy needs to be minimized.

Cross-entropy is often used as cost function for neural networks learning classification tasks.

Cross-entropy

Dividing cross-entropy by the total number $N = n_1 + n_2 + \dots + n_K$ of data points:

$$\text{cross-entropy} \propto - \sum_{k=1}^K \frac{n_k}{N} \log p_k = - \sum_{k=1}^K q_k \log p_k$$

where $q_k = \frac{n_k}{N}$ gives the probabilities of class k , given by data.

It can be shown that the cross-entropy is minimum when $p_k = q_k$ for all k .

In other words, the probabilities given by the model (i.e., classifier) is same as those implied by data when cross-entropy is minimum.

That is, cross-entropy can be used as the cost function of a classifier. The optimal weights and biases obtained by minimizing the cross entropy learns a classifier that fits the model of data.

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

Given an input \mathbf{x} , the activation of the neuron gives $P(y = 1|\mathbf{x})$:

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

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

That is, the activation function of the neuron is given by the sigmoid (or logistic function). Note that the output y of the neuron is not equal to the activation.

The activation of neuron gives the probability of the input belonging to class "1".

And

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

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.

The activation function give the probability that the output class label $y = 1$. That is,

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

The output y of the neuron is given by:

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

Decision boundary is given by $u = 0$.

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

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

GD for logistic regression neuron

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

where $u = \mathbf{x}^T \mathbf{w} + b$ 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 $\frac{\partial J}{\partial u}$:

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

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

Gradient learning equations:

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

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

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

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

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

SGD for logistic regression neuron

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

Set learning rate α

Initialize \mathbf{w} and b

Iterate until convergence:

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

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

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

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

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

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

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 α

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

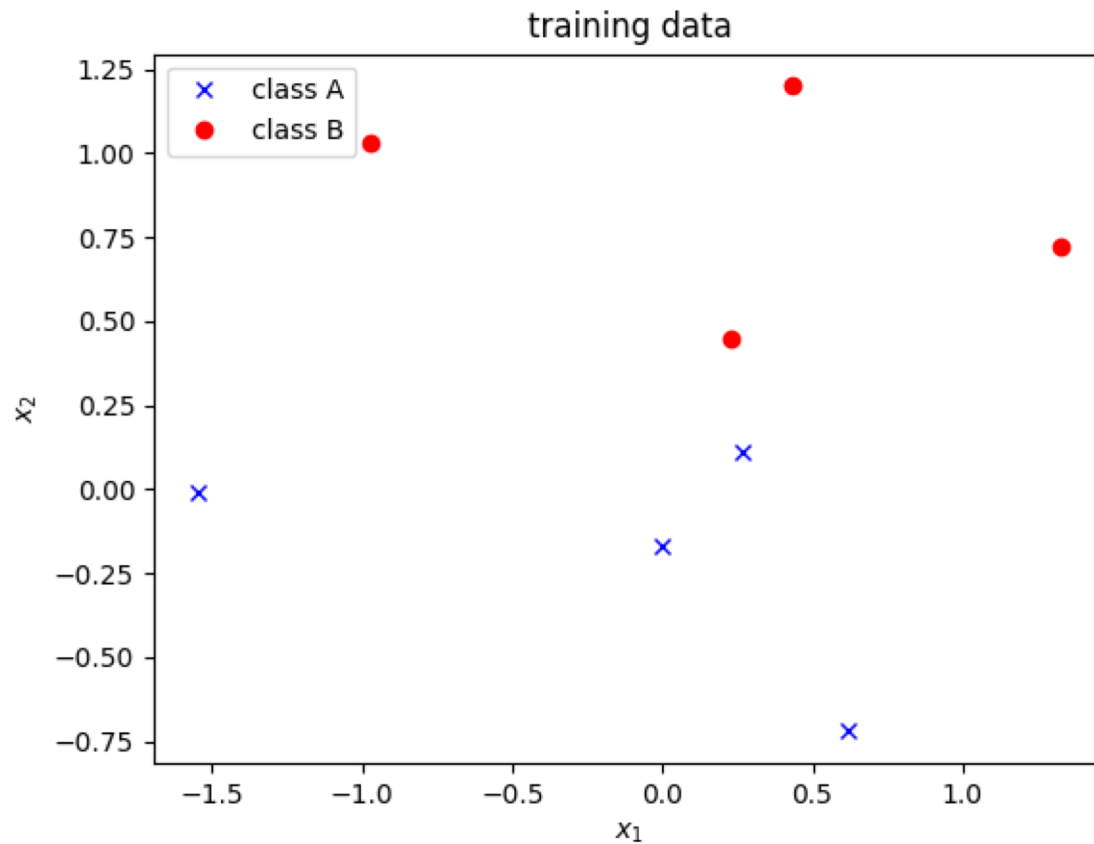
Example 2

Train a logistic regression neuron to perform the following classification:

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

User a learning factor $\alpha = 0.4$.

Example 2



Example 2

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 2

Iteration 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 2

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

$$= \begin{pmatrix} 0.77 \\ 0.02 \end{pmatrix} + 0.4 \begin{pmatrix} 1.33 & -1.55 & 0.62 & 0.27 & 0 & 0.43 & -0.97 & 0.23 \\ 0.72 & -0.01 & -0.72 & 0.11 & -0.17 & 1.2 & 1.03 & 0.45 \end{pmatrix} \begin{pmatrix} 0 \\ 1 \\ 1 \\ 1 \\ 1 \\ 0 \\ 0 \\ 0 \end{pmatrix} - \begin{pmatrix} 0.74 \\ 0.23 \\ 0.61 \\ 0.55 \\ 0.5 \\ 0.59 \\ 0.33 \\ 0.55 \end{pmatrix}$$

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

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

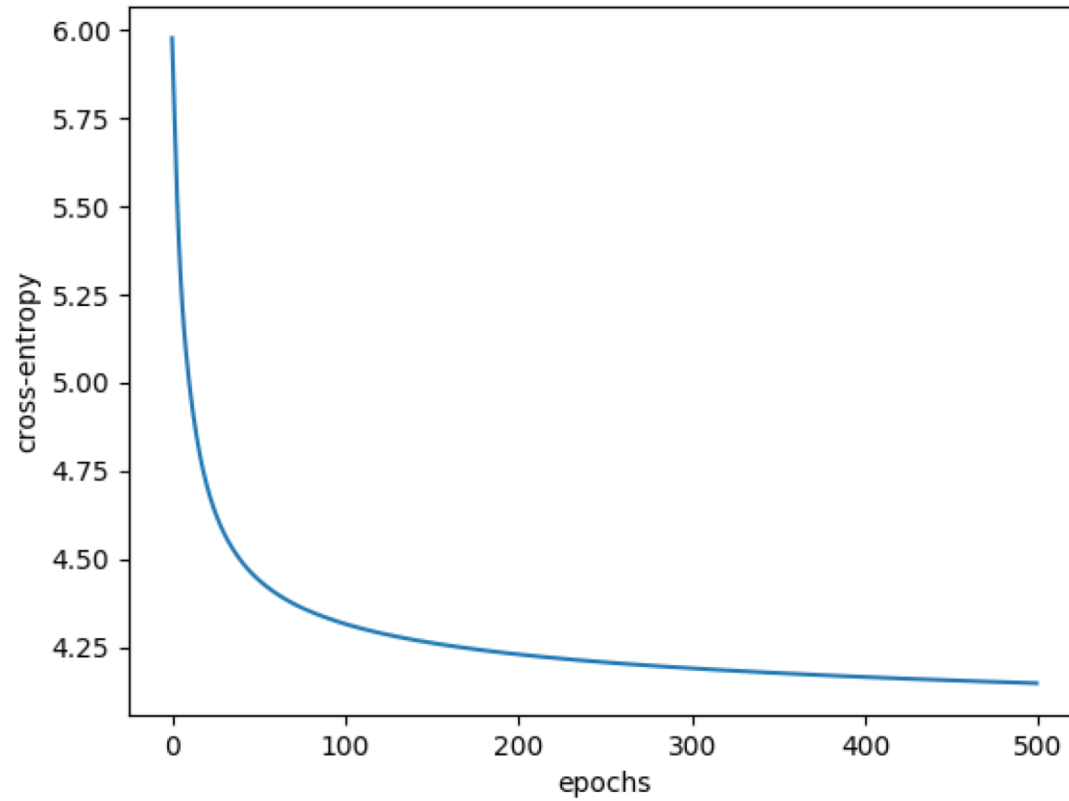
$$= 0.0 + 0.4(1 \quad 1 \quad 1 \quad 1 \quad 1 \quad 1 \quad 1 \quad 1) \begin{pmatrix} 0 \\ 1 \\ 1 \\ 1 \\ 1 \\ 0 \\ 0 \\ 0 \end{pmatrix} - \begin{pmatrix} 0.74 \\ 0.23 \\ 0.61 \\ 0.55 \\ 0.5 \\ 0.59 \\ 0.33 \\ 0.55 \end{pmatrix} = -0.09$$

Example 2

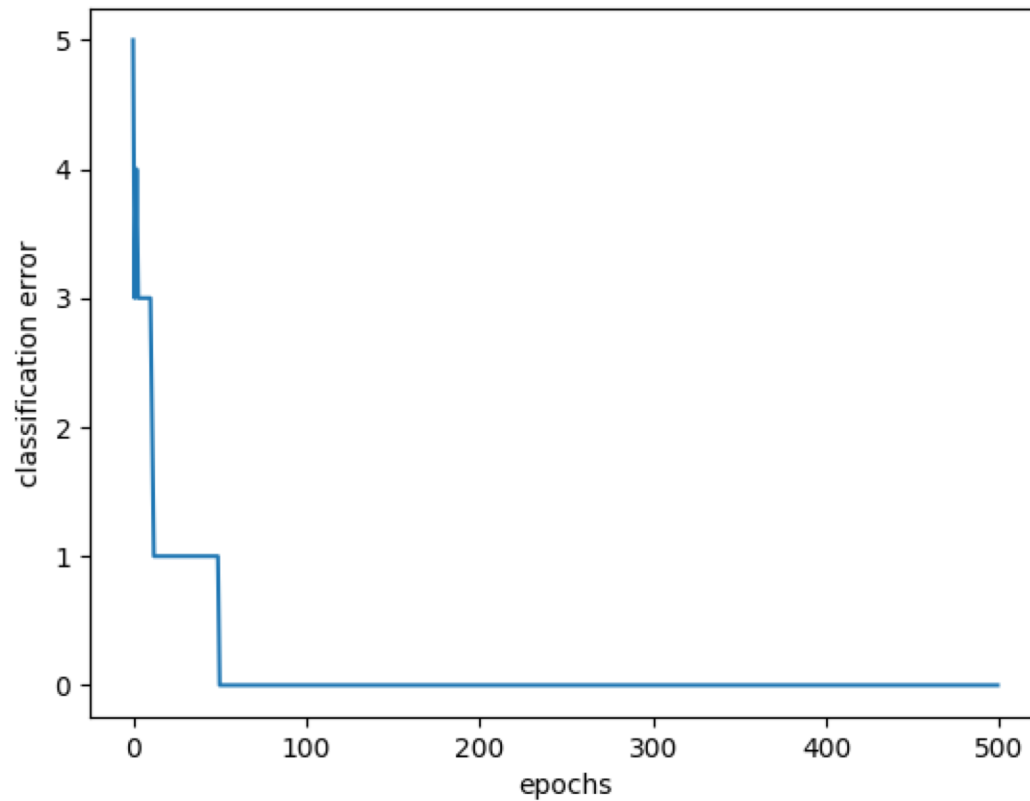
$$\mathbf{y} = 1(f(\mathbf{u}) > 0.5) = 1 \left(\begin{pmatrix} 0.74 \\ 0.23 \\ 0.61 \\ 0.55 \\ 0.5 \\ 0.59 \\ 0.33 \\ 0.55 \end{pmatrix} > 0.5 \right) = \begin{pmatrix} 1 \\ 0 \\ 1 \\ 1 \\ 0 \\ 1 \\ 0 \\ 1 \end{pmatrix}$$

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

Example 2



Example 2

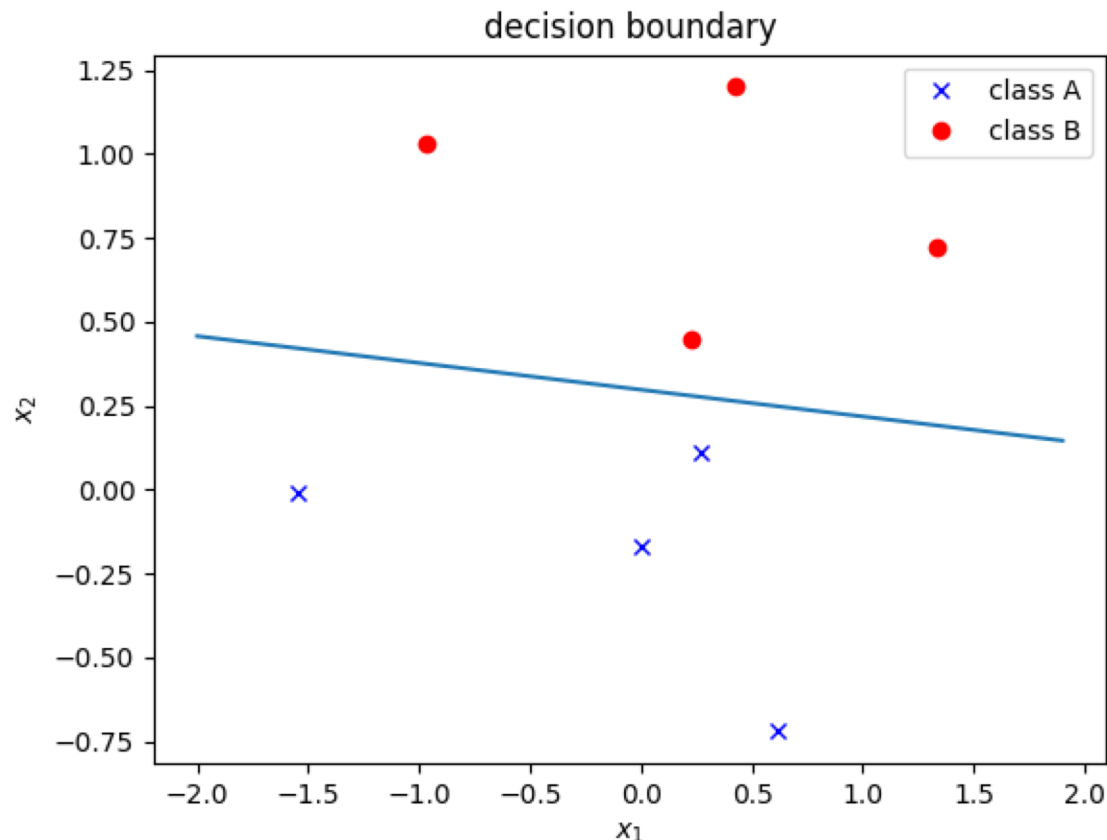


At convergence, $\mathbf{w} = \begin{pmatrix} -1.20 \\ -15.02 \end{pmatrix}$, $b = 4.47$

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

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

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

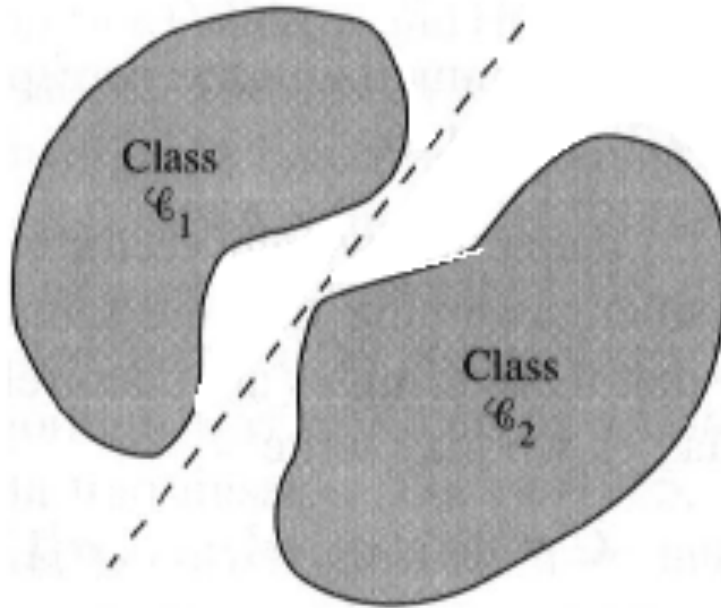


Limitations of Discrete Perceptron and Logistic neuron

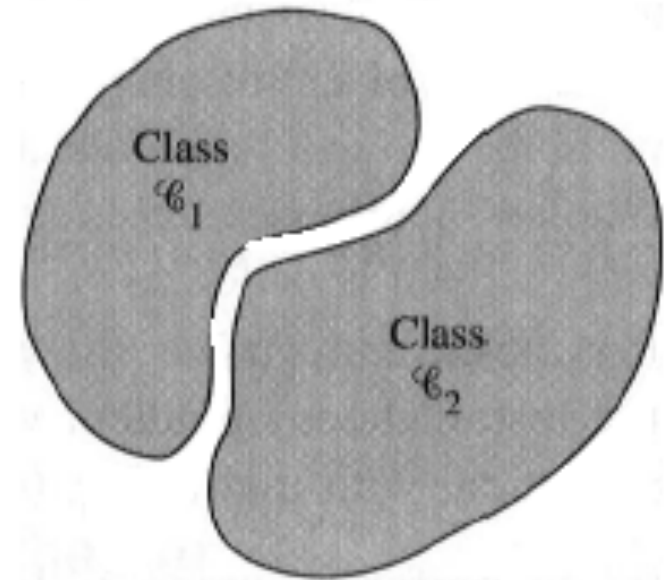
As long as an ANN consists of a *linear combiner* followed by a *non-linear element*, then regardless of the form of non-linearity used, a perceptron can perform pattern classification *only* on *linearly separable* patterns.

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 Discrete Perceptron and Logistic neuron



(a) Linearly
Separable Pattern



(b) Non-Linearly
Separable Pattern

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

Summary of Chapter 3

- Discrete perceptron and logistic regression implements *linear* binary classification (i.e., two class classification) and create linear decision boundaries
- $\mathbf{u} = \mathbf{X}\mathbf{w} + b\mathbf{1}_P$

Discrete perceptron:

$$\mathbf{y} = 1(\mathbf{u} > 0)$$

Learning:

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

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

Logistic regression:

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

$$\mathbf{y} = 1(f(\mathbf{u}) > 0.5)$$

Learning:

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

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

Summary: Neurons

Role	Neuron
Regression (one dimensional)	Linear neuron
	Perceptron
Classification (two class)	Discrete perceptron
	Logistic regression neuron

Summary: GD for a neuron

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

neuron	$f(\mathbf{u}), y$	$\nabla_{\mathbf{u}} J$
Discrete perceptron	$y = 1(f(\mathbf{u}) > 0)$	$-(d - y)$
Logistic regression neuron	$f(\mathbf{u}) = \frac{1}{1 + e^{-u}}$ $y = 1(f(\mathbf{u}) > 0.5)$	$-(d - f(\mathbf{u}))$
Linear neuron	$y = f(\mathbf{u}) = u$	$-(d - y)$
Perceptron	$y = f(\mathbf{u}) = \frac{1}{1 + e^{-u}}$	$-(d - y) \cdot f'(\mathbf{u})$

Summary: SGD for a neuron

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

neuron	$f(u_p), y_p$	$\nabla_{u_p} J$
Discrete perceptron	$y_p = 1(u_p > 0)$	$-(d_p - y_p)$
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	$y_p = f(u_p) = u_p$	$-(d_p - y_p)$
Perceptron	$y_p = f(u_p) = \frac{1}{1 + e^{-u_p}}$	$-(d_p - y_p) \cdot f'(u_p)$