# Introduction to Artificial Intelligence
## Introduction Single-Layer Perceptron

Andres Mendez-Vazquez

March 11, 2019

# Outline

Cinvestav

# Outline

Cinvestav

# History

## At the beginning of Neural Networks (1943 - 1958)

- McCulloch and Pitts (1943) for introducing the idea of neural networks as computing machines.
- Hebb (1949) for postulating the first rule for self-organized learning.
- Rosenblatt (1958) for proposing the perceptron as the first model for learning with a teacher (i.e., supervised learning).

# History

## At the beginning of Neural Networks (1943 - 1958)

- McCulloch and Pitts (1943) for introducing the idea of neural networks as computing machines.
- Hebb (1949) for postulating the first rule for self-organized learning.
- Rosenblatt (1958) for proposing the perceptron as the first model for learning with a teacher (i.e., supervised learning).

## In this Chapter, we are interested in the perceptron

The perceptron is the simplest form of a neural network used for the classification of patterns said to be linearly separable (i.e., patterns that lie on opposite sides of a hyperplane).

Cinvestav

# History

## At the beginning of Neural Networks (1943 - 1958)

- McCulloch and Pitts (1943) for introducing the idea of neural networks as computing machines.
- Hebb (1949) for postulating the first rule for self-organized learning.
- Rosenblatt (1958) for proposing the perceptron as the first model for learning with a teacher (i.e., supervised learning).

In this chapter, we are interested in the perceptron

The perceptron is the simplest form of a neural network used for the classification of patterns said to be linearly separable (i.e., patterns that lie on opposite sides of a hyperplane).

# History

## At the beginning of Neural Networks (1943 - 1958)

- McCulloch and Pitts (1943) for introducing the idea of neural networks as computing machines.
- Hebb (1949) for postulating the first rule for self-organized learning.
- Rosenblatt (1958) for proposing the perceptron as the first model for learning with a teacher (i.e., supervised learning).

## In this chapter, we are interested in the perceptron

The perceptron is the simplest form of a neural network used for the classifica tion of patterns said to be linearly separable (i.e., patterns that lie on opposite sides of a hyperplane).

# In addition

## Something Notable

- The single neuron also forms the basis of an adaptive filter.
- A functional block that is basic to the ever-expanding subject of signal processing.

# In addition

## Something Notable

- The single neuron also forms the basis of an adaptive filter.
- A functional block that is basic to the ever-expanding subject of signal processing.

## Furthermore

The development of adaptive filtering owes much to the classic paper of Widrow and Hoff (1960) for pioneering the so-called least-mean-square (LMS) algorithm, also known as the delta rule.

# In addition

## Something Notable

- The single neuron also forms the basis of an adaptive filter.
- A functional block that is basic to the ever-expanding subject of signal processing.

## Furthermore

The development of adaptive filtering owes much to the classic paper of Widrow and Hoff (1960) for pioneering the so-called least-mean-square (LMS) algorithm, also known as the delta rule.

# Outline

Cinvestav

# Adapting Filtering Problem

## Consider a dynamical system



INPUTS
$x_1(i)$
$x_2(i)$
$\vdots$
$x_m(i)$

Unknown
dynamical
system

Output
$d(i)$

Cinvestav

# Signal-Flow Graph of Adaptive Model

## We have the following equivalence

# Outline

Cinvestav

# Description of the Behavior of the System

We have the data set

$$\mathcal{T} = \{(\boldsymbol{x}(i), d(i)) \mid i = 1, 2, ..., n, ...\} \tag{1}$$

Where

$$\boldsymbol{x}(i) = (x_1(i), x_2(i), ..., x_m(i))^T \tag{2}$$

# Description of the Behavior of the System

## We have the data set

$$\mathcal{T} = \{(\boldsymbol{x}\,(i)\,, d\,(i))\,|i = 1, 2, ..., n, ...\} \tag{1}$$

## Where

$$\boldsymbol{x}\,(i) = (x_1\,(i)\,, x_2\,(i)\,..., x_m\,(i))^T \tag{2}$$

# The Stimulus $\boldsymbol{x}(i)$

## The stimulus $\boldsymbol{x}(i)$ can arise from

The $m$ elements of $\boldsymbol{x}(i)$ originate at different points in space (spatial)



$$\boldsymbol{x}(i) = \begin{pmatrix} x_1(i) \\ x_2(i) \\ x_3(i) \end{pmatrix}$$

# The Stimulus $x(i)$

## The stimulus $x(i)$ can arise from

The $m$ elements of $x(i)$ represent the set of present and $(m-1)$ past values of some excitation that are uniformly spaced in time (temporal).

# Problem

<div style="background: green; color: white; padding: 4px">Quite important</div>

How do we design a multiple input-single output model of the unknown dynamical system?

It is more

We want to build this around a single neuron!!!

Cinvestav

# Problem

> **Quite important**
>
> How do we design a multiple input-single output model of the unknown dynamical system?

> **It is more**
>
> We want to build this around a single neuron!!!

# Thus, we have the following...

# Which steps do you need for the algorithm?

## First

The algorithms starts from an arbitrary setting of the neuron's synaptic weight.

## Second

Adjustments, with respect to changes on the environment, are made on a continuous basis.

# Which steps do you need for the algorithm?

## First

The algorithms starts from an arbitrary setting of the neuron's synaptic weight.

## Second

Adjustments, with respect to changes on the environment, are made on a continuous basis.

- Time is incorporated to the algorithm.

# Which steps do you need for the algorithm?

## First

The algorithms starts from an arbitrary setting of the neuron's synaptic weight.

## Second

Adjustments, with respect to changes on the environment, are made on a continuous basis.

- Time is incorporated to the algorithm.

## Third

Computation of adjustments to synaptic weights are completed inside a time interval that is one sampling period long.

# Signal-Flow Graph of Adaptive Model

## We have the following equivalence

# Thus, This Neural Model $\approx$ Adaptive Filter with two continous processes

## Filtering processes

1. An output, denoted by $y(i)$, that is produced in response to the $m$ elements of the stimulus vector $\boldsymbol{x}(i)$.

2. An error signal, $e(i)$, that is obtained by comparing the output $y(i)$ to the corresponding desired output $d(i)$ produced by the unknown system.

## Adaptive Process

It involves the automatic adjustment of the synaptic weights of the neuron in accordance with the error signal $e(i)$

## Remark

The combination of these two processes working together constitutes a feedback loop acting around the neuron.

# Thus, This Neural Model $\approx$ Adaptive Filter with two continous processes

## Filtering processes

1. An output, denoted by $y(i)$, that is produced in response to the $m$ elements of the stimulus vector $\boldsymbol{x}(i)$.

2. An error signal, $e(i)$, that is obtained by comparing the output $y(i)$ to the corresponding desired output $d(i)$ produced by the unknown system.

## Adaptive Process

It involves the automatic adjustment of the synaptic weights of the neuron in accordance with the error signal $e(i)$

## Remark

The combination of these two processes working together constitutes a feedback loop acting around the neuron.

# Thus, This Neural Model $\approx$ Adaptive Filter with two continous processes

## Filtering processes

1. An output, denoted by $y(i)$, that is produced in response to the $m$ elements of the stimulus vector $\boldsymbol{x}(i)$.

2. An error signal, $e(i)$, that is obtained by comparing the output $y(i)$ to the corresponding desired output $d(i)$ produced by the unknown system.

## Adaptive Process

It involves the automatic adjustment of the synaptic weights of the neuron in accordance with the error signal $e(i)$

## Remark

The combination of these two processes working together constitutes a **feedback loop** acting around the neuron.

# Thus

<div style="background:green;color:white;padding:4px">The output $y(i)$ is exactly the same as the induced local field $v(i)$</div>

$$y\left(i\right) = v\left(i\right) = \sum_{i=1}^{m} w_k\left(i\right) x_k\left(i\right) \tag{3}$$

In matrix form, we have - remember we only have a neuron, so we do not have neuron $k$

$$y\left(i\right) = x^T\left(i\right) w\left(i\right) \tag{4}$$

Error

$$e\left(i\right) = d\left(i\right) - y\left(i\right) \tag{5}$$

# Thus

The output $y(i)$ is exactly the same as the induced local field $v(i)$

$$y(i) = v(i) = \sum_{i=1}^{m} w_k(i) x_k(i) \qquad (3)$$

In matrix form, we have - remember we only have a neuron, so we do not have neuron $k$

$$y(i) = \boldsymbol{x}^T(i)\,\boldsymbol{w}(i) \qquad (4)$$

Error

$$e(i) = d(i) - y(i) \qquad (5)$$

# Thus

The output $y(i)$ is exactly the same as the induced local field $v(i)$

$$y(i) = v(i) = \sum_{i=1}^{m} w_k(i) x_k(i) \tag{3}$$

In matrix form, we have - remember we only have a neuron, so we do not have neuron $k$

$$y(i) = \boldsymbol{x}^T(i) \boldsymbol{w}(i) \tag{4}$$

Error

$$e(i) = d(i) - y(i) \tag{5}$$

# Outline

Cinvestav

# Consider

## A continous differentiable function $J\left(\boldsymbol{w}\right)$

We want to find an optimal solution $\boldsymbol{w}^*$ such that

$$J\left(\boldsymbol{w}^*\right) \leq J\left(\boldsymbol{w}\right), \ \forall \boldsymbol{w} \tag{6}$$

We want to

Minimize the cost function $J(w)$ with respect to the weight vector $w$.

For This

$$\nabla J(w) = 0 \tag{7}$$

# Consider

### A continous differentiable function $J(\boldsymbol{w})$

We want to find an optimal solution $\boldsymbol{w}^*$ such that

$$J(\boldsymbol{w}^*) \leq J(\boldsymbol{w}), \ \forall \boldsymbol{w} \tag{6}$$

### We want to

Minimize the cost function $J(\boldsymbol{w})$ with respect to the weight vector $\boldsymbol{w}$.

### For this

$$\nabla J(\boldsymbol{w}) = 0 \tag{7}$$

# Consider

**A continous differentiable function $J(\boldsymbol{w})$**

We want to find an optimal solution $\boldsymbol{w}^*$ such that

$$J(\boldsymbol{w}^*) \leq J(\boldsymbol{w}), \ \forall \boldsymbol{w} \tag{6}$$

**We want to**

Minimize the cost function $J(\boldsymbol{w})$ with respect to the weight vector $\boldsymbol{w}$.

**For this**

$$\nabla J(\boldsymbol{w}) = 0 \tag{7}$$

# Where

$\nabla$ is the gradient operator

$$\nabla = \left[ \frac{\partial}{\partial w_1}, \frac{\partial}{\partial w_2}, ..., \frac{\partial}{\partial w_m} \right]^T \tag{8}$$

Thus

$$\nabla J(w) = \left[ \frac{\partial J(w)}{\partial w_1}, \frac{\partial J(w)}{\partial w_2}, ..., \frac{\partial J(w)}{\partial w_m} \right]^T \tag{9}$$

# Where

$\nabla$ is the gradient operator

$$\nabla = \left[ \frac{\partial}{\partial w_1}, \frac{\partial}{\partial w_2}, ..., \frac{\partial}{\partial w_m} \right]^T \tag{8}$$

Thus

$$\nabla J\left(\boldsymbol{w}\right) = \left[ \frac{\partial J\left(\boldsymbol{w}\right)}{\partial w_1}, \frac{\partial J\left(\boldsymbol{w}\right)}{\partial w_2}, ..., \frac{\partial J\left(\boldsymbol{w}\right)}{\partial w_m} \right]^T \tag{9}$$

# Thus

**Starting with an initial guess denoted by $\boldsymbol{w}(0)$,**

Then, generate a sequence of weight vectors $\boldsymbol{w}(1), \boldsymbol{w}(2), ...$

Such that you can reduce $J(\boldsymbol{w})$ at each iteration

$$J(\boldsymbol{w}(n+1)) < J(\boldsymbol{w}(n)) \qquad (10)$$

Where, $\boldsymbol{w}(n)$ is the old value and $\boldsymbol{w}(n+1)$ is the new value

# Thus

Starting with an initial guess denoted by $\boldsymbol{w}(0)$,

Then, generate a sequence of weight vectors $\boldsymbol{w}(1), \boldsymbol{w}(2), ...$

Such that you can reduce $J(\boldsymbol{w})$ at each iteration

$$J(\boldsymbol{w}(n+1)) < J(\boldsymbol{w}(n)) \qquad (10)$$

Where: $\boldsymbol{w}(n)$ is the old value and $\boldsymbol{w}(n+1)$ is the new value.

# The Three Main Methods for Unconstrained Optimization

## We will look at

1. Steepest Descent.
2. Newton's Method
3. Gauss-Newton Method

# Outline

Cinvestav

# Steepest Descent

$$\boldsymbol{w}(n+1) \;=\; \boldsymbol{w}(n) - \eta \nabla J(\boldsymbol{w}(n))$$

How we prove that $J(\boldsymbol{w}(n+1)) < J(\boldsymbol{w}(n))$?

We use the first-order Taylor series expansion around $\boldsymbol{w}(n)$

$$J(\boldsymbol{w}(n+1)) \approx J(\boldsymbol{w}(n)) + \nabla J^T(\boldsymbol{w}(n))\Delta \boldsymbol{w}(n) \qquad (11)$$

Remark  This is quite true when the step size is quite small!!! In addition, $\Delta \boldsymbol{w}(n) = \boldsymbol{w}(n+1) - \boldsymbol{w}(n)$

# Steepest Descent

In the method of steepest descent, we have a cost function $J(\boldsymbol{w})$ where

$$\boldsymbol{w}(n+1) \;=\; \boldsymbol{w}(n) - \eta \nabla J(\boldsymbol{w}(\boldsymbol{n}))$$

How, we prove that $J(\boldsymbol{w}(n+1)) < J(\boldsymbol{w}(n))$?

We use the first-order Taylor series expansion around $\boldsymbol{w}(n)$

$$J(\boldsymbol{w}(n+1)) \approx J(\boldsymbol{w}(n)) + \nabla J^T(\boldsymbol{w}(\boldsymbol{n}))\,\Delta\boldsymbol{w}(n) \qquad (11)$$

Remark:  This is quite true when the step size is quite small!!! In addition, $\Delta\boldsymbol{w}(n) = \boldsymbol{w}(n+1) - \boldsymbol{w}(n)$

# Why? Look at the case in $\mathbb{R}$

### The equation of the tangent line to the curve $y = J(w(n))$

$$L(w(n)) = J'(w(n))[w(n+1) - w(n)] + J(w(n)) \qquad (12)$$

Example

# Why? Look at the case in $\mathbb{R}$

## The equation of the tangent line to the curve $y = J(w(n))$

$$L(w(n)) = J'(w(n))[w(n+1) - w(n)] + J(w(n)) \qquad (12)$$

## Example

# Thus, we have that in $\mathbb{R}$

## Remember Something quite Classic



$$\tan \theta = \frac{J\left(w\left(n+1\right)\right) - J\left(w\left(n\right)\right)}{w\left(n+1\right) - w\left(n\right)}$$

# Thus, we have that in $\mathbb{R}$

## Remember Something quite Classic



$$\tan\theta = \frac{J\left(w\left(n+1\right)\right) - J\left(w\left(n\right)\right)}{w\left(n+1\right) - w\left(n\right)}$$

$$\tan\theta\left(w\left(n+1\right) - w\left(n\right)\right) = J\left(w\left(n+1\right)\right) - J\left(w\left(n\right)\right)$$

# Thus, we have that in $\mathbb{R}$

## Remember Something quite Classic



$J\left(w\left(n+1\right)\right) - J\left(w\left(n\right)\right)$

$\theta$

$w\left(n+1\right) - w\left(n\right)$

$$\tan\theta = \frac{J\left(w\left(n+1\right)\right) - J\left(w\left(n\right)\right)}{w\left(n+1\right) - w\left(n\right)}$$

$$\tan\theta \left(w\left(n+1\right) - w\left(n\right)\right) = J\left(w\left(n+1\right)\right) - J\left(w\left(n\right)\right)$$

$$J'\left(w\left(n\right)\right)\left(w\left(n+1\right) - w\left(n\right)\right) = J\left(w\left(n+1\right)\right) - J\left(w\left(n\right)\right)$$

# Thus, we have that

$$J\left(w\left(n\right)\right) \approx J\left(w\left(n\right)\right) + J'\left(w\left(n\right)\right)\left[w\left(n+1\right) - w\left(n\right)\right] \qquad (13)$$

# Now, for Many Variables

# Now, for Many Variables

$$H = \left\{ \boldsymbol{x} | \boldsymbol{a}^T \boldsymbol{x} = b \right\} \qquad (14)$$

Given $\boldsymbol{x} \in H$ and $\boldsymbol{x}_0 \in H$

$$b = \boldsymbol{a}^T \boldsymbol{x} = \boldsymbol{a}^T \boldsymbol{x}_0$$

Thus, we have that

$$H = \left\{ \boldsymbol{x} | \boldsymbol{a}^T \left( \boldsymbol{x} - \boldsymbol{x}_0 \right) = 0 \right\}$$

Cinvestav

# Now, for Many Variables

An hyperplane in $\mathbb{R}^n$ is a set of the form
$$H = \left\{ \boldsymbol{x} \middle| \boldsymbol{a}^T \boldsymbol{x} = b \right\} \tag{14}$$

Given $\boldsymbol{x} \in H$ and $\boldsymbol{x}_0 \in H$
$$b = \boldsymbol{a}^T \boldsymbol{x} = \boldsymbol{a}^T \boldsymbol{x}_0$$

Thus, we have that
$$H = \left\{ \boldsymbol{x} \middle| \boldsymbol{a}^T \left( \boldsymbol{x} - \boldsymbol{x}_0 \right) = 0 \right\}$$

# Thus, we have the following definition

## Definition (Differentiability)

Assume that $J$ is defined in a disk $D$ containing $\boldsymbol{w}(n)$. We say that $J$ is differentiable at $\boldsymbol{w}(n)$ if:

1. $\frac{\partial J(\boldsymbol{w}(n))}{\partial w_i}$ exist for all $i = 1, \ldots, n$.
2. $J$ is locally linear at $\boldsymbol{w}(n)$.

# Thus, we have the following definition

---

### Definition (Differentiability)

Assume that $J$ is defined in a disk $D$ containing $\boldsymbol{w}\left(n\right)$. We say that $J$ is differentiable at $\boldsymbol{w}\left(n\right)$ if:

1. $\frac{\partial J(\boldsymbol{w}(n))}{\partial w_i}$ exist for all $i = 1, ..., n$.

2. $J$ is locally linear at $\boldsymbol{w}\left(n\right)$.

---

# Thus, we have the following definition

## Definition (Differentiability)

Assume that $J$ is defined in a disk $D$ containing $\boldsymbol{w}(n)$. We say that $J$ is differentiable at $\boldsymbol{w}(n)$ if:

1. $\frac{\partial J(\boldsymbol{w}(n))}{\partial w_i}$ exist for all $i = 1, ..., n$.
2. $J$ is locally linear at $\boldsymbol{w}(n)$.

Thus, given $J\left(\boldsymbol{w}\left(n\right)\right)$

## We know that we have the following operator

$$\nabla = \left(\frac{\partial}{\partial w_1}, \frac{\partial}{\partial w_2}, ..., \frac{\partial}{\partial w_m}\right) \tag{15}$$

Thus, given $J\left(\boldsymbol{w}\left(n\right)\right)$

We know that we have the following operator

$$\nabla = \left(\frac{\partial}{\partial w_1}, \frac{\partial}{\partial w_2}, ..., \frac{\partial}{\partial w_m}\right) \tag{15}$$

Thus, we have

$$\nabla J\left(\boldsymbol{w}\left(n\right)\right) = \left(\frac{\partial J\left(\boldsymbol{w}\left(n\right)\right)}{\partial w_1}, \frac{\partial J\left(\boldsymbol{w}\left(n\right)\right)}{\partial w_2}, ..., \frac{\partial J\left(\boldsymbol{w}\left(n\right)\right)}{\partial w_m}\right)$$

$$= \sum_{i=1}^{m} \hat{w}_i \frac{\partial J\left(\boldsymbol{w}\left(n\right)\right)}{\partial w_i}$$

Where: $\hat{w}_i^T = (1, 0, ..., 0) \in \mathbb{R}$

# Now

Given a curve function $r(t)$ that lies on the level set $J(\boldsymbol{w}(n)) = c$ (When is in $\mathbb{R}^3$)

# Level Set

## Definition

$$\{(w_1, w_2, ..., w_m) \in \mathbb{R}^m | J(w_1, w_2, ..., w_m) = c\} \qquad (16)$$

Remark: In a normal Calculus course we will use $x$ and $f$ instead of $w$ and $J$.

## Where

---

**Any curve has the following parametrization**

$$r : [a, b] \to \mathbb{R}^m$$
$$r(t) = (w_1(t), ..., w_m(t))$$

With $r(n+1) = (w_1(n+1), ..., w_m(n+1))$

---

We can write the parametrized version of it

$$z(t) = J(w_1(t), w_2(t), ..., w_m(t)) = c \tag{17}$$

Differentiating with respect to $t$ and using the chain rule for multiple variables

$$\frac{dz(t)}{dt} = \sum_{i=1}^{m} \frac{\partial J(w(t))}{\partial w_i} \cdot \frac{dw_i(t)}{dt} = 0 \tag{18}$$

# Where

## Any curve has the following parametrization

$$r : [a, b] \to \mathbb{R}^m$$

$$r(t) = (w_1(t), ..., w_m(t))$$

With $r(n+1) = (w_1(n+1), ..., w_m(n+1))$

## We can write the parametrized version of it

$$z(t) = J(w_1(t), w_2(t), ..., w_m(t)) = c \qquad (17)$$

Differentiating with respect to $t$ and using the chain rule for multiple variables

$$\frac{dz(t)}{dt} = \sum_{i=1}^{m} \frac{\partial J(w(t))}{\partial w_i} \cdot \frac{dw_i(t)}{dt} = 0 \qquad (18)$$

# Where

## Any curve has the following parametrization

$$r : [a, b] \to \mathbb{R}^m$$
$$r(t) = (w_1(t), ..., w_m(t))$$

With $r(n+1) = (w_1(n+1), ..., w_m(n+1))$

## We can write the parametrized version of it

$$z(t) = J(w_1(t), w_2(t), ..., w_m(t)) = c \tag{17}$$

## Differentiating with respect to $t$ and using the chain rule for multiple variables

$$\frac{dz(t)}{dt} = \sum_{i=1}^{m} \frac{\partial J(\boldsymbol{w}(t))}{\partial w_i} \cdot \frac{dw_i(t)}{dt} = 0 \tag{18}$$

# Note

## First

Given $y = f(\boldsymbol{u}) = (f_1(\boldsymbol{u}), ..., f_l(\boldsymbol{u}))$ and
$\boldsymbol{u} = g(\boldsymbol{x}) = (g_1(\boldsymbol{x}), ..., g_m(\boldsymbol{x}))$.

We have then that

$$\frac{\partial(f_1, f_2, ..., f_l)}{\partial(x_1, x_2, ..., x_k)} = \frac{\partial(f_1, f_2, ..., f_l)}{\partial(g_1, g_2, ..., g_m)} \cdot \frac{\partial(g_1, g_2, ..., g_m)}{\partial(x_1, x_2, ..., x_k)} \qquad (19)$$

Thus

$$\frac{\partial(f_1, f_2, ..., f_l)}{\partial x_i} = \frac{\partial(f_1, f_2, ..., f_l)}{\partial(g_1, g_2, ..., g_m)} \cdot \frac{\partial(g_1, g_2, ..., g_m)}{\partial x_i}$$

$$= \sum_{k=1}^{m} \frac{\partial(f_1, f_2, ..., f_l)}{\partial g_k} \frac{\partial g_k}{\partial x_i}$$

# Note

## First

Given $y = f(\boldsymbol{u}) = (f_1(\boldsymbol{u}), ..., f_l(\boldsymbol{u}))$ and
$\boldsymbol{u} = g(\boldsymbol{x}) = (g_1(\boldsymbol{x}), ..., g_m(\boldsymbol{x}))$.

## We have then that

$$\frac{\partial(f_1, f_2, ..., f_l)}{\partial(x_1, x_2, ..., x_k)} = \frac{\partial(f_1, f_2, ..., f_l)}{\partial(g_1, g_2, ..., g_m)} \cdot \frac{\partial(g_1, g_2, ..., g_m)}{\partial(x_1, x_2, ..., x_k)} \qquad (19)$$

Times

$$\frac{\partial(f_1, f_2, ..., f_l)}{\partial x_i} = \frac{\partial(f_1, f_2, ..., f_l)}{\partial(g_1, g_2, ..., g_m)} \cdot \frac{\partial(g_1, g_2, ..., g_m)}{\partial x_i}$$

$$= \sum_{k=1}^{m} \frac{\partial(f_1, f_2, ..., f_l)}{\partial g_k} \frac{\partial g_k}{\partial x_i}$$

# Note

## First

Given $y = f(\boldsymbol{u}) = (f_1(\boldsymbol{u}), ..., f_l(\boldsymbol{u}))$ and
$\boldsymbol{u} = g(\boldsymbol{x}) = (g_1(\boldsymbol{x}), ..., g_m(\boldsymbol{x}))$.

## We have then that

$$\frac{\partial(f_1, f_2, ..., f_l)}{\partial(x_1, x_2, ..., x_k)} = \frac{\partial(f_1, f_2, ..., f_l)}{\partial(g_1, g_2, ..., g_m)} \cdot \frac{\partial(g_1, g_2, ..., g_m)}{\partial(x_1, x_2, ..., x_k)} \qquad (19)$$

## Thus

$$\frac{\partial(f_1, f_2, ..., f_l)}{\partial x_i} = \frac{\partial(f_1, f_2, ..., f_l)}{\partial(g_1, g_2, ..., g_m)} \cdot \frac{\partial(g_1, g_2, ..., g_m)}{\partial x_i}$$

$$= \sum_{k=1}^{m} \frac{\partial(f_1, f_2, ..., f_l)}{\partial g_k} \frac{\partial g_k}{\partial x_i}$$

# Thus

**Evaluating at $t = n$**

$$\sum_{i=1}^{m} \frac{\partial J\left(\boldsymbol{w}\left(n\right)\right)}{\partial w_i} \cdot \frac{dw_i(n)}{dt} = 0$$

We have that

$$\nabla J\left(w\left(n\right)\right) \cdot r'\left(n\right) = 0 \tag{20}$$

This proves than for every level set the gradient is perpendicular to the tangent to any curve that lies on the level set

In particular to the point $w\left(n\right)$.

# Thus

## Evaluating at $t = n$

$$\sum_{i=1}^{m} \frac{\partial J(\boldsymbol{w}(n))}{\partial w_i} \cdot \frac{dw_i(n)}{dt} = 0$$

## We have that

$$\nabla J(\boldsymbol{w}(n)) \cdot r'(n) = 0 \tag{20}$$

This proves that for every level set the gradient is perpendicular to the tangent to any curve that lies on the level set

In particular to the point $\boldsymbol{w}(n)$.

# Thus

$$\sum_{i=1}^{m} \frac{\partial J\left(\boldsymbol{w}\left(n\right)\right)}{\partial w_i} \cdot \frac{dw_i(n)}{dt} = 0$$

### We have that

$$\nabla J\left(\boldsymbol{w}\left(n\right)\right) \cdot r'\left(n\right) = 0 \tag{20}$$

This proves that for every level set the gradient is perpendicular to the tangent to any curve that lies on the level set

In particular to the point $\boldsymbol{w}\left(n\right)$.

# Now the tangent plane to the surface can be described generally

## Thus

$$L\left(\boldsymbol{w}\left(n+1\right)\right) = J\left(\boldsymbol{w}\left(n\right)\right) + \nabla J^{T}\left(\boldsymbol{w}\left(\boldsymbol{n}\right)\right)\left[\boldsymbol{w}\left(n+1\right) - \boldsymbol{w}\left(n\right)\right] \quad (21)$$

This looks like

# Now the tangent plane to the surface can be described generally

## Thus

$$L\left(\boldsymbol{w}\left(n+1\right)\right) = J\left(\boldsymbol{w}\left(n\right)\right) + \nabla J^T\left(\boldsymbol{w}\left(\boldsymbol{n}\right)\right)\left[\boldsymbol{w}\left(n+1\right) - \boldsymbol{w}\left(n\right)\right] \quad (21)$$

## This looks like

# Proving the fact about the Steepest Descent

## We want the following

$$J\left(\boldsymbol{w}\left(n+1\right)\right) < J\left(\boldsymbol{w}\left(n\right)\right)$$

# Proving the fact about the Steepest Descent

## We want the following

$$J\left(\boldsymbol{w}\left(n+1\right)\right) < J\left(\boldsymbol{w}\left(n\right)\right)$$

## Using the first-order Taylor approximation

$$J\left(\boldsymbol{w}\left(n+1\right)\right) - J\left(\boldsymbol{w}\left(n\right)\right) \approx \nabla J^{T}\left(\boldsymbol{w}\left(\boldsymbol{n}\right)\right)\Delta\boldsymbol{w}\left(n\right)$$

So, we ask the following

$$\Delta\boldsymbol{w}\left(n\right) \approx -\eta\nabla J\left(\boldsymbol{w}\left(n\right)\right) \text{ with } \eta > 0$$

# Proving the fact about the Steepest Descent

$$J\left(\boldsymbol{w}\left(n+1\right)\right) < J\left(\boldsymbol{w}\left(n\right)\right)$$

**Using the first-order Taylor approximation**

$$J\left(\boldsymbol{w}\left(n+1\right)\right) - J\left(\boldsymbol{w}\left(n\right)\right) \approx \nabla J^{T}\left(\boldsymbol{w}\left(\boldsymbol{n}\right)\right)\Delta\boldsymbol{w}\left(n\right)$$

**So, we ask the following**

$$\Delta\boldsymbol{w}\left(n\right) \approx -\eta\nabla J\left(\boldsymbol{w}\left(\boldsymbol{n}\right)\right) \text{ with } \eta > 0$$

# Then

## We have that

$$J\left(\boldsymbol{w}\left(n+1\right)\right) - J\left(\boldsymbol{w}\left(n\right)\right) \approx -\eta \nabla J^{T}\left(\boldsymbol{w}\left(\boldsymbol{n}\right)\right) \nabla J\left(\boldsymbol{w}\left(\boldsymbol{n}\right)\right) = -\eta \left\|\nabla J\left(\boldsymbol{w}\left(\boldsymbol{n}\right)\right)\right\|^{2}$$

# Then

## We have that

$$J\left(\boldsymbol{w}\left(n+1\right)\right) - J\left(\boldsymbol{w}\left(n\right)\right) \approx -\eta \nabla J^{T}\left(\boldsymbol{w}\left(\boldsymbol{n}\right)\right) \nabla J\left(\boldsymbol{w}\left(\boldsymbol{n}\right)\right) = -\eta \left\|\nabla J\left(\boldsymbol{w}\left(\boldsymbol{n}\right)\right)\right\|^{2}$$

## Thus

$$J\left(\boldsymbol{w}\left(n+1\right)\right) - J\left(\boldsymbol{w}\left(n\right)\right) < 0$$

## Or

$$J\left(\boldsymbol{w}\left(n+1\right)\right) < J\left(\boldsymbol{w}\left(n\right)\right)$$

# Then

## We have that

$$J\left(\boldsymbol{w}\left(n+1\right)\right) - J\left(\boldsymbol{w}\left(n\right)\right) \approx -\eta \nabla J^{T}\left(\boldsymbol{w}\left(\boldsymbol{n}\right)\right) \nabla J\left(\boldsymbol{w}\left(\boldsymbol{n}\right)\right) = -\eta \left\|\nabla J\left(\boldsymbol{w}\left(\boldsymbol{n}\right)\right)\right\|^{2}$$

## Thus

$$J\left(\boldsymbol{w}\left(n+1\right)\right) - J\left(\boldsymbol{w}\left(n\right)\right) < 0$$

## Or

$$J\left(\boldsymbol{w}\left(n+1\right)\right) < J\left(\boldsymbol{w}\left(n\right)\right)$$

# Outline

Cinvestav

# Newton's Method

## Here

The basic idea of Newton's method is to minimize the quadratic approximation of the cost function $J(\boldsymbol{w})$ around the current point $\boldsymbol{w}(n)$.

Using a second-order Taylor series expansion of the cost function around the point $\boldsymbol{w}(n)$

$$\Delta J(\boldsymbol{w}(n)) = J(\boldsymbol{w}(n+1)) - J(\boldsymbol{w}(n))$$

$$\approx \nabla J^T(\boldsymbol{w}(n))\Delta\boldsymbol{w}(n) + \frac{1}{2}\Delta\boldsymbol{w}^T(n)H(n)\Delta\boldsymbol{w}(n)$$

Where given that $\boldsymbol{w}(n)$ is a vector with dimension $m$

$$H = \nabla^2 J(\boldsymbol{w}) = \begin{pmatrix} \frac{\partial^2 J(\boldsymbol{w})}{\partial w_1^2} & \frac{\partial^2 J(\boldsymbol{w})}{\partial w_1 \partial w_2} & \cdots & \frac{\partial^2 J(\boldsymbol{w})}{\partial w_1 \partial w_m} \\ \frac{\partial^2 J(\boldsymbol{w})}{\partial w_2 \partial w_1} & \frac{\partial^2 J(\boldsymbol{w})}{\partial w_2^2} & \cdots & \frac{\partial^2 J(\boldsymbol{w})}{\partial w_2 \partial w_m} \\ \vdots & \vdots & & \vdots \\ \frac{\partial^2 J(\boldsymbol{w})}{\partial w_m \partial w_1} & \frac{\partial^2 J(\boldsymbol{w})}{\partial w_m \partial w_2} & \cdots & \frac{\partial^2 J(\boldsymbol{w})}{\partial w_m^2} \end{pmatrix}$$

# Newton's Method

## Using a second-order Taylor series expansion of the cost function around the point $\boldsymbol{w}(n)$

$$\Delta J(\boldsymbol{w}(n)) = J(\boldsymbol{w}(n+1)) - J(\boldsymbol{w}(n))$$

$$\approx \nabla J^T(\boldsymbol{w(n)})\,\Delta\boldsymbol{w}(n) + \frac{1}{2}\Delta\boldsymbol{w}^T(n)\,H(n)\,\Delta\boldsymbol{w}(n)$$

Where given that $w(n)$ is a vector with dimension $m$

$$H = \nabla^2 J(\boldsymbol{w}) = \begin{pmatrix} \frac{\partial^2 J(\boldsymbol{w})}{\partial w_1^2} & \frac{\partial^2 J(\boldsymbol{w})}{\partial w_1 \partial w_2} & \cdots & \frac{\partial^2 J(\boldsymbol{w})}{\partial w_1 \partial w_m} \\ \frac{\partial^2 J(\boldsymbol{w})}{\partial w_2 \partial w_1} & \frac{\partial^2 J(\boldsymbol{w})}{\partial w_2^2} & \cdots & \frac{\partial^2 J(\boldsymbol{w})}{\partial w_2 \partial w_m} \\ \vdots & \vdots & & \vdots \\ \frac{\partial^2 J(\boldsymbol{w})}{\partial w_m \partial w_1} & \frac{\partial^2 J(\boldsymbol{w})}{\partial w_m \partial w_2} & \cdots & \frac{\partial^2 J(\boldsymbol{w})}{\partial w_m^2} \end{pmatrix}$$

# Newton's Method

## Here

The basic idea of Newton's method is to minimize the quadratic approximation of the cost function $J(\boldsymbol{w})$ around the current point $\boldsymbol{w}(n)$.

## Using a second-order Taylor series expansion of the cost function around the point $\boldsymbol{w}(n)$

$$\Delta J(\boldsymbol{w}(n)) = J(\boldsymbol{w}(n+1)) - J(\boldsymbol{w}(n))$$

$$\approx \nabla J^T(\boldsymbol{w(n)})\Delta\boldsymbol{w}(n) + \frac{1}{2}\Delta\boldsymbol{w}^T(n)H(n)\Delta\boldsymbol{w}(n)$$

## Where given that $\boldsymbol{w}(n)$ is a vector with dimension $m$

$$H = \nabla^2 J(\boldsymbol{w}) = \begin{pmatrix} \frac{\partial^2 J(\boldsymbol{w})}{\partial w_1^2} & \frac{\partial^2 J(\boldsymbol{w})}{\partial w_1 \partial w_2} & \cdots & \frac{\partial^2 J(\boldsymbol{w})}{\partial w_1 \partial w_m} \\ \frac{\partial^2 J(\boldsymbol{w})}{\partial w_2 \partial w_1} & \frac{\partial^2 J(\boldsymbol{w})}{\partial w_2^2} & \cdots & \frac{\partial^2 J(\boldsymbol{w})}{\partial w_2 \partial w_m} \\ \vdots & \vdots & & \vdots \\ \frac{\partial^2 J(\boldsymbol{w})}{\partial w_m \partial w_1} & \frac{\partial^2 J(\boldsymbol{w})}{\partial w_m \partial w_2} & \cdots & \frac{\partial^2 J(\boldsymbol{w})}{\partial w_m^2} \end{pmatrix}$$

# Now, we want to minimize $J(\boldsymbol{w}(n+1))$

## Do you have any idea?

Look again

$$J(\boldsymbol{w}(n)) + \nabla J^T(\boldsymbol{w(n)}) \Delta \boldsymbol{w}(n) + \frac{1}{2}\Delta \boldsymbol{w}^T(n) H(n) \Delta \boldsymbol{w}(n) \qquad (22)$$

Derive with respect to $\Delta w(n)$

$$\nabla J(w(n)) + H(n)\Delta w(n) = 0 \qquad (23)$$

Thus

$$\Delta w(n) = -H^{-1}(n)\nabla J(w(n))$$

# Now, we want to minimize $J(\boldsymbol{w}(n+1))$

## Do you have any idea?

Look again

$$J(\boldsymbol{w}(n)) + \nabla J^T(\boldsymbol{w}(\boldsymbol{n}))\Delta\boldsymbol{w}(n) + \frac{1}{2}\Delta\boldsymbol{w}^T(n)H(n)\Delta\boldsymbol{w}(n) \tag{22}$$

## Derive with respect to $\Delta\boldsymbol{w}(n)$

$$\nabla J(\boldsymbol{w}(\boldsymbol{n})) + H(n)\Delta\boldsymbol{w}(n) = 0 \tag{23}$$

Thus

$$\Delta\boldsymbol{w}(n) = -H^{-1}(n)\nabla J(\boldsymbol{w}(\boldsymbol{n}))$$

# Now, we want to minimize $J(\boldsymbol{w}(n+1))$

**Do you have any idea?**

Look again

$$J(\boldsymbol{w}(n)) + \nabla J^T(\boldsymbol{w(n)})\,\Delta\boldsymbol{w}(n) + \frac{1}{2}\Delta\boldsymbol{w}^T(n)\,H(n)\,\Delta\boldsymbol{w}(n) \qquad (22)$$

**Derive with respect to $\Delta\boldsymbol{w}(n)$**

$$\nabla J(\boldsymbol{w(n)}) + H(n)\,\Delta\boldsymbol{w}(n) = 0 \qquad (23)$$

**Thus**

$$\Delta\boldsymbol{w}(n) = -H^{-1}(n)\,\nabla J(\boldsymbol{w(n)})$$

# The Final Method

**Define the following**

$$J\left(\boldsymbol{w}\left(n+1\right)\right) - J\left(\boldsymbol{w}\left(n\right)\right) = -H^{-1}\left(n\right)\nabla J\left(\boldsymbol{w}\left(\boldsymbol{n}\right)\right)$$

**Then**

$$J\left(w\left(n+1\right)\right) = J\left(w\left(n\right)\right) - H^{-1}\left(n\right)\nabla J\left(w\left(n\right)\right)$$

# The Final Method

**Define the following**

$$J\left(\boldsymbol{w}\left(n+1\right)\right) - J\left(\boldsymbol{w}\left(n\right)\right) = -H^{-1}\left(n\right)\nabla J\left(\boldsymbol{w}\left(\boldsymbol{n}\right)\right)$$

**Then**

$$J\left(\boldsymbol{w}\left(n+1\right)\right) = J\left(\boldsymbol{w}\left(n\right)\right) - H^{-1}\left(n\right)\nabla J\left(\boldsymbol{w}\left(\boldsymbol{n}\right)\right)$$

# Outline

Cinvestav

# We have then an error

> **Something Notable**
>
> $$J\left(\boldsymbol{w}\right) = \frac{1}{2}\sum_{i=1}^{n} e^2\left(i\right)$$

Thus using the first order Taylor expansion, we linearize

$$e_l\left(i, w\right) = e\left(i\right) + \left[\frac{\partial e\left(i\right)}{\partial w}\right]^T \left[w - w\left(n\right)\right]$$

In matrix form

$$e_l\left(n, w\right) = e\left(n\right) + J\left(n\right)\left[w - w\left(n\right)\right]$$

# We have then an error

**Something Notable**

$$J\left(\boldsymbol{w}\right) = \frac{1}{2}\sum_{i=1}^{n} e^2\left(i\right)$$

**Thus using the first order Taylor expansion, we linearize**

$$e_l\left(i, \boldsymbol{w}\right) = e\left(i\right) + \left[\frac{\partial e\left(i\right)}{\partial \boldsymbol{w}}\right]^T \left[\boldsymbol{w} - \boldsymbol{w}\left(n\right)\right]$$

In matrix form

$$e_l\left(n, w\right) = e\left(n\right) + J\left(n\right)\left[w - w\left(n\right)\right]$$

# We have then an error

$$J\left(\boldsymbol{w}\right) = \frac{1}{2}\sum_{i=1}^{n} e^2\left(i\right)$$

Thus using the first order Taylor expansion, we linearize

$$e_l\left(i, \boldsymbol{w}\right) = e\left(i\right) + \left[\frac{\partial e\left(i\right)}{\partial \boldsymbol{w}}\right]^T \left[\boldsymbol{w} - \boldsymbol{w}\left(n\right)\right]$$

In matrix form

$$\boldsymbol{e}_l\left(n, \boldsymbol{w}\right) = \boldsymbol{e}\left(n\right) + \mathsf{J}\left(n\right)\left[\boldsymbol{w} - \boldsymbol{w}\left(n\right)\right]$$

# Where

## The error vector is equal to

$$\boldsymbol{e}\left(n\right) = \left[e\left(1\right), e\left(2\right), ..., e\left(n\right)\right]^{T} \tag{24}$$

Thus, we get the famous Jacobian once we derive $\frac{\partial e(i)}{\partial u_{j}}$

$$J\left(n\right) = \begin{pmatrix} \frac{\partial e(1)}{\partial u_{1}} & \frac{\partial e(1)}{\partial u_{2}} & \cdots & \frac{\partial e(1)}{\partial u_{m}} \\ \frac{\partial e(2)}{\partial u_{1}} & \frac{\partial e(2)}{\partial u_{2}} & \cdots & \frac{\partial e(2)}{\partial u_{m}} \\ \vdots & \vdots & & \vdots \\ \frac{\partial e(n)}{\partial u_{1}} & \frac{\partial e(n)}{\partial u_{2}} & \cdots & \frac{\partial e(n)}{\partial u_{m}} \end{pmatrix}$$

Cinvestav

# Where

The error vector is equal to

$$\boldsymbol{e}\left(n\right) = \left[e\left(1\right), e\left(2\right), ..., e\left(n\right)\right]^{T} \tag{24}$$

Thus, we get the famous Jacobian once we derive $\frac{\partial e(i)}{\partial \boldsymbol{w}}$

$$\mathsf{J}\left(n\right) = \begin{pmatrix} \frac{\partial e(1)}{\partial w_1} & \frac{\partial e(1)}{\partial w_2} & \cdots & \frac{\partial e(1)}{\partial w_m} \\ \frac{\partial e(2)}{\partial w_1} & \frac{\partial e(2)}{\partial w_2} & \cdots & \frac{\partial e(2)}{\partial w_m} \\ \vdots & \vdots & & \vdots \\ \frac{\partial e(n)}{\partial w_1} & \frac{\partial e(n)}{\partial w_2} & \cdots & \frac{\partial e(n)}{\partial w_m} \end{pmatrix}$$

# Where

## We want the following

$$\boldsymbol{w}\left(n+1\right)=\underset{\boldsymbol{w}}{arg\min}\left\{\frac{1}{2}\left\|\boldsymbol{e}_l\left(n,\boldsymbol{w}\right)\right\|^2\right\}$$

## Ideas

What if we expand out the equation?

# Where

## We want the following

$$\boldsymbol{w}\left(n+1\right)=arg\min_{\boldsymbol{w}}\left\{\frac{1}{2}\left\|\boldsymbol{e}_l\left(n,\boldsymbol{w}\right)\right\|^2\right\}$$

## Ideas

What if we expand out the equation?

# Expanded Version

> **We get**
>
> $$\frac{1}{2} \left\| \boldsymbol{e}_l \left( n, \boldsymbol{w} \right) \right\|^2 = \frac{1}{2} \left\| \boldsymbol{e} \left( n \right) \right\|^2 + \boldsymbol{e}^T \left( n \right) \boldsymbol{J} \left( n \right) \left( \boldsymbol{w} - \boldsymbol{w} \left( n \right) \right) + ...$$
> $$\frac{1}{2} \left( \boldsymbol{w} - \boldsymbol{w} \left( n \right) \right)^T \boldsymbol{J}^T \left( n \right) \boldsymbol{J} \left( n \right) \left( \boldsymbol{w} - \boldsymbol{w} \left( n \right) \right)$$

# Now, doing the Differential, we get

$$\mathsf{J}^T(n)\,\boldsymbol{e}(n) + \mathsf{J}^T(n)\,\mathsf{J}(n)\,[\boldsymbol{w} - \boldsymbol{w}(n)] = 0$$

We get finally

$$w(n+1) = w(n) - \left(\mathsf{J}^T(n)\mathsf{J}(n)\right)^{-1}\mathsf{J}^T(n)e(n) \qquad (25)$$

# Now, doing the Differential, we get

$$\mathsf{J}^T\left(n\right)\boldsymbol{e}\left(n\right) + \mathsf{J}^T\left(n\right)\mathsf{J}\left(n\right)\left[\boldsymbol{w} - \boldsymbol{w}\left(n\right)\right] = 0$$

We get finally

$$\boldsymbol{w}\left(n+1\right) = \boldsymbol{w}\left(n\right) - \left(\mathsf{J}^T\left(n\right)\mathsf{J}\left(n\right)\right)^{-1}\mathsf{J}^T\left(n\right)e\left(n\right) \tag{25}$$

# Remarks

## We have that

- The Newton's method that requires knowledge of the Hessian matrix of the cost function.
- The Gauss-Newton method only requires the Jacobian matrix of the error vector $e(n)$.

## However

The Gauss-Newton iteration to be computable, the matrix product $J^T(n) J(n)$ must be nonsingular!!!

# Remarks

## We have that

- The Newton's method that requires knowledge of the Hessian matrix of the cost function.
- The Gauss-Newton method only requires the Jacobian matrix of the error vector $\boldsymbol{e}(n)$.

## However

The Gauss-Newton iteration to be computable, the matrix product $\mathsf{J}^T(n)\,\mathsf{J}(n)$ must be nonsingular!!!

Cinvestav

# Outline

Cinvestav

# Introduction

## A linear least-squares filter has two distinctive characteristics

- First, the single neuron around which it is built is linear.
- The cost function $J(\boldsymbol{w})$ used to design the filter consists of the sum of error squares.

Thus, expressing the error

$$e(n) = d(n) - (x(1),...,x(n))^T w(n)$$

Short Version - error is linear in the weight vector $w(n)$

$$e(n) = d(n) - X(n) w(n)$$

- Where $d(n)$ is a $n \times 1$ desired response vector.
- Where $X(n)$ is the $n \times m$ data matrix.

# Introduction

<div style="background:green;color:white">A linear least-squares filter has two distinctive characteristics</div>

- First, the single neuron around which it is built is linear.
- The cost function $J(\boldsymbol{w})$ used to design the filter consists of the sum of error squares.

<div style="background:red;color:white">Thus, expressing the error</div>

$$e(n) = d(n) - (\boldsymbol{x}(1),...,\boldsymbol{x}(n))^T \boldsymbol{w}(n)$$

Short Version - error is linear in the weight vector $w(n)$

$$e(n) = d(n) - X(n) w(n)$$

- Where $d(n)$ is a $n \times 1$ desired response vector.
- Where $X(n)$ is the $n \times m$ data matrix.

# Introduction

**A linear least-squares filter has two distinctive characteristics**

- First, the single neuron around which it is built is linear.
- The cost function $J(\boldsymbol{w})$ used to design the filter consists of the sum of error squares.

**Thus, expressing the error**

$$e(n) = d(n) - (\boldsymbol{x}(1), ..., \boldsymbol{x}(n))^T \boldsymbol{w}(n)$$

**Short Version - error is linear in the weight vector $\boldsymbol{w}(n)$**

$$e(n) = d(n) - \boldsymbol{X}(n)\boldsymbol{w}(n)$$

- Where $d(n)$ is a $n \times 1$ desired response vector.
- Where $\boldsymbol{X}(n)$ is the $n \times m$ data matrix.

# Now, differentiate $e(n)$ with respect to $\boldsymbol{w}(n)$

> **Thus**
>
> $$\nabla e(n) = -\boldsymbol{X}^T(n)$$

Correspondingly, the Jacobian of $e(n)$ is

$$J(n) = -X(n)$$

Let us to use the Gaussian-Newton

$$w(n+1) = w(n) - \left(J^T(n)J(n)\right)^{-1}J^T(n)e(n)$$

# Now, differentiate $e(n)$ with respect to $\boldsymbol{w}(n)$

**Thus**

$$\nabla e(n) = -\boldsymbol{X}^T(n)$$

**Correspondingly, the Jacobian of $e(n)$ is**

$$\mathsf{J}(n) = -\boldsymbol{X}(n)$$

Let us to use the Gaussian-Newton

$$\boldsymbol{w}(n+1) = \boldsymbol{w}(n) - \left(\mathsf{J}^T(n)\mathsf{J}(n)\right)^{-1}\mathsf{J}^T(n)e(n)$$

# Now, differentiate $e(n)$ with respect to $\boldsymbol{w}(n)$

**Thus**

$$\nabla e(n) = -\boldsymbol{X}^T(n)$$

**Correspondingly, the Jacobian of $e(n)$ is**

$$\mathsf{J}(n) = -\boldsymbol{X}(n)$$

**Let us to use the Gaussian-Newton**

$$\boldsymbol{w}(n+1) = \boldsymbol{w}(n) - \left(\mathsf{J}^T(n)\,\mathsf{J}(n)\right)^{-1}\mathsf{J}^T(n)\,e(n)$$

Cinvestav

# Thus

## We have the following

$$\boldsymbol{w}\left(n+1\right)=\boldsymbol{w}\left(n\right)-\left(-\boldsymbol{X}^{T}\left(n\right)\cdot-\boldsymbol{X}\left(n\right)\right)^{-1}\cdot-\boldsymbol{X}^{T}\left(n\right)\left[d\left(n\right)-\boldsymbol{X}\left(n\right)\boldsymbol{w}\left(n\right)\right]$$

# Thus

$$\boldsymbol{w}\left(n+1\right) = \boldsymbol{w}\left(n\right) - \left(-\boldsymbol{X}^{T}\left(n\right) \cdot -\boldsymbol{X}\left(n\right)\right)^{-1} \cdot -\boldsymbol{X}^{T}\left(n\right)\left[d\left(n\right) - \boldsymbol{X}\left(n\right)\boldsymbol{w}\left(n\right)\right]$$

## We have then

$$\boldsymbol{w}\left(n+1\right) = \boldsymbol{w}\left(n\right) + \left(\boldsymbol{X}^{T}\left(n\right)\boldsymbol{X}\left(n\right)\right)^{-1}\boldsymbol{X}^{T}\left(n\right)d\left(n\right) - ...$$
$$\left(\boldsymbol{X}^{T}\left(n\right)\boldsymbol{X}\left(n\right)\right)^{-1}\boldsymbol{X}^{T}\left(n\right)\boldsymbol{X}\left(n\right)\boldsymbol{w}\left(n\right)$$

Thus, we have

$w(n+1) = w(n) + \left(X^T(n)X(n)\right)^{-1}X^T(n)d(n) - w(n)$

$= \left(X^T(n)X(n)\right)^{-1}X^T(n)d(n)$

# Thus

## We have the following

$$\boldsymbol{w}\left(n+1\right)=\boldsymbol{w}\left(n\right)-\left(-\boldsymbol{X}^{T}\left(n\right)\cdot-\boldsymbol{X}\left(n\right)\right)^{-1}\cdot-\boldsymbol{X}^{T}\left(n\right)\left[d\left(n\right)-\boldsymbol{X}\left(n\right)\boldsymbol{w}\left(n\right)\right]$$

## We have then

$$\boldsymbol{w}\left(n+1\right)=\boldsymbol{w}\left(n\right)+\left(\boldsymbol{X}^{T}\left(n\right)\boldsymbol{X}\left(n\right)\right)^{-1}\boldsymbol{X}^{T}\left(n\right)d\left(n\right)-...$$
$$\left(\boldsymbol{X}^{T}\left(n\right)\boldsymbol{X}\left(n\right)\right)^{-1}\boldsymbol{X}^{T}\left(n\right)\boldsymbol{X}\left(n\right)\boldsymbol{w}\left(n\right)$$

## Thus, we have

$$\boldsymbol{w}\left(n+1\right)=\boldsymbol{w}\left(n\right)+\left(\boldsymbol{X}^{T}\left(n\right)\boldsymbol{X}\left(n\right)\right)^{-1}\boldsymbol{X}^{T}\left(n\right)d\left(n\right)-\boldsymbol{w}\left(n\right)$$
$$=\left(\boldsymbol{X}^{T}\left(n\right)\boldsymbol{X}\left(n\right)\right)^{-1}\boldsymbol{X}^{T}\left(n\right)d\left(n\right)$$

# Outline

Cinvestav

# Again Our Error Cost function

## We have

$$J\left(\boldsymbol{w}\right) = \frac{1}{2}e^2\left(n\right)$$

where $e\left(n\right)$ is the error signal measured at time $n$.

Again differentiating against the vector $w$

$$\frac{\partial J\left(w\right)}{\partial w} = e\left(n\right)\frac{\partial e\left(n\right)}{\partial w}$$

LMS algorithm operates with a linear neuron so we may express the error signal as

$$e\left(n\right) = d\left(n\right) - x^T\left(n\right)w\left(n\right) \tag{26}$$

# Again Our Error Cost function

## We have

$$J\left(\boldsymbol{w}\right) = \frac{1}{2}e^2\left(n\right)$$

where $e\left(n\right)$ is the error signal measured at time $n$.

## Again differentiating against the vector $\boldsymbol{w}$

$$\frac{\partial J\left(\boldsymbol{w}\right)}{\partial \boldsymbol{w}} = e\left(n\right)\frac{\partial e\left(n\right)}{\partial \boldsymbol{w}}$$

LMS algorithm operates with a linear neuron so we may express the error signal as

$$e\left(n\right) = d\left(n\right) - x^T\left(n\right)w\left(n\right) \tag{26}$$

# Again Our Error Cost function

**We have**

$$J(\boldsymbol{w}) = \frac{1}{2}e^2(n)$$

where $e(n)$ is the error signal measured at time $n$.

**Again differentiating against the vector $\boldsymbol{w}$**

$$\frac{\partial J(\boldsymbol{w})}{\partial \boldsymbol{w}} = e(n)\frac{\partial e(n)}{\partial \boldsymbol{w}}$$

**LMS algorithm operates with a linear neuron so we may express the error signal as**

$$e(n) = d(n) - \boldsymbol{x}^T(n)\,\boldsymbol{w}(n) \tag{26}$$

# We have

> **Something Notable**
>
> $$\frac{\partial e\left(n\right)}{\partial \boldsymbol{w}} = -\boldsymbol{x}\left(n\right)$$

Then

$$\frac{\partial J\left(w\right)}{\partial w} = -x\left(n\right)e\left(n\right)$$

Using this as an estimate for the gradient vector, we have for the gradient descent

$$w\left(n+1\right) = w\left(n\right) + \eta x\left(n\right)e\left(n\right) \tag{27}$$

# We have

## Something Notable

$$\frac{\partial e(n)}{\partial \boldsymbol{w}} = -\boldsymbol{x}(n)$$

## Then

$$\frac{\partial J(\boldsymbol{w})}{\partial \boldsymbol{w}} = -\boldsymbol{x}(n)\, e(n)$$

Using this as an estimate for the gradient vector, we have for the gradient descent

$$\boldsymbol{w}(n+1) = \boldsymbol{w}(n) + \eta \boldsymbol{x}(n)\, e(n) \tag{27}$$

# We have

$$\frac{\partial e(n)}{\partial \boldsymbol{w}} = -\boldsymbol{x}(n)$$

**Then**

$$\frac{\partial J(\boldsymbol{w})}{\partial \boldsymbol{w}} = -\boldsymbol{x}(n)\, e(n)$$

**Using this as an estimate for the gradient vector, we have for the gradient descent**

$$\widehat{\boldsymbol{w}}(n+1) = \widehat{\boldsymbol{w}}(n) + \eta \boldsymbol{x}(n)\, e(n) \tag{27}$$

Cinvestav

# Remarks

## The feedback loop around the weight vector low-pass filter

- It behaves like a low-pass filter.
  - It passes the low frequency component of the error signal and attenuating its high frequency component.

# Remarks

## The feedback loop around the weight vector low-pass filter

- It behaves like a low-pass filter.
- It passes the low frequency component of the error signal and attenuating its high frequency component.

## Low-Pass filter

# Remarks

## The feedback loop around the weight vector low-pass filter

- It behaves like a low-pass filter.
- It passes the low frequency component of the error signal and attenuating its high frequency component.

## Low-Pass filter

# Actually

> **Thus**
>
> The average time constant of this filtering action is inversely proportional to the learning-rate parameter $\eta$.

# Actually

<div style="background-color:green">Thus</div>

The average time constant of this filtering action is inversely proportional to the learning-rate parameter $\eta$.

<div style="background-color:red">Thus</div>

Assigning a small value to $\eta$, the adaptive process progresses slowly.

# Actually

**Thus**

The average time constant of this filtering action is inversely proportional to the learning-rate parameter $\eta$.

**Thus**

Assigning a small value to $\eta$, the adaptive process progresses slowly.

**Thus**

- More of the past data are remembered by the LMS algorithm.
- Thus, LMS is a more accurate filter

# Actually

**Thus**

The average time constant of this filtering action is inversely proportional to the learning-rate parameter $\eta$.

**Thus**

Assigning a small value to $\eta$, the adaptive process progresses slowly.

**Thus**

- More of the past data are remembered by the LMS algorithm.
- Thus, LMS is a more accurate filter.

# Outline

Cinvestav

# Convergence of the LMS

## This convergence depends on the following points

- The statistical characteristics of the input vector $\boldsymbol{x}(n)$.
- The learning-rate parameter $\eta$.

## Something Notable

However instead using $E[\tilde{w}(n)]$ as $n \to \infty$, we use $E[e^2(n)] \to$ constant as $n \to \infty$

# Convergence of the LMS

## This convergence depends on the following points

- The statistical characteristics of the input vector $\boldsymbol{x}(n)$.
- The learning-rate parameter $\eta$.

## Something Notable

However instead using $E[\widehat{\boldsymbol{w}}(n)]$ as $n \to \infty$, we use $E[e^2(n)] \to$ constant as $n \to \infty$

# To make this analysis practical

## We take the following assumptions

- The successive input vectors $x(1), x(2), ..$ are statistically independent of each other.

- At time step $n$, the input vector $x(n)$ is statistically independent of all previous samples of the desired response, namely $d(1), d(2), ..., d(n-1)$.

- At time step $n$, the desired response $d(n)$ is dependent on $x(n)$, but statistically independent of all previous values of the desired response.

- The input vector $x(n)$ and desired response $d(n)$ are drawn from Gaussiandistributed populations.

# To make this analysis practical

## We take the following assumptions

- The successive input vectors $\boldsymbol{x}(1), \boldsymbol{x}(2), ..$ are statistically independent of each other.
- At time step $n$, the input vector $\boldsymbol{x}(n)$ is statistically independent of all previous samples of the desired response, namely $d(1), d(2), ..., d(n-1)$.
- At time step $n$, the desired response $d(n)$ is dependent on $x(n)$, but statistically independent of all previous values of the desired response.
- The input vector $x(n)$ and desired response $d(n)$ are drawn from Gaussiandistributed populations.

# To make this analysis practical

## We take the following assumptions

- The successive input vectors $\boldsymbol{x}(1), \boldsymbol{x}(2), ..$ are statistically independent of each other.
- At time step $n$, the input vector $\boldsymbol{x}(n)$ is statistically independent of all previous samples of the desired response, namely $d(1), d(2), ..., d(n-1)$.
- At time step $n$, the desired response $d(n)$ is dependent on $x(n)$, but statistically independent of all previous values of the desired response.
- The input vector $x(n)$ and desired response $d(n)$ are drawn from Gaussiandistributed populations.

# To make this analysis practical

## We take the following assumptions

- The successive input vectors $x(1), x(2), ..$ are statistically independent of each other.
- At time step $n$, the input vector $x(n)$ is statistically independent of all previous samples of the desired response, namely $d(1), d(2), ..., d(n-1)$.
- At time step $n$, the desired response $d(n)$ is dependent on $x(n)$, but statistically independent of all previous values of the desired response.
- The input vector $x(n)$ and desired response $d(n)$ are drawn from Gaussiandistributed populations.

# We get the following

## The LMS is convergent in the mean square provided that $\eta$ satisfies

$$0 < \eta < \frac{2}{\lambda_{\max}} \tag{28}$$

Because $\lambda_{\max}$ is the largest eigenvalue of the correlation sample $R_x$

This can be difficult in reality... then we use the trace instead

$$0 < \eta < \frac{2}{\text{trace}\,[R_x]} \tag{29}$$

However, each diagonal element of $R_x$ is equal the mean-squared value of the corresponding of the sensor input

We can re-state the previous condition as

$$0 < \eta < \frac{2}{\text{sum of the mean-square values of the sensor input}} \tag{30}$$

# We get the following

**The LMS is convergent in the mean square provided that $\eta$ satisfies**

$$0 < \eta < \frac{2}{\lambda_{\max}} \tag{28}$$

**Because $\lambda_{\max}$ is the largest eigenvalue of the correlation sample $\boldsymbol{R_x}$**

This can be difficult in reality.... then we use the trace instead

$$0 < \eta < \frac{2}{\text{trace}\,[\boldsymbol{R_x}]} \tag{29}$$

However, each diagonal element of $\boldsymbol{R_x}$ is equal the mean-squared value of the corresponding of the sensor input

We can re-state the previous condition as

$$0 < \eta < \frac{2}{\text{sum of the mean-square values of the sensor input}} \tag{30}$$

# We get the following

The LMS is convergent in the mean square provided that $\eta$ satisfies

$$0 < \eta < \frac{2}{\lambda_{\max}} \tag{28}$$

Because $\lambda_{\max}$ is the largest eigenvalue of the correlation sample $\boldsymbol{R_x}$

This can be difficult in reality.... then we use the trace instead

$$0 < \eta < \frac{2}{\text{trace}\,[\boldsymbol{R_x}]} \tag{29}$$

However, each diagonal element of $\boldsymbol{R_x}$ is equal the mean-squared value of the corresponding of the sensor input

We can re-state the previous condition as

$$0 < \eta < \frac{2}{\text{sum of the mean-square values of the sensor input}} \tag{30}$$

# Virtues and Limitations of the LMS Algorithm

## Virtues

- An important virtue of the LMS algorithm is its simplicity.
- The model is independent and robust to the error (small disturbances = small estimation error).

# Virtues and Limitations of the LMS Algorithm

## Virtues

- An important virtue of the LMS algorithm is its simplicity.
- The model is independent and robust to the error (small disturbances = small estimation error).

# Virtues and Limitations of the LMS Algorithm

## Virtues

- An important virtue of the LMS algorithm is its simplicity.
- The model is independent and robust to the error (small disturbances = small estimation error).

## Not only that, the LMS algorithm is optimal in accordance with the minimax criterion

If you do not know what you are up against, plan for the worst and optimize.

## Primary Limitation

- The slow rate of convergence and sensitivity to variations in the eigenstructure of the input.
- The LMS algorithms requires about 10 times the dimensionality of the input space for convergence.

# Virtues and Limitations of the LMS Algorithm

## Virtues

- An important virtue of the LMS algorithm is its simplicity.
- The model is independent and robust to the error (small disturbances = small estimation error).

## Not only that, the LMS algorithm is optimal in accordance with the minimax criterion

If you do not know what you are up against, plan for the worst and optimize.

## Primary Limitation

- The slow rate of convergence and sensitivity to variations in the eigenstructure of the input.
- The LMS algorithms requires about 10 times the dimensionality of the input space for convergence.

# Virtues and Limitations of the LMS Algorithm

## Virtues

- An important virtue of the LMS algorithm is its simplicity.
- The model is independent and robust to the error (small disturbances = small estimation error).

## Not only that, the LMS algorithm is optimal in accordance with the minimax criterion

If you do not know what you are up against, plan for the worst and optimize.

## Primary Limitation

- The slow rate of convergence and sensitivity to variations in the eigenstructure of the input.
- The LMS algorithms requires about 10 times the dimensionality of the input space for convergence.

# More of this in...

## Simon Haykin

Simon Haykin - Adaptive Filter Theory (3rd Edition)

# Exercises

## We have from NN by Haykin

3.1, 3.2, 3.3, 3.4, 3.5, 3.7 and 3.8

# Outline

Cinvestav

# Objective

## Goal

Correctly classify a series of samples (External applied stimuli) $x_1, x_2, x_3, ..., x_m$ into one of two classes, $C_1$ and $C_2$.

### Output of each input

1. Class $C_1$ output y +1.
2. Class $C_2$ output y -1.

# Objective

**Goal**

Correctly classify a series of samples (External applied stimuli) $x_1, x_2, x_3, ..., x_m$ into one of two classes, $C_1$ and $C_2$.

**Output of each input**

1. Class $C_1$ output y +1.
2. Class $C_2$ output y -1.

# History

## Frank Rosenblatt

The perceptron algorithm was invented in 1957 at the Cornell Aeronautical Laboratory by Frank Rosenblatt.

## Something Notable

Frank Rosenblatt was a Psychologist!!! Working at a militar R&D!!!

## Frank Rosenblatt

He helped to develop the Mark I Perceptron - a new machine based in the connectivity of neural networks!!!

## Some problems with it

- The most important is the impossibility to use the perceptron with a single neuron to solve the XOR problem

# History

## Frank Rosenblatt

The perceptron algorithm was invented in 1957 at the Cornell Aeronautical Laboratory by Frank Rosenblatt.

## Something Notable

Frank Rosenblatt was a Psychologist!!! Working at a militar R&D!!!

## Frank Rosenblatt

He helped to develop the Mark I Perceptron - a new machine based in the connectivity of neural networks!!!
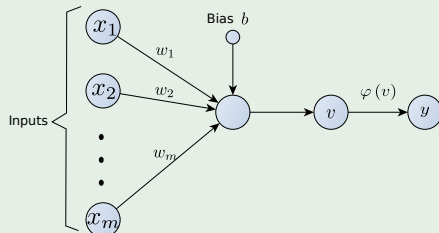
## Some problems with it

- The most important is the impossibility to use the perceptron with a single neuron to solve the XOR problem

# History

## Frank Rosenblatt

The perceptron algorithm was invented in 1957 at the Cornell Aeronautical Laboratory by Frank Rosenblatt.

## Something Notable

Frank Rosenblatt was a Psychologist!!! Working at a militar R&D!!!

## Frank Rosenblatt

He helped to develop the Mark I Perceptron - a new machine based in the connectivity of neural networks!!!

## Some problems with it

- The most important is the impossibility to use the perceptron with a single neuron to solve the XOR problem

# History

## Frank Rosenblatt

The perceptron algorithm was invented in 1957 at the Cornell Aeronautical Laboratory by Frank Rosenblatt.

## Something Notable

Frank Rosenblatt was a Psychologist!!! Working at a militar R&D!!!

## Frank Rosenblatt

He helped to develop the Mark I Perceptron - a new machine based in the connectivity of neural networks!!!

## Some problems with it

- The most important is the impossibility to use the perceptron with a single neuron to solve the XOR problem

# Outline

Cinvestav

# Perceptron: Local Field of a Neuron

## Signal-Flow

# Perceptron: Local Field of a Neuron

Bias $b$

Inputs

$x_1$, $w_1$

$x_2$, $w_2$

$x_m$, $w_m$

$v$, $\varphi(v)$, $y$

## Induced local field of a neuron

$$v = \sum_{i=1}^{m} w_i x_i + b \qquad (31)$$

# Outline

Cinvestav

# Perceptron: One Neuron Structure

## Based in the previous induced local field

In the simplest form of the perceptron there are two decision regions separated by an **hyperplane:**

$$\sum_{i=1}^{m} w_i x_i + b = 0 \tag{32}$$
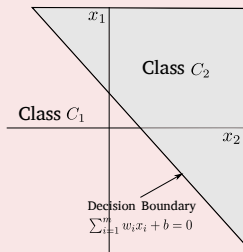
Example with two signals

# Perceptron: One Neuron Structure

## Based in the previous induced local field

In the simplest form of the perceptron there are two decision regions separated by an **hyperplane:**

$$\sum_{i=1}^{m} w_i x_i + b = 0 \tag{32}$$

## Example with two signals

# Outline

Cinvestav

# Deriving the Algorithm

First, you put signals together

$$x(n) = [1, x_1(n), x_2(n), ..., x_m(n)]^T \qquad (33)$$

# Deriving the Algorithm

$$x(n) = [1, x_1(n), x_2(n), ..., x_m(n)]^T \tag{33}$$

**Weights**

$$v(n) = \sum_{i=0}^{m} w_i(n) x_i(n) = \boldsymbol{w}^T(n) \boldsymbol{x}(n) \tag{34}$$

# Deriving the Algorithm

$$x(n) = [1, x_1(n), x_2(n), ..., x_m(n)]^T \tag{33}$$

**Weights**

$$v(n) = \sum_{i=0}^{m} w_i(n) x_i(n) = \boldsymbol{w}^T(n) \boldsymbol{x}(n) \tag{34}$$

**Note IMPORTANT - Perceptron works only if $C_1$ and $C_2$ are linearly separable**



Decision Boundary

Class $C_1$

Class $C_2$

# Rule for Linear Separable Classes

## There must exist a vector $\boldsymbol{w}$

1. $\boldsymbol{w}^T \boldsymbol{x} > 0$ for every input vector $\boldsymbol{x}$ belonging to class $C_1$.

2. $\boldsymbol{w}^T \boldsymbol{x} \leq 0$ for every input vector $\boldsymbol{x}$ belonging to class $C_2$.

# Rule for Linear Separable Classes

1. $\boldsymbol{w}^T \boldsymbol{x} > 0$ for every input vector $\boldsymbol{x}$ belonging to class $C_1$.
2. $\boldsymbol{w}^T \boldsymbol{x} \leq 0$ for every input vector $\boldsymbol{x}$ belonging to class $C_2$.

What is the derivative of $\frac{dv(n)}{dw}$?

$$\frac{dv(n)}{dw} = x(n) \qquad (35)$$

# Rule for Linear Separable Classes

> **There must exist a vector $\boldsymbol{w}$**
>
> 1. $\boldsymbol{w}^T \boldsymbol{x} > 0$ for every input vector $\boldsymbol{x}$ belonging to class $C_1$.
> 2. $\boldsymbol{w}^T \boldsymbol{x} \leq 0$ for every input vector $\boldsymbol{x}$ belonging to class $C_2$.

> **What is the derivative of $\frac{dv(n)}{d\boldsymbol{w}}$?**
>
> $$\frac{dv(n)}{d\boldsymbol{w}} = \boldsymbol{x}(n) \tag{35}$$

# Finally

## No correction is necessary

1. $\boldsymbol{w}(n+1) = \boldsymbol{w}(n)$ if $\boldsymbol{w}^T \boldsymbol{x}(n) > 0$ and $\boldsymbol{x}(n)$ belongs to class $C_1$.

2. $\boldsymbol{w}(n+1) = \boldsymbol{w}(n)$ if and $\boldsymbol{w}^T \boldsymbol{x}(n) \leq 0$ and $\boldsymbol{x}(n) > 0$ belongs to class $C_2$.

# Finally

## No correction is necessary

1. $\boldsymbol{w}(n+1) = \boldsymbol{w}(n)$ if $\boldsymbol{w}^T \boldsymbol{x}(n) > 0$ and $\boldsymbol{x}(n)$ belongs to class $C_1$.

2. $\boldsymbol{w}(n+1) = \boldsymbol{w}(n)$ if and $\boldsymbol{w}^T \boldsymbol{x}(n) \leq 0$ and $\boldsymbol{x}(n) > 0$ belongs to class $C_2$.

# Finally

## No correction is necessary

1. $\boldsymbol{w}(n+1) = \boldsymbol{w}(n)$ if $\boldsymbol{w}^T \boldsymbol{x}(n) > 0$ and $\boldsymbol{x}(n)$ belongs to class $C_1$.
2. $\boldsymbol{w}(n+1) = \boldsymbol{w}(n)$ if and $\boldsymbol{w}^T \boldsymbol{x}(n) \leq 0$ and $\boldsymbol{x}(n) > 0$ belongs to class $C_2$.

## Correction is necessary

1. $\boldsymbol{w}(n+1) = \boldsymbol{w}(n) - \eta(n)\,\boldsymbol{x}(n)$ if $\boldsymbol{w}^T(n)\,\boldsymbol{x}(n) > 0$ and $\boldsymbol{x}(n)$ belongs to class $C_2$.
2. $\boldsymbol{w}(n+1) = \boldsymbol{w}(n) + \eta(n)\,\boldsymbol{x}(n)$ if and $\boldsymbol{w}^T(n)\,\boldsymbol{x}(n) \leq 0$ and $\boldsymbol{x}(n)$ belongs to class $C_1$.

Where $\eta(n)$ is a learning parameter adjusting the learning rate.

# Finally

## No correction is necessary

1. $\boldsymbol{w}(n+1) = \boldsymbol{w}(n)$ if $\boldsymbol{w}^T\boldsymbol{x}(n) > 0$ and $\boldsymbol{x}(n)$ belongs to class $C_1$.
2. $\boldsymbol{w}(n+1) = \boldsymbol{w}(n)$ if and $\boldsymbol{w}^T\boldsymbol{x}(n) \leq 0$ and $\boldsymbol{x}(n) > 0$ belongs to class $C_2$.

## Correction is necessary

1. $\boldsymbol{w}(n+1) = \boldsymbol{w}(n) - \eta(n)\boldsymbol{x}(n)$ if $\boldsymbol{w}^T(n)\boldsymbol{x}(n) > 0$ and $\boldsymbol{x}(n)$ belongs to class $C_2$.
2. $\boldsymbol{w}(n+1) = \boldsymbol{w}(n) + \eta(n)\boldsymbol{x}(n)$ if and $\boldsymbol{w}^T(n)\boldsymbol{x}(n) \leq 0$ and $\boldsymbol{x}(n)$ belongs to class $C_1$.

Where $\eta(n)$ is a learning parameter adjusting the learning rate.

# Finally

## No correction is necessary

1. $\boldsymbol{w}(n+1) = \boldsymbol{w}(n)$ if $\boldsymbol{w}^T \boldsymbol{x}(n) > 0$ and $\boldsymbol{x}(n)$ belongs to class $C_1$.

2. $\boldsymbol{w}(n+1) = \boldsymbol{w}(n)$ if and $\boldsymbol{w}^T \boldsymbol{x}(n) \leq 0$ and $\boldsymbol{x}(n) > 0$ belongs to class $C_2$.

## Correction is necessary

1. $\boldsymbol{w}(n+1) = \boldsymbol{w}(n) - \eta(n)\,\boldsymbol{x}(n)$ if $\boldsymbol{w}^T(n)\,\boldsymbol{x}(n) > 0$ and $\boldsymbol{x}(n)$ belongs to class $C_2$.

2. $\boldsymbol{w}(n+1) = \boldsymbol{w}(n) + \eta(n)\,\boldsymbol{x}(n)$ if and $\boldsymbol{w}^T(n)\,\boldsymbol{x}(n) \leq 0$ and $\boldsymbol{x}(n)$ belongs to class $C_1$.

Where $\eta(n)$ is a learning parameter adjusting the learning rate.

# A little bit on the Geometry

## For Example, $\boldsymbol{w}(n+1) = \boldsymbol{w}(n) - \eta(n)\boldsymbol{x}(n)$

# Outline

Cinvestav

# Under Linear Separability - Convergence happens!!!

## If we assume

Linear Separabilty for the classes $C_1$ and $C_2$.

# Under Linear Separability - Convergence happens!!!

## If we assume

Linear Separabilty for the classes $C_1$ and $C_2$.

## Rosenblatt - 1962

- Let the subsets of training vectors $C_1$ and $C_2$ be linearly separable.
- Let the inputs presented to the perceptron originate from these two subsets.
- The perceptron converges after some $n_0$ iterations, in the sense that is a solution vector for

$$w(n_0) = w(n_0 + 1) = w(n_0 + 2) = \dots \qquad (36)$$

is a solution vector for $n_0 \leq n_{max}$.

# Under Linear Separability - Convergence happens!!!

## If we assume

Linear Separabilty for the classes $C_1$ and $C_2$.

## Rosenblatt - 1962

- Let the subsets of training vectors $C_1$ and $C_2$ be linearly separable.
- Let the inputs presented to the perceptron originate from these two subsets.
- The perceptron converges after some $n_0$ iterations, in the sense that is a solution vector for

$$w(n_0) = w(n_0 + 1) = w(n_0 + 2) = \ldots \tag{36}$$

is a solution vector for $n_0 \leq n_{max}$

# Under Linear Separability - Convergence happens!!!

## If we assume

Linear Separabilty for the classes $C_1$ and $C_2$.

## Rosenblatt - 1962

- Let the subsets of training vectors $C_1$ and $C_2$ be linearly separable.
- Let the inputs presented to the perceptron originate from these two subsets.
- The perceptron converges after some $n_0$ iterations, in the sense that is a solution vector for

$$\boldsymbol{w}(n_0) = \boldsymbol{w}(n_0 + 1) = \boldsymbol{w}(n_0 + 2) = ... \qquad (36)$$

is a solution vector for $n_0 \leq n_{max}$

# Outline

Cinvestav

## Initialization

$$\boldsymbol{w}\left(0\right) = 0 \tag{37}$$

# Proof I - First a Lower Bound for $\|\boldsymbol{w}\left(n+1\right)\|^2$

## Initialization

$$\boldsymbol{w}\left(0\right) = 0 \tag{37}$$

## Now assume for time $n = 1, 2, 3, \ldots$

$$\boldsymbol{w}^T\left(n\right)\boldsymbol{x}\left(n\right) < 0 \tag{38}$$

with $x(n)$ belongs to class $C_1$.

PERCEPTRON INCORRECTLY CLASSIFY THE VECTORS

$x\left(1\right), x\left(2\right), \ldots$

# Proof I - First a Lower Bound for $\|\boldsymbol{w}(n+1)\|^2$

**Initialization**

$$\boldsymbol{w}(0) = 0 \tag{37}$$

**Now assume for time $n = 1, 2, 3, ...$**

$$\boldsymbol{w}^T(n)\,\boldsymbol{x}(n) < 0 \tag{38}$$

with $\boldsymbol{x}(n)$ belongs to class $C_1$.

PERCEPTRON INCORRECTLY CLASSIFY THE VECTORS
$\boldsymbol{x}(1), \boldsymbol{x}(2), ...$

Apply the correction formula

$$\boldsymbol{w}(n+1) = \boldsymbol{w}(n) + \boldsymbol{x}(n) \tag{39}$$

# Proof I - First a Lower Bound for $\|\boldsymbol{w}\left(n+1\right)\|^2$

## Initialization

$$\boldsymbol{w}\left(0\right) = 0 \tag{37}$$

## Now assume for time $n = 1, 2, 3, ...$

$$\boldsymbol{w}^T\left(n\right)\boldsymbol{x}\left(n\right) < 0 \tag{38}$$

with $\boldsymbol{x}(n)$ belongs to class $C_1$.

> PERCEPTRON INCORRECTLY CLASSIFY THE VECTORS
> $\boldsymbol{x}\left(1\right), \boldsymbol{x}\left(2\right), ...$

## Apply the correction formula

$$w\left(n+1\right) = w\left(n\right) + x\left(n\right) \tag{39}$$

# Proof I - First a Lower Bound for $\|\boldsymbol{w}(n+1)\|^2$

$$\boldsymbol{w}(0) = 0 \tag{37}$$

Now assume for time $n = 1, 2, 3, ...$

$$\boldsymbol{w}^T(n)\boldsymbol{x}(n) < 0 \tag{38}$$

with $\boldsymbol{x}(n)$ belongs to class $C_1$.

PERCEPTRON INCORRECTLY CLASSIFY THE VECTORS
$$\boldsymbol{x}(1), \boldsymbol{x}(2), ...$$

Apply the correction formula

$$\boldsymbol{w}(n+1) = \boldsymbol{w}(n) + \boldsymbol{x}(n) \tag{39}$$

# Proof II

Apply the correction iteratively

$$\boldsymbol{w}(n+1) = \boldsymbol{x}(1) + \boldsymbol{x}(2) + ... + \boldsymbol{x}(n) \tag{40}$$

We know that there is a solution $w_0$ (Linear Separability)

$$\alpha = \min_{x(n) \in C} \ w_0^T x(n) \tag{41}$$

Then, we have

$$w_0^T w(n+1) = w_0^T x(1) + w_0^T x(2) + ... + w_0^T x(n) \tag{42}$$

# Proof II

The box below with green header.

**Apply the correction iteratively**

$$\boldsymbol{w}\left(n+1\right) = \boldsymbol{x}\left(1\right) + \boldsymbol{x}\left(2\right) + ... + \boldsymbol{x}\left(n\right) \quad (40)$$

**We know that there is a solution $\boldsymbol{w}_0$(Linear Separability)**

$$\alpha = \min_{\boldsymbol{x}(n)\in C_1} \boldsymbol{w}_0^T \boldsymbol{x}\left(n\right) \quad (41)$$

Then, we have

$$w_0^T w\left(n+1\right) = w_0^T x\left(1\right) + w_0^T x\left(2\right) + ... + w_0^T x\left(n\right) \quad (42)$$

footer

# Proof II

**Apply the correction iteratively**

$$\boldsymbol{w}\left(n+1\right) = \boldsymbol{x}\left(1\right) + \boldsymbol{x}\left(2\right) + ... + \boldsymbol{x}\left(n\right) \tag{40}$$

**We know that there is a solution $\boldsymbol{w}_0$(Linear Separability)**

$$\alpha = \min_{\boldsymbol{x}(n) \in C_1} \boldsymbol{w}_0^T \boldsymbol{x}\left(n\right) \tag{41}$$

**Then, we have**

$$\boldsymbol{w_0^T} \boldsymbol{w}\left(n+1\right) = \boldsymbol{w_0^T} \boldsymbol{x}\left(1\right) + \boldsymbol{w}_0^T \boldsymbol{x}\left(2\right) + ... + \boldsymbol{w_0^T} \boldsymbol{x}\left(n\right) \tag{42}$$

# Proof III

## Apply the correction iteratively

$$w(n+1) = x(1) + x(2) + ... + x(n) \tag{43}$$

We know that there is a solution $w_0$ (Linear Separability)

$$\alpha = \min_{x(n) \in C_1} w_0^T x(n) \tag{44}$$

Then, we have

$$w_0^T w(n+1) = w_0^T x(1) + w_0^T x(2) + ... + w_0^T x(n) \tag{45}$$

# Proof III

Apply the correction iteratively

$$\boldsymbol{w}\left(n+1\right) = \boldsymbol{x}\left(1\right) + \boldsymbol{x}\left(2\right) + ... + \boldsymbol{x}\left(n\right) \tag{43}$$

We know that there is a solution $\boldsymbol{w}_0$(Linear Separability)

$$\alpha = \min_{\boldsymbol{x}(n) \in C_1} \boldsymbol{w}_0^T \boldsymbol{x}\left(n\right) \tag{44}$$

Then, we have

$$w_0^T w\left(n+1\right) = w_0^T x\left(1\right) + w_0^T x\left(2\right) + ... + w_0^T x\left(n\right) \tag{45}$$

# Proof III

**Apply the correction iteratively**

$$\boldsymbol{w}(n+1) = \boldsymbol{x}(1) + \boldsymbol{x}(2) + ... + \boldsymbol{x}(n) \tag{43}$$

**We know that there is a solution $\boldsymbol{w}_0$(Linear Separability)**

$$\alpha = \min_{\boldsymbol{x}(n) \in C_1} \boldsymbol{w}_0^T \boldsymbol{x}(n) \tag{44}$$

**Then, we have**

$$\boldsymbol{w_0^T} \boldsymbol{w}(n+1) = \boldsymbol{w_0^T} \boldsymbol{x}(1) + \boldsymbol{w_0^T} \boldsymbol{x}(2) + ... + \boldsymbol{w_0^T} \boldsymbol{x}(n) \tag{45}$$

Cinvestav

# Proof IV

Thus we use the $\alpha$

$$\boldsymbol{w_0^T} \boldsymbol{w}(n+1) \geq n\alpha \tag{46}$$

Thus using the Cauchy-Schwartz Inequality

$$\left\| \boldsymbol{w_0^T} \right\|^2 \left\| \boldsymbol{w}(n+1) \right\|^2 \geq \left[ \boldsymbol{w_0^T} \boldsymbol{w}(n+1) \right]^2 \tag{47}$$

$\|\cdot\|$ is the Euclidean distance

Thus

$$\left\| \boldsymbol{w_0^T} \right\|^2 \left\| \boldsymbol{w}(n+1) \right\|^2 \geq n^2 \alpha^2$$

$$\left\| \boldsymbol{w}(n+1) \right\|^2 \geq \frac{n^2 \alpha^2}{\left\| \boldsymbol{w_0^T} \right\|^2}$$

# Proof IV

**Thus we use the $\alpha$**

$$\boldsymbol{w_0^T}\boldsymbol{w}\left(n+1\right) \geq n\alpha \tag{46}$$

**Thus using the Cauchy-Schwartz Inequality**

$$\left\|\boldsymbol{w_0^T}\right\|^2 \left\|\boldsymbol{w}\left(n+1\right)\right\|^2 \geq \left[\boldsymbol{w_0^T}\boldsymbol{w}\left(n+1\right)\right]^2 \tag{47}$$

$\|\cdot\|$ is the Euclidean distance.

Thus

$$\left\|w_0^T\right\|^2 \left\|w\left(n+1\right)\right\|^2 \geq n^2\alpha^2$$

$$\left\|w\left(n+1\right)\right\|^2 \geq \frac{n^2\alpha^2}{\left\|w_0^T\right\|^2}$$

# Proof IV

**Thus we use the $\alpha$**

$$\boldsymbol{w_0^T} \boldsymbol{w}\,(n+1) \geq n\alpha \qquad (46)$$

**Thus using the Cauchy-Schwartz Inequality**

$$\left\| \boldsymbol{w_0^T} \right\|^2 \|\boldsymbol{w}\,(n+1)\|^2 \geq \left[ \boldsymbol{w_0^T} \boldsymbol{w}\,(n+1) \right]^2 \qquad (47)$$

$\|\cdot\|$ is the Euclidean distance.

**Thus**

$$
\begin{aligned}
\left\| \boldsymbol{w}_0^T \right\|^2 \|\boldsymbol{w}\,(n+1)\|^2 &\geq n^2\alpha^2 \\
\|\boldsymbol{w}\,(n+1)\|^2 &\geq \frac{n^2\alpha^2}{\left\| \boldsymbol{w}_0^T \right\|^2}
\end{aligned}
$$

# Proof V - Now a Upper Bound for $\|\boldsymbol{w}(n+1)\|^2$

## Now rewrite equation 39

$$\boldsymbol{w}(k+1) = \boldsymbol{w}(k) + \boldsymbol{x}(k) \tag{48}$$

for $k = 1, 2, ..., n$ and $\boldsymbol{x}(k) \in C_1$

Squaring the Euclidean norm of both sides

$$\|\boldsymbol{w}(k+1)\|^2 = \|\boldsymbol{w}(k)\|^2 + \|\boldsymbol{x}(k)\|^2 + 2\boldsymbol{w}^T(k)\boldsymbol{x}(k) \tag{49}$$

Now taking that $\boldsymbol{w}^T(k)\boldsymbol{x}(k) < 0$

$$\|\boldsymbol{w}(k+1)\|^2 \leq \|\boldsymbol{w}(k)\|^2 + \|\boldsymbol{x}(k)\|^2$$
$$\|\boldsymbol{w}(k+1)\|^2 - \|\boldsymbol{w}(k)\|^2 \leq \|\boldsymbol{x}(k)\|^2$$

# Proof V - Now a Upper Bound for $\|\boldsymbol{w}\left(n+1\right)\|^2$

## Now rewrite equation 39

$$\boldsymbol{w}\left(k+1\right) = \boldsymbol{w}\left(k\right) + \boldsymbol{x}\left(k\right) \tag{48}$$

for $k = 1, 2, ..., n$ and $\boldsymbol{x}\left(k\right) \in C_1$

## Squaring the Euclidean norm of both sides

$$\|\boldsymbol{w}\left(k+1\right)\|^2 = \|\boldsymbol{w}\left(k\right)\|^2 + \|\boldsymbol{x}\left(k\right)\|^2 + 2\boldsymbol{w}^T\left(k\right)\boldsymbol{x}\left(k\right) \tag{49}$$

# Proof V - Now a Upper Bound for $\|\boldsymbol{w}(n+1)\|^2$

## Now rewrite equation 39

$$\boldsymbol{w}(k+1) = \boldsymbol{w}(k) + \boldsymbol{x}(k) \tag{48}$$

for $k = 1, 2, ..., n$ and $\boldsymbol{x}(k) \in C_1$

## Squaring the Euclidean norm of both sides

$$\|\boldsymbol{w}(k+1)\|^2 = \|\boldsymbol{w}(k)\|^2 + \|\boldsymbol{x}(k)\|^2 + 2\boldsymbol{w}^T(k)\boldsymbol{x}(k) \tag{49}$$

## Now taking that $\boldsymbol{w}^T(k)\boldsymbol{x}(k) < 0$

$$
\begin{aligned}
\|\boldsymbol{w}(k+1)\|^2 &\leq \|\boldsymbol{w}(k)\|^2 + \|\boldsymbol{x}(k)\|^2 \\
\|\boldsymbol{w}(k+1)\|^2 - \|\boldsymbol{w}(k)\|^2 &\leq \|\boldsymbol{x}(k)\|^2
\end{aligned}
$$

# Proof VI

## Use the telescopic sum

$$\sum_{k=0}^{n} \left[ \|\boldsymbol{w}(k+1)\|^2 - \|\boldsymbol{w}(k)\|^2 \right] \leq \sum_{k=0}^{n} \|\boldsymbol{x}(k)\|^2 \qquad (50)$$

Assume

$$w(0) = 0$$
$$x(0) = 0$$

Thus

$$\|w(n+1)\|^2 \leq \sum_{k=1}^{n} \|x(k)\|^2$$

# Proof VI

$$\sum_{k=0}^{n}\left[\|\boldsymbol{w}(k+1)\|^2 - \|\boldsymbol{w}(k)\|^2\right] \leq \sum_{k=0}^{n}\|\boldsymbol{x}(k)\|^2 \tag{50}$$

**Assume**

$$
\begin{aligned}
\boldsymbol{w}(0) &= \boldsymbol{0} \\
\boldsymbol{x}(0) &= \boldsymbol{0}
\end{aligned}
$$

Thus

$$\|w(n+1)\|^2 \leq \sum_{k=1}^{n} \|x(k)\|^2$$

Cinvestav

87 / 102

# Proof VI

## Use the telescopic sum

$$\sum_{k=0}^{n} \left[ \|\boldsymbol{w}(k+1)\|^2 - \|\boldsymbol{w}(k)\|^2 \right] \leq \sum_{k=0}^{n} \|\boldsymbol{x}(k)\|^2 \tag{50}$$

## Assume

$$\begin{aligned} \boldsymbol{w}(0) &= \boldsymbol{0} \\ \boldsymbol{x}(0) &= \boldsymbol{0} \end{aligned}$$

## Thus

$$\|\boldsymbol{w}(n+1)\|^2 \leq \sum_{k=1}^{n} \|x(k)\|^2$$

Cinvestav

# Proof VII

**Then, we can define a positive number**

$$\beta = \max_{\boldsymbol{x}(k) \in C_1} \|\boldsymbol{x}(k)\|^2 \tag{51}$$

Thus

$$\|x(k+1)\|^2 \leq \sum_{i=1}^{n} \|x(k)\|^2 \leq n\beta$$

Thus, we satisfies the equations only when exists a $n_{max}$ (Using Our Sandwich )

$$\frac{n_{max}^2 \alpha^2}{\|w_0\|^2} = n_{max}\beta \tag{52}$$

# Proof VII

## Then, we can define a positive number

$$\beta = \max_{\boldsymbol{x}(k) \in C_1} \|\boldsymbol{x}(k)\|^2 \tag{51}$$

## Thus

$$\|\boldsymbol{w}(k+1)\|^2 \leq \sum_{k=1}^{n} \|x(k)\|^2 \leq n\beta$$

Thus, we satisfies the equations only when exists a $n_{max}$ (Using Our Sandwich )

$$\frac{n_{max}^2 \alpha^2}{\|w_0\|^2} = n_{max}\beta \tag{52}$$

# Proof VII

Then, we can define a positive number

$$\beta = \max_{\boldsymbol{x}(k) \in C_1} \|\boldsymbol{x}(k)\|^2 \qquad (51)$$

Thus

$$\|\boldsymbol{w}(k+1)\|^2 \leq \sum_{k=1}^{n} \|x(k)\|^2 \leq n\beta$$

Thus, we satisfies the equations only when exists a $n_{max}$ (Using Our Sandwich )

$$\frac{n_{max}^2 \alpha^2}{\|\boldsymbol{w}_0\|^2} = n_{max}\beta \qquad (52)$$

Cinvestav

# Proof VIII

## Solving

$$n_{max} = \frac{\beta \left\| \boldsymbol{w}_0 \right\|^2}{\alpha^2} \tag{53}$$

## Thus

For $\eta(n) = 1$ for all $n$, $w(0) = 0$ and a solution vector $w_0$.

- The rule for adaptying the synaptic weights of the perceptron must terminate after at most $n_{max}$ steps.

## In addition

Because $w_0$ the solution is not unique.

# Proof VIII

## Solving

$$n_{max} = \frac{\beta \left\| \boldsymbol{w}_0 \right\|^2}{\alpha^2} \tag{53}$$

## Thus

For $\eta(n) = 1$ for all n, $\boldsymbol{w}(0) = \boldsymbol{0}$ and a solution vector $\boldsymbol{w}_0$:

- The rule for adaptying the synaptic weights of the perceptron must terminate after at most $n_{max}$ steps.

In addition

Because $w_0$ the solution is not unique.

# Proof VIII

## Solving

$$n_{max} = \frac{\beta \left\| \boldsymbol{w}_0 \right\|^2}{\alpha^2} \tag{53}$$

## Thus

For $\eta(n) = 1$ for all n, $\boldsymbol{w}(0) = \boldsymbol{0}$ and a solution vector $\boldsymbol{w}_0$:

- The rule for adaptying the synaptic weights of the perceptron must terminate after at most $n_{max}$ steps.

## In addition

**Because $w_0$ the solution is not unique.**

# Outline

Cinvestav

# Algorithm Using Error-Correcting

Now, if we use the $\frac{1}{2}e_k(n)^2$

We can actually simplify the rules and the final algorithm!!!

# Algorithm Using Error-Correcting

## Now, if we use the $\frac{1}{2}e_k(n)^2$

We can actually simplify the rules and the final algorithm!!!

## Thus, we have the following Delta Value

$$\Delta \boldsymbol{w}(n) = \eta\left(\left(d_j - y_j(n)\right)\right)\boldsymbol{x}(n) \tag{54}$$

# Outline

Cinvestav

# We could use the previous Rule

## In order to generate an algorithm

However, you need classes that are linearly separable!!!

Therefore, we can use a more generale Gradient Descent Rule

To obtain an algorithm to the best separation hyperplane!!!

# We could use the previous Rule

## In order to generate an algorithm

However, you need classes that are linearly separable!!!

## Therefore, we can use a more generals Gradient Descent Rule

To obtain an algorithm to the best separation hyperplane!!!

# Gradient Descent Algorithm

# Gradient Descent Algorithm

## Algorithm - Off-line/Batch Learning

1. Set $n = 0$.
2. Set $d_j = \begin{cases} +1 & \text{if } \boldsymbol{x}_j\left(n\right) \in \text{Class 1} \\ -1 & \text{if } \boldsymbol{x}_j\left(n\right) \in \text{Class 2} \end{cases}$ for all $j = 1, 2, ..., m$.
3. Initialize the weights, $\boldsymbol{w}^T = \left(w_1\left(n\right), w_2\left(n\right), ..., w_n\left(n\right)\right)$.
   - Weights may be initialized to 0 or to a small random value.
4. Initialize Dummy outputs so you can enter loop $\boldsymbol{y}^t = \langle y_1\left(n\right)., y_2\left(n\right), ..., y_m\left(n\right)\rangle$
5. Initialize Stopping error $\epsilon > 0$.
6. Initialize learning error $\eta$.
7. While $\frac{1}{m} \sum_{j=1}^{m} \|d_j - y_j\left(n\right)\| > \epsilon$
   - For each sample $(x_j, d_j)$ for $j = 1, ..., m$:
     - ⋆ Calculate output $y_j = \varphi\left(\boldsymbol{w}^T\left(n\right) \cdot \boldsymbol{x}_j\right)$
     - ⋆ Update weights $w_i\left(n+1\right) = w_i\left(n\right) + \eta\left(d_j - y_j\left(n\right)\right)x_{ij}$.
   - $n = n + 1$

# Nevertheless

We have the following problem

$$\epsilon > 0 \tag{55}$$

Thus

Convergence to the best linear separation is a tweaking business!!!

# Nevertheless

## We have the following problem

$$\epsilon > 0 \tag{55}$$

## Thus...

Convergence to the best linear separation is a tweaking business!!!

# Outline

Cinvestav

# However, if we limit our features!!!

**The Winnow Algorithm!!!**

It converges even with no-linear separability.

**Feature Vector**

A Boolean-valued features $X = \{0,1\}^d$

**Weight Vector $w$**

1. $w^t = (w_1, w_2, \ldots, w_n)$ for all $w_i \in \mathbb{R}$
2. For all $i$, $w_i \geq 0$.

# However, if we limit our features!!!

## The Winnow Algorithm!!!

It converges even with no-linear separability.

## Feature Vector

A Boolean-valued features $X = \{0, 1\}^d$

## Weight Vector $w$

1. $w^t = (w_1, w_2, \dots, w_n)$ for all $w_i \in \mathbb{R}$
2. For all $i$, $w_i \geq 0$.

# However, if we limit our features!!!

## The Winnow Algorithm!!!

It converges even with no-linear separability.

## Feature Vector

A Boolean-valued features $X = \{0,1\}^d$

## Weight Vector $\boldsymbol{w}$

1. $\boldsymbol{w}^t = (w_1, w_2, ..., w_p)$ for all $w_i \in \mathbb{R}$
2. For all $i$, $w_i \geq 0$.

# Classification Scheme

## We use a specific $\theta$

1. $w^T x \geq \theta \Rightarrow$ positive classification Class 1 if $x \in$ Class 1
2. $w^T x < \theta \Rightarrow$ negative classification Class 2 if $x \in$ Class 2

## Rule

We use two possible Rules for training!!! With a learning rate of $\alpha > 1$.

## Rule 1

- When misclassifying a positive training example $x \in$ Class 1 i.e. $w^T x < \theta$

$$\forall x_i = 1 : w_i \leftarrow \alpha w_i \qquad (56)$$

# Classification Scheme

## We use a specific $\theta$

**1** $\boldsymbol{w}^T\boldsymbol{x} \geq \theta \Rightarrow$ positive classification Class 1 if $\boldsymbol{x} \in$ Class 1

**2** $\boldsymbol{w}^T\boldsymbol{x} < \theta \Rightarrow$ negative classification Class 2 if $\boldsymbol{x} \in$ Class 2

## Rule

We use two possible Rules for training!!! With a learning rate of $\alpha > 1$.

## Rule 1

- When misclassifying a positive training example $x \in$ Class 1 i.e. $w^T x < \theta$

$$\forall x_i = 1 : w_i \leftarrow \alpha w_i \tag{56}$$

# Classification Scheme

## We use a specific $\theta$

① $\boldsymbol{w}^T \boldsymbol{x} \geq \theta \Rightarrow$ positive classification Class 1 if $\boldsymbol{x} \in$ Class 1

② $\boldsymbol{w}^T \boldsymbol{x} < \theta \Rightarrow$ negative classification Class 2 if $\boldsymbol{x} \in$ Class 2

## Rule

We use two possible Rules for training!!! With a learning rate of $\alpha > 1$.

## Rule 1

- When misclassifying a positive training example $\boldsymbol{x} \in$ Class 1 i.e. $\boldsymbol{w}^T \boldsymbol{x} < \theta$

$$\forall x_i = 1 : w_i \leftarrow \alpha w_i \tag{56}$$

# Classification Scheme

## Rule 2

- When misclassifying a negative training example $\boldsymbol{x} \in$ Class 2 i.e. $\boldsymbol{w}^T \boldsymbol{x} \geq \theta$

$$\forall x_i = 1 : w_i \leftarrow \frac{w_i}{\alpha} \tag{57}$$

## Rule 3

- If samples are correctly classified do nothing!!!

# Classification Scheme

## Rule 2

- When misclassifying a negative training example $\boldsymbol{x} \in$ Class 2 i.e. $\boldsymbol{w}^T \boldsymbol{x} \geq \theta$

$$\forall x_i = 1 : w_i \leftarrow \frac{w_i}{\alpha} \qquad (57)$$

## Rule 3

- If samples are correctly classified do nothing!!!

# Properties of Winnow

- If there are many irrelevant variables Winnow is better than the Perceptron.

Drawback

- Sensitive to the learning rate $\alpha$

# Properties of Winnow

## Property

- If there are many irrelevant variables Winnow is better than the Perceptron.

## Drawback

- Sensitive to the learning rate $\alpha$.

# The Pocket Algorithm

## A variant of the Perceptron Algorithm

It was suggested by Geman et al. in

- "Perceptron based learning algorithms." IEEE Transactions on Neural Networks, Vol. 1(2), pp. 179–191, 1990.
- It converges to an optimal solution even if the linear separability is not fulfilled.

# The Pocket Algorithm

## A variant of the Perceptron Algorithm

It was suggested by Geman et al. in

- "Perceptron based learning algorithms," IEEE Transactions on Neural Networks,Vol. 1(2), pp. 179–191, 1990.

# The Pocket Algorithm

## A variant of the Perceptron Algorithm

It was suggested by Geman et al. in

- "Perceptron based learning algorithms," IEEE Transactions on Neural Networks, Vol. 1(2), pp. 179–191, 1990.
- It converges to an optimal solution even if the linear separability is not fulfilled.

# Finally

## It consists of the following steps

1. Initialize the weight vector $w(0)$ in a random way.

2. Define a storage pocket vector $w_s$ and a history counter $h_s$ to zero for the same pocket vector

3. At the $t^{th}$ iteration step compute the update $w(n+1)$ using the Perceptron rule.

4. Use the update weight to find the number $h$ of samples correctly classified

5. If at any moment $h > h_s$ replace $w_s$ with $w(n+1)$ and $h_s$ with $h$

6. Keep iterating to 3 until convergence.

7. Return $w_s$.

# Finally

## It consists of the following steps

1. Initialize the weight vector $\boldsymbol{w}(0)$ in a random way.
2. Define a storage pocket vector $\boldsymbol{w}_s$ and a history counter $h_s$ to zero for the same pocket vector.
3. At the $i^{th}$ iteration step compute the update $w(n+1)$ using the Perceptron rule.
4. Use the update weight to find the number $h$ of samples correctly classified
5. If at any moment $h > h_s$ replace $w_s$ with $w(n+1)$ and $h_s$ with $h$
6. Keep iterating to 3 until convergence.
7. Return $w_s$.

# Finally

## It consists of the following steps

1. Initialize the weight vector $\boldsymbol{w}(0)$ in a random way.

2. Define a storage pocket vector $\boldsymbol{w}_s$ and a history counter $h_s$ to zero for the same pocket vector.

3. At the $i^{th}$ iteration step compute the update $\boldsymbol{w}(n+1)$ using the Perceptron rule.

4. Use the update weight to find the number $h$ of samples correctly classified

5. If at any moment $h > h_s$ replace $w_s$ with $w(n+1)$ and $h_s$ with $h$

6. Keep iterating to 3 until convergence.

7. Return $w_s$.

# Finally

## It consists of the following steps

1. Initialize the weight vector $\boldsymbol{w}(0)$ in a random way.
2. Define a storage pocket vector $\boldsymbol{w}_s$ and a history counter $h_s$ to zero for the same pocket vector.
3. At the $i^{th}$ iteration step compute the update $\boldsymbol{w}(n+1)$ using the Perceptron rule.
4. Use the update weight to find the number $h$ of samples correctly classified.
5. If at any moment $h > h_s$ replace $w_s$ with $w(n+1)$ and $h_s$ with $h$
6. Keep iterating to 3 until convergence.
7. Return $w_s$

# Finally

## It consists of the following steps

1. Initialize the weight vector $\boldsymbol{w}(0)$ in a random way.

2. Define a storage pocket vector $\boldsymbol{w}_s$ and a history counter $h_s$ to zero for the same pocket vector.

3. At the $i^{th}$ iteration step compute the update $\boldsymbol{w}(n+1)$ using the Perceptron rule.

4. Use the update weight to find the number $h$ of samples correctly classified.

5. If at any moment $h > h_s$ replace $\boldsymbol{w}_s$ with $\boldsymbol{w}(n+1)$ and $h_s$ with $h$

6. Keep iterating to 3 until convergence.

7. Return $w_s$.

# Finally

## It consists of the following steps

1. Initialize the weight vector $\boldsymbol{w}(0)$ in a random way.
2. Define a storage pocket vector $\boldsymbol{w}_s$ and a history counter $h_s$ to zero for the same pocket vector.
3. At the $i^{th}$ iteration step compute the update $\boldsymbol{w}(n+1)$ using the Perceptron rule.
4. Use the update weight to find the number $h$ of samples correctly classified.
5. If at any moment $h > h_s$ replace $\boldsymbol{w}_s$ with $\boldsymbol{w}(n+1)$ and $h_s$ with $h$
6. Keep iterating to 3 until convergence.

# Finally

## It consists of the following steps

1. Initialize the weight vector $\boldsymbol{w}(0)$ in a random way.

2. Define a storage pocket vector $\boldsymbol{w}_s$ and a history counter $h_s$ to zero for the same pocket vector.

3. At the $i^{th}$ iteration step compute the update $\boldsymbol{w}(n+1)$ using the Perceptron rule.

4. Use the update weight to find the number $h$ of samples correctly classified.

5. If at any moment $h > h_s$ replace $\boldsymbol{w}_s$ with $\boldsymbol{w}(n+1)$ and $h_s$ with $h$

6. Keep iterating to 3 until convergence.

7. Return $\boldsymbol{w}_s$.