

# Introduction to Neural Networks and Deep Learning

## Multilayer Perceptron

Andres Mendez-Vazquez

September 28, 2019

# Outline

## 1 Multi-Layer Perceptron

- The XOR Problem
- Architecture
- The Forward and Backward Propagation
- The Quadratic Learning Error Function
- Hidden-to-Output Weights
- Input-to-Hidden Weights
- Total Training Error
- About Stopping Criteria
- Final Basic Batch Algorithm

## 2 Implementing Using Matrix Operations

- Using Matrix Operations to Simplify the Pseudo-Code
- Generating the Output  $z_k$
- Generating the Weights from Hidden to Output Layer
- Generating the Weights from Input to Hidden Layer

## 3 Policies for Multilayer Perceptron

- Maximizing information content
- Activation Functions
- Target Values
- Normalizing the inputs
- Virtues and limitations of Back-Propagation Algorithm

## 4 The Universal Approximation Theorem

- Introduction
- Topology
- Compactness
- About Density in a Topology
- Hausdorff Space
- Measure
- Discriminatory Functions
- Universal Representation Theorem

# Do you remember?

The Perceptron has the following problem

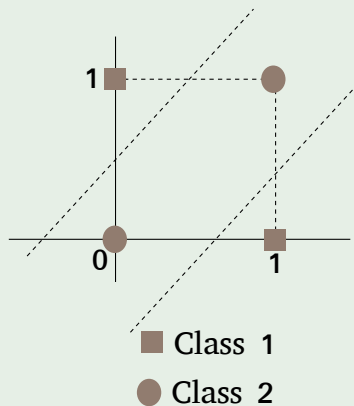
Given that the perceptron is a linear classifier

It is clear that

It will never be able to classify stuff that is not linearly separable

# Example: XOR Problem

## The Problem



# The Perceptron cannot solve it

Because

The perceptron is a linear classifier!!!

Thus

Something needs to be done!!!

Maybe

Add an extra layer!!!

## A little bit of history

### It was first cited by Vapnik

Vapnik cites (Bryson, A.E.; W.F. Denham; S.E. Dreyfus. Optimal programming problems with inequality constraints. I: Necessary conditions for extremal solutions. AIAA J. 1, 11 (1963) 2544-2550 [1]) as the first publication of the backpropagation algorithm in his book "Support Vector Machines."

### It was first used by

Arthur E. Bryson and Yu-Chi Ho described it as a multi-stage dynamic system optimization method in 1969.

### However

It was not until 1974 and later, when applied in the context of neural networks and through the work of Paul Werbos, David E. Rumelhart, Geoffrey E. Hinton and Ronald J. Williams that it gained recognition [2, 3, 4].

# Then

## Something Notable

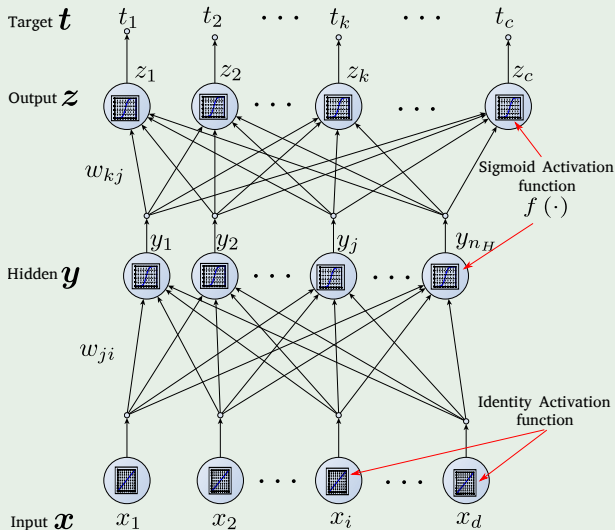
It led to a “renaissance” in the field of artificial neural network research.

## Nevertheless

During the 2000s it fell out of favor but has returned again in the 2010s, now able to train much larger networks using huge modern computing power such as GPUs.

# Multi-Layer Perceptron (MLP) [5]

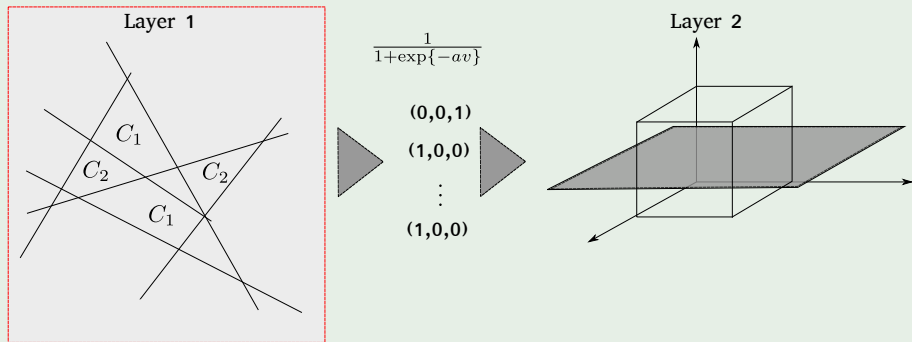
## Multi-Layer Architecture





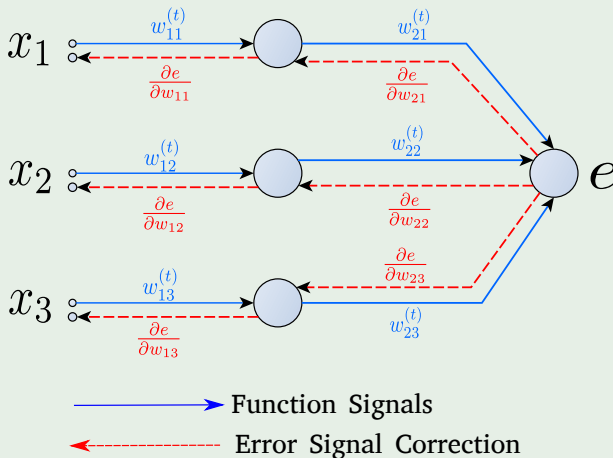
# What do we believe happens in the MLP

We have something like this



# Information Flow

We have the following information flow



## Based in the Previous Idea

People noticed that you require a way to get the info

- In order to build the error signal...

This can be done by simply

- Evaluating the function composition to get the error

$$\mathbf{x} \longrightarrow f_n \circ f_{n-1} \circ \cdots \circ f_1 (\mathbf{x}) = y \longrightarrow e = t - y$$

# Now, if we want to use the gradient descent

## We have a small problem

- If you have the derivative of the cost function by weights that are deep into the network from the cost function

$$\frac{\partial F(\mathbf{x})}{\partial w} = \frac{\partial (t - f_n \circ f_{n-1} \circ \dots \circ f_1(\mathbf{x}))^2}{\partial w}$$

## Therefore, we need to use the chain rule of derivatives

- To reach those functions to build the gradient descent for weights deep into the network

# Chain Rule

## Definition

- Given a composition of differentiable functions

$$F(\mathbf{x}) = f_n \circ f_{n-1} \circ \cdots \circ f_1(\mathbf{x})$$

Then

$$\frac{\partial F(\mathbf{x})}{\partial \mathbf{x}} = \frac{\partial f_n}{\partial f_{n-1}} \times \frac{\partial f_{n-1}}{\partial f_{n-2}} \times \frac{\partial f_{n-2}}{\partial f_{n-3}} \times \cdots \times \frac{\partial f_1}{\partial \mathbf{x}}$$

## Therefore

If we derivate with respect to weights at each level, we are actually back-propagating the error

$$\begin{aligned}J(\mathbf{W}) &= \frac{1}{2} (t - z)^2 \\z &= f(\mathbf{w}_1^T \mathbf{y}) \\ \mathbf{y} &= g(\mathbf{w}_2^T \mathbf{x})\end{aligned}$$

Therefore, for the first layer

$$\frac{\partial J(\mathbf{W})}{\partial \mathbf{w}_1} = \frac{\partial J(\mathbf{W})}{\partial z} \times \frac{\partial f(\mathbf{w}_1^T \mathbf{y})}{\partial \mathbf{w}_1^T \mathbf{y}} \times \frac{\partial \mathbf{w}_1^T \mathbf{y}}{\partial \mathbf{w}_1}$$

What about the second layer?

- We go to the blackboard!!!

# The Quadratic Learning Error Function

Cost Function our well know error at **pattern**  $m$

$$J(m) = \frac{1}{2} e_k^2(m) \quad (1)$$

Delta Rule or Widrow-Hoff Rule

$$\Delta w_{kj}(m) = -\eta e_k(m) x_j(m) \quad (2)$$

Actually this is know as Gradient Descent

$$w_{kj}(m+1) = w_{kj}(m) + \Delta w_{kj}(m) \quad (3)$$

# Back-propagation

## Setup

Let  $t_k$  be the  $k$ -th target (or desired) output and  $z_k$  be the  $k$ -th computed output with  $k = 1, \dots, d$  and  $\mathbf{w}$  represents all the weights of the network

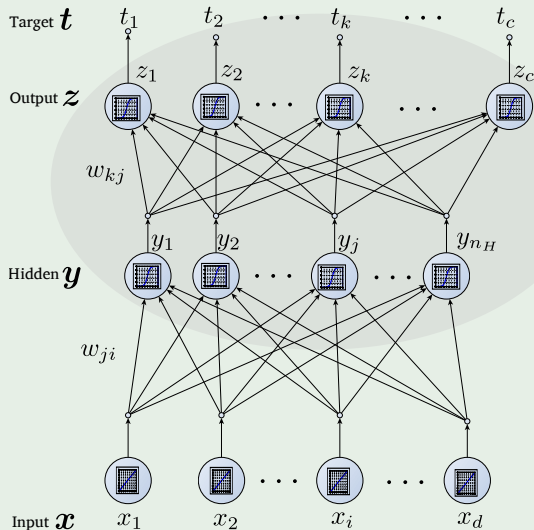
## Training Error for a single Pattern or Sample!!!

$$J(\mathbf{w}) = \frac{1}{2} \sum_{k=1}^c (t_k - z_k)^2 = \frac{1}{2} \|\mathbf{t} - \mathbf{z}\|^2 \quad (4)$$



# Multilayer Architecture

## Multilayer Architecture: hidden-to-output weights



# Observation about the activation function

Hidden Output is equal to

$$y_j = f \left( \sum_{i=1}^d w_{ji} x_i \right)$$

Output is equal to

$$z_k = f \left( \sum_{j=1}^{y_{n_H}} w_{kj} y_j \right)$$

# Hidden-to-Output Weights

$net_k$

- It describes how the overall error changes with the activation of the unit's net:

$$net_k = \sum_{j=1}^{y_{n_H}} w_{kj} y_j = \mathbf{w}_k^T \cdot \mathbf{y} \quad (5)$$

Which is composed with an activation function  $f$

$$z_k = f(net_k) \quad (6)$$

Thus

$$\frac{\partial z_k}{\partial net_k} = f'(net_k) \quad (7)$$

Now, we have the cost function  $J$

Thus, we can apply the chain rule to the cost function

$$\frac{\partial J(z_k)}{\partial w_{kj}} = \frac{\partial J(z_k)}{\partial net_k} \cdot \frac{\partial net_k}{\partial w_{kj}} = -\delta_k \cdot \frac{\partial net_k}{\partial w_{kj}} \quad (8)$$

Still, we need to apply the same for  $\frac{\partial J(z_k)}{\partial net_k}$

$$\frac{\partial J(z_k)}{\partial net_k} = \frac{\partial J(z_k)}{\partial z_k} \cdot \frac{\partial z_k}{\partial net_k} = -(t_k - z_k) f'(net_k) \quad (9)$$

Then

We create a new variable

$$\delta_k = -\frac{\partial J(\mathbf{z}_k)}{\partial net_k} = (t_k - z_k) f'(net_k)$$

Not only that, but we need

$$\frac{\partial net_k}{\partial w_{kj}}$$

Since  $net_k = \mathbf{w}_k^T \cdot \mathbf{y}$  therefore:

$$\frac{\partial net_k}{\partial w_{kj}} = y_j \quad (10)$$

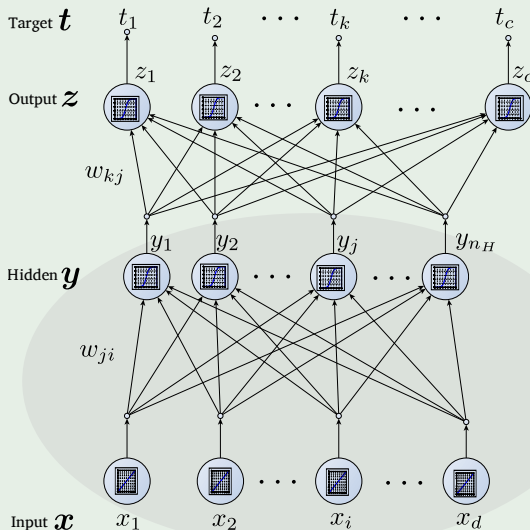
## Finally

The weight update (or learning rule) for the hidden-to-output weights is:

$$\Delta w_{kj} = \eta \delta_k y_j = \eta (t_k - z_k) f'(net_k) y_j \quad (11)$$

# Multi-Layer Architecture

## Going deeper into the network



## Input-to-Hidden Weights

Chain rule on the Input-to-Hidden weights

$$\frac{\partial J}{\partial w_{ji}} = \frac{\partial J}{\partial y_j} \cdot \frac{\partial y_j}{\partial net_j} \cdot \frac{\partial net_j}{\partial w_{ji}} \quad (12)$$

Thus

$$\begin{aligned} \frac{\partial J}{\partial y_j} &= \frac{\partial}{\partial y_j} \left[ \frac{1}{2} \sum_{k=1}^c (t_k - z_k)^2 \right] \\ &= - \sum_{k=1}^c (t_k - z_k) \frac{\partial z_k}{\partial y_j} \\ &= - \sum_{k=1}^c (t_k - z_k) \frac{\partial z_k}{\partial net_k} \cdot \frac{\partial net_k}{\partial y_j} \\ &= - \sum_{k=1}^c (t_k - z_k) \frac{\partial f(net_k)}{\partial net_k} \cdot w_{kj} \end{aligned}$$



# Input-to-Hidden Weights

Finally

$$\frac{\partial J}{\partial y_j} = - \sum_{k=1}^c (t_k - z_k) f'(net_k) \cdot w_{kj} \quad (13)$$

Remember

$$\delta_k = - \frac{\partial J}{\partial net_k} = (t_k - z_k) f'(net_k) \quad (14)$$

What is  $\frac{\partial y_j}{\partial net_j}$ ?

First

$$net_j = \sum_{i=1}^d w_{ji} x_i = \mathbf{w}_j^T \cdot \mathbf{x} \quad (15)$$

Then

$$y_j = f(net_j)$$

Then

$$\frac{\partial y_j}{\partial net_j} = \frac{\partial f(net_j)}{\partial net_j} = f'(net_j)$$

Then, we can define  $\delta_j$

By defining the sensitivity for a hidden unit:

$$\delta_j = f'(net_j) \sum_{k=1}^c w_{kj} \delta_k \quad (16)$$

Which means that:

“The sensitivity at a hidden unit is simply the sum of the individual sensitivities at the output units weighted by the **hidden-to-output weights**  $w_{kj}$ ; all multiplied by  $f'(net_j)$ ”

What about  $\frac{\partial net_j}{\partial w_{ji}}$ ?

We have that

$$\frac{\partial net_j}{\partial w_{ji}} = \frac{\partial \mathbf{w}_j^T \cdot \mathbf{x}}{\partial w_{ji}} = \frac{\partial \sum_{i=1}^d w_{ji} x_i}{\partial w_{ji}} = x_i$$

## Finally

The learning rule for the input-to-hidden weights is:

$$\Delta w_{ji} = \eta x_i \delta_j = \eta \left[ \sum_{k=1}^c w_{kj} \delta_k \right] f'(net_j) x_i \quad (17)$$

Basically, the entire training process has the following steps

### Initialization

Assuming that no prior information is available, pick the synaptic weights and thresholds

### Forward Computation

Compute the induced function signals of the network by proceeding forward through the network, layer by layer.

### Backward Computation

Compute the local gradients of the network.

### Finally

Adjust the weights!!!

## However, you have a problem

### Problems with Hidden Layers

- 1 Increase complexity of Training
- 2 It is necessary to think about “Long and Narrow” network vs “Short and Fat” network.

### Intuition for a One Hidden Layer

- 1 For every input case of region, that region can be delimited by hyperplanes on all sides using hidden units on the first hidden layer.
- 2 A hidden unit in the second layer than ANDs them together to bound the region.

### Advantages

It has been proven that an MLP with one hidden layer can learn any nonlinear function of the input.

## Now, Calculating Total Change

We have for that

The Total Training Error is the sum over the errors of  $N$  individual patterns

The Total Training Error

$$J = \sum_{p=1}^N J_p = \frac{1}{2} \sum_{p=1}^N \sum_{k=1}^d (t_k^p - z_k^p)^2 = \frac{1}{2} \sum_{p=1}^n \|\mathbf{t}^p - \mathbf{z}^p\|^2 \quad (18)$$



# About the Total Training Error

## Remarks

- A weight update may reduce the error on the single pattern being presented but can increase the error on the full training set.
- However, given a large number of such individual updates, the total error of equation (18) decreases.

## Now, we want the training to stop

### Therefore

It is necessary to have a way to stop when the change of the weights is enough!!!

### A simple way to stop the training

- The algorithm terminates when the change in the criterion function  $J(\mathbf{w})$  is smaller than some preset value  $\Theta$ .

$$\Delta J(\mathbf{w}) = |J(\mathbf{w}(t+1)) - J(\mathbf{w}(t))| \quad (19)$$

- There are other stopping criteria that lead to better performance than this one.

## Other Stopping Criteria

### Norm of the Gradient

The back-propagation algorithm is considered to have converged when the Euclidean norm of the gradient vector reaches a sufficiently small gradient threshold.

$$\|\nabla_w J(m)\| < \Theta \quad (20)$$

### Rate of change in the average error per epoch

The back-propagation algorithm is considered to have converged when the absolute rate of change in the average squared error per epoch is sufficiently small.

$$\left| \frac{1}{N} \sum_{p=1}^N J_p \right| < \Theta \quad (21)$$

# About the Stopping Criteria

## Observations

- ① Before training starts, the error on the training set is high.
  - ▶ Through the learning process, the error becomes smaller.
- ② The error per pattern depends on the amount of training data and the expressive power (such as the number of weights) in the network.
- ③ The average error on an independent test set is always higher than on the training set, and it can decrease as well as increase.
- ④ A validation set is used in order to decide when to stop training.
  - ▶ We do not want to over-fit the network and decrease the power of the classifier generalization “we stop training at a minimum of the error on the validation set”

# Some More Terminology

## Epoch

As with other types of backpropagation, 'learning' is a supervised process that occurs with each cycle or 'epoch' through a forward activation flow of outputs, and the backwards error propagation of weight adjustments.

## In our case

I am using the batch sum of all correcting weights to define that epoch.

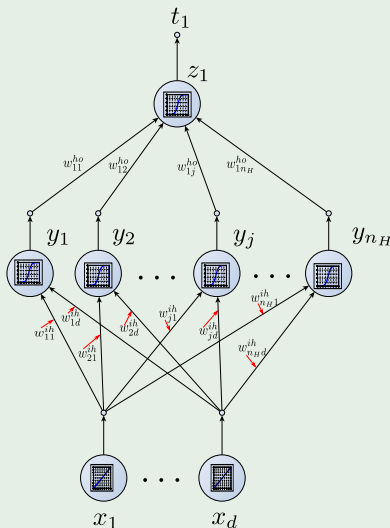
# Final Basic Batch Algorithm

Perceptron( $\mathbf{X}$ )

- 1 Initialize random  $\mathbf{w}$ , number of hidden units  $n_H$ , number of outputs  $z$ , stopping criterion  $\Theta$ , learning rate  $\eta$ , epoch  
 $m = 0$
- 2 do
- 3      $m = m + 1$
- 4     for  $s = 1$  to  $N$
- 5          $\mathbf{x}(m) = \mathbf{X}(:, s)$
- 6         for  $k = 1$  to  $c$
- 7              $\delta_k = (t_k - z_k) f'(\mathbf{w}_k^T \cdot \mathbf{y})$
- 8             for  $j = 1$  to  $n_H$
- 9                  $net_j = \mathbf{w}_j^T \cdot \mathbf{x}; y_j = f(net_j)$
- 10                  $w_{kj}(m) = w_{kj}(m) + \eta \delta_k y_j(m)$
- 11             for  $j = 1$  to  $n_H$
- 12                  $\delta_j = f'(net_j) \sum_{k=1}^c w_{kj} \delta_k$
- 13                 for  $i = 1$  to  $d$
- 14                      $w_{ji}(m) = w_{ji}(m) + \eta \delta_j x_i(m)$
- 15     until  $\|\nabla_{\mathbf{w}} J(m)\| < \Theta$
- 16 return  $\mathbf{w}(m)$

## Example of Architecture to be used

Given the following Architecture and assuming  $N$  samples



## Generating the output $z_k$

Given the input

$$\mathbf{X} = \begin{bmatrix} \mathbf{x}_1 & \mathbf{x}_2 & \cdots & \mathbf{x}_N \end{bmatrix} \quad (22)$$

Where

$\mathbf{x}_i$  is a vector of features

$$\mathbf{x}_i = \begin{pmatrix} x_{1i} \\ x_{2i} \\ \vdots \\ x_{di} \end{pmatrix} \quad (23)$$



Therefore

We must have the following matrix for the input to hidden inputs

$$W_{IH} = \begin{pmatrix} w_{11} & w_{12} & \cdots & w_{1d} \\ w_{21} & w_{22} & \cdots & w_{2d} \\ \vdots & \vdots & \ddots & \vdots \\ w_{n_H 1} & w_{n_H 2} & \cdots & w_{n_H d} \end{pmatrix} = \begin{pmatrix} w_1^T \\ w_2^T \\ \vdots \\ w_{n_H}^T \end{pmatrix} \quad (24)$$

Given that  $w_j = \begin{pmatrix} w_{j1} \\ w_{j2} \\ \vdots \\ w_{jd} \end{pmatrix}$

Thus

We can create the  $net_j$  for all the inputs by simply

$$net_j = W_{IH} X = \begin{pmatrix} w_1^T x_1 & w_1^T x_2 & \cdots & w_1^T x_N \\ w_2^T x_1 & w_2^T x_2 & \cdots & w_2^T x_N \\ \vdots & \vdots & \ddots & \vdots \\ w_{n_H}^T x_1 & w_{n_H}^T x_2 & \cdots & w_{n_H}^T x_N \end{pmatrix} \quad (25)$$

Now, we need to generate the  $\mathbf{y}_k$

We apply the activation function element by element in  $\mathbf{net}_j$

$$\mathbf{y}_1 = \begin{pmatrix} f(\mathbf{w}_1^T \mathbf{x}_1) & f(\mathbf{w}_1^T \mathbf{x}_2) & \cdots & f(\mathbf{w}_1^T \mathbf{x}_N) \\ f(\mathbf{w}_2^T \mathbf{x}_1) & f(\mathbf{w}_2^T \mathbf{x}_2) & \cdots & f(\mathbf{w}_2^T \mathbf{x}_N) \\ \vdots & \vdots & \ddots & \vdots \\ f(\mathbf{w}_{n_H}^T \mathbf{x}_1) & f(\mathbf{w}_{n_H}^T \mathbf{x}_2) & \cdots & f(\mathbf{w}_{n_H}^T \mathbf{x}_N) \end{pmatrix} \quad (26)$$

### IMPORTANT about overflows!!!

- Be careful about the numeric stability of the activation function.
- In the case of python, we can use the ones provided by `scipy.special`

## However, We can create a Sigmoid function

It is possible to use the following pseudo-code

Sigmoid( $x$ )

- ① if try  $\left\{ \frac{1.0}{1.0 + \exp\{-\alpha x\}} \right\}$  catch  $\{OVERFLOW\}$   $\triangleleft$  I will use a  
 $\triangleleft$  try-and-catch to catch  
 $\triangleleft$  the overflow
- ② if  $x < 0$
- ③ return 0
- ④ else
- ⑤ return 1
- ⑥ else
- ⑦ return  $\frac{1.0}{1.0 + \exp\{-\alpha x\}}$   $\triangleleft$  1.0 refers to the floating point (Rationals  
 $\triangleleft$  trying to represent Reals)

For this, we get  $net_k$

For this, we obtain the  $W_{HO}$

$$W_{HO} = \begin{pmatrix} w_{11}^o & w_{12}^o & \cdots & w_{1n_H}^o \end{pmatrix} = \left( w_o^T \right) \quad (27)$$

Thus

$$net_k = \begin{pmatrix} w_{11}^o & w_{12}^o & \cdots & w_{1n_H}^o \end{pmatrix} \begin{pmatrix} f\left(w_1^T x_1\right) & f\left(w_1^T x_2\right) & \cdots & f\left(w_1^T x_N\right) \\ f\left(w_2^T x_1\right) & f\left(w_2^T x_2\right) & \cdots & f\left(w_2^T x_N\right) \\ \vdots & \vdots & \ddots & \vdots \\ \underbrace{f\left(w_{n_H}^T x_1\right)}_{y_{k1}} & \underbrace{f\left(w_{n_H}^T x_2\right)}_{y_{k2}} & \cdots & \underbrace{f\left(w_{n_H}^T x_N\right)}_{y_{kN}} \end{pmatrix} \quad (28)$$

In matrix notation

$$net_k = \begin{pmatrix} w_o^T y_{k1} & w_o^T y_{k2} & \cdots & w_o^T y_{kN} \end{pmatrix} \quad (29)$$

Now, we have

Thus, we have  $z_k$  (In our case  $k = 1$ , but it could be a range of values)

$$z_k = \left( f(w_o^T y_{k1}) \quad f(w_o^T y_{k2}) \quad \cdots \quad f(w_o^T y_{kN}) \right) \quad (30)$$

Thus, we generate a vector of differences

$$d = t - z_k = \left( t_1 - f(w_o^T y_{k1}) \quad t_2 - f(w_o^T y_{k2}) \quad \cdots \quad t_N - f(w_o^T y_{kN}) \right) \quad (31)$$

where  $t = \left( t_1 \quad t_2 \quad \cdots \quad t_N \right)$  is a row vector of desired outputs for each sample.

Now, we multiply element wise

We have the following vector of derivatives of *net*

$$\mathbf{D}_f = \left( \eta f'(\mathbf{w}_o^T \mathbf{y}_{k1}) \quad \eta f'(\mathbf{w}_o^T \mathbf{y}_{k2}) \quad \cdots \quad \eta f'(\mathbf{w}_o^T \mathbf{y}_{kN}) \right) \quad (32)$$

where  $\eta$  is the step rate.

Finally, by element wise multiplication (Hadamard Product)

$$\mathbf{d} = \left( \eta [t_1 - f(\mathbf{w}_o^T \mathbf{y}_{k1})] f'(\mathbf{w}_o^T \mathbf{y}_{k1}) \quad \eta [t_2 - f(\mathbf{w}_o^T \mathbf{y}_{k2})] f'(\mathbf{w}_o^T \mathbf{y}_{k2}) \quad \cdots \right. \\ \left. \eta [t_N - f(\mathbf{w}_o^T \mathbf{y}_{kN})] f'(\mathbf{w}_o^T \mathbf{y}_{kN}) \right)$$

## Tile $d$

Tile downward

$$d_{tile} = n_H \text{ rows } \left\{ \begin{pmatrix} d \\ d \\ \vdots \\ d \end{pmatrix} \right. \quad (33)$$

Finally, we multiply element wise against  $y_1$  (Hadamard Product)

$$\Delta w_{1j}^{temp} = y_1 \circ d_{tile} \quad (34)$$

We obtain the total  $\Delta \mathbf{w}_{1j}$

We sum along the rows of  $\Delta \mathbf{w}_{1j}^{temp}$

$$\Delta \mathbf{w}_{1j} = \begin{pmatrix} \eta [t_1 - f(\mathbf{w}_o^T \mathbf{y}_{k1})] f'(\mathbf{w}_o^T \mathbf{y}_{k1}) y_{11} + \eta [t_1 - f(\mathbf{w}_o^T \mathbf{y}_{k1})] f'(\mathbf{w}_o^T \mathbf{y}_{k1}) y_{1N} \\ \vdots \\ \eta [t_1 - f(\mathbf{w}_o^T \mathbf{y}_{k1})] f'(\mathbf{w}_o^T \mathbf{y}_{k1}) y_{n_H 1} + \eta [t_1 - f(\mathbf{w}_o^T \mathbf{y}_{k1})] f'(\mathbf{w}_o^T \mathbf{y}_{k1}) y_{n_H N} \end{pmatrix} \quad (35)$$

where  $y_{hm} = f(\mathbf{w}_h^T \mathbf{x}_m)$  with  $h = 1, 2, \dots, n_H$  and  $m = 1, 2, \dots, N$ .



Finally, we update the first weights

We have then

$$\mathbf{W}_{HO}(t+1) = \mathbf{W}_{HO}(t) + \Delta \mathbf{w}_{1j}^T(t) \quad (36)$$

# First

We multiply element wise the  $\mathbf{W}_{HO}$  and  $\Delta \mathbf{w}_{1j}$

$$\mathbf{T} = \Delta \mathbf{w}_{1j}^T \circ \mathbf{W}_{HO}^T \quad (37)$$

Now, we obtain the element wise derivative of  $net_j$

$$Dnet_j = \begin{pmatrix} f'(\mathbf{w}_1^T \mathbf{x}_1) & f'(\mathbf{w}_1^T \mathbf{x}_2) & \cdots & f'(\mathbf{w}_1^T \mathbf{x}_N) \\ f'(\mathbf{w}_2^T \mathbf{x}_1) & f'(\mathbf{w}_2^T \mathbf{x}_2) & \cdots & f'(\mathbf{w}_2^T \mathbf{x}_N) \\ \vdots & \vdots & \ddots & \vdots \\ f'(\mathbf{w}_{n_H}^T \mathbf{x}_1) & f'(\mathbf{w}_{n_H}^T \mathbf{x}_2) & \cdots & f'(\mathbf{w}_{n_H}^T \mathbf{x}_N) \end{pmatrix} \quad (38)$$

Thus

We tile to the right  $T$

$$\mathbf{T}_{tile} = \underbrace{\begin{pmatrix} \mathbf{T} & \mathbf{T} & \cdots & \mathbf{T} \end{pmatrix}}_{N \text{ Columns}} \quad (39)$$

Now, we multiply element wise together with  $\eta$

$$\mathbf{P}_t = \eta (\mathbf{Dnet}_j \circ \mathbf{T}_{tile}) \quad (40)$$

where  $\eta$  is constant multiplied against the result the Hadamar Product  
(Result a  $n_H \times N$  matrix)

Finally

We get use the transpose of  $\mathbf{X}$  which is a  $N \times d$  matrix

$$\mathbf{X}^T = \begin{pmatrix} \mathbf{x}_1^T \\ \mathbf{x}_2^T \\ \vdots \\ \mathbf{x}_N^T \end{pmatrix} \quad (41)$$

Finally, we get a  $n_H \times d$  matrix

$$\Delta \mathbf{w}_{ij} = \mathbf{P}_t \mathbf{X}^T \quad (42)$$

Thus, given  $\mathbf{W}_{IH}$

$$\mathbf{W}_{IH}(t+1) = \mathbf{W}_{HO}(t) + \Delta \mathbf{w}_{ij}^T(t) \quad (43)$$

# Maximizing information content

## Two ways of achieving this, LeCun 1993

- The use of an example that results in the largest training error.
- The use of an example that is radically different from all those previously used.

## For this

Randomized the samples presented to the multilayer perceptron when not doing batch training.

## Or use an emphasizing scheme

By using the error identify the difficult vs. easy patterns:

- Use them to train the neural network

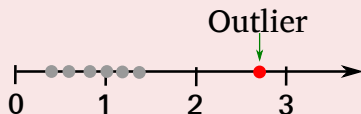
## However

### Be careful about emphasizing scheme

- The distribution of examples within an epoch presented to the network is distorted.
- The presence of an outlier or a mislabeled example can have a catastrophic consequence on the performance of the algorithm.

### Definition of Outlier

An outlier is an observation that lies outside the overall pattern of a distribution (Moore and McCabe 1999).



# Activation Functions [5]

We say that

An activation function  $f(v)$  is antisymmetric if  $f(-v) = -f(v)$

It seems to be

That the multilayer perceptron learns faster using an antisymmetric function.

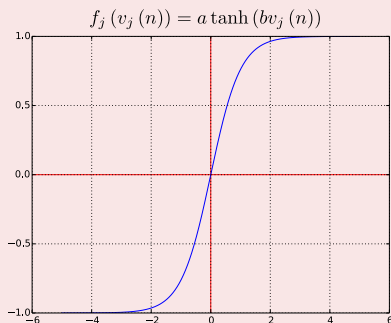
- For example, the hyperbolic tangent function

# Hyperbolic tangent function

Another commonly used form of sigmoid non linearity is the hyperbolic tangent function

$$f_j(v_j(n)) = a \tanh(bv_j(n)) \quad (44)$$

## Example





# The differential of the hyperbolic tangent

We have

$$\begin{aligned}f_j(v_j(n)) &= ab \operatorname{sech}^2(bv_j(n)) \\ &= ab(1 - \tanh^2(bv_j(n)))\end{aligned}$$

BTW

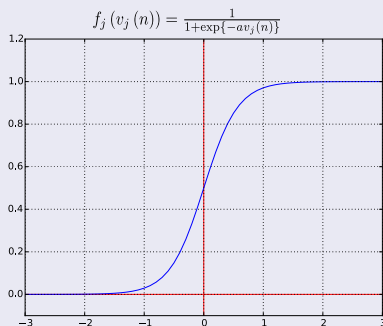
I leave to you to figure out the outputs.

# Logistic Function

This non-linear function has the following definition for a neuron  $j$

$$f_j(v_j(n)) = \frac{1}{1 + \exp\{-av_j(n)\}} \quad a > 0 \text{ and } -\infty < v_j(n) < \infty \quad (45)$$

## Example



# The differential of the sigmoid function

Now, if we differentiate it, we have

$$\begin{aligned} f'_j(v_j(n)) &= \left[ \frac{1}{1 + \exp\{-av_j(n)\}} \right] \left[ 1 - \frac{1}{1 + \exp\{-av_j(n)\}} \right] \\ &= \frac{\exp\{-av_j(n)\}}{(1 + \exp\{-av_j(n)\})^2} \end{aligned}$$

## The outputs finish as

### For the output neurons

$$\begin{aligned}\delta_k &= (t_k - z_k) f'(\text{net}_k) \\ &= (t_k - f_k(v_k(n))) f_k(v_k(n)) (1 - f_k(v_k(n)))\end{aligned}$$

### For the hidden neurons

$$\delta_j = f_j(v_j(n)) (1 - f_j(v_j(n))) \sum_{k=1}^c w_{kj} \delta_k$$

# Problems

## It is clear from our first class

- The problems of the sigmoid activation function, the vanishing gradient

## Not only that, we want activation functions

- That accelerate the learning
- That can control the previous problem
- An any other possible properties

## Something Notable

- This activation function was first introduced to a dynamical network by Hahnloser et al. [6].

## Finally eleven years after

- It was proved that the ReLU function could accelerate the training of deep networks [7].

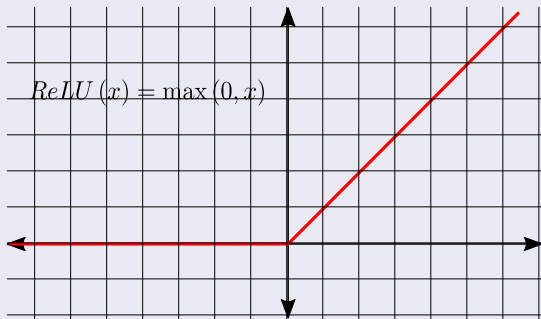
# Definition

## Rectified Linear Unit (ReLU)

$$\text{ReLU}: \mathbb{R} \longrightarrow \mathbb{R}$$

$$\text{ReLU}(x) = \max(0, x), \quad x \in \mathbb{R}$$

## Plotting ReLU



# Noisy ReLUs

## Definition

$$f(x) = \max(0, x + Y) \text{ with } Y \sim N(0, \sigma(x))$$

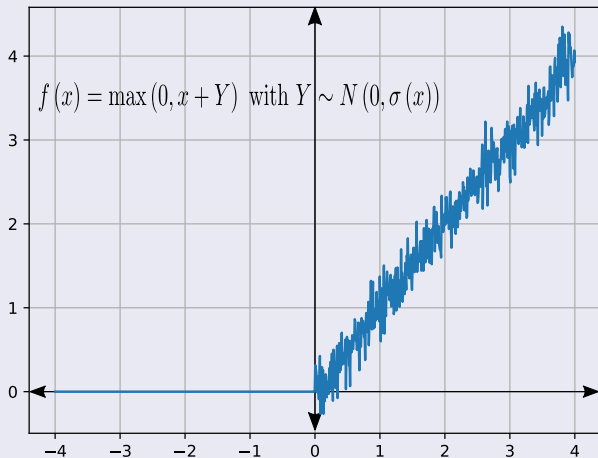
## They have been used

- In restricted Boltzmann machines for computer-vision tasks.



We have

Plot

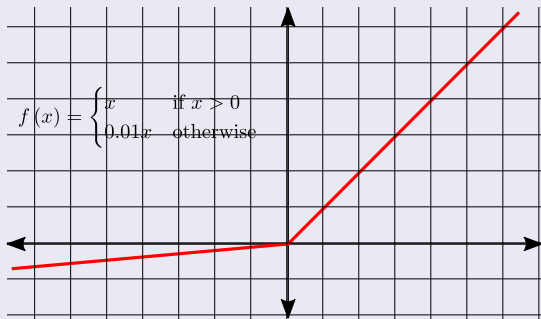


# Leaky ReLU

## Definition

$$f(x) = \begin{cases} x & \text{if } x > 0 \\ 0.01x & \text{otherwise} \end{cases}$$

## Plotting



# Furthermore

## The Parametric ReLU

$$f(x) = \begin{cases} x & \text{if } x > 0 \\ ax & \text{otherwise} \end{cases}$$

- with  $a \leq 1$

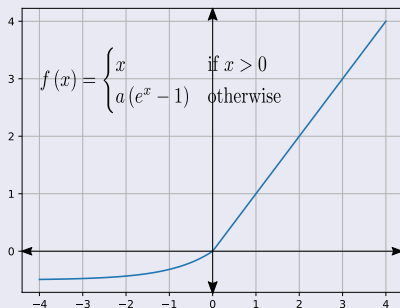
# ELU

## Definition

$$f(x) = \begin{cases} x & \text{if } x > 0 \\ a(e^x - 1) & \text{otherwise} \end{cases}$$

- where  $a$  is a hyper-parameter to be tuned, and  $a \geq 0$ .

## Plotting



# Target Values

## Important

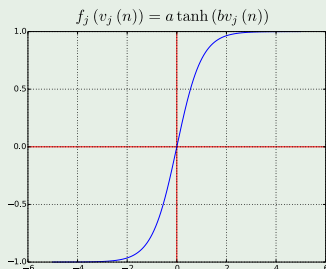
It is important that the target values be chosen within the range of the sigmoid activation function.

## Specifically

The desired response for neuron in the output layer of the multilayer perceptron should be offset by some amount  $\epsilon$

For example

Given the  $a$  limiting value



We have then

- If we have a limiting value  $+a$ , we set  $t = a - \epsilon$ .
- If we have a limiting value  $-a$ , we set  $t = -a + \epsilon$ .

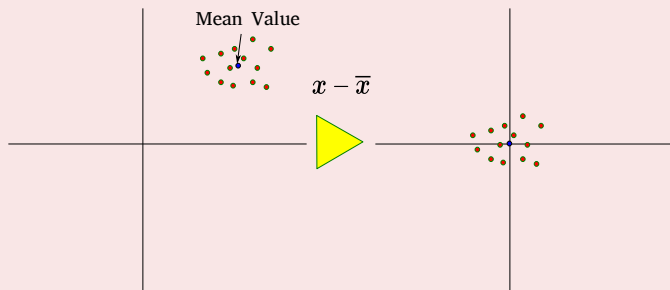
# Normalizing the inputs

## Something Important (LeCun, 1993)

Each input variable should be preprocessed so that:

- The mean value, averaged over the entire training set, is close to zero.
- Or it is small compared to its standard deviation.

## Example

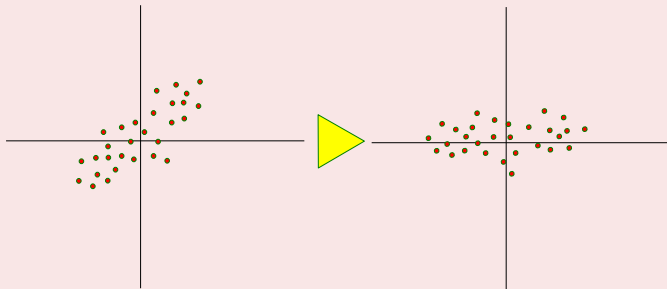


# The normalization must include two other measures

## Uncorrelated

We can use the principal component analysis

### Example





## In addition

### Quite interesting

- The decorrelated input variables should be scaled so that their covariances are approximately equal.

### Why

- Ensuring that the different synaptic weights learn at approximately the same speed.

# There are other heuristics

As

- Initialization
- Learning form hints
- Learning rates
- etc

## In addition

### In section 4.15, Simon Haykin

We have the following techniques:

- Network growing
  - ▶ You start with a small network and add neurons and layers to accomplish the learning task.
- Network pruning
  - ▶ Start with a large network, then prune weights that are not necessary in an orderly fashion.

# Virtues and limitations of Back-Propagation Algorithm

## Something Notable

The back-propagation algorithm has emerged as the most popular algorithm for the training of multilayer perceptrons.

## It has two distinct properties

- It is simple to compute locally.
- It performs stochastic gradient descent in weight space when doing pattern-by-pattern training

# Connectionism

## Back-propagation

It is an example of a connectionist paradigm that relies on local computations to discover the processing capabilities of neural networks.

## This form of restriction

It is known as the locality constraint

# Why this is advocated in Artificial Neural Networks

## First

Artificial neural networks that perform local computations are often held up as metaphors for biological neural networks.

## Second

The use of local computations permits a graceful degradation in performance due to hardware errors, and therefore provides the basis for a fault-tolerant network design.

## Third

Local computations favor the use of parallel architectures as an efficient method for the implementation of artificial neural networks.

However, all this has been seriously questioned [8, 9, 10]

## First

- The reciprocal synaptic connections between the neurons of a multilayer perceptron may assume weights that are excitatory or inhibitory.
- In the real nervous system, neurons usually appear to be the one or the other.

## Second

In a multilayer perceptron, hormonal and other types of global communications are ignored.

However, all this has been seriously questioned [8, 9, 10]

### Third

- In back-propagation learning, a synaptic weight is modified by a presynaptic activity and an error (learning) signal independent of postsynaptic activity.
- There is evidence from neurobiology to suggest otherwise.

### Fourth

- In a neurobiological sense, the implementation of back-propagation learning requires the rapid transmission of information backward along an axon.
- It appears highly unlikely that such an operation actually takes place in the brain.



However, all this has been seriously questioned [8, 9, 10]

## Fifth

- Back-propagation learning implies the existence of a "teacher," which in the context of the brain would presumably be another set of neurons with novel properties.
- The existence of such neurons is biologically implausible.

# Computational Efficiency

## Something Notable

The computational complexity of an algorithm is usually measured in terms of the number of multiplications, additions, and storage involved in its implementation.

- This is the electrical engineering approach!!!

Taking in account the total number of synapses,  $W$  including biases

We have  $\Delta w_{kj} = \eta \delta_k y_j = \eta (t_k - z_k) f'(net_k) y_j$  (Backward Pass)

We have that for this step

- 1 We need to calculate  $net_k$  linear in the number of weights.
- 2 We need to calculate  $y_j = f(net_j)$  which is linear in the number of weights.

# Computational Efficiency

## Now the Forward Pass

$$\Delta w_{ji} = \eta x_i \delta_j = \eta f'(net_j) \left[ \sum_{k=1}^c w_{kj} \delta_k \right] x_i$$

## We have that for this step

$\left[ \sum_{k=1}^c w_{kj} \delta_k \right]$  takes, because of the previous calculations of  $\delta_k$ 's, linear on the number of weights

## Clearly all this takes to have memory

In addition the calculation of the derivatives of the activation functions, but assuming a constant time.

We have that

The Complexity of the multi-layer perceptron is

$O(W)$  Complexity

# Exercises

We have from NN by Haykin

4.2, 4.3, 4.6, 4.8, 4.16, 4.17, 3.7

# Introduction

## Representation of functions

The main result in multi-layer perceptron is its power of representation.

## Furthermore

After all, it is quite striking if we can represent continuous functions of the form  $f : \mathbb{R}^n \mapsto \mathbb{R}$  as a finite sum of simple functions.

# Therefore

## Our main goal

We want to know under which conditions the sum of the form:

$$G(\mathbf{x}) = \sum_{j=1}^N \alpha_j f(\mathbf{w}^T \mathbf{x} + \theta_j) \quad (46)$$

can represent continuous functions in a specific domain.

# Setup of the problem

## Definition of $I_n$

It is an  $n$ -dimensional unit cube  $[0, 1]^n$

In addition, we have the following set of functions

$$C(I_n) = \{f : I_n \rightarrow \mathbb{R} \mid f \text{ is a continuous function}\} \quad (47)$$



## Now, some important definitions

### Definition (Topological Space)

A topological space is then a set  $X$  together with a collection of subsets of  $X$ , called **open sets** and satisfying the following axioms:

- 1 The empty set and  $X$  itself are open.
- 2 Any union of open sets is open.
- 3 The intersection of any finite number of open sets is open.

### Examples

Given the set  $\{1, 2, 3, 4\}$  we have that the following set is a topology  $\{\emptyset, \{1\}, \{1, 2\}, \{1, 2, 3, 4\}\}$ .

### Remark

- This is quite axiomatic... because any set in the collection of  $X$  is open...

Then

We are interested in defining open and close sets in metric spaces

- After all this will allow to define the concept of closed set

### Definition

- A subset  $U$  of a metric space  $(M, d)$  is called open if, given any point  $x \in U$ , there exists a real number  $\epsilon > 0$  such that, given any point  $y \in M$  with  $d(x, y) < \epsilon$ ,  $y$  also belongs to  $U$ .

### Therefore

- A set  $V \subset M$  is closed if  $M - V$  is open.

## Theorem

A compact set is closed and bounded.

## Thus

$I_n$  is a compact set in  $\mathbb{R}^n$ .

## Definition (Continuous functions)

A function  $f : X \rightarrow Y$  where the pre-image of every open set in  $Y$  is open in  $X$ .

Thus

We have the following statement

Let  $K$  be a nonempty subset of  $\mathbb{R}^n$ , where  $n > 1$ . Then If  $K$  is compact, then every continuous real-valued function defined on  $K$  is **bounded**.

### Definition (Supremum Norm)

Let  $X$  be a topological space and let  $F$  be the space of all bounded complex-valued continuous functions defined on  $X$ . The supremum norm is the norm defined on  $F$  by

$$\|f\| = \sup_{x \in X} |f(x)| \quad (48)$$

# Limit Points

## Definition

If  $X$  is a topological space and  $p$  is a point in  $X$ , a neighborhood of  $p$  is a subset  $V$  of  $X$  that includes an open set  $U$  containing  $p$ ,  $p \in U \subseteq V$ .

- This is also equivalent to  $p \in X$  being in the interior of  $V$ .

## Example in a metric space

In a metric space  $(X, d)$ , a set  $V$  is a neighborhood of a point  $p$  if there exists an open ball with center at  $p$  and radius  $r > 0$ , such that

$$B_r(p) = B(p; r) = \{x \in X \mid d(x, p) < r\} \quad (49)$$

is contained in  $V$ .

# Limit Points

## Definition of a Limit Point

Let  $S$  be a subset of a topological space  $X$ . A point  $x \in X$  is a limit point of  $S$  if every neighborhood of  $x$  contains at least one point of  $S$  different from  $x$  itself.

## Example in $\mathbb{R}$

Which are the limit points of the set  $\left\{\frac{1}{n}\right\}_{n=1}^{\infty}$ ?

# This allows to define the idea of density

## Something Notable

A subset  $A$  of a topological space  $X$  is dense in  $X$ , if for any point  $x \in X$ , any neighborhood of  $x$  contains at least one point from  $A$ .

## Classic Example

The real numbers with the usual topology have the rational numbers as a countable dense subset.

- Why do you believe the floating-point numbers are rational?

## In addition

Also the irrational numbers.

# From this, you have the idea of closure

## Definition

The closure of a set  $S$  is the set of all points of closure of  $S$ , that is, the set  $S$  together with all of its limit points.

## Example

The closure of the following set  $(0, 1) \cup \{2\}$

## Meaning

Not all points in the closure are limit points.



# Hausdorff Space

## Definition of Separation

Points  $x$  and  $y$  in a topological space  $X$  can be separated by neighborhoods if there exists a neighborhood  $U$  of  $x$  and a neighborhood  $V$  of  $y$  such that  $U$  and  $V$  are disjoint.

## Definition

$X$  is a Hausdorff space if any two distinct points of  $X$  can be separated by neighborhoods.

We have then

Look at what we have

- 1  $C(I_n)$  is compact
- 2 The continuous functions there are bounded!!!

# Now, the Measure Concept

## Definition of $\sigma$ -algebra

Let  $\mathcal{A} \subset \mathcal{P}(X)$ , we say that  $\mathcal{A}$  to be an algebra if

- 1  $\emptyset, X \in \mathcal{A}$ .
- 2  $A, B \in \mathcal{A}$  then  $A \cup B \in \mathcal{A}$ .
- 3  $A \in \mathcal{A}$  then  $A^c \in \mathcal{A}$ .

## Definition

An algebra  $\mathcal{A}$  in  $\mathcal{P}(X)$  is said to be a  $\sigma$ -algebra, if for any sequence  $\{A_n\}$  of elements in  $\mathcal{A}$ , we have  $\bigcup_{n=1}^{\infty} A_n \in \mathcal{A}$

## Example

In  $X = [0, 1)$ , the class  $\mathcal{A}_0$  consisting of  $\emptyset$ , and all finite unions  $A = \bigcup_{i=1}^n [a_i, b_i)$  with  $0 \leq a_i < b_i \leq a_{i+1} \leq 1$  is an algebra.

## Now, the Measure Concept

### Definition of additivity

Let  $\mu : \mathcal{A} \rightarrow [0, +\infty]$  be such that  $\mu(\emptyset) = 0$ , we say that  $\mu$  is  $\sigma$ -additive if for any  $\{A_i\}_{i \in I} \subset \mathcal{A}$  (Where  $I$  can be finite or infinite countable) of mutually disjoint sets such that  $\cup_{i \in I} A_i \in \mathcal{A}$ , we have that

$$\mu\left(\cup_{i \in I} A_i\right) = \sum_{i \in I} \mu(A_i) \quad (50)$$

### Definition of Measurability

Let  $\mathcal{A}$  be a  $\sigma$ -algebra of subsets of  $X$ , we say that the pair  $(X, \mathcal{A})$  is a measurable space where a  $\sigma$ -additive function  $\mu : \mathcal{A} \rightarrow [0, +\infty]$  is called a measure on  $(X, \mathcal{A})$ .

# A Borel Measure

## Definition

The Borel  $\sigma$ -algebra is defined to be the  $\sigma$ -algebra generated by the open sets (or equivalently, by the closed sets).

## Definition of a Borel Measure

If  $\mathcal{F}$  is the Borel  $\sigma$ -algebra on some topological space, then a measure  $\mu : \mathcal{F} \rightarrow \mathbb{R}$  is said to be a Borel measure (or Borel probability measure). For a Borel measure, all continuous functions are measurable.

## Definition of a signed Borel Measure

A signed Borel measure  $\mu : \mathcal{B}(X) \rightarrow \mathbb{R}$  is a measure such that

- 1  $\mu(\emptyset) = 0$ .
- 2  $\mu$  is  $\sigma$ -additive.
- 3  $\sup_{A \in \mathcal{B}(X)} |\mu(A)| < \infty$

# A Borel Measure

## Regularity

A measure  $\mu$  is Borel regular measure:

- 1 For every Borel set  $B \subseteq \mathbb{R}^n$  and  $A \subseteq \mathbb{R}^n$ ,  
 $\mu(A) = \mu(A \cap B) + \mu(A - B)$ .
- 2 For every  $A \subseteq \mathbb{R}^n$ , there exists a Borel set  $B \subseteq \mathbb{R}^n$  such that  $A \subseteq B$  and  $\mu(A) = \mu(B)$ .

# Discriminatory Functions

## Definition

Given the set  $M(I_n)$  of signed regular Borel measures, a function  $f$  is discriminatory if for a measure  $\mu \in M(I_n)$

$$\int_{I_n} f(\mathbf{w}^T \mathbf{x} + \theta) d\mu = 0 \quad (51)$$

for all  $\mathbf{w} \in \mathbb{R}^n$  and  $\theta \in \mathbb{R}$  implies that  $\mu = 0$

## Definition

We say that  $f$  is sigmoidal if

$$f(t) \rightarrow \begin{cases} 1 & \text{as } t \rightarrow +\infty \\ 0 & \text{as } t \rightarrow -\infty \end{cases}$$

# The Important Theorem

## Theorem 1

Let  $f$  be any continuous discriminatory function. Then finite sums of the form

$$G(\mathbf{x}) = \sum_{j=1}^N \alpha_j f(\mathbf{w}_j^T \mathbf{x} + \theta_j), \quad (52)$$

where  $\mathbf{w}_j \in \mathbb{R}^n$  and  $\alpha_j, \theta_j \in \mathbb{R}$  are fixed, are dense in  $C(I_n)$

## Meaning

Basically given any function  $g \in C(I_n)$  and any neighborhood  $V$  of  $g$ , you have a  $G \in V$ .



## Furthermore

### In other words

Given any  $g \in C(I_n)$  and  $\epsilon > 0$ , there is a sum,  $G(x)$ , of the above form, for which

$$|G(x) - g(x)| < \epsilon \quad \forall x \in I_n \quad (53)$$

# Proof

Let  $S \subset C(I_n)$  be the set of functions of the form  $G(x)$

First,  $S$  is a linear subspace of  $C(I_n)$

## Definition

A subset  $V$  of  $\mathbb{R}^n$  is called a linear subspace of  $\mathbb{R}^n$  if  $V$  contains the zero vector, and is closed under vector addition and scaling. That is, for  $X, Y \in V$  and  $c \in \mathbb{R}$ , we have  $X + Y \in V$  and  $cX \in V$ .

We claim that the closure of  $S$  is all of  $C(I_n)$

Assume that the closure of  $S$  is not all of  $C(I_n)$

# Proof

Then

- The closure of  $S$ , say  $R$ , is a closed proper subspace of  $C(I)$

We use the Hahn-Banach Theorem

- If  $p : V \rightarrow \mathbb{R}$  is a sublinear function
  - ▶  $p(x + y) \leq p(x) + p(y)$
  - ▶  $p(\alpha x) = \alpha p(x)$

And  $\varphi : U \rightarrow \mathbb{R}$  is a linear functional on a linear subspace  $U \subseteq V$

- It is dominated by  $p$  on  $U$ , i.e.  $\varphi(x) \leq p(x) \forall x \in U$ .

# Hahn-Banach Theorem

Then

There exists a **linear extension**  $\psi : V \rightarrow \mathbb{R}$  of  $\varphi$  to the whole space  $V$ , i.e., there exists a linear functional  $\psi$  such that

- ①  $\psi(x) = \varphi(x) \quad \forall x \in U.$
- ②  $\psi(x) \leq p(x) \quad \forall x \in V.$

## Proof

It is possible to construct sublinear function defined as follow

We define the following linear functional

$$T(f) = \begin{cases} f & \text{if } f \in C(I_n) - R \\ 0 & \text{if } f \in R \end{cases} \quad (54)$$

Then

We have there is a bounded linear functional (Using  $T$  as  $p$  and  $\varphi$ , and  $V \in C(I_n)$  and  $U = R$ ) called  $L \neq 0$  with  $L(R) = L(S) = 0$ .

## Proof

Now, we use the Riesz Representation Theorem

Let  $X$  be a locally compact Hausdorff space. For any positive linear functional  $\psi$  on  $C(X)$ , there is a unique regular Borel measure  $\mu$  on  $X$  such that

$$\psi = \int_X f(x) d\mu(x) \quad (55)$$

for all  $f$  in  $C(X)$

We can then do the following

$$L(h) = \int_{I_n} h(\mathbf{x}) d\mu(\mathbf{x}) \quad (56)$$

Where?

For some  $\mu \in M(I_n)$ , for all  $h \in C(I_n)$

# Proof

In particular

Given that  $f(\mathbf{w}^T \mathbf{x} + \theta)$  is in  $R$  for all  $\mathbf{w}$  and  $\theta$

We must have that

$$\int_{I_n} f(\mathbf{w}^T \mathbf{x} + \theta) d\mu(\mathbf{x}) = 0 \quad (57)$$

for all  $\mathbf{w}$  and  $\theta$

But we assumed that  $f$  is discriminatory!!!

Then...  $\mu = 0$  contradicting the fact that  $L \neq 0$ !!! We have a contradiction!!!

# Proof

Finally

The subspace  $S$  of sums of the form  $G$  is dense!!!



## Now, we deal with the sigmoidal function

### Lemma 1

Any bounded, measurable sigmoidal function,  $f$ , is discriminatory. In particular, any continuous sigmoidal function is discriminatory.

### Proof

I will leave this to you... it is possible I will get a question from this proof for the first midterm.

# We have the theorem finally!!!

## Universal Representation Theorem for the multi-layer perceptron

Let  $f$  be any continuous sigmoid function. Then finite sums of the form

$$G(\mathbf{x}) = \sum_{j=1}^N \alpha_j f(\mathbf{w}^T \mathbf{x} + \theta_j) \quad (58)$$

are dense in  $C(I_n)$ .

## In other words

Given any  $g \in C(I_n)$  and  $\epsilon > 0$ , there is a sum  $G(\mathbf{x})$  of the above form, for which

$$|G(\mathbf{x}) - g(\mathbf{x})| < \epsilon \quad \forall \mathbf{x} \in I_n \quad (59)$$

## Simple

Combine the theorem and lemma 1... and because the continuous sigmoid satisfy the conditions of the lemma... we have our representation!!!



A. E. Bryson, W. F. Denham, and S. E. Dreyfus, “Optimal programming problems with inequality constraints,” *AIAA journal*, vol. 1, no. 11, pp. 2544–2550, 1963.



P. J. Werbos *et al.*, “Backpropagation through time: what it does and how to do it,” *Proceedings of the IEEE*, vol. 78, no. 10, pp. 1550–1560, 1990.



D. E. Rumelhart and D. Zipser, “Feature discovery by competitive learning,” *Cognitive science*, vol. 9, no. 1, pp. 75–112, 1985.



D. E. Rumelhart, G. E. Hinton, and R. J. Williams, “Neurocomputing: Foundations of research,” ch. Learning Representations by Back-propagating Errors, pp. 696–699, Cambridge, MA, USA: MIT Press, 1988.



S. Haykin, *Neural Networks and Learning Machines*.

No. v. 10 in Neural networks and learning machines, Prentice Hall, 2009.



R. H. Hahnloser, R. Sarpeshkar, M. A. Mahowald, R. J. Douglas, and H. S. Seung, "Digital selection and analogue amplification coexist in a cortex-inspired silicon circuit," *Nature*, vol. 405, no. 6789, p. 947, 2000.



X. Glorot, A. Bordes, and Y. Bengio, "Deep sparse rectifier neural networks," in *Proceedings of the fourteenth international conference on artificial intelligence and statistics*, pp. 315–323, 2011.



G. M. Shepherd, "The dendritic spine: a multifunctional integrative unit," *Journal of neurophysiology*, vol. 75, no. 6, pp. 2197–2210, 1996.



F. Crick, "The recent excitement about neural networks.," *Nature*, vol. 337, no. 6203, pp. 129–132, 1989.



D. G. Stork, "Is backpropagation biologically plausible," in *International Joint Conference on Neural Networks*, vol. 2, pp. 241–246, IEEE Washington, DC, 1989.