

Introducción a Redes Neuronales

Dr. Mauricio Toledo-Acosta

Diplomado Ciencia de Datos con Python

Table of Contents

Introducción a
Redes
Neuronales

Introduction

The
Rosenblatt
Perceptron

The Perceptron
Learning Algorithm
Limitations
Combining
Perceptrons

Gradient
Based
Learning

Loss Functions
Optimizers

Fully
Connected
Networks

1 Introduction

2 The Rosenblatt Perceptron

- The Perceptron Learning Algorithm
- Limitations
- Combining Perceptrons

3 Gradient Based Learning

- Loss Functions
- Optimizers

4 Fully Connected Networks

Introduction

Introducción a
Redes
Neuronales

Introduction

The
Rosenblatt
Perceptron

The Perceptron
Learning Algorithm
Limitations
Combining
Perceptrons

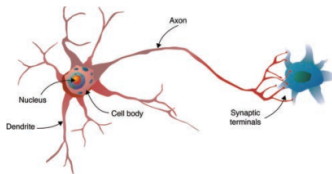
Gradient
Based
Learning

Loss Functions
Optimizers

Fully
Connected
Networks

Neural Networks

Artificial neural networks are machine learning techniques that simulate the mechanism of learning in biological organisms. The human nervous system contains cells, referred to as neurons. The neurons are connected to one another by axons and dendrites, and the connecting regions between axons and dendrites are called synapses. The strengths of synaptic connections often change in response to external stimuli. This change is how learning takes place in living organisms.



El modelo lineal de clasificación

Introducción a
Redes
Neuronales

Introduction

The
Rosenblatt
Perceptron

The Perceptron
Learning Algorithm
Limitations
Combining
Perceptrons

Gradient
Based
Learning

Loss Functions
Optimizers

Fully
Connected
Networks

- The Perceptrón algorithm (Rosenblatt, 1961) played an important role in Machine Learning history. It was first simulated in a computer IBM 704 at Cornell in 1957. By the early 60s, a dedicated hardware was designed to implement the learning algorithm.

El modelo lineal de clasificación

Introducción a
Redes
Neuronales

Introduction

The
Rosenblatt
Perceptron

The Perceptron
Learning Algorithm
Limitations
Combining
Perceptrons

Gradient
Based
Learning

Loss Functions
Optimizers

Fully
Connected
Networks

- The Perceptrón algorithm (Rosenblatt, 1961) played an important role in Machine Learning history. It was first simulated in a computer IBM 704 at Cornell in 1957. By the early 60s, a dedicated hardware was designed to implement the learning algorithm.
- It was criticized by Marvin Minsky, who showed the limitations of the perceptron algorithm when dealing with a non linear separable set of points. This caused a void in the neural computation research lasting until the mid 80s.

Table of Contents

Introducción a
Redes
Neuronales

Introduction

The
Rosenblatt
Perceptron

The Perceptron
Learning Algorithm
Limitations
Combining
Perceptrons

Gradient
Based
Learning

Loss Functions
Optimizers

Fully
Connected
Networks

1 Introduction

2 The Rosenblatt Perceptron

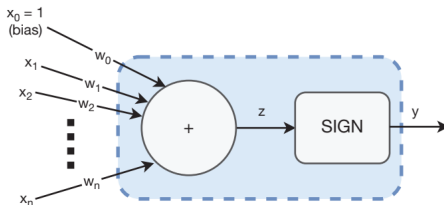
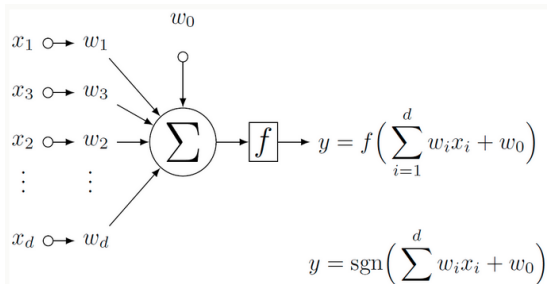
- The Perceptron Learning Algorithm
- Limitations
- Combining Perceptrons

3 Gradient Based Learning

- Loss Functions
- Optimizers

4 Fully Connected Networks

The Rosenblatt Perceptron



The Algorithm

We iterate over the training set

$$(x_1, y_1), (x_2, y_2) \cdots, (x_N, y_N).$$

Introducción a
Redes
Neuronales

Introduction

The
Rosenblatt
Perceptron

**The Perceptron
Learning Algorithm**

Limitations

Combining
Perceptrons

Gradient
Based
Learning

Loss Functions

Optimizers

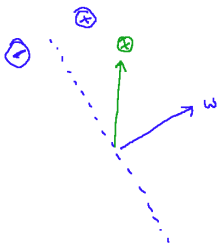
Fully
Connected
Networks

The Algorithm

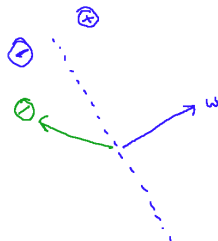
We iterate over the training set

$$(x_1, y_1), (x_2, y_2) \cdots, (x_N, y_N).$$

- If the point is well classified:



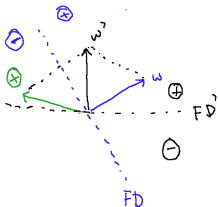
(a) If $y = +1$



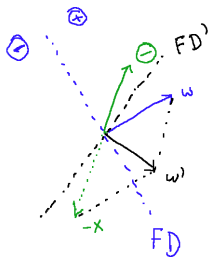
(b) If $y = -1$

The Algorithm

- If the point is not well classified:



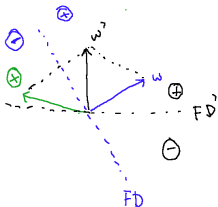
(a) If $y = +1$



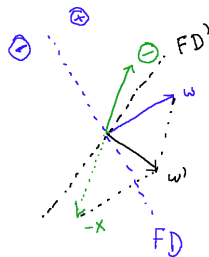
(b) If $y = -1$

The Algorithm

- If the point is not well classified:



(a) If $y = +1$



(b) If $y = -1$

We update the vector w with $w + yx$

The Algorithm

For each epoch, the algorithm is summarized as follows:

Data: A : training set of points, $X = \{x_1, \dots, x_N\} \subset \mathbb{R}^{D+1}$,
 $Y = \{y_1, \dots, y_N\}$: set of labels
 $\eta > 0$: learning rate.

Result: w : weight vector defining the decision frontier.

```
1 Function Perceptron( $X, y, \eta$ ):  
2    $w = 0$ ;  
3   converged = False;  
4   while converged == False do  
5     for  $i \in \{1, \dots, N\}$  do  
6       if  $y_i \langle w, x \rangle \leq 0$  then  
7          $w = w + y_i \eta x_i$ ;  
8         converged = False  
9   return  $w$ ;
```

The Algorithm

For each epoch, the algorithm is summarized as follows:

Data: A : training set of points, $X = \{x_1, \dots, x_N\} \subset \mathbb{R}^{D+1}$,
 $Y = \{y_1, \dots, y_N\}$: set of labels
 $\eta > 0$: learning rate.

Result: w : weight vector defining the decision frontier.

```
1 Function Perceptron( $X, y, \eta$ ):  
2    $w = 0$ ;  
3   converged = False;  
4   while converged == False do  
5     for  $i \in \{1, \dots, N\}$  do  
6       if  $y_i \langle w, x \rangle \leq 0$  then  
7          $w = w + y_i \eta x_i$ ;  
8         converged = False  
9   return  $w$ ;
```

w is a linear combination of x_1, \dots, x_N .

Example

Introducción a
Redes
Neuronales

Introduction

The
Rosenblatt
Perceptron

**The Perceptron
Learning Algorithm**

Limitations

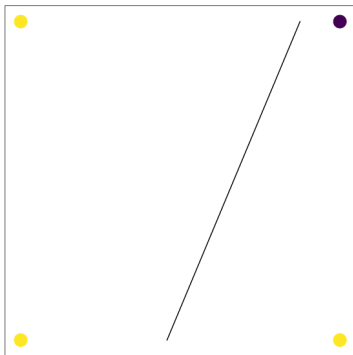
Combining
Perceptrons

Gradient
Based
Learning

Loss Functions

Optimizers

Fully
Connected
Networks



Example

Introducción a
Redes
Neuronales

Introduction

The
Rosenblatt
Perceptron

**The Perceptron
Learning Algorithm**

Limitations

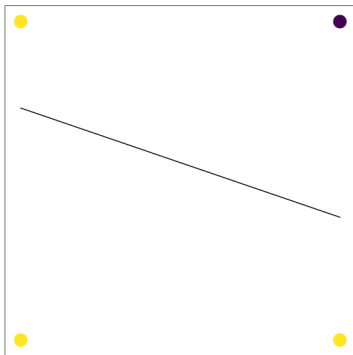
Combining
Perceptrons

Gradient
Based
Learning

Loss Functions

Optimizers

Fully
Connected
Networks



Example

Introducción a
Redes
Neuronales

Introduction

The
Rosenblatt
Perceptron

**The Perceptron
Learning Algorithm**

Limitations

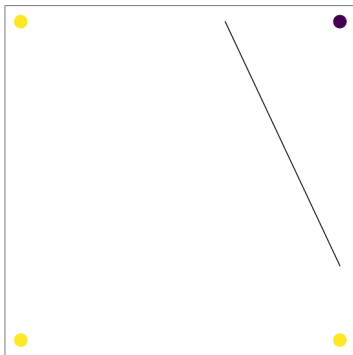
Combining
Perceptrons

Gradient
Based
Learning

Loss Functions

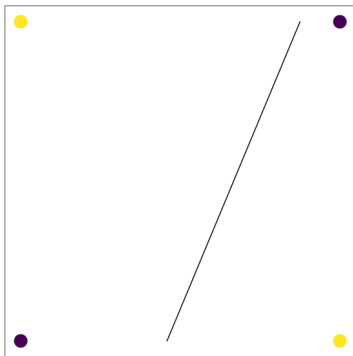
Optimizers

Fully
Connected
Networks



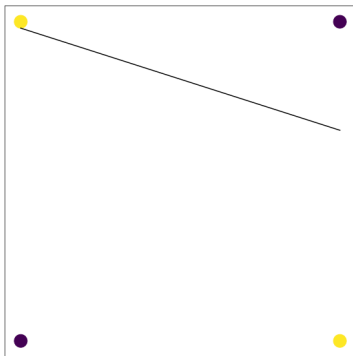
Limitations of the Perceptron

What happens if a straight line cannot separate the data points?



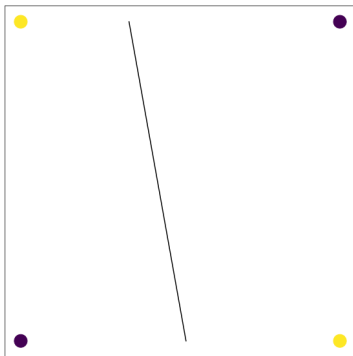
Limitations of the Perceptron

What happens if a straight line cannot separate the data points?



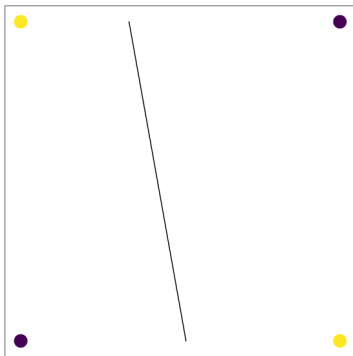
Limitations of the Perceptron

What happens if a straight line cannot separate the data points?



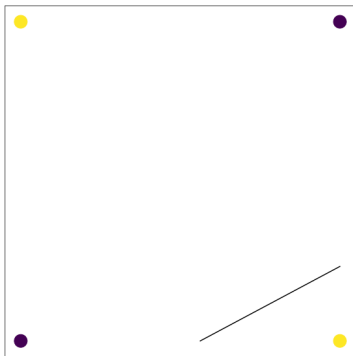
Limitations of the Perceptron

What happens if a straight line cannot separate the data points?



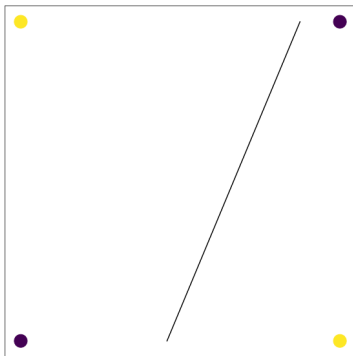
Limitations of the Perceptron

What happens if a straight line cannot separate the data points?



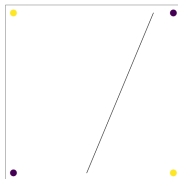
Limitations of the Perceptron

What happens if a straight line cannot separate the data points?



Limitations of the Perceptron

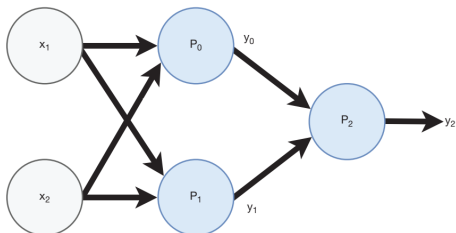
What happens if a straight line cannot separate the data points?



Alternatives:

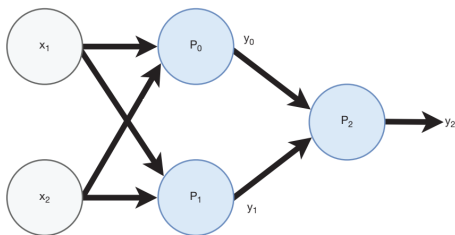
- We change the model of a neuron or
- We combine multiple of them to solve this limitation.

Combining Perceptrons



$$\begin{array}{c|c} -1 & (-1, -1) \\ +1 & (1, -1) \\ +1 & (-1, 1) \\ -1 & (1, 1) \end{array}$$

Combining Perceptrons



This neural network is one of the simplest examples of a fully connected feedforward network. **Fully connected** means that the output of each neuron in one layer is connected to all neurons in the next layer. **Feedforward** means that there are no backward connections. A **multilevel neural network** has an input layer, one or more hidden layers, and an output layer. The input layer does not contain neurons but contains only the inputs themselves.

Table of Contents

Introducción a
Redes
Neuronales

Introduction

The
Rosenblatt
Perceptron

The Perceptron
Learning Algorithm
Limitations
Combining
Perceptrons

Gradient
Based
Learning

Loss Functions
Optimizers

Fully
Connected
Networks

1 Introduction

2 The Rosenblatt Perceptron

- The Perceptron Learning Algorithm
- Limitations
- Combining Perceptrons

3 Gradient Based Learning

- Loss Functions
- Optimizers

4 Fully Connected Networks

How to find the minimum of a function

- Let $f(x)$ be a function, we want to find the minimum of f .

Introducción a
Redes
Neuronales

Introduction

The
Rosenblatt
Perceptron

The Perceptron
Learning Algorithm

Limitations

Combining
Perceptrons

Gradient
Based
Learning

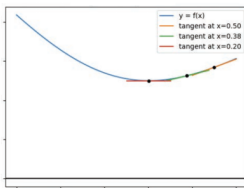
Loss Functions

Optimizers

Fully
Connected
Networks

How to find the minimum of a function

- Let $f(x)$ be a function, we want to find the minimum of f .
- The derivative at the point x that minimizes the value of f is 0.



How to find the minimum of a function

Introducción a
Redes
Neuronales

Introduction

The
Rosenblatt
Perceptron

The Perceptron
Learning Algorithm

Limitations

Combining
Perceptrons

Gradient
Based
Learning

Loss Functions

Optimizers

Fully
Connected
Networks

- Let $f(x)$ be a function, we want to find the minimum of f .
- The derivative at the point x that minimizes the value of f is 0.
- Given an initial value x , the sign of the derivative $f'(x)$ indicates in what direction to adjust x to reduce the value of $f(x)$.

How to find the minimum of a function

- Let $f(x)$ be a function, we want to find the minimum of f .
- The derivative at the point x that minimizes the value of f is 0.
- Given an initial value x , the sign of the derivative $f'(x)$ indicates in what direction to adjust x to reduce the value of $f(x)$.
- Gradient descent uses the value of the derivative to decide how much to adjust x .

$$x_{n+1} = x_n - \eta f'(x_n).$$

How to find the minimum of a function

- Let $f(x)$ be a function, we want to find the minimum of f .
- The derivative at the point x that minimizes the value of f is 0.
- Given an initial value x , the sign of the derivative $f'(x)$ indicates in what direction to adjust x to reduce the value of $f(x)$.
- Gradient descent uses the value of the derivative to decide how much to adjust x .

$$x_{n+1} = x_n - \eta f'(x_n).$$

- If the learning rate η is too large, gradient descent can overshoot the solution and fail to converge.

How to find the minimum of a function

- Let $f(x)$ be a function, we want to find the minimum of f .
- The derivative at the point x that minimizes the value of f is 0.
- Given an initial value x , the sign of the derivative $f'(x)$ indicates in what direction to adjust x to reduce the value of $f(x)$.
- Gradient descent uses the value of the derivative to decide how much to adjust x .

$$x_{n+1} = x_n - \eta f'(x_n).$$

- If the learning rate η is too large, gradient descent can overshoot the solution and fail to converge.
- The algorithm is not guaranteed to find the global minimum because it can get stuck in a local minimum.

Learning Rate

Introducción a
Redes
Neuronales

Introduction

The
Rosenblatt
Perceptron

The Perceptron
Learning Algorithm

Limitations

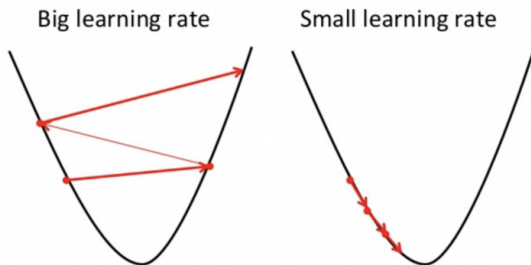
Combining
Perceptrons

Gradient
Based
Learning

Loss Functions

Optimizers

Fully
Connected
Networks



The high-dimensional case

- Now, consider the function $f(x_0, x_1)$.

Introducción a
Redes
Neuronales

Introduction

The
Rosenblatt
Perceptron

The Perceptron
Learning Algorithm

Limitations

Combining
Perceptrons

Gradient
Based
Learning

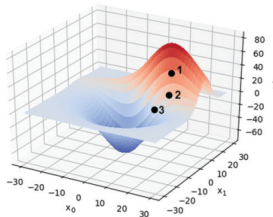
Loss Functions

Optimizers

Fully
Connected
Networks

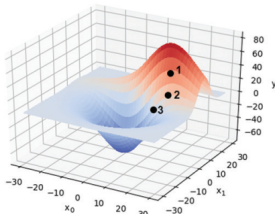
The high-dimensional case

- Now, consider the function $f(x_0, x_1)$.
- The gradient is a vector consisting of partial derivatives and indicates the direction in the input space that results in the steepest ascent for the value of f .



The high-dimensional case

- Now, consider the function $f(x_0, x_1)$.
- The gradient is a vector consisting of partial derivatives and indicates the direction in the input space that results in the steepest ascent for the value of f .



- If we are at the point $\mathbf{x}_n = (x_0^{(n)}, x_1^{(n)})$ and want to minimize y , then we choose our next point as

$$\mathbf{x}_{n+1} = \mathbf{x}_n - \eta \nabla f(\mathbf{x}_n)$$

Gradient-based Learning

Introducción a
Redes
Neuronales

Introduction

The
Rosenblatt
Perceptron

The Perceptron
Learning Algorithm
Limitations
Combining
Perceptrons

Gradient
Based
Learning

Loss Functions
Optimizers

Fully
Connected
Networks

When applying gradient descent to our neural network, we consider input values \mathbf{x} to be constants, with our goal being to adjust the weights \mathbf{w} .

- Which function we want to minimize? A loss function.

Gradient-based Learning

Introducción a
Redes
Neuronales

Introduction

The
Rosenblatt
Perceptron

The Perceptron
Learning Algorithm
Limitations
Combining
Perceptrons

Gradient
Based
Learning

Loss Functions
Optimizers

Fully
Connected
Networks

When applying gradient descent to our neural network, we consider input values \mathbf{x} to be constants, with our goal being to adjust the weights \mathbf{w} .

- Which function we want to minimize? A loss function.
- In the context of an optimization algorithm, the function used to evaluate a candidate solution is referred to as the **objective function**. With neural networks, we seek to minimize the error. As such, the objective function is often referred to as a cost function or a **loss function**.

Gradient-based Learning

- In the case of the Perceptron, the loss function is given by

$$L^{(0/1)}(\mathbf{w}) = (y_i - \text{sign}\langle \mathbf{w}, \mathbf{x}_i \rangle)^2$$

Introducción a
Redes
Neuronales

Introduction

The
Rosenblatt
Perceptron

The Perceptron
Learning Algorithm

Limitations

Combining
Perceptrons

Gradient
Based
Learning

Loss Functions

Optimizers

Fully
Connected
Networks

Gradient-based Learning

Introducción a
Redes
Neuronales

Introduction

The
Rosenblatt
Perceptron

The Perceptron
Learning Algorithm
Limitations
Combining
Perceptrons

Gradient
Based
Learning

Loss Functions
Optimizers

Fully
Connected
Networks

- In the case of the Perceptron, the loss function is given by

$$L^{(0/1)}(\mathbf{w}) = \frac{1}{2}(y_i - \text{sign}\langle \mathbf{w}, \mathbf{x}_i \rangle)^2 = 1 - y_i \text{sign}\langle \mathbf{w}, \mathbf{x}_i \rangle$$

- This function is not smooth, we use the smooth surrogate loss function

$$L(\mathbf{w}) = \max\{-y_i \text{sign}\langle \mathbf{w}, \mathbf{x}_i \rangle, 0\}.$$

Gradient-based Learning

Introducción a
Redes
Neuronales

Introduction

The
Rosenblatt
Perceptron

The Perceptron
Learning Algorithm
Limitations
Combining
Perceptrons

Gradient
Based
Learning

Loss Functions
Optimizers

Fully
Connected
Networks

- In the case of the Perceptron, the loss function is given by
- This function is not smooth, we use the smooth surrogate loss function

$$L(\mathbf{w}) = \max\{-y_i \text{sign}\langle \mathbf{w}, \mathbf{x}_i \rangle, 0\}.$$

- Applying gradient descent

$$\begin{aligned}\mathbf{w}_{n+1} &= \mathbf{w}_n - \eta \nabla L(\mathbf{w}) \\ &= \begin{cases} \mathbf{w}_n + \eta y_i \mathbf{x}_i, & \text{well classified} \\ \mathbf{w}_n, & \text{misclassified} \end{cases}\end{aligned}$$

Loss Functions

Introducción a
Redes
Neuronales

Introduction

The
Rosenblatt
Perceptron

The Perceptron
Learning Algorithm

Limitations

Combining
Perceptrons

Gradient
Based
Learning

Loss Functions

Optimizers

Fully
Connected
Networks

Other examples of loss functions

► Keras Losses

Loss Functions: MSE

Introducción a
Redes
Neuronales

Introduction

The
Rosenblatt
Perceptron

The Perceptron
Learning Algorithm

Limitations

Combining
Perceptrons

Gradient
Based
Learning

Loss Functions

Optimizers

Fully
Connected
Networks

Mean Square Error (MSE)

$$L(y, t) = (t - y)^2.$$

- L2 loss.

Loss Functions: MSE

Introducción a
Redes
Neuronales

Introduction

The
Rosenblatt
Perceptron

The Perceptron
Learning Algorithm

Limitations

Combining
Perceptrons

Gradient
Based
Learning

Loss Functions

Optimizers

Fully
Connected
Networks

Mean Square Error (MSE)

$$L(y, t) = (t - y)^2.$$

- L2 loss.
- Good for regression tasks.

Loss Functions: MSE

Introducción a
Redes
Neuronales

Introduction

The
Rosenblatt
Perceptron

The Perceptron
Learning Algorithm

Limitations

Combining
Perceptrons

Gradient
Based
Learning

Loss Functions

Optimizers

Fully
Connected
Networks

Mean Square Error (MSE)

$$L(y, t) = (t - y)^2.$$

- L2 loss.
- Good for regression tasks.
- Trivial derivative for gradient descent.

Loss Functions: MAE

Introducción a
Redes
Neuronales

Introduction

The
Rosenblatt
Perceptron

The Perceptron
Learning Algorithm

Limitations

Combining
Perceptrons

Gradient
Based
Learning

Loss Functions

Optimizers

Fully
Connected
Networks

Mean Absolute Error (MAE)

$$L = |t - y|.$$

- L1 loss.
- More robust to outliers than mse.
- Good for regression tasks.
- Discontinuity in its derivative.

Loss Functions: MAE

Introducción a
Redes
Neuronales

Introduction

The
Rosenblatt
Perceptron

The Perceptron
Learning Algorithm

Limitations

Combining
Perceptrons

Gradient
Based
Learning

Loss Functions

Optimizers

Fully
Connected
Networks

Mean Absolute Error (MAE)

$$L = |t - y|.$$

- L1 loss.
- More robust to outliers than mse.
- Good for regression tasks.
- Discontinuity in its derivative.

Loss Functions: Hinge

Introducción a
Redes
Neuronales

Introduction

The
Rosenblatt
Perceptron

The Perceptron
Learning Algorithm

Limitations

Combining
Perceptrons

Gradient
Based
Learning

Loss Functions

Optimizers

Fully
Connected
Networks

Hinge Loss

$$L = \max\{-y_i\hat{y}_i, 0\}.$$

- Used in SVMs and Perceptron.
- Penalizes errors, but also correct predictions of low confidence (probabilities).
- Good for binary classification tasks.

Loss Functions: Categorical cross entropy

Introducción a
Redes
Neuronales

Introduction

The
Rosenblatt
Perceptron

The Perceptron
Learning Algorithm
Limitations
Combining
Perceptrons

Gradient
Based
Learning

Loss Functions
Optimizers

Fully
Connected
Networks

Categorical cross entropy

$$L = \sum_i^n y_i \log(\hat{y}_i).$$

- Good for multi-class classification problems.
- Considers y to be a one-hot encoding vector in n classes.

Optimizers: Gradient Descent

Introducción a
Redes
Neuronales

Introduction

The
Rosenblatt
Perceptron

The Perceptron
Learning Algorithm

Limitations

Combining
Perceptrons

Gradient
Based
Learning

Loss Functions

Optimizers

Fully
Connected
Networks

Gradient descent is the **most basic** but most used optimization algorithm.

Optimizers: Gradient Descent

Introducción a
Redes
Neuronales

Introduction

The
Rosenblatt
Perceptron

The Perceptron
Learning Algorithm

Limitations

Combining
Perceptrons

Gradient
Based
Learning

Loss Functions

Optimizers

Fully
Connected
Networks

Gradient descent is the **most basic** but most used optimization algorithm. It's used heavily in linear regression and classification algorithms.

Optimizers: Gradient Descent

Introducción a
Redes
Neuronales

Introduction

The
Rosenblatt
Perceptron

The Perceptron
Learning Algorithm

Limitations
Combining
Perceptrons

Gradient
Based
Learning

Loss Functions

Optimizers

Fully
Connected
Networks

Gradient descent is the **most basic** but most used optimization algorithm. It's used heavily in linear regression and classification algorithms. Gradient descent is a **first-order optimization algorithm** which is dependent on the first order derivative of a loss function.

Optimizers: Gradient Descent

Introducción a
Redes
Neuronales

Introduction

The
Rosenblatt
Perceptron

The Perceptron
Learning Algorithm

Limitations

Combining
Perceptrons

Gradient
Based
Learning

Loss Functions

Optimizers

Fully
Connected
Networks

Advantages:

- Easy computation.
- Easy to implement.
- Easy to understand.

Optimizers: Gradient Descent

Introducción a
Redes
Neuronales

Introduction

The
Rosenblatt
Perceptron

The Perceptron
Learning Algorithm

Limitations

Combining
Perceptrons

Gradient
Based
Learning

Loss Functions

Optimizers

Fully
Connected
Networks

Advantages:

- Easy computation.
- Easy to implement.
- Easy to understand.

Disadvantages:

- May trap at local minima.
- Weights are changed after calculating gradient on the whole dataset. May take a long time to converge.
- Requires large memory to calculate gradient on the whole dataset.

Optimizers

Introducción a
Redes
Neuronales

Introduction

The
Rosenblatt
Perceptron

The Perceptron
Learning Algorithm

Limitations

Combining
Perceptrons

Gradient
Based
Learning

Loss Functions

Optimizers

Fully
Connected
Networks

- Gradient Descent
- Stochastic Gradient Descent
- Stochastic Gradient descent with momentum
- Mini-Batch Gradient Descent
- Adagrad
- RMSProp
- AdaDelta
- Adam

► Keras Optimizers

Stochastic Gradient Descent

Introducción a
Redes
Neuronales

Introduction

The
Rosenblatt
Perceptron

The Perceptron
Learning Algorithm

Limitations

Combining
Perceptrons

Gradient
Based
Learning

Loss Functions

Optimizers

Fully
Connected
Networks

- Instead of taking the whole dataset for each iteration in each iteration, we randomly shuffle the data and take a batch.

Stochastic Gradient Descent

Introducción a
Redes
Neuronales

Introduction

The
Rosenblatt
Perceptron

The Perceptron
Learning Algorithm

Limitations

Combining
Perceptrons

Gradient
Based
Learning

Loss Functions

Optimizers

Fully
Connected
Networks

- Instead of taking the whole dataset for each iteration in each iteration, we randomly shuffle the data and take a batch.
- The path took by the algorithm is full of noise as compared to the gradient descent algorithm.

Stochastic Gradient Descent

Introducción a
Redes
Neuronales

Introduction

The
Rosenblatt
Perceptron

The Perceptron
Learning Algorithm

Limitations

Combining
Perceptrons

Gradient
Based
Learning

Loss Functions

Optimizers

Fully
Connected
Networks

- Instead of taking the whole dataset for each iteration in each iteration, we randomly shuffle the data and take a batch.
- The path took by the algorithm is full of noise as compared to the gradient descent algorithm.
- Due to an increase in the number of iterations, the overall computation time increases. Still, the computation cost is still less than that of the gradient descent optimizer.

Stochastic Gradient Descent

Introducción a
Redes
Neuronales

Introduction

The
Rosenblatt
Perceptron

The Perceptron
Learning Algorithm

Limitations

Combining
Perceptrons

Gradient
Based
Learning

Loss Functions

Optimizers

Fully
Connected
Networks

- Instead of taking the whole dataset for each iteration in each iteration, we randomly shuffle the data and take a batch.
- The path took by the algorithm is full of noise as compared to the gradient descent algorithm.
- Due to an increase in the number of iterations, the overall computation time increases. Still, the computation cost is still less than that of the gradient descent optimizer.
- If the data is enormous and computational time is an essential factor, stochastic gradient descent should be preferred over batch gradient descent algorithm.

Adagrad: Adaptive gradient descent

- The adaptive gradient descent algorithm uses different learning rates for each iteration. The change in learning rate depends upon the difference in the parameters during training.

Introducción a
Redes
Neuronales

Introduction

The
Rosenblatt
Perceptron

The Perceptron
Learning Algorithm

Limitations

Combining
Perceptrons

Gradient
Based
Learning

Loss Functions

Optimizers

Fully
Connected
Networks

Adagrad: Adaptive gradient descent

Introducción a
Redes
Neuronales

Introduction

The
Rosenblatt
Perceptron

The Perceptron
Learning Algorithm

Limitations

Combining
Perceptrons

Gradient
Based
Learning

Loss Functions

Optimizers

Fully
Connected
Networks

- The adaptive gradient descent algorithm uses different learning rates for each iteration. The change in learning rate depends upon the difference in the parameters during training.
- The more the weights change, the least the learning rate changes.

Adagrad: Adaptive gradient descent

Introducción a
Redes
Neuronales

Introduction

The
Rosenblatt
Perceptron

The Perceptron
Learning Algorithm
Limitations
Combining
Perceptrons

Gradient
Based
Learning

Loss Functions
Optimizers

Fully
Connected
Networks

- The adaptive gradient descent algorithm uses different learning rates for each iteration. The change in learning rate depends upon the difference in the parameters during training.
- The more the weights change, the least the learning rate changes.
- The benefit of using Adagrad is that it abolishes the need to modify the learning rate manually. It is more reliable than gradient descent algorithmss, as it reaches convergence at a higher speed.

Adagrad: Adaptive gradient descent

Introducción a
Redes
Neuronales

Introduction

The
Rosenblatt
Perceptron

The Perceptron
Learning Algorithm
Limitations
Combining
Perceptrons

Gradient
Based
Learning
Loss Functions
Optimizers

Fully
Connected
Networks

- The adaptive gradient descent algorithm uses different learning rates for each iteration. The change in learning rate depends upon the difference in the parameters during training.
- The more the weights change, the least the learning rate changes.
- The benefit of using Adagrad is that it abolishes the need to modify the learning rate manually. It is more reliable than gradient descent algorithmss, as it reaches convergence at a higher speed.
- One downside of AdaGrad optimizer is that it decreases the learning rate aggressively and monotonically. There might be a point when the learning rate becomes extremely small.

Adam: ADaptive Moment estimation

- It is an extension of stochastic gradient descent.

Introducción a
Redes
Neuronales

Introduction

The
Rosenblatt
Perceptron

The Perceptron
Learning Algorithm

Limitations

Combining
Perceptrons

Gradient
Based
Learning

Loss Functions

Optimizers

Fully
Connected
Networks

Adam: ADAptive Moment estimation

Introducción a
Redes
Neuronales

Introduction

The
Rosenblatt
Perceptron

The Perceptron
Learning Algorithm

Limitations

Combining
Perceptrons

Gradient
Based
Learning

Loss Functions

Optimizers

Fully
Connected
Networks

- It is an extension of stochastic gradient descent.
- Adam optimizer updates the learning rate for each network weight individually.

Adam: ADaptive Moment estimation

Introducción a
Redes
Neuronales

Introduction

The
Rosenblatt
Perceptron

The Perceptron
Learning Algorithm

Limitations

Combining
Perceptrons

Gradient
Based
Learning

Loss Functions

Optimizers

Fully
Connected
Networks

- It is an extension of stochastic gradient descent.
- Adam optimizer updates the learning rate for each network weight individually.
- The Adam optimizers inherit the features of both Adagrad and RMSProp algorithms.

Adam: ADaptive Moment estimation

Introducción a
Redes
Neuronales

Introduction

The
Rosenblatt
Perceptron

The Perceptron
Learning Algorithm

Limitations

Combining
Perceptrons

Gradient
Based
Learning

Loss Functions

Optimizers

Fully
Connected
Networks

- It is an extension of stochastic gradient descent.
- Adam optimizer updates the learning rate for each network weight individually.
- The Adam optimizers inherit the features of both Adagrad and RMSProp algorithms.
- Instead of adapting learning rates based upon the first moment (mean), it also uses the second moment of the gradients (variance).

Adam: ADaptive Moment estimation

Introducción a
Redes
Neuronales

Introduction

The
Rosenblatt
Perceptron

The Perceptron
Learning Algorithm

Limitations

Combining
Perceptrons

Gradient
Based
Learning

Loss Functions

Optimizers

Fully
Connected
Networks

- It is an extension of stochastic gradient descent.
- Adam optimizer updates the learning rate for each network weight individually.
- The Adam optimizers inherit the features of both Adagrad and RMSProp algorithms.
- Instead of adapting learning rates based upon the first moment (mean), it also uses the second moment of the gradients (variance).
- It is often used as a default optimization algorithm. It has faster running time, low memory requirements, and requires less tuning than any other optimization algorithm.

Adam: ADaptive Moment estimation

Introducción a
Redes
Neuronales

Introduction

The
Rosenblatt
Perceptron

The Perceptron
Learning Algorithm

Limitations

Combining
Perceptrons

Gradient
Based
Learning

Loss Functions

Optimizers

Fully
Connected
Networks

- It is an extension of stochastic gradient descent.
- Adam optimizer updates the learning rate for each network weight individually.
- The Adam optimizers inherit the features of both Adagrad and RMSProp algorithms.
- Instead of adapting learning rates based upon the first moment (mean), it also uses the second moment of the gradients (variance).
- It is often used as a default optimization algorithm. It has faster running time, low memory requirements, and requires less tuning than any other optimization algorithm.
- It tends to focus on faster computation time, it might not generalize the data well enough.

Optimizers

Introducción a
Redes
Neuronales

Introduction

The
Rosenblatt
Perceptron

The Perceptron
Learning Algorithm

Limitations

Combining
Perceptrons

Gradient
Based
Learning

Loss Functions

Optimizers

Fully
Connected
Networks

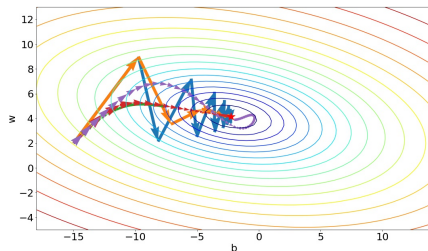
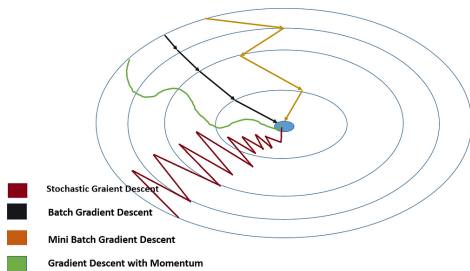


Table of Contents

Introducción a
Redes
Neuronales

Introduction

The
Rosenblatt
Perceptron

The Perceptron
Learning Algorithm

Limitations

Combining
Perceptrons

Gradient
Based
Learning

Loss Functions

Optimizers

Fully
Connected
Networks

1 Introduction

2 The Rosenblatt Perceptron

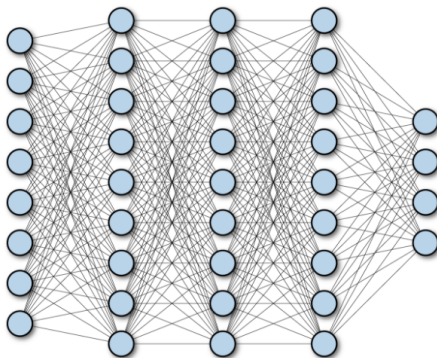
- The Perceptron Learning Algorithm
- Limitations
- Combining Perceptrons

3 Gradient Based Learning

- Loss Functions
- Optimizers

4 Fully Connected Networks

Fully Connected Networks



A fully connected neural network consists of a series of fully connected layers that connect every neuron in one layer to every neuron in the other layer.

Advantages and Disadvantages

Introducción a Redes Neuronales

Introduction

The Rosenblatt Perceptron

The Perceptron Learning Algorithm

Limitations

Combining Perceptrons

Gradient Based Learning

Loss Functions

Optimizers

Fully Connected Networks

- The major advantage of fully connected networks is that they are no special assumptions needed to be made about the input.

Advantages and Disadvantages

Introducción a
Redes
Neuronales

Introduction

The
Rosenblatt
Perceptron

The Perceptron
Learning Algorithm

Limitations

Combining
Perceptrons

Gradient
Based
Learning

Loss Functions

Optimizers

Fully
Connected
Networks

- The major advantage of fully connected networks is that they are no special assumptions needed to be made about the input.
- While being structure agnostic makes fully connected networks very broadly applicable, such networks do tend to have weaker performance than special-purpose networks tuned to the structure of a problem space.

The Algorithm

Introducción a
Redes
Neuronales

Introduction

The
Rosenblatt
Perceptron

The Perceptron
Learning Algorithm

Limitations

Combining
Perceptrons

Gradient
Based
Learning

Loss Functions

Optimizers

Fully
Connected
Networks

The algorithm consists of three simple steps:

- First, present one or more training examples to the neural network: **Feed-Forward**.

The Algorithm

Introducción a
Redes
Neuronales

Introduction

The
Rosenblatt
Perceptron

The Perceptron
Learning Algorithm

Limitations

Combining
Perceptrons

Gradient
Based
Learning

Loss Functions

Optimizers

Fully
Connected
Networks

The algorithm consists of three simple steps:

- First, present one or more training examples to the neural network: **Feed-Forward**.
- Second, compare the output of the neural network to the desired value: **Loss function**.

The Algorithm

Introducción a
Redes
Neuronales

Introduction

The
Rosenblatt
Perceptron

The Perceptron
Learning Algorithm

Limitations

Combining
Perceptrons

Gradient
Based
Learning

Loss Functions

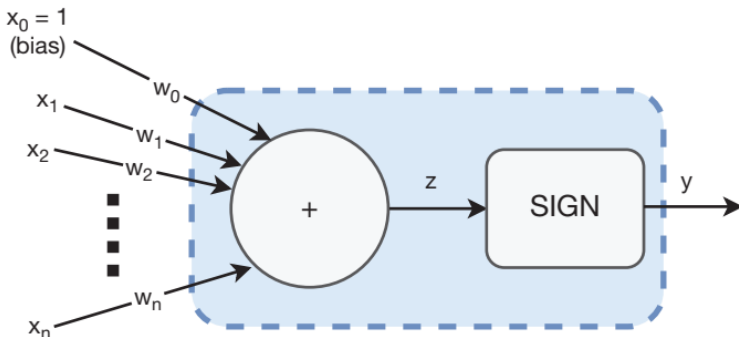
Optimizers

Fully
Connected
Networks

The algorithm consists of three simple steps:

- First, present one or more training examples to the neural network: **Feed-Forward**.
- Second, compare the output of the neural network to the desired value: **Loss function**.
- Finally, adjust the weights to make the output get closer to the desired value using gradient descent: **Back propagation**.

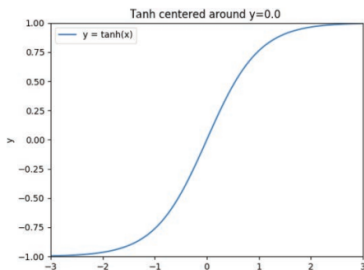
Feed-Forward



The sign function is an activation function of the neuron. It is not smooth because of the discontinuity in 0.

Other Activation Functions

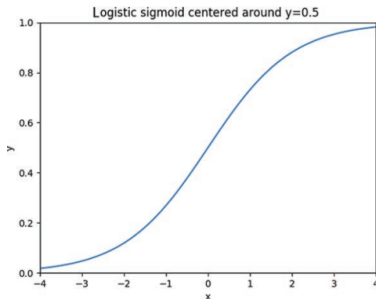
$$\tanh(x) = \frac{e^{2x} - 1}{e^{2x} + 1}$$



$$\tanh'(x) = 1 - \tanh^2(x)$$

Other Activation Functions

$$S(x) = \frac{e^x}{e^x + 1}$$



$$S'(x) = S(x)(1 - S(x))$$

Choice of Activation Functions

Introducción a
Redes
Neuronales

Introduction

The
Rosenblatt
Perceptron

The Perceptron
Learning Algorithm
Limitations
Combining
Perceptrons

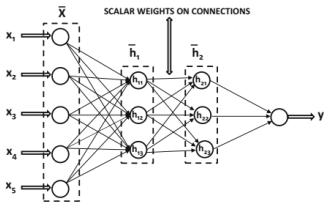
Gradient
Based
Learning

Loss Functions
Optimizers

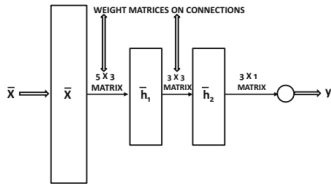
Fully
Connected
Networks

There exist a large number of activation functions. Two popular choices are tanh and the logistic sigmoid function. When picking between the two, choose tanh for hidden layers and logistic sigmoid for the output layer.

Vector notation



(c) Scalar notation and architecture



(d) Vector notation and architecture

$$\bar{h}_1 = \Phi_1(W_1^T \bar{x})$$

$$\bar{h}_2 = \Phi_2(W_2^T \bar{h}_1)$$

$$\bar{o} = \Phi_3(W_3^T \bar{h}_2)$$

Φ_j son las funciones de activación y W_j , las matrices de pesos.

Feed-Forward: An Example

Introducción a Redes Neuronales

Introduction

The Rosenblatt Perceptron

The Perceptron Learning Algorithm

Limitations

Combining Perceptrons

Gradient Based Learning

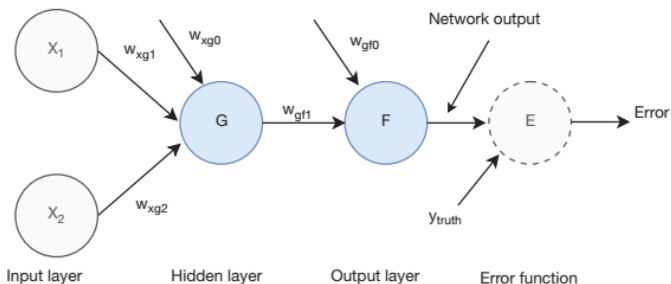
Loss Functions

Optimizers

Fully Connected Networks

...

Example: Back-propagation



This neural network implements

$$\hat{y} = S(w_{gf0} + w_{gf1} \tanh(w_{xg0} + w_{xg1}x_1 + w_{xg2}x_2))$$

We use the MSE:

$$e(f) = \frac{(y - f)^2}{2}.$$

Example: Back-propagation

- The error is given by

$$\text{Error} = \frac{1}{2} (y - S(w_{gf0} + w_{gf1} \tanh(w_{xg0} + w_{xg1}x_1 + w_{xg2}x_2)))$$

- We write it as:

$$e(f) = \frac{1}{2}(y - f)^2$$

$$f(z_f) = S(z_f)$$

$$z_f(w_{gf0}, w_{gf1}, g) = w_{gf0} + w_{gf1}g$$

$$g(z_g) = \tanh(z_g)$$

$$z_g(w_{xg0}, w_{xg1}, w_{xg2}) = w_{xg0} + w_{xg1}x_1 + w_{xg2}x_2$$

Example: Back-propagation

$$\frac{\partial e}{\partial w_{gf0}} = -(y - f) \cdot S'(z_f)$$

$$\frac{\partial e}{\partial w_{gf1}} = -(y - f) \cdot S'(z_f) \cdot g$$

$$\frac{\partial e}{\partial w_{xg0}} = -(y - f) \cdot S'(z_f) \cdot w_{gf1} \cdot \tanh'(z_g)$$

$$\frac{\partial e}{\partial w_{xg1}} = -(y - f) \cdot S'(z_f) \cdot w_{gf1} \cdot \tanh'(z_g) \cdot x_1$$

$$\frac{\partial e}{\partial w_{xg2}} = -(y - f) \cdot S'(z_f) \cdot w_{gf1} \cdot \tanh'(z_g) \cdot x_2$$

Example: Back-propagation

Finally, we update the weights, via gradient descent

$$w_{gf0} \leftarrow w_{gf0} + \eta \frac{\partial e}{\partial w_{gf0}}$$

$$w_{gf1} \leftarrow w_{gf1} + \eta \frac{\partial e}{\partial w_{gf1}}$$

$$w_{xg0} \leftarrow w_{xg0} + \eta \frac{\partial e}{\partial w_{xg0}}$$

$$w_{xg1} \leftarrow w_{xg1} + \eta \frac{\partial e}{\partial w_{xg1}}$$

$$w_{xg2} \leftarrow w_{xg2} + \eta \frac{\partial e}{\partial w_{xg2}}$$