

# The perceptron algorithm

Christos Dimitrakakis

October 7, 2025

# Outline

## The Perceptron

- Introduction

- The algorithm

## Gradient methods

- Gradients for optimisation

- The perceptron as a gradient algorithm

## Lab and Assignment

## The Perceptron

- Introduction
- The algorithm

## Gradient methods

- Gradients for optimisation
- The perceptron as a gradient algorithm

## Lab and Assignment

# Guessing gender from height

- ▶ Feature space  $\mathcal{X} \subset \mathbb{R}$ : e.g. height
- ▶ Label space  $\mathcal{Y} = \{-1, 1\}$ : e.g. gender
- ▶ Can we find some  $\theta_1 \in \mathbb{R}$  and a direction  $\theta_0 \in \{-1, +1\}$  so as to separate the genders?

# Guessing gender from height

- ▶ Feature space  $\mathcal{X} \subset \mathbb{R}$ : e.g. height
- ▶ Label space  $\mathcal{Y} = \{-1, 1\}$ : e.g. gender
- ▶ Can we find some  $\theta_1 \in \mathbb{R}$  and a direction  $\theta_0 \in \{-1, +1\}$  so as to separate the genders?

## Online learning: At time $t$

- ▶ We choose a separator  $\theta_0^t, \theta_1^t$
- ▶ We observe a new datapoint  $x_t, y_t$
- ▶ We make a mistake at time  $t$  if:

$$\theta^t x_t - \theta_0^t \leq 0.$$

- ▶ If we stop making mistakes, then we are classifying everything perfectly.

# Guessing gender from height

- ▶ Feature space  $\mathcal{X} \subset \mathbb{R}$ : e.g. height
- ▶ Label space  $\mathcal{Y} = \{-1, 1\}$ : e.g. gender
- ▶ Can we find some  $\theta_1 \in \mathbb{R}$  and a direction  $\theta_0 \in \{-1, +1\}$  so as to separate the genders?

## Online learning: At time $t$

- ▶ We choose a separator  $\theta_0^t, \theta_1^t$
- ▶ We observe a new datapoint  $x_t, y_t$
- ▶ We make a mistake at time  $t$  if:

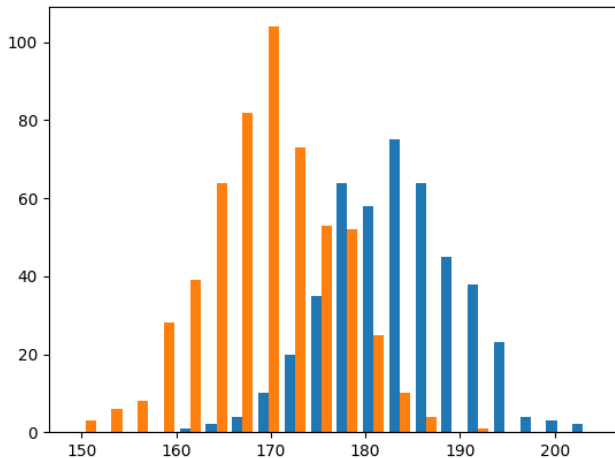
$$\theta^t x_t - \theta_0^t \leq 0.$$

- ▶ If we stop making mistakes, then we are classifying everything perfectly.

Can you find a threshold that makes a small number of mistakes?

`./src/Perceptron/perceptron_simple.py`

## Non-separable classes



- In general, we cannot perfectly classify everything
- But we can estimate  $\mathbb{P}(y | x)$  ... more on this later.

## A more complex example

- ▶ Feature space  $\mathcal{X} \subset \mathbb{R}^n$ : e.g. height and weight for  $n = 2$
- ▶ Label space  $\mathcal{Y} = \{-1, 1\}$ : e.g. gender
- ▶ Can we find some line so as to separate the genders?

-./src/Perceptron/show\_class\_data\_labels.py

- ▶ Is there an algorithm for doing so?



# A linear classifier

## The separating hyperplane

We now have parameters  $\theta_0 \in \mathbb{R}$  and  $\theta \in \mathbb{R}^n$  defining a **hyperplane**  
 $f(x) = 0$  in  $\mathbb{R}^n$

$$f(x) = \theta_0 + \theta^\top x = \theta_0 + \sum_{i=1}^n \theta_i x_i.$$

# A linear classifier

## The separating hyperplane

We now have parameters  $\theta_0 \in \mathbb{R}$  and  $\theta \in \mathbb{R}^n$  defining a **hyperplane**  $f(x) = 0$  in  $\mathbb{R}^n$

$$f(x) = \theta_0 + \theta^\top x = \theta_0 + \sum_{i=1}^n \theta_i x_i.$$

If we augment  $x$  with an additional component  $x_0 = 1$ , we can write

$$f(x) = \theta^\top x = \sum_{i=0}^n \theta_i x_i.$$

# A linear classifier

## The separating hyperplane

We now have parameters  $\theta_0 \in \mathbb{R}$  and  $\theta \in \mathbb{R}^n$  defining a **hyperplane**  $f(x) = 0$  in  $\mathbb{R}^n$

$$f(x) = \theta_0 + \theta^\top x = \theta_0 + \sum_{i=1}^n \theta_i x_i.$$

If we augment  $x$  with an additional component  $x_0 = 1$ , we can write

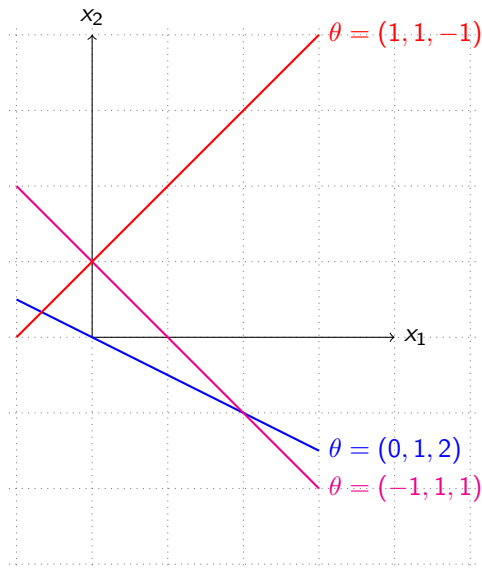
$$f(x) = \theta^\top x = \sum_{i=0}^n \theta_i x_i.$$

## The classifier

The **perceptron decision rule** is  $\pi(x) = \text{sign}(f(x))$

- ▶ If  $f(x) \geq 0$ , we assign class +1
- ▶ If  $f(x) < 0$ , we assign class -1

# Hyperplanes in 2 dimensions (lines)



- These lines are the solution to  $f(x) = 0$ .
- $\theta_1, \theta_2$  can be thought of as the line to each hyperplane

# The Perceptron



Figure: Pitts



Figure: Rosenblatt

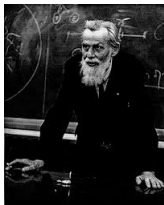


Figure: McCulloch

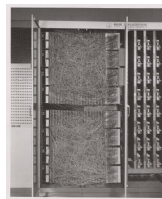


Figure: Perceptron Mark I

# The perceptron algorithm

## Input

- ▶ Feature space  $X \subset \mathbb{R}^n$ .
- ▶ Label space  $Y = \{-1, 1\}$ .
- ▶ Data  $(x_t, y_t)$ ,  $t \in [T]$ , with  $x_t \in X, y_t \in Y$ .

# The perceptron algorithm

## Input

- ▶ Feature space  $X \subset \mathbb{R}^n$ .
- ▶ Label space  $Y = \{-1, 1\}$ .
- ▶ Data  $(x_t, y_t)$ ,  $t \in [T]$ , with  $x_t \in X, y_t \in Y$ .

## Algorithm

- ▶  $\theta^0 \sim \text{Normal}^n(0, I)$ . % Initialise parameters
- ▶ For  $t = 1, \dots, T$ 
  - ▶  $a_t = \text{sgn}(\theta^t \cdot x_t)$ . % Classify example
  - ▶ If  $a_t \neq y_t$ 
    - ▶  $\theta^t = \theta^{t-1} + y_t x_t$  % Move hyperplane
  - ▶ Else
    - ▶  $\theta^t = \theta^{t-1}$  % Do nothing for correct examples
  - ▶ EndIf
- ▶ Return  $\theta^T$

# Perceptron examples

## Example 1: One-dimensional data

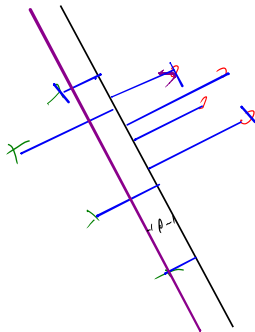
- ▶ Done on the board
- ▶ Shows how the algorithm works.
- ▶ Demonstrates the idea of a margin

## Example 2: Two-dimensional data

- ▶ See in-class programming exercise



# Margins and the perceptron theorem



- ▶ The **hyperplane**  $\theta^*$  separates the examples
- ▶ The **margin**  $\rho$  is the minimum distance  $\rho$  between  $\theta^*$  and any point.

## Theorem (Perceptron theorem)

The number of mistakes is bounded by  $\rho^{-2}$ , where  $\|x_t\| \leq 1$ ,  
 $\rho \leq y_t(x_t^\top \theta^*)$  for some **margin**  $\rho$  and **hyperplane**  $\theta^*$  with  $\|\theta^*\| = 1$ .

# Simple proof

- ▶ Scale data:  $\|x\| \leq 1$

## Simple proof

- ▶ Scale data:  $\|x\| \leq 1$
- ▶ Separating plane:  $y_t(x_t \cdot \theta^*) \geq \rho \forall t, \|\theta^*\| = 1.$

## Simple proof

- ▶ Scale data:  $\|x\| \leq 1$
- ▶ Separating plane:  $y_t(x_t \cdot \theta^*) \geq \rho \forall t, \|\theta^*\| = 1.$
- ▶ When we make an update:  $y_t(x_t \cdot \theta^t) \leq 0.$

## Simple proof

- ▶ Scale data:  $\|x\| \leq 1$
- ▶ Separating plane:  $y_t(x_t \cdot \theta^*) \geq \rho \forall t, \|\theta^*\| = 1.$
- ▶ When we make an update:  $y_t(x_t \cdot \theta^t) \leq 0.$
- ▶ At each mistake,  $\theta^t \cdot \theta^*$  grows by **at least  $\rho$** .

## Simple proof

- ▶ Scale data:  $\|x\| \leq 1$
- ▶ Separating plane:  $y_t(x_t \cdot \theta^*) \geq \rho \forall t, \|\theta^*\| = 1.$
- ▶ When we make an update:  $y_t(x_t \cdot \theta^t) \leq 0.$
- ▶ At each mistake,  $\theta^t \cdot \theta^*$  grows by **at least  $\rho$** .

## Simple proof

- ▶ Scale data:  $\|x\| \leq 1$
- ▶ Separating plane:  $y_t(x_t \cdot \theta^*) \geq \rho \forall t, \|\theta^*\| = 1.$
- ▶ When we make an update:  $y_t(x_t \cdot \theta^t) \leq 0.$
- ▶ At each mistake,  $\theta^t \cdot \theta^*$  grows by **at least  $\rho$** .

$$\theta^{t+1} \cdot \theta^* = (\theta^t + y_t x_t) \cdot \theta^* = \theta^t \cdot \theta^* + y_t(x_t \cdot \theta^*) \geq \theta^t \cdot \theta^* + \rho$$

## Simple proof

- ▶ Scale data:  $\|x\| \leq 1$
- ▶ Separating plane:  $y_t(x_t \cdot \theta^*) \geq \rho \forall t, \|\theta^*\| = 1.$
- ▶ When we make an update:  $y_t(x_t \cdot \theta^t) \leq 0.$
- ▶ At each mistake,  $\theta^t \cdot \theta^*$  grows by **at least  $\rho$** .

$$\theta^{t+1} \cdot \theta^* = (\theta^t + y_t x_t) \cdot \theta^* = \theta^t \cdot \theta^* + y_t(x_t \cdot \theta^*) \geq \theta^t \cdot \theta^* + \rho$$

- ▶ At each mistake,  $\theta \cdot \theta$  grows by **at most 1**.

$$\theta^{t+1} \cdot \theta^{t+1} = (\theta^t + y_t x_t) \cdot (\theta^t + y_t x_t) = \theta^t \cdot \theta^t + 2y_t(\theta^t \cdot x_t) + y_t^2(x_t \cdot x_t) \leq \theta^t \cdot \theta^t + 1$$



## Simple proof

- ▶ Scale data:  $\|x\| \leq 1$
- ▶ Separating plane:  $y_t(x_t \cdot \theta^*) \geq \rho \forall t, \|\theta^*\| = 1.$
- ▶ When we make an update:  $y_t(x_t \cdot \theta^t) \leq 0.$
- ▶ At each mistake,  $\theta^t \cdot \theta^*$  grows by **at least  $\rho$** .

$$\theta^{t+1} \cdot \theta^* = (\theta^t + yx_t) \cdot \theta^* = \theta^t \cdot \theta^* + y(x_t \cdot \theta^*) \geq \theta^t \cdot \theta^* + \rho$$

- ▶ At each mistake,  $\theta \cdot \theta$  grows by **at most 1**.

$$\theta^{t+1} \cdot \theta^{t+1} = (\theta^t + yx_t) \cdot (\theta^t + yx_t) = \theta^t \cdot \theta^t + 2y(\theta^t \cdot x_t) + y^2(x_t \cdot x_t) \leq \theta^t \cdot \theta^t + 1$$

## Putting it together

After  $M$  mistakes:

- ▶  $\theta^t \cdot \theta^* \geq M\rho$

## Simple proof

- ▶ Scale data:  $\|x\| \leq 1$
- ▶ Separating plane:  $y_t(x_t \cdot \theta^*) \geq \rho \forall t, \|\theta^*\| = 1.$
- ▶ When we make an update:  $y_t(x_t \cdot \theta^t) \leq 0.$
- ▶ At each mistake,  $\theta^t \cdot \theta^*$  grows by **at least  $\rho$** .

$$\theta^{t+1} \cdot \theta^* = (\theta^t + yx_t) \cdot \theta^* = \theta^t \cdot \theta^* + y(x_t \cdot \theta^*) \geq \theta^t \cdot \theta^* + \rho$$

- ▶ At each mistake,  $\theta^t \cdot \theta^t$  grows by **at most 1**.

$$\theta^{t+1} \cdot \theta^{t+1} = (\theta^t + yx_t) \cdot (\theta^t + yx_t) = \theta^t \cdot \theta^t + 2y(\theta^t \cdot x_t) + y^2(x_t \cdot x_t) \leq \theta^t \cdot \theta^t + 1$$

## Putting it together

After  $M$  mistakes:

- ▶  $\theta^t \cdot \theta^* \geq M\rho$
- ▶  $\theta^t \cdot \theta^t \leq M$

## Simple proof

- ▶ Scale data:  $\|x\| \leq 1$
- ▶ Separating plane:  $y_t(x_t \cdot \theta^*) \geq \rho \forall t, \|\theta^*\| = 1.$
- ▶ When we make an update:  $y_t(x_t \cdot \theta^t) \leq 0.$
- ▶ At each mistake,  $\theta^t \cdot \theta^*$  grows by **at least  $\rho$** .

$$\theta^{t+1} \cdot \theta^* = (\theta^t + y_t x_t) \cdot \theta^* = \theta^t \cdot \theta^* + y_t(x_t \cdot \theta^*) \geq \theta^t \cdot \theta^* + \rho$$

- ▶ At each mistake,  $\theta \cdot \theta$  grows by **at most 1**.

$$\theta^{t+1} \cdot \theta^{t+1} = (\theta^t + y_t x_t) \cdot (\theta^t + y_t x_t) = \theta^t \cdot \theta^t + 2y_t(\theta^t \cdot x_t) + y_t^2(x_t \cdot x_t) \leq \theta^t \cdot \theta^t + 1$$

## Putting it together

After  $M$  mistakes:

- ▶  $\theta^t \cdot \theta^* \geq M\rho$
- ▶  $\theta^t \cdot \theta^t \leq M$
- ▶ So  $M\rho \leq \theta^t \cdot \theta^* \leq \|\theta^t\| \cdot \|\theta^*\| = \|\theta^t\| = \sqrt{\theta^t \cdot \theta^t} \leq \sqrt{M}.$

## Simple proof

- ▶ Scale data:  $\|x\| \leq 1$
- ▶ Separating plane:  $y_t(x_t \cdot \theta^*) \geq \rho \forall t, \|\theta^*\| = 1.$
- ▶ When we make an update:  $y_t(x_t \cdot \theta^t) \leq 0.$
- ▶ At each mistake,  $\theta^t \cdot \theta^*$  grows by **at least  $\rho$** .

$$\theta^{t+1} \cdot \theta^* = (\theta^t + y_t x_t) \cdot \theta^* = \theta^t \cdot \theta^* + y_t(x_t \cdot \theta^*) \geq \theta^t \cdot \theta^* + \rho$$

- ▶ At each mistake,  $\theta^t \cdot \theta^t$  grows by **at most 1**.

$$\theta^{t+1} \cdot \theta^{t+1} = (\theta^t + y_t x_t) \cdot (\theta^t + y_t x_t) = \theta^t \cdot \theta^t + 2y_t(\theta^t \cdot x_t) + y_t^2(x_t \cdot x_t) \leq \theta^t \cdot \theta^t + 1$$

## Putting it together

After  $M$  mistakes:

- ▶  $\theta^t \cdot \theta^* \geq M\rho$
- ▶  $\theta^t \cdot \theta^t \leq M$
- ▶ So  $M\rho \leq \theta^t \cdot \theta^* \leq \|\theta^t\| \cdot \|\theta^*\| = \|\theta^t\| = \sqrt{\theta^t \cdot \theta^t} \leq \sqrt{M}.$
- ▶ Thus,  $M \leq \rho^{-2}.$

# Promise of the perceptron

## NEW NAVY DEVICE LEARNS BY DOING

Psychologist Shows Embryo  
of Computer Designed to  
Read and Grow Wiser

WASHINGTON, July 7 (UPI)—The Navy revealed the embryo of an electronic computer today that it expects will be able to walk, talk, see, write, reproduce itself and be conscious of its existence.

The embryo—the Weather Bureau's \$2,000,000 "704" computer—learned to differentiate between right and left after fifty attempts in the Navy's demonstration for newsmen.

The service said it would use this principle to build the first of its Perceptron thinking machines that will be able to read and write. It is expected to be finished in about a year at a cost of \$100,000.

Dr. Frank Rosenblatt, designer of the Perceptron, conducted the demonstration. He said the machine would be the first device to think as the human brain. As do human be-

ings, Perceptron will make mistakes at first, but will grow wiser as it gains experience, he said.

Dr. Rosenblatt, a research psychologist at the Cornell Aeronautical Laboratory, Buffalo, said Perceptrons might be fired to the planets as mechanical space explorers.

### Without Human Controls

The Navy said the perceptron would be the first non-living mechanism "capable of receiving, recognizing and identifying its surroundings without any human training or control."

The "brain" is designed to remember images and information it has perceived itself. Ordinary computers remember only what is fed into them on punch cards or magnetic tape.

Later Perceptrons will be able to recognize people and call out their names and instantly translate speech in one language to speech or writing in another language, it was predicted.

Mr. Rosenblatt said in principle it would be possible to build brains that could reproduce themselves on an assembly line and which would be conscious of their existence.

## 1958 New York Times...

In today's demonstration, the "704" was fed two cards, one with squares marked on the left side and the other with squares on the right side.

### Learns by Doing

In the first fifty trials, the machine made no distinction between them. It then started registering a "Q" for the left squares and "O" for the right squares.

Dr. Rosenblatt said he could explain why the machine learned only in highly technical terms. But he said the computer had undergone a "self-induced change in the wiring diagram."

The first Perceptron will have about 1,000 electronic "association cells" receiving electrical impulses from an eye-like scanning device with 400 photo-cells. The human brain has 10,000,000,000 responsive cells, including 100,000,000 connections with the eyes.

# Promise versus reality

## Focus on classification

- ▶ Rosenblatt only consider classification problems
- ▶ Many problems in learning and AI are not simply classification problems
- ▶ Classification requires labels. These are not always easily available.

## Separable representation assumption

- ▶ Rosenblatt assumed that there was a representation available that would allow us to distinguish classes.
- ▶ However, it is not clear *a priori* how to obtain such a data representation from the data. Progress followed roughly these steps:
  - ▶ Hand-crafted features
  - ▶ Random features
  - ▶ Multi-layer perceptrons, hand-crafted architectures, and backpropagation
  - ▶ Attention mechanisms

## The Perceptron

Introduction

The algorithm

## Gradient methods

Gradients for optimisation

The perceptron as a gradient algorithm

## Lab and Assignment

# The gradient descent method: one dimension

- ▶ Function to minimise  $f : \mathbb{R} \rightarrow \mathbb{R}$ .
- ▶ Derivative  $\frac{d}{d\theta} f(\theta)$



# The gradient descent method: one dimension

- ▶ Function to minimise  $f : \mathbb{R} \rightarrow \mathbb{R}$ .
- ▶ Derivative  $\frac{d}{d\theta} f(\theta)$

## Gradient descent algorithm

- ▶ Input: initial value  $\theta^0$ , **learning rate** schedule  $\alpha_t$
- ▶ For  $t = 1, \dots, T$ 
  - ▶  $\theta^{t+1} = \theta^t - \alpha_t \frac{d}{d\theta} f(\theta^t)$
- ▶ Return  $\theta^T$

# The gradient descent method: one dimension

- ▶ Function to minimise  $f : \mathbb{R} \rightarrow \mathbb{R}$ .
- ▶ Derivative  $\frac{d}{d\theta} f(\theta)$

## Gradient descent algorithm

- ▶ Input: initial value  $\theta^0$ , **learning rate** schedule  $\alpha_t$
- ▶ For  $t = 1, \dots, T$ 
  - ▶  $\theta^{t+1} = \theta^t - \alpha_t \frac{d}{d\theta} f(\theta^t)$
- ▶ Return  $\theta^T$

## Properties

- ▶ If  $\sum_t \alpha_t = \infty$  and  $\sum_t \alpha_t^2 < \infty$ , it finds a local minimum  $\theta^T$ , i.e. there is  $\epsilon > 0$  so that

$$f(\theta^T) < f(\theta), \forall \theta : \|\theta^T - \theta\| < \epsilon.$$

# Gradient methods for expected value

Estimate the expected value

$x_t \sim P$  with  $\mathbb{E}_P[x_t] = \mu$ .

# Gradient methods for expected value

Estimate the expected value

$x_t \sim P$  with  $\mathbb{E}_P[x_t] = \mu$ .

Objective: mean squared error

Here  $\ell(x, \theta) = (x - \theta)^2$ .

$$\min_{\theta} \mathbb{E}_P[(x_t - \theta)^2].$$

# Gradient methods for expected value

## Estimate the expected value

$x_t \sim P$  with  $\mathbb{E}_P[x_t] = \mu$ .

## Objective: mean squared error

Here  $\ell(x, \theta) = (x - \theta)^2$ .

$$\min_{\theta} \mathbb{E}_P[(x_t - \theta)^2].$$

## Exact gradient update

If we know  $P$ , then we can calculate

$$\theta^{t+1} = \theta^t - \alpha_t \frac{d}{d\theta} \mathbb{E}_P[(x - \theta^t)^2] \quad (1)$$

$$\frac{d}{d\theta} \mathbb{E}_P[(x - \theta^t)^2] = 2 \mathbb{E}_P[x] - \theta^t \quad (2)$$

# Gradient for mean estimation

- ▶ Let us show this in detail

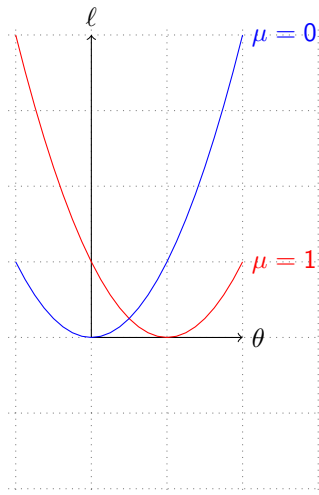
$$\begin{aligned}\frac{d}{d\theta} \mathbb{E}_P[(x - \theta)^2] &= \int_{-\infty}^{\infty} dP(x) \frac{d}{d\theta} (x - \theta)^2 \\ &= \int_{-\infty}^{\infty} dP(x) 2(x - \theta) \\ &= 2 \mathbb{E}_P[x] - 2\theta.\end{aligned}$$

- ▶ If we set the derivative to zero, then we find the optimal solution:

$$\theta^* = \mathbb{E}_P[x]$$

- ▶ How can we do this if we only have data  $x_t \sim P$ ?

# Mean-squared error cost function



Here we see a plot of  $\ell(\mu, \theta) = (\theta - \mu)^2$ .

# Gradient descent and sampling

## Theorem (Sampling)

For any bounded random variable  $f$ ,

$$\mathbb{E}_P[f] = \int_X dP(x) f(x) = \lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T f(x_t) = \mathbb{E}_P \left[ \frac{1}{T} \sum_{t=1}^T f(x_t) \right], \quad x_t \sim P$$

## Sampled derivatives

$$\frac{d}{d\theta} \mathbb{E}_P[f(\theta, x)] = \mathbb{E}_P \left[ \frac{d}{d\theta} f(\theta, x) \right] \approx \frac{1}{T} \sum_{t=1}^T \frac{d}{d\theta} f(\theta, x_t)$$



# Gradient descent and sampling

## Theorem (Sampling)

For any bounded random variable  $f$ ,

$$\mathbb{E}_P[f] = \int_X dP(x) f(x) = \lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T f(x_t) = \mathbb{E}_P \left[ \frac{1}{T} \sum_{t=1}^T f(x_t) \right], \quad x_t \sim P$$

## Sampled derivatives

$$\frac{d}{d\theta} \mathbb{E}_P[f(\theta, x)] = \mathbb{E}_P \left[ \frac{d}{d\theta} f(\theta, x) \right] \approx \frac{1}{T} \sum_{t=1}^T \frac{d}{d\theta} f(\theta, x_t)$$

- Batch updates (over all samples)

$$\theta^{k+1} = \theta^k - 2\alpha_k \frac{d}{d\theta} \sum_{t=1}^T f(\theta, x_t) / T.$$

# Gradient descent and sampling

## Theorem (Sampling)

For any bounded random variable  $f$ ,

$$\mathbb{E}_P[f] = \int_X dP(x) f(x) = \lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T f(x_t) = \mathbb{E}_P \left[ \frac{1}{T} \sum_{t=1}^T f(x_t) \right], \quad x_t \sim P$$

## Sampled derivatives

$$\frac{d}{d\theta} \mathbb{E}_P[f(\theta, x)] = \mathbb{E}_P \left[ \frac{d}{d\theta} f(\theta, x) \right] \approx \frac{1}{T} \sum_{t=1}^T \frac{d}{d\theta} f(\theta, x_t)$$

- Batch updates (over all samples)

$$\theta^{k+1} = \theta^k - 2\alpha_k \frac{d}{d\theta} \sum_{t=1}^T f(\theta, x_t) / T.$$

- Online updates (after each sample)

$$\theta^{t+1} = \theta^t - 2\alpha_t \frac{d}{d\theta} f(\theta, x_t)$$

# The gradient method

- ▶ Function to minimise  $f : \mathbb{R}^n \rightarrow \mathbb{R}$ .
- ▶ Derivative  $\nabla_{\theta} f(\theta) = \left( \frac{\partial f(\theta)}{\partial \theta_1}, \dots, \frac{\partial f(\theta)}{\partial \theta_n} \right)$ , where  $\frac{\partial f}{\partial \theta_n}$  denotes the **partial** derivative, i.e. varying one argument and keeping the others fixed.

## Gradient descent algorithm

- ▶ Input: initial value  $\theta^0$ , learning rate schedule  $\alpha_t$
- ▶ For  $t = 1, \dots, T$ 
  - ▶  $\theta^{t+1} = \theta^t - \alpha_t \nabla_{\theta} f(\theta^t)$
- ▶ Return  $\theta^T$

## Properties

- ▶ If  $\sum_t \alpha_t = \infty$  and  $\sum_t \alpha_t^2 < \infty$ , it finds a local minimum  $\theta^T$ , i.e. there is  $\epsilon > 0$  so that

$$f(\theta^T) < f(\theta), \forall \theta : \|\theta^T - \theta\| < \epsilon.$$

# Stochastic gradient method

This is the same as the gradient method, but with added noise:

- ▶  $\theta^{t+1} = \theta^t - \alpha_t [\nabla_{\theta} f(\theta^t) + \omega_t]$
- ▶  $\mathbb{E}[\omega_t] = 0$  is sufficient for convergence.

# Stochastic gradient method

This is the same as the gradient method, but with added noise:

- ▶  $\theta^{t+1} = \theta^t - \alpha_t [\nabla_{\theta} f(\theta^t) + \omega_t]$
- ▶  $\mathbb{E}[\omega_t] = 0$  is sufficient for convergence.

## Example (When the cost is an expectation)

In machine learning, the cost is frequently an expectation of some function  $\ell$ ,

$$f(\theta) = \int_{\mathcal{X}} dP(x) \ell(x, \theta)$$

This can be approximated with a sample

$$f(\theta) \approx \frac{1}{T} \sum_t \ell(x_t, \theta)$$

The same holds for the gradient:

$$\nabla_{\theta} f(\theta) = \int_{\mathcal{X}} dP(x) \nabla_{\theta} \ell(x, \theta) \approx \frac{1}{T} \sum_t \nabla_{\theta} \ell(x_t, \theta)$$

# Perceptron algorithm as gradient descent

## Target error function

$$\mathbb{E}_{\mathbf{P}}^{\theta}[\ell] = \int_{\mathcal{X}} d\mathbf{P}(x) \sum_y \mathbf{P}(y|x) \ell(x, y, \theta)$$

Minimises the error on the true distribution.

# Perceptron algorithm as gradient descent

## Target error function

$$\mathbb{E}_{\mathbf{P}}^{\theta}[\ell] = \int_{\mathcal{X}} d\mathbf{P}(x) \sum_y \mathbf{P}(y|x) \ell(x, y, \theta)$$

Minimises the error on the true distribution.

## Empirical error function

$$\mathbb{E}_{\mathbf{D}}^{\theta}[\ell] = \frac{1}{T} \sum_{t=1}^T \ell(x_t, y_t, \theta), \quad \mathbf{D} = (x_t, y_t)_{t=1}^T, \quad x_t, y_t \sim P.$$

Minimises the error on the empirical distribution.

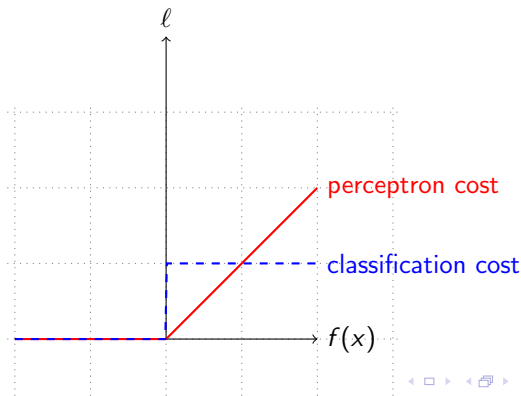
# Cost functions and the chain rule

## Perceptron cost function

The cost of each example

$$\ell(x, y, \theta) = \overbrace{\mathbb{I}\{y(x^\top \theta) < 0\}}^{\text{misclassified?}} \overbrace{[-y(x^\top \theta)]}^{\text{margin of error}} \quad (3)$$

where the **indicator function**  $\mathbb{I}\{A\}$  is 1 when  $A$  is true and 0 otherwise.





## Derivative of the perceptron cost function

The total loss over the data is defined as  $L(D, \theta) = \sum_{(x,y) \in D} \ell(x, y, \theta)$ . Taking the derivative, we have

$$\nabla_{\theta} L(D, \theta) = \nabla_{\theta} \sum_{(x,y) \in D} \ell(x, y, \theta) = \sum_{(x,y) \in D} \nabla_{\theta} \ell(x, y, \theta)$$

### Reminder: The chain rule

Let  $z = g(y)$ ,  $y = f(x)$  so that  $z = g(f(x))$ . Then  $\frac{dz}{dx} = \frac{dz}{dy} \frac{dy}{dx}$

---

<sup>1</sup>Note that for  $z = 0$ , there is no derivative but the subdifferential  $[0, 1]$ .

# Derivative of the perceptron cost function

The total loss over the data is defined as  $L(D, \theta) = \sum_{(x,y) \in D} \ell(x, y, \theta)$ .  
Taking the derivative, we have

$$\nabla_{\theta} L(D, \theta) = \nabla_{\theta} \sum_{(x,y) \in D} \ell(x, y, \theta) = \sum_{(x,y) \in D} \nabla_{\theta} \ell(x, y, \theta)$$

## Reminder: The chain rule

Let  $z = g(y)$ ,  $y = f(x)$  so that  $z = g(f(x))$ . Then  $\frac{dz}{dx} = \frac{dz}{dy} \frac{dy}{dx}$

## Applying the chain rule to calculate the gradient

- Set  $z = \theta^{\top} xy$ .

---

<sup>1</sup>Note that for  $z = 0$ , there is no derivative but the subdifferential  $[0, 1]$ .

# Derivative of the perceptron cost function

The total loss over the data is defined as  $L(D, \theta) = \sum_{(x,y) \in D} \ell(x, y, \theta)$ . Taking the derivative, we have

$$\nabla_{\theta} L(D, \theta) = \nabla_{\theta} \sum_{(x,y) \in D} \ell(x, y, \theta) = \sum_{(x,y) \in D} \nabla_{\theta} \ell(x, y, \theta)$$

## Reminder: The chain rule

Let  $z = g(y)$ ,  $y = f(x)$  so that  $z = g(f(x))$ . Then  $\frac{dz}{dx} = \frac{dz}{dy} \frac{dy}{dx}$

## Applying the chain rule to calculate the gradient

- ▶ Set  $z = \theta^{\top} xy$ .
- ▶  $\nabla_{\theta} \ell(x, y, \theta) = -\frac{d}{dz}(\mathbb{I}\{z < 0\} z) \nabla_{\theta} z$ .

---

<sup>1</sup>Note that for  $z = 0$ , there is no derivative but the subdifferential  $[0, 1]$ .

# Derivative of the perceptron cost function

The total loss over the data is defined as  $L(D, \theta) = \sum_{(x,y) \in D} \ell(x, y, \theta)$ . Taking the derivative, we have

$$\nabla_{\theta} L(D, \theta) = \nabla_{\theta} \sum_{(x,y) \in D} \ell(x, y, \theta) = \sum_{(x,y) \in D} \nabla_{\theta} \ell(x, y, \theta)$$

## Reminder: The chain rule

Let  $z = g(y)$ ,  $y = f(x)$  so that  $z = g(f(x))$ . Then  $\frac{dz}{dx} = \frac{dz}{dy} \frac{dy}{dx}$

## Applying the chain rule to calculate the gradient

- ▶ Set  $z = \theta^{\top} xy$ .
- ▶  $\nabla_{\theta} \ell(x, y, \theta) = -\frac{d}{dz}(\mathbb{I}\{z < 0\} z) \nabla_{\theta} z$ .
- ▶  $\frac{d}{dz}(\mathbb{I}\{z < 0\} z) = \frac{d}{dz} \mathbb{I}\{z < 0\} \times z + \mathbb{I}\{z < 0\} \times 1 = \mathbb{I}\{z < 0\}$ .<sup>1</sup>

---

<sup>1</sup>Note that for  $z = 0$ , there is no derivative but the subdifferential  $[0, 1]$ .

## Derivative of the perceptron cost function

The total loss over the data is defined as  $L(D, \theta) = \sum_{(x,y) \in D} \ell(x, y, \theta)$ . Taking the derivative, we have

$$\nabla_{\theta} L(D, \theta) = \nabla_{\theta} \sum_{(x,y) \in D} \ell(x, y, \theta) = \sum_{(x,y) \in D} \nabla_{\theta} \ell(x, y, \theta)$$

### Reminder: The chain rule

Let  $z = g(y)$ ,  $y = f(x)$  so that  $z = g(f(x))$ . Then  $\frac{dz}{dx} = \frac{dz}{dy} \frac{dy}{dx}$

### Applying the chain rule to calculate the gradient

- ▶ Set  $z = \theta^{\top} xy$ .
- ▶  $\nabla_{\theta} \ell(x, y, \theta) = -\frac{d}{dz}(\mathbb{I}\{z < 0\} z) \nabla_{\theta} z$ .
- ▶  $\frac{d}{dz}(\mathbb{I}\{z < 0\} z) = \frac{d}{dz} \mathbb{I}\{z < 0\} \times z + \mathbb{I}\{z < 0\} \times 1 = \mathbb{I}\{z < 0\}$ .<sup>1</sup>
- ▶  $\frac{\partial}{\partial \theta_i} [y(x_t^{\top} \theta)] = y x_{t,i}$  (gradient of Perceptron's output)

---

<sup>1</sup>Note that for  $z = 0$ , there is no derivative but the subdifferential  $[0, 1]$ .

# Derivative of the perceptron cost function

The total loss over the data is defined as  $L(D, \theta) = \sum_{(x,y) \in D} \ell(x, y, \theta)$ .  
Taking the derivative, we have

$$\nabla_{\theta} L(D, \theta) = \nabla_{\theta} \sum_{(x,y) \in D} \ell(x, y, \theta) = \sum_{(x,y) \in D} \nabla_{\theta} \ell(x, y, \theta)$$

## Reminder: The chain rule

Let  $z = g(y)$ ,  $y = f(x)$  so that  $z = g(f(x))$ . Then  $\frac{dz}{dx} = \frac{dz}{dy} \frac{dy}{dx}$

## Applying the chain rule to calculate the gradient

- ▶ Set  $z = \theta^{\top} xy$ .
- ▶  $\nabla_{\theta} \ell(x, y, \theta) = -\frac{d}{dz}(\mathbb{I}\{z < 0\} z) \nabla_{\theta} z$ .
- ▶  $\frac{d}{dz}(\mathbb{I}\{z < 0\} z) = \frac{d}{dz} \mathbb{I}\{z < 0\} \times z + \mathbb{I}\{z < 0\} \times 1 = \mathbb{I}\{z < 0\}$ .<sup>1</sup>
- ▶  $\frac{\partial}{\partial \theta_i} [y(x_t^{\top} \theta)] = yx_{t,i}$  (gradient of Perceptron's output)
- ▶ Gradient update:  
 $\theta^{t+1} = \theta^t - \nabla_{\theta} \ell(x_t, y_t, \theta_t) = \theta^t + yx_t \mathbb{I}\{yx_t^{\top} \theta^t < 0\}$  with  $\alpha_t = 1$ .

---

<sup>1</sup>Note that for  $z = 0$ , there is no derivative but the subdifferential  $[0, 1]$ .

# Derivative of the perceptron cost function

The total loss over the data is defined as  $L(D, \theta) = \sum_{(x,y) \in D} \ell(x, y, \theta)$ .  
Taking the derivative, we have

$$\nabla_{\theta} L(D, \theta) = \nabla_{\theta} \sum_{(x,y) \in D} \ell(x, y, \theta) = \sum_{(x,y) \in D} \nabla_{\theta} \ell(x, y, \theta)$$

## Reminder: The chain rule

Let  $z = g(y)$ ,  $y = f(x)$  so that  $z = g(f(x))$ . Then  $\frac{dz}{dx} = \frac{dz}{dy} \frac{dy}{dx}$

## Applying the chain rule to calculate the gradient

- ▶ Set  $z = \theta^{\top} xy$ .
- ▶  $\nabla_{\theta} \ell(x, y, \theta) = -\frac{d}{dz}(\mathbb{I}\{z < 0\} z) \nabla_{\theta} z$ .
- ▶  $\frac{d}{dz}(\mathbb{I}\{z < 0\} z) = \frac{d}{dz} \mathbb{I}\{z < 0\} \times z + \mathbb{I}\{z < 0\} \times 1 = \mathbb{I}\{z < 0\}$ .<sup>1</sup>
- ▶  $\frac{\partial}{\partial \theta_i} [y(x_t^{\top} \theta)] = yx_{t,i}$  (gradient of Perceptron's output)
- ▶ Gradient update:  
 $\theta^{t+1} = \theta^t - \nabla_{\theta} \ell(x_t, y_t, \theta_t) = \theta^t + yx_t \mathbb{I}\{yx_t^{\top} \theta^t < 0\}$  with  $\alpha_t = 1$ .

---

<sup>1</sup>Note that for  $z = 0$ , there is no derivative but the subdifferential  $[0, 1]$ .

# Derivative of the perceptron cost function

The total loss over the data is defined as  $L(D, \theta) = \sum_{(x,y) \in D} \ell(x, y, \theta)$ .  
Taking the derivative, we have

$$\nabla_{\theta} L(D, \theta) = \nabla_{\theta} \sum_{(x,y) \in D} \ell(x, y, \theta) = \sum_{(x,y) \in D} \nabla_{\theta} \ell(x, y, \theta)$$

## Reminder: The chain rule


Let  $z = g(y)$ ,  $y = f(x)$  so that  $z = g(f(x))$ . Then  $\frac{dz}{dx} = \frac{dz}{dy} \frac{dy}{dx}$

## Applying the chain rule to calculate the gradient

- ▶ Set  $z = \theta^{\top} xy$ .
- ▶  $\nabla_{\theta} \ell(x, y, \theta) = -\frac{d}{dz}(\mathbb{I}\{z < 0\} z) \nabla_{\theta} z$ .
- ▶  $\frac{d}{dz}(\mathbb{I}\{z < 0\} z) = \frac{d}{dz} \mathbb{I}\{z < 0\} \times z + \mathbb{I}\{z < 0\} \times 1 = \mathbb{I}\{z < 0\}$ .<sup>1</sup>
- ▶  $\frac{\partial}{\partial \theta_i} [y(x_t^{\top} \theta)] = yx_{t,i}$  (gradient of Perceptron's output)
- ▶ Gradient update:  
 $\theta^{t+1} = \theta^t - \nabla_{\theta} \ell(x_t, y_t, \theta_t) = \theta^t + yx_t \mathbb{I}\{yx_t^{\top} \theta^t < 0\}$  with  $\alpha_t = 1$ .

The classification error cost function is **not** differentiable :(

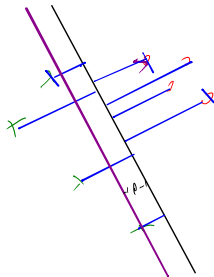
---

<sup>1</sup>Note that for  $z = 0$ , there is no derivative but the subdifferential  $[0, 1]$ . 



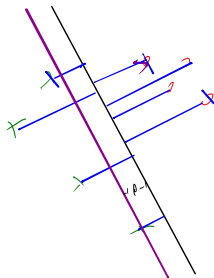
## Margins and confidences

We can think of the output of the network as a measure of confidence



# Margins and confidences

We can think of the output of the network as a measure of confidence



By applying the **logit** function, we can bound a real number  $x$  to  $[0, 1]$ :

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

# Logistic regression

Output as a measure of confidence, given the parameter  $\theta$

$$P_{\theta}(y = 1|x) = \frac{1}{1 + \exp(-x_t^{\top} \theta)}$$

The original output  $x_t^{\top} \theta$  is now passed through the logit function.

# Logistic regression

Output as a measure of confidence, given the parameter  $\theta$

$$P_{\theta}(y = 1|x) = \frac{1}{1 + \exp(-x_t^{\top} \theta)}$$

The original output  $x_t^{\top} \theta$  is now passed through the logit function.

## Negative Log likelihood

$$\ell(x_t, y_t, \theta) = -\ln P_{\theta}(y_t|x_t) = \ln(1 + \exp(-y_t x_t^{\top} \theta))$$

$$\begin{aligned}\nabla_{\theta} \ell(x_t, y_t, \theta) &= \frac{1}{1 + \exp(-y_t x_t^{\top} \theta)} \nabla_{\theta} [1 + \exp(-y_t x_t^{\top} \theta)] \\ &= \frac{1}{1 + \exp(-y_t x_t^{\top} \theta)} \exp(-y_t x_t^{\top} \theta) [\nabla_{\theta} (-y_t x_t^{\top} \theta)] \\ &= -\frac{1}{1 + \exp(x_t^{\top} \theta)} (x_t)_i^n e\end{aligned}$$

$$\blacktriangleright \mathbb{E}_P(\ell) = \int_X dP(x) \sum_{y \in Y} P(y|x) P_{\theta}(y_t + x_t)$$

## The Perceptron

Introduction

The algorithm

## Gradient methods

Gradients for optimisation

The perceptron as a gradient algorithm

## Lab and Assignment

# The Perceptron and Gradients

`./src/Perceptron/Perceptron_gd.ipynb`

- ▶ Perceptron implementation to fill in
- ▶ Gradient descent implementation
- ▶ Experiment on the learning rate with sklearn