

Machine Learning II

Introduction

Souhaib Ben Taieb

February 8, 2022

University of Mons

Table of contents

About the course

The supervised learning problem

The perceptron learning model

S-INFO-075: Machine Learning II

- Everything in **English** (lectures, labs, communications, etc)
- **Instructor**
 - Prof. Souhaib BEN TAIEB
 - De Vinci Building, second floor, room 2.15
 - Email: souhaib.bentaieb@umons.ac.be
- **Teaching assistant**
 - Sukanya Patra (PhD student)
 - De Vinci Building, ground floor
 - Email: sukanya.patra@umons.ac.be
- **Course Webpage**
 - <https://github.com/bsouhaib/ML2-2022>
 - Lecture notes, project details, etc.
- **Moodle**
 - <https://moodle.umons.ac.be/course/view.php?id=2786>
 - Forum for asking questions, etc.
 - **No email please — use the Moodle forum**

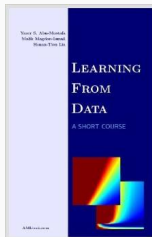
Prerequisites

- Machine learning I (S-INFO-256)
- Probability and Statistics
- Multivariate calculus
- Linear algebra
- Optimization (linear and non-linear)

Key reference

**Yaser S. Abu-Mostafa, Malik Magdon-Ismail, Hsuan-Tien Lin
(2012) Learning from Data. AMLBook.**

<https://work.caltech.edu/telecourse.html>



Assessment

- Exam (E) (*open book*): **60%**
- Project (P): **20%**
- Four assignments (A): **20%** (5% each)
- Final mark:
 - If $E \geq 45\%$ and $P \geq 45\%$ and $A \geq 45\%$:
 - Final mark = $E \times 0.6 + P \times 0.2 + A \times 0.2$
 - Otherwise:
 - Final mark = $\min(E, P, A)$

Task	Due Date	Value
Final exam	Official exam period	60%
Project	TBA	20%
Assignments 1–4	TBA	20%

Main topics

- Theory of learning
- Linear models (classification and regression)
- (Stochastic) Gradient descent
- (Deep) Neural networks and backpropagation
- (?) Support Vector Machines
- (?) Recommender systems, text mining, ...



Table of contents

About the course

The supervised learning problem

The perceptron learning model

Supervised learning

Components of learning

Metaphor: Credit approval

Applicant information:

age	23 years
gender	male
annual salary	\$30,000
years in residence	1 year
years in job	1 year
current debt	\$15,000
...	...

Approve credit?

Components of learning

Components of learning

Formalization:

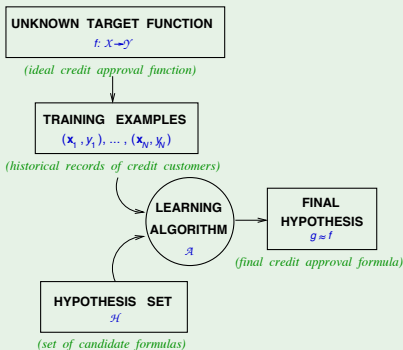
- Input: \mathbf{x} (*customer application*)
- Output: y (*good/bad customer?*)
- Target function: $f : \mathcal{X} \rightarrow \mathcal{Y}$ (*ideal credit approval formula*)
- Data: $(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_N, y_N)$ (*historical records*)



- Hypothesis: $g : \mathcal{X} \rightarrow \mathcal{Y}$ (*formula to be used*)

Note: $\mathbf{x} \in \mathcal{X}$ and $y \in \mathcal{Y}$.

The learning process



The learning algorithm \mathcal{A} picks $g \approx f$ from a hypothesis set \mathcal{H} using the training examples (data).

The learning model

Solution components

The 2 solution components of the learning problem:

- The Hypothesis Set

$$\mathcal{H} = \{h\} \quad g \in \mathcal{H}$$

- The Learning Algorithm

Together, they are referred to as the *learning model*.

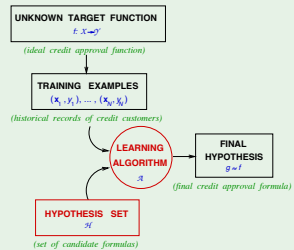


Table of contents

About the course

The supervised learning problem

The perceptron learning model

The perceptron learning model

Let $\mathcal{X} = \mathbb{R}^d$ be the input space, and let $\mathcal{Y} = \{+1, -1\}$ be the output space, denoting a binary (yes/no) decision.

To specify the perceptron **learning model**, we need to define:

1. The perceptron **hypothesis set**
2. The perceptron **learning algorithm**

For $s \in \mathbb{R}$, we define the sign function as

$$\text{sign}(s) = \begin{cases} -1 & \text{if } s < 0, \\ 1 & \text{if } s > 0. \end{cases}$$

Note: for the moment, $\text{sign}(0)$ is ignored (technicality).

The perceptron hypothesis set

A simple hypothesis set - the 'perceptron'

For input $\mathbf{x} = (x_1, \dots, x_d)$ 'attributes of a customer'

Approve credit if $\sum_{i=1}^d w_i x_i > \text{threshold},$

Deny credit if $\sum_{i=1}^d w_i x_i < \text{threshold}.$

This linear formula $h \in \mathcal{H}$ can be written as

$$h(\mathbf{x}) = \text{sign} \left(\left(\sum_{i=1}^d w_i x_i \right) - \text{threshold} \right)$$

The perceptron hypothesis set

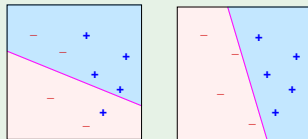
$$h(\mathbf{x}) = \text{sign} \left(\left(\sum_{i=1}^d \mathbf{w}_i x_i \right) + \mathbf{w}_0 \right)$$

Introduce an artificial coordinate $x_0 = 1$:

$$h(\mathbf{x}) = \text{sign} \left(\sum_{i=0}^d \mathbf{w}_i x_i \right)$$

In vector form, the perceptron implements

$$h(\mathbf{x}) = \text{sign}(\mathbf{w}^T \mathbf{x})$$



'linearly separable' data

Note

The L_2 norm of $\mathbf{z} = (z_0, z_1, \dots, z_d)^T$ is given by

$$\|\mathbf{z}\|_2 = \sqrt{\sum_{i=0}^d z_i^2}.$$

The dot product of two vectors \mathbf{w} and \mathbf{x} is defined by

$$\mathbf{w} \cdot \mathbf{x} = \mathbf{w}^T \mathbf{x} = \|\mathbf{w}\|_2 \|\mathbf{x}\|_2 \cos(\theta),$$

where $\|\cdot\|_2$ is the L_2 norm, θ is the angle between \mathbf{w} and \mathbf{x} .

- If the angle between \mathbf{w} and \mathbf{x} is less than 90 degrees, the dot product will be positive, as $\cos(\theta)$ will be positive.
- If the angle between \mathbf{w} and \mathbf{x} is greater than 90 degrees, the dot product will be negative, as $\cos(\theta)$ will be negative.
- If \mathbf{w} and \mathbf{x} are perpendicular (at 90 degrees to each other), the result of the dot product will be zero, because $\cos(\theta)$ will be zero.

The perceptron learning algorithm

A simple learning algorithm - PLA

The perceptron implements

$$h(\mathbf{x}) = \text{sign}(\mathbf{w}^T \mathbf{x})$$

Given the training set:

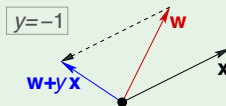
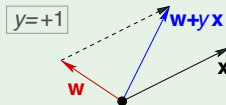
$$(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_N, y_N)$$

pick a **misclassified** point:

$$\text{sign}(\mathbf{w}^T \mathbf{x}_n) \neq y_n$$

and update the weight vector:

$$\mathbf{w} \leftarrow \mathbf{w} + y_n \mathbf{x}_n$$



The perceptron learning algorithm

Data: $\{(\mathbf{x}_n, y_n)\}_{n=1}^N$

initialise weights at $t = 0$ to $\mathbf{w}(0)$;

for $t = 0, 1, 2, \dots$ **do**

 select a misclassified point (\mathbf{x}_n, y_n) ;

 update the weights: $\mathbf{w}(t + 1) = \mathbf{w}(t) + y_n \mathbf{x}_n$;

 iterate to the next step until all points are well classified;

end

Return the final weights $\mathbf{w}(t + 1)$.

PLA is guaranteed to converge if data is **linearly separable** (see labs).

Exercise I

Problem 1.2 Consider the perceptron in two dimensions: $h(\mathbf{x}) = \text{sign}(\mathbf{w}^T \mathbf{x})$ where $\mathbf{w} = [w_0, w_1, w_2]^T$ and $\mathbf{x} = [1, x_1, x_2]^T$. Technically, \mathbf{x} has three coordinates, but we call this perceptron two-dimensional because the first coordinate is fixed at 1.

- (a) Show that the regions on the plane where $h(\mathbf{x}) = +1$ and $h(\mathbf{x}) = -1$ are separated by a line. If we express this line by the equation $x_2 = ax_1 + b$, what are the slope a and intercept b in terms of w_0, w_1, w_2 ?
- (b) Draw a picture for the cases $\mathbf{w} = [1, 2, 3]^T$ and $\mathbf{w} = -[1, 2, 3]^T$.

In more than two dimensions, the $+1$ and -1 regions are separated by a *hyperplane*, the generalization of a line.

(Source: Abu-Mostafa et al. Learning from data. AMLbook)

Solution to Exercise I

- (a) We have $\mathbf{w}^T \mathbf{x} > 0$ if $h(\mathbf{x}) = +1$, and $\mathbf{w}^T \mathbf{x} < 0$ if $h(\mathbf{x}) = -1$.
These two regions are separated by the line $\mathbf{w}^T \mathbf{x} = 0$. This can also be written as $w_0 + w_1 x_1 + w_2 x_2 = 0$. In other words, if $w_2 \neq 0$,
 $a = -\frac{w_1}{w_2}$ and $b = -\frac{w_0}{w_2}$.
- (b) See board.

Exercise II

Consider the following dataset

$\mathbf{x}_1 = (3, 1)$, $\mathbf{x}_2 = (1, -3)$, $\mathbf{x}_3 = (-1, 3)$, $\mathbf{x}_4 = (2.5, -1)$ and
 $y_1 = 1, y_2 = -1, y_3 = 1, y_4 = 1$.

1. Plot the data set in $[-1, 3] \times [-3, 3]$
2. Is the data linearly separable?
3. Run the perceptron algorithm with $\mathbf{w}(0) = (-3, 1, 1)^T$.

Exercise III

The weight update rule, at time step $t + 1$,

$$\mathbf{w}(t + 1) \leftarrow \mathbf{w}(t) + y(t)\mathbf{x}(t),$$

has the nice interpretation that it moves in the direction of classifying $\mathbf{x}(t)$ correctly.

1. Show that $y(t)\mathbf{w}^T(t)\mathbf{x}(t) < 0$.
[Hint: $\mathbf{x}(t)$ is misclassified by $\mathbf{w}(t)$]
2. Show that $y(t)\mathbf{w}^T(t + 1)\mathbf{x}(t) > y(t)\mathbf{w}^T(t)\mathbf{x}(t)$.
[Hint: Use the update rule].

Solution to Exercise III

1. Let $s(t) = \mathbf{w}^T(t)\mathbf{x}(t)$. If $\mathbf{x}(t)$ is misclassified by $\mathbf{w}(t)$, then we have $\text{sign}(s(t)) = +1$ and $y(t) = -1$, or $\text{sign}(s(t)) = -1$ and $y(t) = +1$. In other words, we have

$$y(t)s(t) < 0 \equiv y(t)\mathbf{w}^T(t)\mathbf{x}(t) < 0.$$

2. We have

$$y(t)\mathbf{w}^T(t+1)\mathbf{x}(t) = y(t)[\mathbf{w}^T(t) + y(t)\mathbf{x}^T(t)]\mathbf{x}(t) \quad (1)$$

$$= y(t)\mathbf{w}^T(t)\mathbf{x}(t) + [y(t)]^2\mathbf{x}^T(t)\mathbf{x}(t) \quad (2)$$

$$> y(t)\mathbf{w}^T(t)\mathbf{x}(t), \quad (3)$$

since $[y(t)]^2 > 0$ and $\mathbf{x}^T(t)\mathbf{x}(t) > 0$ (since $x_0(t) = 1$).