



MONASH University

ETC3555

Statistical Machine Learning

The learning problem

25 July 2018

Outline

1 (Statistical) Machine learning

2 The supervised learning problem

3 The perceptron

(Statistical) Machine learning

- The essence of machine learning
 - A pattern exists
 - We cannot pin it down mathematically
 - We have data on it
- Learning examples
 - Spam Detection
 - Product Recommendation
 - Credit Card Fraud Detection
 - Medical Diagnosis
- Other views of “learning from data”:
Statistics, Data Mining, Data Science,
etc.

Exercise I

Which of the following problems are best suited for Machine Learning?

- 1** Classifying numbers into primes and non-primes.
 - 2** Detecting potential fraud in credit card charges.
 - 3** Determining the time it would take a falling object to hit the ground.
 - 4** Determining the optimal cycle for traffic lights in a busy intersection.
- a)** (2) and (4)
 - b)** (1) and (2)
 - c)** (1), (2), and (3)
 - d)** (3)

Different learning problems

- Supervised learning
 - (input, correct output)
- Unsupervised learning
 - (input)
- Reinforcement learning
 - (input, some output, grade for this output)

Exercise II

For each of the following tasks, identify which type of learning is involved (supervised, unsupervised or reinforced) and the training data to be used. If a task can fit more than one type, explain how and describe the training data for each type.

- Recommending a book to a user in an online bookstore
- Playing tic-tac-toe
- Categorizing movies into different types
- Learning to play music

Outline

- 1 (Statistical) Machine learning
- 2 The supervised learning problem**
- 3 The perceptron

Components of learning

Components of learning

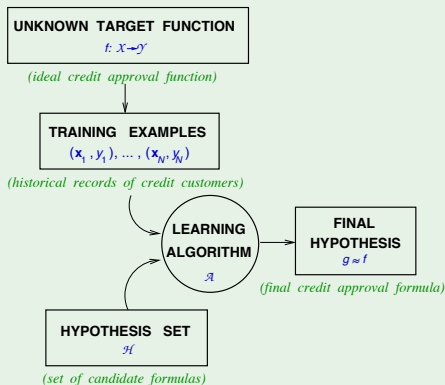
Formalization:

- 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*)



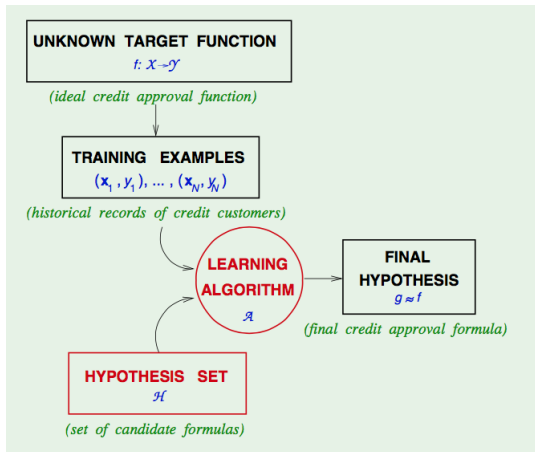
- Hypothesis: $g : \mathcal{X} \rightarrow \mathcal{Y}$ (*formula to be used*)

The learning process



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

The learning model



“The hypothesis set and learning algorithm are referred to informally as the *learning model*.”

Outline

- 1 (Statistical) Machine learning
- 2 The supervised learning problem
- 3 The perceptron**

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

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

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

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

The perceptron hypothesis set

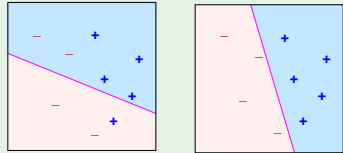
$$h(\mathbf{x}) = \text{sign} \left(\left(\sum_{i=1}^d \mathbf{w}_i x_i \right) + \mathbf{w}_0 \right)$$

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 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 $\cos(\theta)$ is the angle between \mathbf{w} and \mathbf{x} .

- If the angle between \mathbf{w} and \mathbf{x} are 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} are 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:

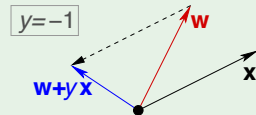
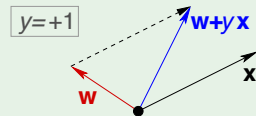
$$(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_N, y_N)$$

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

Iterations of PLA

- One iteration of the PLA:

$$\mathbf{w} \leftarrow \mathbf{w} + y\mathbf{x}$$

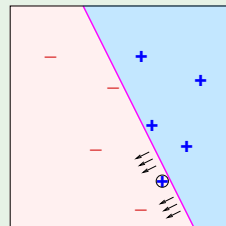
where (\mathbf{x}, y) is a misclassified training point.

- At iteration $t = 1, 2, 3, \dots$, pick a misclassified point from

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

and run a PLA iteration on it.

- That's it!



→ Guarantee to converge if linearly separable data

Exercise III

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.

Solution

$h(\mathbf{x}) = +1 \implies \mathbf{w}^T \mathbf{x} > 0$, and $h(\mathbf{x}) = -1 \implies \mathbf{w}^T \mathbf{x} < 0$. 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, $a = -\frac{w_1}{w_2}$ and $b = -\frac{w_0}{w_2}$.

Exercise IV

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 V

The weight update rule $w(t+1) \leftarrow w(t) + y(t)x(t)$ has the nice interpretation that it moves in the direction of classifying $x(t)$ correctly.

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

- 1 Let $\hat{y}(t) = w^T(t)x(t)$. If $x(t)$ is misclassified by $w(t)$, then we have $\text{sign}(\hat{y}(t)) = 1$ and $y(t) = -1$, or $\text{sign}(\hat{y}(t)) = -1$ and $y(t) = 1$. This is equivalent to $y(t)\hat{y}(t) < 0$ or equivalently $y(t)w^T(t)x(t) < 0$.
- 2 $y(t)w^T(t+1)x(t) = y(t)w^T(t)x(t) + [y(t)]^2x^T(t)x(t) > y(t)w^T(t)x(t)$ since $x^T(t)x(t) > 0$. Note that the first coordinate of $x(t)$ is 1.