

[FIN414 Machine Learning Algorithms]

Supervised Learning

Regression

Jae Yun JUN KIM and **Yves RAKOTONDRATSIMBA**

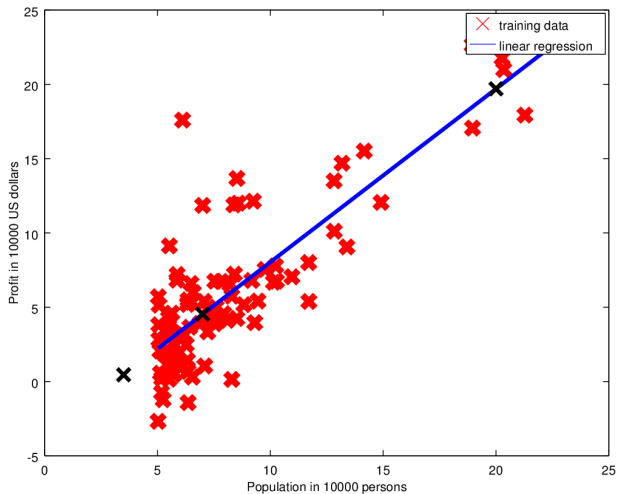
ECE Paris, Graduate School of Engineering
37, quai de Grenelle 75015 Paris, France.

September 21, 2016

Supervised Learning

Linear regression

Motivation



Notation

N = number of training examples

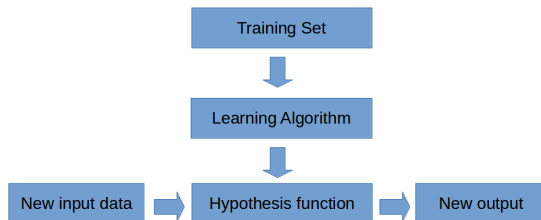
x = input variables / features

y = output variable(s) / “target” variable(s)

(x, y) - training example

For example, i^{th} example: $(x^{(i)}, y^{(i)})$.

Learning Scheme



Hypothesis function:

$$h(x) = h_{\beta}(x) = \beta_0 + \beta_1 x_1 + \cdots + \beta_p x_p$$

For conciseness, define $x_0 = 1$. Then

$$h_{\beta}(x) = \sum_{j=0}^p \beta_j x_j = \beta^T x$$

where p is the number of features, and β are the **model parameters**.

Goal of regression problems

The goal of regression problems is to find the parameter values (i.e., β) that minimize the error between y and $h_{\beta}(x)$ for all the training examples. One possible definition of error can be the following

$$J(\beta) = \frac{1}{2} \sum_{i=1}^N \left(h_{\beta} \left(x^{(i)} \right) - y^{(i)} \right)^2.$$

Hence, a *regression problem* can be mathematically defined as

$$\min_{\beta} J(\beta) = \min_{\beta} \frac{1}{2} \sum_{i=1}^N \left(h_{\beta} \left(x^{(i)} \right) - y^{(i)} \right)^2.$$

But, how can this problem be solved?

Method 1: [Batch gradient descent](#)

Method 2: [Stochastic gradient descent](#)

Method 3: [Closed-form solution](#)

Method 1: Batch Gradient Descent

The idea is

1. Initialize β with some values (say $\beta = \vec{0}$).
2. Keep changing β to reduce $J(\beta)$ with the following update rule:

$$\beta_j := \beta_j - \alpha \frac{\partial}{\partial \beta_j} J(\beta)$$

But, what is $\frac{\partial}{\partial \beta_j} J(\beta)$?

For i fixed (i.e., for a given training example),

$$\begin{aligned} \frac{\partial}{\partial \beta_j} J(\beta) &= \frac{\partial}{\partial \beta_j} \frac{1}{2} (h_{\beta}(x) - y)^2 \\ &= (h_{\beta}(x) - y) \frac{\partial}{\partial \beta_j} (h_{\beta}(x) - y) \\ &= (h_{\beta}(x) - y) \frac{\partial}{\partial \beta_j} (\beta_0 x_0 + \cdots + \beta_j x_j + \beta_p x_p - y) \\ &= (h_{\beta}(x) - y) x_j \end{aligned}$$

Method 1: Batch Gradient Descent

Hence, the update rule with a single training example is

$$\beta_j := \beta_j - \alpha \left(h_{\beta}(x^{(i)}) - y^{(i)} \right) x_j^{(i)}$$

But, we need to do this taking into account for all the training examples,

$$\beta_j := \beta_j - \alpha \sum_{i=1}^N \left(h_{\beta}(x^{(i)}) - y^{(i)} \right) x_j^{(i)}$$

This update rule is known as the **Batch Gradient Descent**.

Method 2: Stochastic Gradient Descent

But, when $N \gg 1$, the *Batch Gradient Descent* method may be very inefficient.

Alternatively, one might be able to use the **Stochastic Gradient Descent**, which can be formulated as follows:

```

For  $i = 1$  to  $N$  {
  For  $j = 1$  to  $p$  {
     $\beta_j := \beta_j - \alpha (h_{\beta}(x^{(i)}) - y^{(i)}) x_j^{(i)}$ 
  }
}

```

This method can allow us to *approximately but more quickly* find the parameters that minimize $J(\beta)$.

Hence, the choice between the two optimization methods depends on the preference between the *accuracy* and the *efficiency*, respectively.

Method 3: Closed-form solution

Let us define the *gradient* of the function $J(\beta)$ as follows

$$\nabla_{\beta} J = \begin{bmatrix} \frac{\partial J}{\partial \beta_0} \\ \vdots \\ \frac{\partial J}{\partial \beta_p} \end{bmatrix} \in \mathbb{R}^{(p+1) \times 1} \quad (1)$$

Hence, the gradient descent method can be vectorially expressed as

$$\beta := \beta - \alpha \nabla_{\beta} J \quad (2)$$

where $\beta \in \mathbb{R}^{(p+1) \times 1}$.

Method 3: Closed-form solution

In the sequel, some mathematical backgrounds are given to be able to find the closed-form solution for the linear regression problem.

Let $f : \mathbb{R}^{N \times (p+1)} \longrightarrow \mathbb{R}$ (i.e., $f(A) \in \mathbb{R}$, $A \in \mathbb{R}^{N \times (p+1)}$).

The *gradient* of f is defined as

$$\nabla_A f = \begin{bmatrix} \frac{\partial f}{\partial A_{11}} & \cdots & \frac{\partial f}{\partial A_{1(p+1)}} \\ \vdots & \ddots & \vdots \\ \frac{\partial f}{\partial A_{m1}} & \cdots & \frac{\partial f}{\partial A_{N(p+1)}} \end{bmatrix}$$

The *trace* of $A \in \mathbb{R}^{N \times (p+1)}$ is defined as

$$\text{tr } A = \sum_{i=1}^p A_{ii}$$

Method 3: Closed-form solution

Some facts with respect to the *gradient of f* and the *trace of A* are given as follows (which will be used to obtain the closed-form solution)

$$\text{tr } AB = \text{tr } BA$$

$$\text{tr } ABC = \text{tr } CAB = \text{tr } BCA$$

$$\text{If } f(A) = \text{tr } AB, \text{ then } \nabla_A \text{tr } AB = B^T$$

$$\text{If } a \in \mathbb{R}, \text{ then } \text{tr } a = a$$

$$\nabla_A \text{tr } ABA^T C = CAB + C^T AB^T$$

Method 3: Closed-form solution

Let

$$X = \begin{bmatrix} (x^{(1)})^T \\ \vdots \\ (x^{(N)})^T \end{bmatrix}, \quad y = \begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(N)} \end{bmatrix}$$

Then,

$$X\beta = \begin{bmatrix} (x^{(1)})^T \beta \\ \vdots \\ (x^{(N)})^T \beta \end{bmatrix} = \begin{bmatrix} h_{\beta}(x^{(1)}) \\ \vdots \\ h_{\beta}(x^{(N)}) \end{bmatrix}$$

and,

$$X\beta - y = \begin{bmatrix} h_{\beta}(x^{(1)}) - y^{(1)} \\ \vdots \\ h_{\beta}(x^{(N)}) - y^{(N)} \end{bmatrix}$$

Method 3: Closed-form solution

By recalling that $z^T z = \sum_i z_i^2$,

$$(X\beta - y)^T (X\beta - y) = \sum_{i=1}^N \left(h(x^{(i)}) - y \right)^2 \triangleq 2J(\beta)$$

Because the *regression problem* consists of finding the parameters (β) that minimize $J(\beta)$, let us impose that

$$\nabla_{\beta} J(\beta) = \vec{0},$$

and find an expression for β that satisfies this constraint.

Method 3: Closed-form solution

Let us now develop the expression of $\nabla_{\beta} J(\beta)$.

$$\begin{aligned}
 \nabla_{\beta} J(\beta) &= \nabla_{\beta} \frac{1}{2} (X\beta - y)^T (X\beta - y) \\
 &= \frac{1}{2} \nabla_{\beta} \left(\beta^T X^T X \beta - \beta^T X^T y - y^T X \beta + y^T y \right) \\
 &= \frac{1}{2} \nabla_{\beta} \text{tr} \left(\beta^T X^T X \beta - \beta^T X^T y - y^T X \beta + y^T y \right) \\
 &= \frac{1}{2} \left[\nabla_{\beta} \text{tr} \left(\beta \beta^T X^T X \right) - \nabla_{\beta} \text{tr} \left(y \beta^T X^T \right) - \nabla_{\beta} \text{tr} \left(y^T X \beta \right) \right]
 \end{aligned}$$

Using the mathematical facts given on a previous slide, one can show that

$$\nabla_{\beta} \text{tr} (\beta \beta^T X^T X) = X^T X \beta + X^T X \beta$$

$$\nabla_{\beta} \text{tr} (y \beta^T X^T) = X^T y$$

$$\nabla_{\beta} \text{tr} (y^T X \beta) = X^T y$$

Method 3: Closed-form solution

Therefore,

$$\nabla_{\beta} J(\beta) = \frac{1}{2} (X^T X \beta + X^T X \beta - X^T y - X^T y) = X^T X \beta - X^T y$$

Now imposing the condition

$$\nabla_{\beta} J(\beta) = \vec{0}$$

we obtain

$$X^T X \beta = X^T y$$

which are known as the **normal equations**.

Finally, the expression of the *parameters* β that optimize $J(\beta)$ is

$$\beta = (X^T X)^{-1} X^T y$$

Further readings

Lecture Notes (available on campus.ece.fr)

Exercises (available on campus.ece.fr)