



ITESO, Universidad
Jesuita de Guadalajara

Maestría de Ciencia de Datos

Optimización Convexa

Tarea 3: Regresión lineal múltiple

Estudiante: Daniel Nuño

Profesor: Dr. Juan Diego Sanchez Torres

Fecha entrega: 21 de febrero de 2022

Introduction

In most of the real-life problems related to data analysis, when the objective is to explain or predict a given variable, several other explanatory variables are often used. It is pretty common to propose input-output models where the explanatory variables are the inputs and the variable to explain or predict is the output. For this modeling process, the first hypothesis is the output variable admits a representation as a linear combination of the input variables. Since in real-world data, the output does not belong to the input-set span, it is necessary to perform a sort of approximation method. As a solution, the approximate output or model's prediction is defined as an output's projection on the set spanned by the inputs. This straightforward reasoning leads to multiple linear regression, a case of linear regression with multiple input variables.

Problema 1

First, let start with a calculus refresher.

Problem 1: Warm Up

Solve the following calculus exercises:

1. For $f_1(x_1, x_2) = x_1^2 x_2 + x_1 x_2^3 \in \mathbb{R}$, calculate the gradient $\frac{df_1}{dx} = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \frac{\partial f_1}{\partial x_2} \end{bmatrix}$. Similarly, for $f_2(x_1, x_2) = x_1^2 + 2x_2$, where $x_1 = \sin(t)$ and $x_2 = \cos(t)$, calculate the gradient $\frac{df_2}{dt} = \begin{bmatrix} \frac{\partial f_2}{\partial x_1} & \frac{\partial f_2}{\partial x_2} \end{bmatrix} \begin{bmatrix} \frac{\partial x_1(t)}{\partial t} \\ \frac{\partial x_2(t)}{\partial t} \end{bmatrix}$.
2. For $f(x) = Ax$, $f(x) \in \mathbb{R}^M$, $A \in \mathbb{R}^{M \times N}$, and $x \in \mathbb{R}^N$, calculate the gradient $\frac{df}{dx}$.

3. Consider the function $h: \mathbb{R} \rightarrow \mathbb{R}, h(t) = (f \circ g)(t)$ with

$$\begin{aligned} f: \mathbb{R}^2 &\rightarrow \mathbb{R} \\ g: \mathbb{R} &\rightarrow \mathbb{R}^2 \\ f(x) &= \exp(x_1 x_2^2) \\ x &= \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = g(t) = \begin{bmatrix} t \cos t \\ t \sin t \end{bmatrix} \end{aligned}$$

and compute the gradient of h with respect to t

4. Let us consider the linear model $Y = X\theta$, where $\theta \in \mathbb{R}^D$ is a parameter vector, $X \in \mathbb{R}^{N \times D}$ are input features and $Y \in \mathbb{R}^N$ are the corresponding observations. Also, let the functions $L(e) = \|e\|_2^2$ and $e(\theta) = Y - X\theta$. Calculate critical point of $L(e)$, that is the solution to $\frac{\partial L}{\partial \theta} = 0$.

Daniel Nuño

HW 3.1: Regresión lineal múltiple

Problema 1: Calentamiento. Resuelve las siguientes ejercicios de cálculo.

1 Para $f_1(x_1, x_2) = x_1^2 x_2 + x_1 x_2^3 \in \mathbb{R}$, calcula el gradiente.

$$\frac{\partial f_1}{\partial x_1} = 2x_1 x_2 + x_2^3$$

$$\frac{\partial f_1}{\partial x_2} = x_1^2 + 3x_1 x_2^2$$

$$\nabla f_1(x_1, x_2) = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} \\ \frac{\partial f_1}{\partial x_2} \end{bmatrix}$$

Para $f_2(x_1, x_2) = x_1^2 + 2x_2$, donde $x_1 = \sin(t)$, $x_2 = \cos(t)$, calcula el gradiente.

$$\frac{\partial f_2}{\partial t} = \begin{bmatrix} \frac{\partial f_2}{\partial x_1} & \frac{\partial f_2}{\partial x_2} \end{bmatrix} \begin{bmatrix} \frac{\partial x_1(t)}{\partial t} \\ \frac{\partial x_2(t)}{\partial t} \end{bmatrix}$$

$$\frac{\partial f_2}{\partial x_1} = 2x_1 \quad \frac{\partial x_1}{\partial t} = \cos(t)$$

$$\frac{\partial f_2}{\partial x_2} = 2 \quad \frac{\partial x_2}{\partial t} = -\sin(t)$$

$$\nabla \frac{\partial f_2}{\partial t} = [2x_1 \ 2] \begin{bmatrix} \cos(t) \\ \sin(t) \end{bmatrix}$$

2 Para $F(X) = AX$, $F(X) \in \mathbb{R}^M$, $A \in \mathbb{R}^{M \times N}$, and $X \in \mathbb{R}^N$, calcula el gradiente.

$$A = \begin{bmatrix} A_{11} & A_{12} & \dots & A_{1N} \\ A_{21} & 0 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ A_{M1} & 0 & \dots & A_{MN} \end{bmatrix}_{M \times N}$$

$$X = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_N \end{bmatrix}$$

$$N \times 1 = M \times 1$$

$$AX = \begin{bmatrix} A_{11}x_1 + A_{12}x_2 + \dots + A_{1N}x_N \\ \vdots \\ A_{M1}x_1 + A_{M2}x_2 + \dots + A_{MN}x_N \end{bmatrix} = \begin{bmatrix} F_1(x_1, x_2, \dots, x_N) \\ \vdots \\ F_M(x_1, x_2, \dots, x_N) \end{bmatrix}$$

$$\nabla \frac{\partial f}{\partial x} = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \frac{\partial f_1}{\partial x_2} & \dots & \frac{\partial f_1}{\partial x_n} \\ \frac{\partial f_2}{\partial x_1} & \frac{\partial f_2}{\partial x_2} & \dots & \frac{\partial f_2}{\partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1} & \frac{\partial f_m}{\partial x_2} & \dots & \frac{\partial f_m}{\partial x_n} \end{bmatrix}$$

3. Considera la función $h: \mathbb{R} \rightarrow \mathbb{R}$, $h(t) = (f \circ g)(t)$ con
 $f: \mathbb{R}^2 \rightarrow \mathbb{R}$, $g: \mathbb{R} \rightarrow \mathbb{R}^2$

$$f(x) = \exp(x_1 x_2^2)$$

$$x_t = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = g(t) = \begin{bmatrix} t \cos t \\ t \sin t \end{bmatrix}$$

Calcula el gradiente de h con respecto

$$\frac{\partial f}{\partial x_1} = x_2^2 e^{x_1 x_2^2}$$

$$\frac{\partial x_1}{\partial t} = \cos(t) - t \sin(t)$$

$$\frac{\partial f}{\partial x_2} = 2x_1 x_2 e^{x_1 x_2^2}$$

$$\frac{\partial x_2}{\partial t} = \sin(t) + t \cos(t)$$

$$\frac{\partial f}{\partial t} = \begin{bmatrix} x_2^2 e^{x_1 x_2^2} & 2x_1 x_2 e^{x_1 x_2^2} \end{bmatrix} \begin{bmatrix} \cos(t) - t \sin(t) \\ \sin(t) + t \cos(t) \end{bmatrix}$$

$$x_1 = t \cos(t)$$

$$x_2 = t \sin(t)$$

4. Considera el modelo lineal $y = X\theta$, sea $\theta \in \mathbb{R}^D$ es un vector parámetro.
 $X \in \mathbb{R}^{n \times D}$ son entradas y $y \in \mathbb{R}^n$ son las observaciones correspondientes.
 También, sea las funciones

$L(\theta) = \|e\|_2^2$ y $e(\theta) = y - X\theta$. Calcula el punto crítico de $L(\theta)$,
 que es solución a $\frac{\partial L}{\partial \theta} = 0$

$$\frac{\partial L(\theta)}{\partial \theta} = 2X^T(y - X\theta) = 0$$

Propiedad de la norma 2

$$= 2X^T y - 2X^T X \theta = 0$$

$$\|e\|_2^2 = e^T e$$

$$\Rightarrow 2X^T y = 2X^T X \theta$$

$$\frac{d\|e\|_2^2}{d\theta} = 2e$$

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

Problema 2

The following case considers the application of the ML and MAP approaches to the multiple regression.

Problem 2: Multiple Linear Regression

Consider a standard linear regression problem, in which for $i = 1, \dots, n$ the mean of the conditional distribution of y_i is specified given a $k \times 1$ predictor vector \mathbf{x}_i of the form $y_i = \mathbf{x}_i^T \boldsymbol{\theta} + \varepsilon_i$ where $\boldsymbol{\theta}$ is a $k \times 1$ vector, and the ε_i are independent and identically normally distributed random variables $\varepsilon_i \sim N(0, \sigma^2)$.

1. Show that the likelihood function for this problem is

$$L(\mathbf{Y}|\mathbf{X}, \boldsymbol{\theta}, \sigma^2) \propto (\sigma^2)^{-n/2} \exp\left(-\frac{1}{2\sigma^2}(\mathbf{Y} - \mathbf{X}\boldsymbol{\theta})^T(\mathbf{Y} - \mathbf{X}\boldsymbol{\theta})\right)$$

and find the parameter vector $\boldsymbol{\theta}$ that minimizes L .

2. Propose a prior normal distribution for $\boldsymbol{\theta}$, find the posterior distribution for $\boldsymbol{\theta}$.
3. Propose a prior Laplace distribution for $\boldsymbol{\theta}$, find the posterior distribution for $\boldsymbol{\theta}$.
4. (5 pt) Revisit this problem using the least-squares approach, exposing in a very detailed way the conditions in such the probabilistic formulation is equivalent to the Tikhonov (RIDGE) regularization and the LASSO regularization. Finally, explain with an example how the LASSO

regularization can lead to a sparse solution for the regression problem.

Problema 2. Regresión lineal múltiple

Considero una regresión lineal estándar, en cual para $i = 1, 2, \dots, n$ la media de la distribución de probabilidad condicional de $y_i = X_i^T \theta + E_i$ sea θ un vector $k \times 1$, y la E_i son independientes e idénticamente distribuidos. $E_i \sim N(0, \sigma^2)$

1.- Muestra que la ~~probabilidad~~ función de máxima verosimilitud es

$$L(Y|X, \theta, \sigma^2) \propto (\sigma^2)^{-n/2} \exp\left(-\frac{1}{2\sigma^2} (Y - X\theta)^T (Y - X\theta)\right)$$

y encuentre un parámetro que ~~maximice~~ maximice L .

$$L(Y|X, \sigma) \text{ puede ser escrito como } \propto \left(\frac{1}{\sigma^2}\right)^{n/2} \exp\left[-\frac{1}{2\sigma^2} \|Y - X\theta\|_2^2\right]$$

$$\ln(L(Y|X, \sigma)) = -\frac{n}{2} \ln(\sigma^2) - \frac{1}{2\sigma^2} \|Y - X\theta\|_2^2$$

$$\frac{\partial \ln(L(Y|X, \sigma))}{\partial \theta} = \frac{1}{2\sigma^2} \cdot 2X^T (Y - X\theta) = \theta$$

$$\begin{aligned} X^T Y &= X^T X \theta \\ \theta &= (X^T X)^{-1} X^T Y \end{aligned}$$

$$\frac{\partial^2 \ln(L)}{\partial \theta^2} = \frac{1}{\sigma^2} (X^T Y - X^T X \theta) = -\frac{X^T X}{\sigma^2}$$

2. Proponga una distribución normal priori para θ , encuentre la distribución posteriori para θ .

una distribución normal $\theta \sim N(0, \sigma^2)$ donde $e = Y - X\theta$ y

$\theta \sim N(0, \frac{1}{\lambda})$ obtenemos

$$P(e) = \frac{1}{(2\pi)^{n/2} \sigma^n} \exp\left[-\frac{1}{2\sigma^2} \|e\|_2^2\right]$$

$$P(\theta) = \frac{\lambda^n}{(2\pi)^{n/2}} \exp\left[-\frac{\lambda}{2} \|\theta - 0\|_2^2\right]$$

$$P(\theta|D) \propto \frac{1}{(2\pi)^{n/2} \sigma^n} \exp\left[-\frac{1}{2\sigma^2} \|e\|_2^2\right] \cdot \frac{\lambda^n}{(2\pi)^{n/2}} \exp\left[-\frac{\lambda}{2} \|\theta - 0\|_2^2\right]$$

3. proponga una distribución de Laplace priori para θ , encuentre la distribución posteriori para θ .