

Tarea 2 - Optimización

Erick Salvador Alvarez Valencia

CIMAT A.C.,
erick.alvarez@cimat.mx

1. Problema 1

Escribe la expansión de serie de Taylor de segundo orden de la función:

$$f(x) = 100(x_2 - x_1^2)^2 + (1 - x_1)^2$$

donde $x \in R^2$.

La expansión en series de Taylor de orden dos está definida como:

$$f(x) = f(x_0) + Df(x_0)(x - x_0) + \frac{1}{2}(x - x_0)^T D^2 f(x_0)(x - x_0) + o(\|x - x_0\|^2) \quad (1)$$

Por lo tanto necesitamos calcular el gradiente y la Hessiana de la función.

$$\nabla f(x) = \begin{bmatrix} \frac{\partial f(x)}{\partial x_1} \\ \frac{\partial f(x)}{\partial x_2} \end{bmatrix} = \begin{bmatrix} -400(x_1 x_2 - x_1^3) - 2(1 - x_1) \\ 200(x_2 - x_1^2) \end{bmatrix} \quad (2)$$

$$\nabla^2 f(x) = \begin{pmatrix} \frac{\partial^2 f(x)}{\partial x_1^2} & \frac{\partial^2 f(x)}{\partial x_1 x_2} \\ \frac{\partial^2 f(x)}{\partial x_2 x_1} & \frac{\partial^2 f(x)}{\partial x_2^2} \end{pmatrix} = \begin{pmatrix} -400(x_2 - 3x_1^2) + 2 & -400x_1 \\ -400x_1 & 200 \end{pmatrix} \quad (3)$$

Una vez teniendo lo anterior podemos construir la aproximación en series de Taylor de orden 2, pero como no se indicó en qué punto está centrada, asumiremos que está centrada en el punto $x_0 = (x_a, x_b)^T$.

$$f(x) = 100(x_b - x_a^2)^2 + (1 - x_a)^2 + \begin{bmatrix} -400(x_1 x_2 - x_1^3) - 2(1 - x_1) \\ 200(x_2 - x_1^2) \end{bmatrix} \begin{bmatrix} x_1 - x_a \\ x_2 - x_b \end{bmatrix}^T \quad (4)$$

$$+ \frac{1}{2} \begin{bmatrix} x_1 - x_a \\ x_2 - x_b \end{bmatrix}^T \begin{pmatrix} -400(x_2 - 3x_1^2) + 2 & -400x_1 \\ -400x_1 & 200 \end{pmatrix} \begin{bmatrix} x_1 - x_a \\ x_2 - x_b \end{bmatrix} \quad (5)$$

2. Problema 2

Supón que a través de un experimento el valor de una función g es observado en m puntos, x_1, x_2, \dots, x_m , lo que significa que los valores $g(x_1), g(x_2), \dots, g(x_m)$ son conocidos. Queremos aproximar la función $g(\cdot) : R \rightarrow R$ por un polinomio

$$h(x) = a_0 + a_1x + a_2x^2 + \dots + a_nx^n$$

con $n < m$. Propón un método para aproximar $g(\cdot)$ y calcula los valores de a_0, a_1, \dots, a_n .

Genera las observaciones del modelo

$$g(x) = \frac{\sin(x)}{x} + \eta$$

, $x \in [0, 1, 10]$ donde $\eta \sim N(0, 1)$, $0, 1 < x_1 < x_2 < \dots < x_m = 10$.

Queremos generar una aproximación a la función $g(\cdot)$ mediante un polinomio de grado n pero tenemos la restricción que hay m observaciones y $n < m$ lo cual podría generar un sistema sobredeterminado que al final podría no tener solución, para ello usamos el enfoque de mínimos cuadrados. Para ello tenemos una matriz A definida de la siguiente manera:

$$A = \begin{bmatrix} x_1^2 & x_1 & 1 \\ x_2^2 & x_2 & 1 \\ \vdots & \vdots & \vdots \\ x_n^2 & x_n & 1 \end{bmatrix} \quad (6)$$

Y queremos encontrar los coeficientes a, b, c por lo cual tenemos el vector $x = (a, b, c)^T$ y el vector $b = (y_1, y_2, \dots, y_n)^T$. Ahora lo que se quiere es encontrar el error generado por $Ax = b$ y minimizarlo de la siguiente manera:

$$E(a, b, c) = \sum_i ||E_i||^2 = ||E||_2^2 = ||Ax - b||_2^2$$

Para ello podemos encontrar el gradiente del error, pero antes hay que tener en cuenta que:

$$\nabla(c^T x) = c$$

$$\nabla(x^T Ax) = (A + A^T)x$$

Para $c, x \in R^n$ y $A \in R^{n \times n}$.

De esta forma el gradiente del error es:

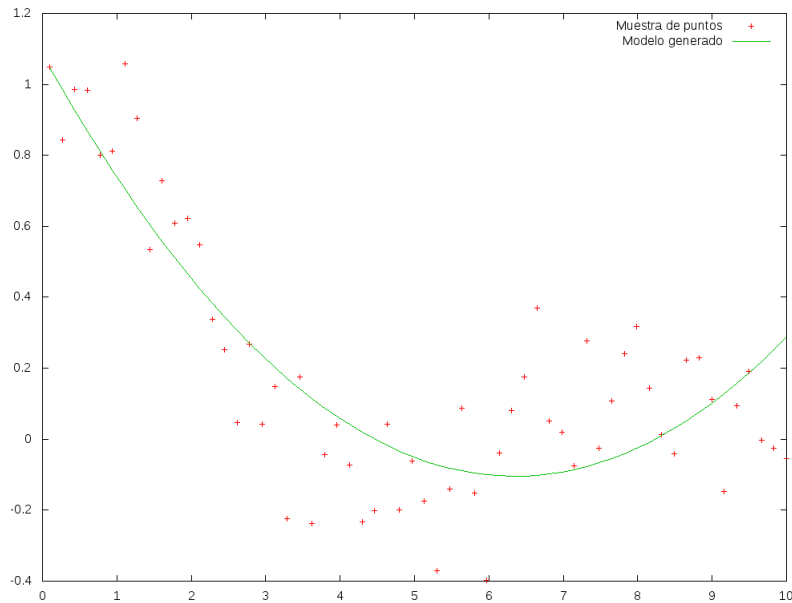
$$\begin{aligned} E(a, b, c) &= ||Ax - b||_2^2 \\ E(a, b, c) &= (Ax - b)^T (Ax - b) \\ E(a, b, c) &= X^T A^T Ax - x^T A^T y - y^T Ax + y^T y \\ E(a, b, c) &= X^T A^T Ax - 2y^T Ax + y^T y \\ \nabla E &= (A^T A + AA^T)x - 2A^T y \end{aligned} \quad (7)$$

Igualamos a cero para encontrar el mínimo.

$$\begin{aligned}(A^T A + AA^T)x - 2A^T y &= 0 \\ 2A^T Ax &= 2A^T y \\ A^T Ax &= A^T y\end{aligned}\tag{8}$$

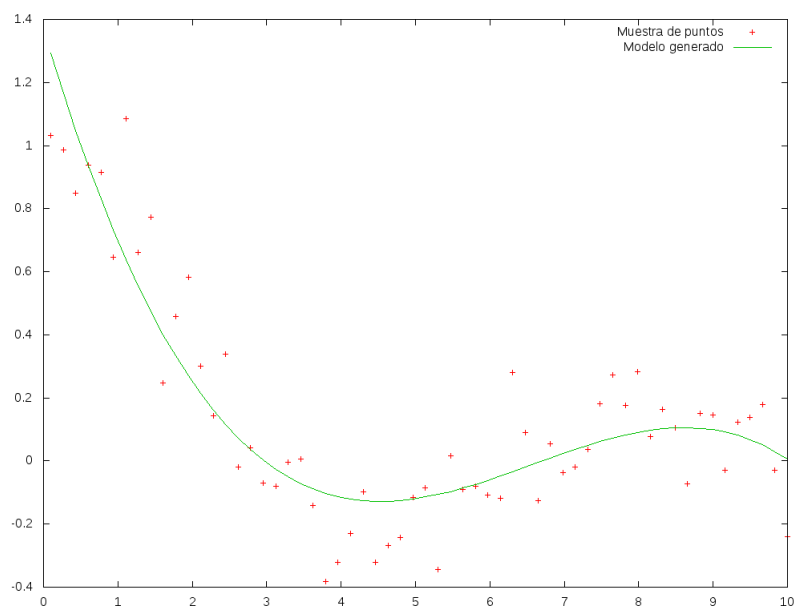
De esta el método de mínimos cuadrados nos dará una aproximación a la función con error mínimo y ahora trabajaremos con un sistema consistente. Aunque para que este sistema funcione hay que asumir que la matriz A es pseudoinversa. Para lo anterior se utilizó un modelo de polinomio cuadrático, pero esto se puede generalizar a polinomios de grados distintos.

Para este ejercicio se realizó un programa que ejecutaba el algoritmo anterior usando 4 modelos de polinomios distintos, desde uno de grado 2 a uno de grado 5. De igual manera se generaron las muestras con la función indicada en la descripción del ejercicio y para el ruido aleatorio se utilizó el Teorema del Límite Central para obtener números aleatorios de una distribución normal ya que C por defecto no tiene funciones que realizan esto anterior. A continuación se muestran los resultados obtenidos por el programa para un conjunto de 50 puntos generados.

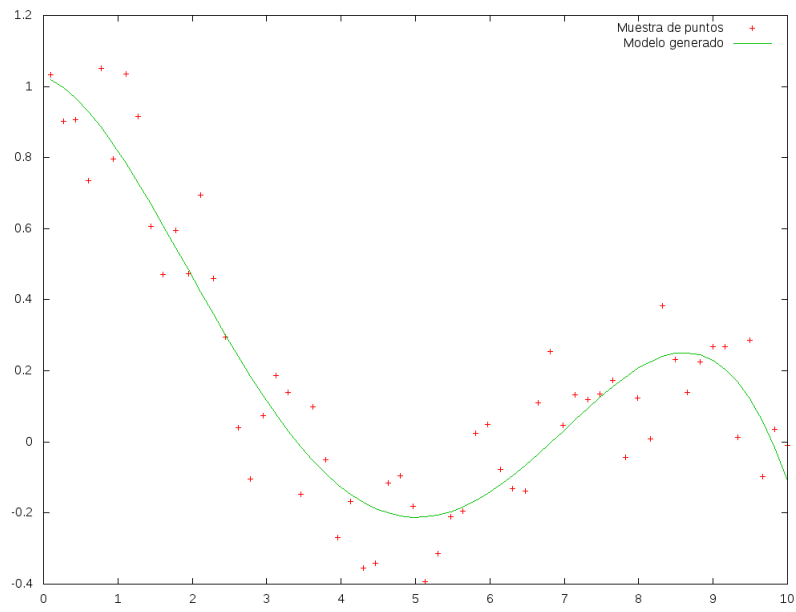


(a) Figura 1. Mínimos cuadrados, generación de un modelo cuadrático.

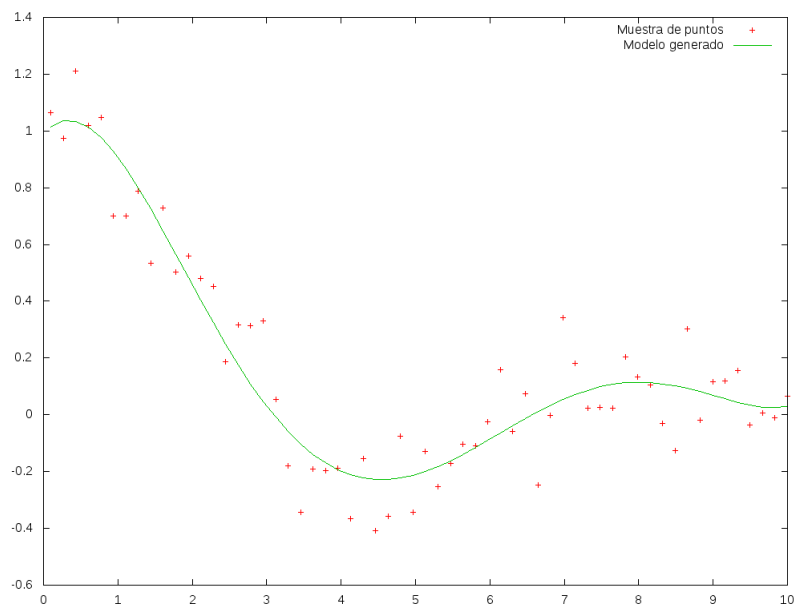
En la Figura 1 podemos apreciar en amarillo el modelo generado por el algoritmo, intencionalmente se pidió que dicho modelo fuera cuadrático, aún así ya empieza a ajustarse a cierto conjunto de puntos.



(b) Figura 1. Mínimos cuadrados, generación de un modelo cúbico.



(c) Figura 1. Mínimos cuadrados, generación de un modelo de grado 4.



(d) Figura 1. Mínimos cuadrados, generación de un modelo de grado 5.

En las últimas tres Figuras podemos apreciar cómo se fue aumentando el grado del polinomio hasta lograr un mejor ajuste hacia el conjunto de puntos. Las gráficas fueron generadas usando GNUPlot.

3. Problema 3

Dada una función continua $f(x)$ en $[a, b]$. Encuentra el polinomio de aproximación de grado n

$$p(x) = a_0 + a_1x + a_2x^2 + \dots + a_nx^n$$

tal que minimiza

$$\int_a^b [f(x) - p(x)]^2 dx$$

Tenemos que nuestro polinomio se define como: $\sum_{k=0}^n a_k x^k$. Ahora, queremos encontrar los coeficientes del anterior polinomio que minimizan la integral propuesta. Así que hacemos la sustitución de $p(x)$ y desarrollamos el binomio.

$$\int_a^b [f(x) - p(x)]^2 dx = \int_a^b f(x)^2 - 2f(x) \sum_{j=0}^n a_j x^j - \left(\sum_{j=0}^n a_j x^j \right)^2 dx \quad (9)$$

Para minimizar lo anterior podemos derivar la integral y para ello usamos la regla de Leibniz que nos permite derivar la función que hay dentro de la integral.

$$\frac{d}{dx} \left(\int_a^b f(x, t) dt \right) = \int_a^b \frac{\partial}{\partial x} f(x, t) dt \quad (10)$$

Aplicamos lo anterior en la integral y separamos.

$$\begin{aligned} \int_a^b \frac{df(x)^2}{dx} - \frac{2f(x) \sum_{j=0}^n a_j x^j}{dx} - \frac{(\sum_{j=0}^n a_j x^j)^2}{dx} \\ - 2 \int_a^b f(x) x^k dx - 2 \int_a^b x^k \sum_{j=0}^n a_j x^j dx \\ - 2 \int_a^b f(x) x^k dx - 2 \sum_{j=0}^n a_j \int_a^b x^j x^k dx \\ - 2 \int_a^b f(x) x^k dx - 2 \sum_{j=0}^n a_j \left(\frac{b^{j+k+1} - a^{j+k+1}}{j+k+1} \right) \end{aligned} \quad (11)$$

Igualamos a cero y despejamos.

$$\sum_{j=0}^n a_j \left(\frac{b^{j+k+1} - a^{j+k+1}}{j+k+1} \right) = \int_a^b f(x) x^k dx \quad (12)$$

Para despejar los coeficientes podemos ver lo anterior como un producto matriz-vector donde la matriz es simétrica y pseudoinvertible.

$$\begin{bmatrix} \frac{b^{1+1+1} - a^{1+1+1}}{2+1+1} & \frac{b^{1+2+1} - a^{1+2+1}}{2+2+1} & \dots & \frac{b^{1+n+1} - a^{1+n+1}}{2+n+1} \\ \frac{b^{2+1+1} - a^{2+1+1}}{2+1+1} & \frac{b^{2+2+1} - a^{2+2+1}}{2+2+1} & \dots & \frac{b^{2+n+1} - a^{2+n+1}}{2+n+1} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{b^{n+1+1} - a^{n+1+1}}{n+1+1} & \frac{b^{n+2+1} - a^{n+2+1}}{n+2+1} & \dots & \frac{b^{n+n+1} - a^{n+n+1}}{n+n+1} \end{bmatrix} \begin{bmatrix} a_0 \\ a_1 \\ \vdots \\ a_n \end{bmatrix} = \begin{bmatrix} \int_a^b f(x) dx \\ \int_a^b f(x) x dx \\ \vdots \\ \int_a^b f(x) x^n dx \end{bmatrix} \quad (13)$$

Finalmente resolvemos el sistema anterior por algún método para encontrar los coeficientes a_k .

4. Problema 4

Para la distribución normal $N(\mu, \sigma^2)$ la cual tiene una función de densidad:

$$f(x|\mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$$

La correspondiente función de densidad para una muestra $\{x_i\}_{i=0}^n$ de n variables aleatorias independientemente distribuidas (la verosimilitud) es:

$$L(\mu, \sigma^2; x) = f(x_1, \dots, x_n | \mu, \sigma^2) = \prod_{i=1}^n f(x_i | \mu, \sigma^2)$$

con $x = x_1, \dots, x_n$ calcula:

$$(\mu^*, \sigma^*) = \arg \max_{\mu, \sigma} L(\mu, \sigma^2; x)$$

Para calcular lo anterior primero encontramos el producto de la función de densidad:

$$L(\mu, \sigma^2; x) = \frac{1}{(\sqrt{2\pi\sigma^2})^n} \exp\left(-\frac{\sum_{i=1}^n (x_i - \mu)^2}{2\sigma^2}\right) \quad (14)$$

Ahora podemos trabajar con la log-verosimilitud ya que el logaritmo conserva la monotonía.

$$l(\mu, \sigma^2; x) = \frac{-\sum_{i=1}^n (x_i - \mu)^2}{2\sigma^2} - \frac{n \log(2\pi\sigma^2)}{2} \quad (15)$$

Encontramos el estimador para la media derivando con respecto a μ e igualando a cero.

$$\begin{aligned} \frac{\partial l}{\partial \mu} &= \frac{2 \sum_{i=1}^n (x_i - \mu)}{2\sigma^2} = 0 \\ \sum_{i=1}^n x_i - n\hat{\mu} &= 0 \\ \hat{\mu} &= \frac{\sum_{i=1}^n x_i}{n} \end{aligned} \quad (16)$$

Ahora se encontrará el estimador para la varianza, por lo cual realizamos el mismo proceso hecho con la media excepto que ahora se derivará con respecto a σ .

$$\begin{aligned} \frac{\partial l}{\partial \sigma} &= \frac{-4\sigma(-\sum_{i=1}^n (x_i - \mu)^2)}{4\sigma^4} - \frac{n(4\pi\sigma)}{4\pi\sigma^2} \\ &\quad \frac{\sum_{i=1}^n (x_i - \mu)^2 - n\hat{\sigma}^2}{\hat{\sigma}^3} = 0 \\ \hat{\sigma}^2 &= \frac{\sum_{i=1}^n (x_i - \mu)^2}{n} \end{aligned} \quad (17)$$

Una vez que se encontraron los dos estimadores se procederá a verificar que son máximos y por lo cual se usa la matriz Hessiana.

$$\frac{\partial^2 l}{\partial \mu^2} = \frac{-n}{\sigma^2} \quad (18)$$

$$\begin{aligned} \frac{\partial^2 l}{\partial \sigma^2} &= \frac{-2n\sigma\sigma^3 - 3\sigma^2[\sum_{i=1}^n (x_i - \mu)^2 - n\sigma^2]}{\sigma^2} \\ \frac{\partial^2 l}{\partial \sigma^2} &= \frac{-2n\sigma^4 - 3\sigma^2 \sum_{i=1}^n (x_i - \mu)^2 + 3n\sigma^4}{\sigma^6} \\ \frac{\partial^2 l}{\partial \sigma^2} &= \frac{-2n}{\sigma^2} - \frac{3 \sum_{i=1}^n (x_i - \mu)^2}{\sigma^4} + \frac{3n}{\sigma^2} \\ \frac{\partial^2 l}{\partial \sigma^2} &= \frac{\sigma^2 n - 3 \sum_{i=1}^n (x_i - \mu)^2}{\sigma^4} \end{aligned} \quad (19)$$

$$\frac{\partial^2 l}{\partial \sigma \mu} = \frac{\partial^2 l}{\partial \mu \sigma} = \frac{-2n \sum_{i=1}^n (x_i - \mu)}{\sigma^4} \quad (20)$$

$$\nabla^2 f = \begin{bmatrix} \frac{-n}{\sigma^2} & \frac{\partial^2 l}{\partial \mu \sigma} = \frac{-2n \sum_{i=1}^n (x_i - \mu)}{\sigma^4} \\ \frac{\partial^2 l}{\partial \mu \sigma} = \frac{-2n \sum_{i=1}^n (x_i - \mu)}{\sigma^4} & \frac{\sigma^2 n - 3 \sum_{i=1}^n (x_i - \mu)^2}{\sigma^4} \end{bmatrix} \quad (21)$$

Una vez que se tiene el Hessiano podemos analizarlo, para lo cual se sustituyen los puntos críticos en él.

Para el primer componente de la matriz: $\frac{-n}{\sigma^2}$ podemos notar que tanto los términos n y σ^2 son positivos y con el signo menos siempre será negativo el término. Ahora para los términos de las derivadas parciales mixtas sustituimos los puntos críticos.

$$\frac{\partial^2 l}{\partial \mu \sigma} = \frac{-2n \sum_{i=1}^n (x_i - \mu)}{\sigma^4} = \frac{-2 \frac{\sum_{i=1}^n (x_i - \mu)^2}{n} \sum_{i=1}^n (x_i - \frac{\sum_{i=1}^n x_i}{n})}{(\frac{\sum_{i=1}^n (x_i - \mu)^2}{n})^2} \quad (22)$$

De lo anterior podemos notar que todo el denominador es una cantidad positiva, y del numerador la única parte que no se aprecia a simple vista si es positiva o negativa es:

$$\sum_{i=1}^n (x_i - \frac{\sum_{i=1}^n x_i}{n}) \quad (23)$$

Hacemos un poco de desarrollo.

$$\sum_{i=1}^n (x_i - \frac{\sum_{i=1}^n x_i}{n}) = \sum_{i=1}^n x_i - \frac{\sum_{i=1}^n \sum_{j=1}^n x_i}{n} = \sum_{i=1}^n x_i - \sum_{i=1}^n x_i = 0 \quad (24)$$

Se convierte en cero por lo cual todo el término se anula, y como las derivadas mixtas son iguales éstas se vuelven cero.

$$\nabla^2 f = \begin{bmatrix} \frac{-n}{\sigma^2} & 0 \\ 0 & \frac{\sigma^2 n - 3 \sum_{i=1}^n (x_i - \mu)^2}{\sigma^4} \end{bmatrix} \quad (25)$$

Únicamente nos falta analizar el último término. Para lo cual hacemos la sustitución con sus puntos críticos.

$$\frac{\sigma^2 n - 3 \sum_{i=1}^n (x_i - \mu)^2}{\sigma^4} = \frac{\frac{\sum_{i=1}^n (x_i - \mu)^2}{n} n - 3 \sum_{i=1}^n (x_i - \frac{\sum_{i=1}^n x_i}{n})^2}{\sigma^4} \quad (26)$$

De (21) podemos apreciar que el denominador es positivo y que las dos primeras n en el numerador se cancelan, quedándonos.

$$\frac{\sum_{i=1}^n (x_i - \mu)^2 - 3 \sum_{i=1}^n (x_i - \mu)^2}{\sigma^4} - \frac{4 \sum_{i=1}^n (x_i - \mu)^2}{\sigma^4} \quad (27)$$

Por lo cual todo el término es completamente negativo.
Finalmente nos queda:

$$\nabla^2 f = \begin{bmatrix} \frac{-n}{\sigma^2} & 0 \\ 0 & \frac{-4 \sum_{i=1}^n (x_i - \mu)^2}{\sigma^4} \end{bmatrix} \quad (28)$$

Como tenemos una matriz diagonal entonces los elementos diagonales de la misma representan a sus eigenvalores y se sabe que un criterio de optimalidad para el máximo es: si los eigenvalores de la matriz son negativos entonces el punto crítico es máximo local, lo cual pasa en nuestro caso.

$$\begin{aligned} \lambda_1 &= \frac{-n}{\sigma^2} < 0 \\ \lambda_2 &= \frac{-4 \sum_{i=1}^n (x_i - \mu)^2}{\sigma^4} < 0 \end{aligned} \quad (29)$$

Por lo tanto los estimadores encontrados son máximo verosímiles.

5. Problema 5

- Sea $f : R \rightarrow R$ es convexa y $a, b \in \text{dom } f$ con $a < b$. Muestra que:

$$a(f(x) - f(b)) + x(f(b) - f(a)) + b(f(a) - f(x)) \geq 0$$

para todo $x \in [a, b]$.

Como x está entre a y b lo podemos definir como: $x = \alpha b + (1 - \alpha)a$. Despejando nos queda: $\alpha = \frac{x-a}{b-a}$. Sustituimos lo anterior en la definición de convexidad y nos queda:

$$f(\alpha b + (1 - \alpha)a) \leq \alpha f(b) + (1 - \alpha)f(a) \quad (30)$$

Para la ecuación anterior se pueden usar los puntos a, b en la definición ya que como f es convexa, lo anterior se debe cumplir con cualquier par de puntos.

$$\begin{aligned} f(x) &\leq \alpha \frac{x-a}{b-a} f(b) + 1 - \frac{x-a}{b-a} f(a) \\ (b-a)f(x) &\leq (x-a)f(b) + (b-a)f(a) - (x-a)f(a) \\ a(f(x) - f(b)) + x(f(b) - f(a)) + b(f(a) - f(x)) &\geq 0 \end{aligned} \quad (31)$$

- Sea $f : R^n \rightarrow R$ una función convexa con $A \in R^{n \times m}$, y $b \in R^n$. Muestra que la función $g : R^m \rightarrow R$ definida por $g(x) = f(Ax + b)$, con $\text{dom } g = \{x | Ax + b \in \text{dom } f\}$, es convexa.

Sabemos por hipótesis que $f(x)$ es convexa, por lo cual se cumple que para algún $\alpha \in (0, 1)$:

$$f(\alpha x + (1 - \alpha)y) \leq \alpha f(x) + (1 - \alpha)f(y) \quad (32)$$

Veamos si g es definida aplicando el mismo enfoque, por lo cual tenemos:

$$\begin{aligned} g(\alpha x + (1 - \alpha)y) &= f(A[\alpha x + (1 - \alpha)y] + b) \\ &= f(\alpha Ax + Ay - \alpha Ay + b) \end{aligned} \quad (33)$$

Ahora a la ecuación anterior sumamos $\alpha b - \alpha b$ y agrupamos.

$$\begin{aligned} &= f(\alpha Ax + Ay - \alpha Ay + b + \alpha b - \alpha b) \\ &= f(\alpha[Ax + b] + [Ay + b] - \alpha[Ay + b]) \\ &= f(\alpha[Ax + b] + (1 - \alpha)[Ay + b]) \\ &\leq \alpha f(Ax + b) + (1 - \alpha)f(Ay + b) \\ &= \alpha g(x) + (1 - \alpha)g(y) \end{aligned} \quad (34)$$

Hemos llegado a la definición de convexidad para funciones por lo cual vemos que g es una función convexa.