

Práctica 1

Regresión lineal

Mario Emilio Jiménez Vizcaíno
A01173359@itesm.mx
Tecnológico de Monterrey
Ingeniería en Tecnologías Computacionales
Monterrey, N.L., México

Jesus Abraham Haros Madrid
A01252642@itesm.mx
Tecnológico de Monterrey
Ingeniería en Tecnologías Computacionales
Monterrey, N.L., México

ABSTRACT

Esta práctica tiene el propósito de demostrar dos metodologías para la creación de modelos de regresión lineal, cuyo objetivo es crear una función a partir de una combinación lineal de sus parámetros para aproximar el resultado real. Estos métodos son la regresión lineal por mínimos cuadrados y el gradiente descendente, ambos implementados dentro de la librería scikit-learn, usada durante la práctica.

1 INTRODUCCIÓN

El aprendizaje automático es la rama de los algoritmos computacionales cuyo objetivo es emular la inteligencia humana a través de "aprender" del entorno. Las técnicas basadas en el aprendizaje automático se han aplicado con éxito en diversos campos como el reconocimiento de patrones, la visión computacional, diferentes ramas de la ingeniería, las finanzas y la biología.[1]

Uno de los modelos más comunes del aprendizaje automático es la regresión, un tipo de análisis predictivo básico y comúnmente utilizado. Se diferencia de los modelos de clasificación porque estima un valor numérico, mientras que los modelos de clasificación identifican a qué categoría pertenece una observación.

Las estimaciones de regresión se utilizan para explicar la relación entre una variable dependiente y una o más variables independientes. La forma más sencilla de este modelo es la ecuación de regresión con una variable dependiente y otra independiente, definida por la fórmula $y = c + b * x$ donde y representa el valor estimado de la variable dependiente, c es una constante, b es el coeficiente de regresión y x es el valor de la variable independiente.

2 CONCEPTOS PREVIOS

- Programación básica en los lenguajes R y Python
- Conocimiento de las librerías scikit-learn, pandas y numpy
- Conocimientos de estadística y de regresión lineal

3 METODOLOGÍA

Para la demostración de los modelos de regresión se utilizaron las implementaciones de estos algoritmos de scikit-learn, específicamente las clases *LinearRegression* y *SGDRegressor*, probadas en dos datasets diferentes: *genero.txt* y *mtcars.txt*.

3.1 Dataset GENERO

Este dataset está compuesto por 10,000 filas de información, que representan la información del género, el peso y la estatura de una persona. El género está representado por una cadena de caracteres

"Male" o "Female", el peso por un número decimal de libras, y la estatura por un número decimal de pulgadas.

3.1.1 Análisis exploratorio.

Inicialmente se utilizó R para hacer un análisis estadístico de este dataset, específicamente las columnas *Weight* y *Height*, usadas como variable independiente y dependiente respectivamente. La función *summary* de R arrojó el siguiente resultado:

Gender	Height	Weight
Length: 10000	Min: 54.26	Min: 64.7
Class: character	1st Qu: 63.51	1st Qu: 135.8
Mode: character	Median: 66.32	Median: 161.2
	Mean: 66.37	Mean: 161.4
	3rd Qu: 69.17	3rd Qu: 187.2
	Max: 79.00	Max: 270.0

También las funciones *boxplot* y *plot* fueron utilizadas para crear dos gráficas que nos dan una idea general sobre cómo se comportan nuestros datos. La figura 1 representa dos boxplots de las variables de GENERO que usaremos.

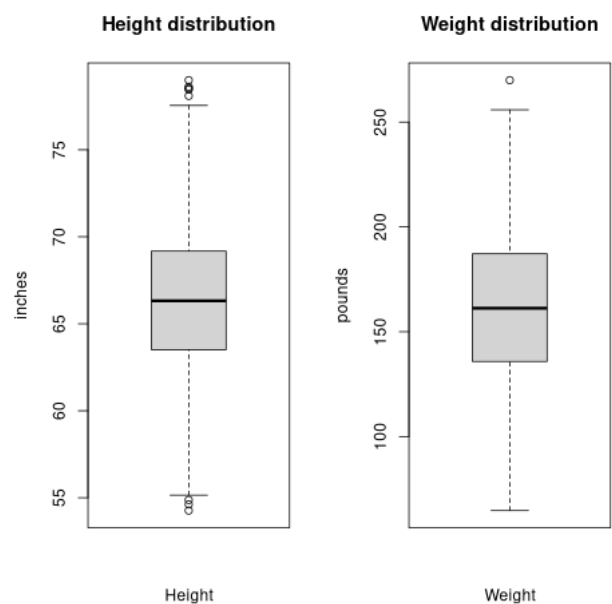


Figure 1: Boxplot del dataset GENERO

Por otro lado, la figura 2 representa los diversos puntos de información graficados en un plano con el peso en el eje horizontal y la altura en el eje vertical.

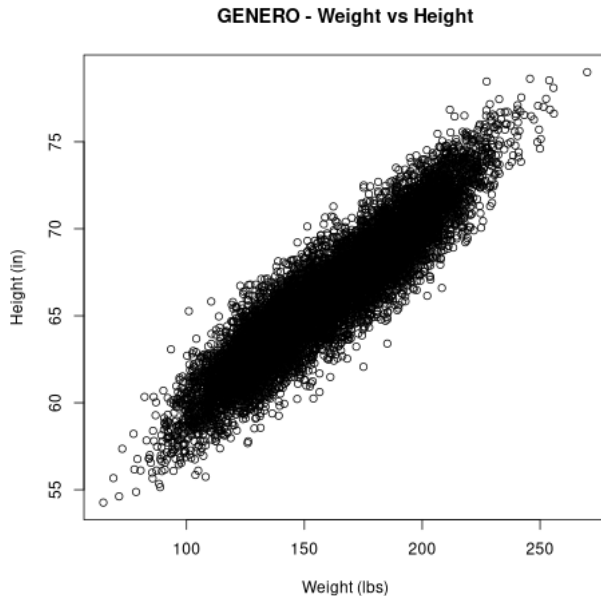


Figure 2: Gráfico de dispersión del dataset GENERO

El código fuente de este análisis puede ser encontrado en el apéndice C.

3.1.2 Regresión lineal.

Como se mencionó anteriormente, se utilizó la implementación del algoritmo de regresión lineal de *scikit-learn*, encontrado en la clase *sklearn.linear_model.LinearRegression*. El código para este ejemplo es simple:

- (1) Usando la librería *pandas* se carga el archivo *genero.txt* a un *Dataframe* usando la función *read_csv*
- (2) Se extraen las columnas de altura y peso del *Dataframe*
- (3) De la librería *sklearn*, se utiliza la función *train_test_split* para dividir tanto la columna de altura y de peso en dos partes cada una: la parte que se usará para entrenar al modelo, con el 80% de los datos, y la parte que se usará para probar el modelo, con el restante 20%
- (4) Se instancia la clase *LinearRegression* y se utiliza el método *fit* para asignarle la columna de datos independiente (en este caso la altura) y la columna de datos dependiente (el peso)
- (5) Se llama al método *predict* de la instancia del modelo de regresión con los datos de prueba de la variable independiente, lo que nos devuelve una columna de datos predichos
- (6) Finalmente, se comparan estos datos predichos con los datos reales (los datos de prueba de la variable dependiente) para sacar un error cuadrático medio, que suele rondar entre 140 y 155

En este caso se decidió normalizar la variable independiente (Height) para que fuera fácilmente comparable con la implementación

de gradiente descendente. El código usado en para demostrar la regresión lineal puede ser encontrado en el apéndice D.

3.1.3 Gradiente descendente.

Para la demostración de este algoritmo se optó por una implementación propia, que exponga las mismas funciones básicas que la clase de *sklearn* pero con un mecanismo interno diferente. La clase de *GradientDescent* se compone por 3 funciones principales:

- El constructor, que inicializa variables como la velocidad de aprendizaje, el número máximo de iteraciones y la precisión final.
- La función *fit*, que recibe dos *Dataframes*, uno de las variables independientes y otro de la dependiente. Este método ejecuta las siguientes instrucciones:
 - (1) Primero normaliza cada variable independiente, esto con el fin de que el algoritmo tenga un mejor funcionamiento. Normalizar las variables se trata de encontrar el promedio (*mu*) y la desviación estándar (*sigma*) de los datos, para así primero restarle el promedio y dividirlos entre la desviación, con el fin de que los datos se comporten como una distribución normal
 - (2) Después se añade una columna de 1's para que en la función de salida tenga una constante no multiplicada por una variable independiente
 - (3) Se crea la variable de los coeficientes e intercepto *theta*, inicializándola con 0's
 - (4) Después, por cada iteración:
 - (a) Se genera una nueva predicción haciendo el producto punto entre las variables independientes y la variable de coeficientes e intercepto
 - (b) A la predicción se le resta la variable dependiente para obtener los errores de la predicción
 - (c) Se multiplica la velocidad de aprendizaje, entre el número de filas en las variables independientes, por el producto punto de las variables independientes por los errores
 - (d) Se resta este resultado a *theta*
 - (e) Se calcula el "costo" de esta predicción usando el error cuadrático medio
 - (f) Se compara el costo con el costo de la iteración previa, y si la diferencia es menor a la precisión del algoritmo, este se detiene
- La función *predict*, que:
 - (1) Recibe variables independientes *x*
 - (2) Las normaliza con los mismos parámetros (promedio y desviación estándar) con los que fueron normalizados los originales
 - (3) Añade la columna de 1's
 - (4) Multiplica estos datos por *theta* usando el producto punto, para así obtener la predicción de la variable dependiente

El código fuente de esta clase se encuentra en el apéndice F y está basado en el artículo de Gunjal, S.[2].

Finalmente, gracias a que esta clase expone las mismas funciones que *LinearRegression* de *sklearn*, se pueden sustituir fácilmente, por lo que el código de esta demostración es muy parecido al de la regresión lineal, lo que cambian son los resultados, que no son conclusivos ya que la selección de datos de entrenamiento y pruebas es

aleatoria. El código fuente de este apartado puede ser encontrado en el apéndice G.

3.2 Dataset MTCARS

Este segundo dataset se compone de 32 filas, cada uno con una variedad de datos sobre el rendimiento y las características del motor de diferentes modelos de automóviles. En esta demostración se utilizarán solamente las columnas de cilindrada ("disp"), peso ("wt") y caballos de fuerza ("hp"); las dos primeras como variables independientes y la tercera como variable dependiente o de salida.

3.2.1 Análisis exploratorio.

Al igual que con el dataset GENERO, se realizó un análisis estadístico inicial con R. La función *summary* arrojó los siguientes resultados sobre las columnas que nos conciernen en este ejemplo:

disp	hp	wt
Min: 71.1	Min: 52.0	Min: 1.513
1st Qu: 120.8	1st Qu: 96.5	1st Qu: 2.581
Median: 196.3	Median: 123.0	Median: 3.325
Mean: 230.7	Mean: 146.7	Mean: 3.217
3rd Qu: 326.0	3rd Qu: 180.0	3rd Qu: 3.610
Max: 472.0	Max: 335.0	Max: 5.424

Para este dataset se usaron las funciones *boxplot* y *pairs* para crear dos tipos de gráficos con los datos: La figura 3 incluye tres boxplots de las variables que incluiremos en la demostración: la cilindrada, el peso y los caballos de fuerza del motor.

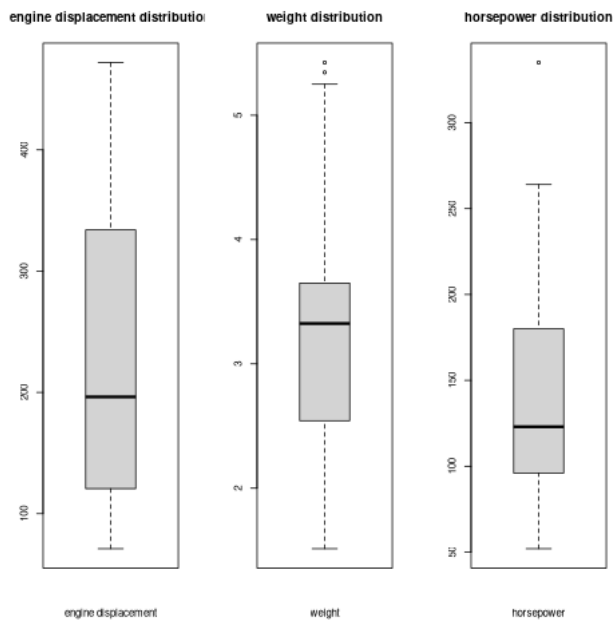


Figure 3: Boxplot del dataset MTCARS

En segundo lugar está la figura 4, que incluye gráficas 6 gráficas de dispersión, comparando las 3 variables con cada combinación posible. De izquierda a derecha, de arriba hacia abajo:

- (1) Caballos de fuerza vs cilindrada

- (2) Peso vs cilindrada
- (3) Cilindrada vs caballos de fuerza
- (4) Peso vs caballos de fuerza
- (5) Cilindrada vs peso
- (6) Caballos de fuerza vs peso

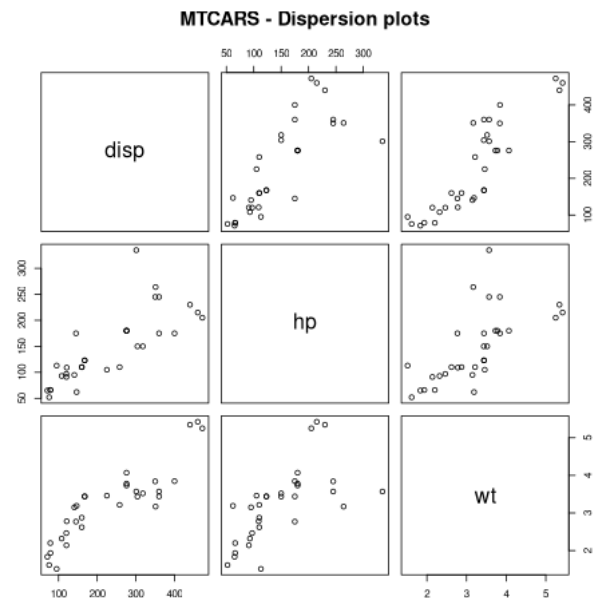


Figure 4: Gráficos de dispersión del dataset MTCARS

El código fuente de este análisis puede ser encontrado en el apéndice C.

3.2.2 Regresión lineal.

Como consecuencia de que este dataset incluye muy pocas instancias o filas, se utilizó la metodología de validación *n-fold cross*, en la que se divide el dataset en pequeños conjuntos, y así se elige uno de estos para realizar las pruebas y los demás dataset para realizar el entrenamiento. Esta selección de conjunto de elementos rota entre todo el dataset para así poder escoger el conjunto que al ser probado tenga el menor error. Para realizar estas iteraciones entre los conjuntos se utilizó la clase *KFold* del paquete *sklearn.model_selection*. Para este ejemplo el código sigue los pasos:

- (1) Primero se lee el dataset utilizando *pandas.read_csv* y se obtiene un *Dataframe*
- (2) Se crea el modelo de validación con *sklearn.model_selection.KFold* y se iteran sobre los conjuntos de datos de entrenamiento y pruebas. Por cada uno de ellos:
 - (a) Se separan las columnas de variables independientes ("disp" y "wt") y dependientes ("hp") de los datos de entrenamiento
 - (b) Se crea un modelo de regresión lineal con la clase *sklearn.linear_model.LinearRegression*[3].
 - (c) Se entrena el modelo con los datos de entrenamiento (independientes y dependientes)

- (d) Se extraen las columnas de los datos de prueba
- (e) Se predice una serie de datos usando las variables independientes de los datos de prueba
- (f) Se comparan los datos predichos con la variable dependiente de los datos de prueba, mediante la fórmula del error cuadrático medio
- (g) Se añade el valor del error a una lista
- (3) Finalmente se promedian los errores para obtener un error promedio

Por defecto, la implementación de *scikit-learn* de *KFold* divide el dataset en 5 partes [3], por lo que se obtienen 5 errores cuadráticos medios y coeficientes de determinación:

- (1) 1084.202, 0.583
- (2) 489.059, 0.742
- (3) 1113.586, 0.818
- (4) 1739.342, 0.463
- (5) 7259.875, 0.102

El promedio de estos errores es de un error cuadrático medio de 2337.21 y un coeficiente de determinación de 0.542.

El código fuente de este análisis puede ser encontrado en el apéndice E.

3.2.3 Gradiente descendente.

Al igual que en el apartado de gradiente descendente del anterior dataset, se utilizó la misma implementación de gradiente descendente, y como consecuencia de que exporta los mismos métodos que *LinearRegression*, el código es muy parecido. Como resultados obtuvimos 5 errores cuadráticos medios y 5 coeficientes de determinación:

- (1) 1050.437, 0.597
- (2) 447.73, 0.764
- (3) 1124.836, 0.816
- (4) 1570.729, 0.515
- (5) 7431.13, 0.082

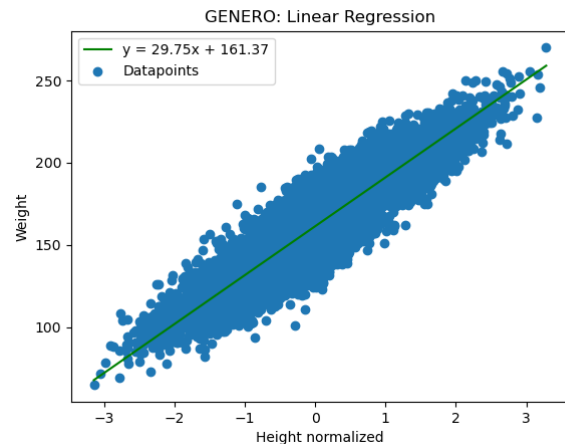
El promedio de estos es 2324.97 para el error cuadrático medio y 0.555 para el coeficiente de determinación.

El código fuente de este análisis puede ser encontrado en el apéndice H.

4 RESULTADOS

4.1 Dataset GENERO

Se elaboró un análisis exploratorio para las 10,000 filas de información del dataset; al observar el resumen de los datos hechos con R se puede apreciar que el comportamiento de los datos es muy semejante a una distribución normal pero agregando algunos outliers, es por esto que el boxplot tiene dicha figura centrada y balanceada; al graficar podemos apreciar que la distribución de los puntos es muy concentrada en el centro, dispersándose poco a poco a los extremos, parecido a la "campana" que se forma al graficar una distribución normal. Al elaborar la regresión lineal se utilizó la metodología de validación simple ya que el dataset era grande, y al evaluar el método nos indica que su **coeficiente de determinación es de 0.85** y su **error cuadrático medio (MSE) de 146**. Al graficar la regresión lineal se aprecia que la línea sigue la tendencia de los datos, concordando con el coeficiente de determinación mencionado anteriormente.

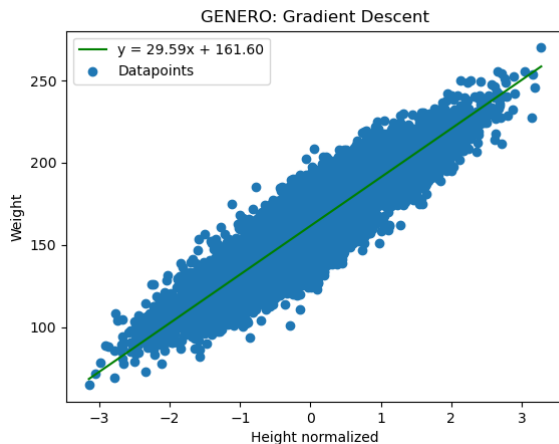


4.2 Dataset MTCARS

Se elaboró un análisis exploratorio para las 32 filas de información del dataset; de primera instancia se aprecia que el dataset no es muy grande, por lo que hacer un modelo de regresión lineal o de predicción resultará en algo muy alejado de la realidad. Al comenzar a explorarlo en R también observamos como se comportan los datos, los diagramas de caja muestran distribuciones diferentes para cada variable, al graficar los datos también quedan demasiado dispersos y no es posible apreciar una tendencia. Cuando se hizo la regresión lineal se utilizó la metodología de n-fold cross validation ya que el dataset era pequeño, el evaluar el método nos indica que **coeficiente de determinación promedio es de 0.542** y su **error cuadrático medio (MSE) promedio de 2337.21**, como era de esperarse, muy alejados de un resultado óptimo. Al graficarlo se muestra una tendencia entre los puntos, pero la dispersión es grande que la correlación termina siendo pequeña. Las gráficas de las validaciones pueden ser encontradas en apéndice A.

4.3 Gradiente descendente

Los resultados con el método del gradiente descendente fueron muy similares al de regresión lineal en términos del comportamiento y de precisión. El error cuadrático de este método arrojaba resultados parecidos al de la regresión lineal, aunque con un costo computacional mayor, por lo que consideramos más conveniente el método de regresión lineal. Como continuación de esta práctica nos gustaría experimentar con datasets más grandes y con diferentes comportamientos de datos para poder comparar cuál es mejor y en qué tipo de casos.



Las gráficas del dataset MTCARS usando gradiente descendente se encuentran en el apéndice B.

5 CONCLUSIONES Y REFLEXIONES

Existen múltiples metodologías para la creación de modelos de regresión lineal, en esta práctica pudimos probar dos: regresión lineal por mínimos cuadrados y gradiente descendente. Sin importar la metodología utilizada, las regresiones lineales tratan de predecir el comportamiento de una función a partir de variables. Durante esta práctica pudimos apreciar que estos métodos se comportan distinto dependiendo de la situación, por lo que no existe un método universal, si no que todo depende del problema a resolver y de las herramientas con las que se cuenta. Nosotros pudimos haber programado el cálculo de la regresión lineal desde cero pero al existir herramientas como scikit learn, podemos desarrollar fácilmente un modelo de regresión lineal para un conjunto de datos; algunas veces no todas las funciones se encuentran descritas en librerías, como la de BGD, por lo que siempre debemos de conocer los conceptos y teoría de las herramientas que empleamos. Si no existiera scikit learn aun así deberíamos poder haber realizado esta práctica, herramientas como numpy, pandas y R la hacen mucho más sencilla, pero el conocer y dominar un tema como este es invaluable. El análisis de los datos resulta muy útil para entender su comportamiento antes de empezar a trabajar con ellos y comprobar si nuestro modelo nos está dando los resultados que esperamos. Quizá no pudimos llegar al 100% de precisión con estos métodos, pero entre más grande es nuestro conjunto de datos, mejor nos aproximamos.

5.1 Reflexión de Abraham

Esta práctica fue todo un reto desde el principio, en la clase había entendido a un muy alto nivel los conceptos, pero al ponerlos en práctica me doy cuenta que hay que ejercitar el conocimiento en entregables como este, ya que lo que había entendido aquel día de clase ahora es mucho más claro al implementarlo. Pensé que esta práctica sería sencilla, pero mis conocimientos de Python, R y estadística no son muchos, es por ello que era difícil conectar conceptos cuando las instrucciones no eran del todo claras, por ejemplo cuando en el análisis pedía generar los boxplots "correspondientes"

no sabía a qué columnas se refería, igual que al implementar métodos de scikit learn, tuve que ir aprendiendo y comprendiendo desde la marcha. Mis conocimientos son mejores ahora, también mi habilidad para investigar mejoró. Ahora hablando específicamente de la práctica, se me hizo muy interesante cómo la regresión lineal nos puede ayudar a predecir y también como el BGD también nos ayuda a aproximarnos, para la siguiente práctica me gustaría poder implementar las cosas con más facilidad y comprender mejor para no revolverme, el resultado se logró y estoy seguro de que lo aprendido me será de mucha utilidad en mi carrera profesional.

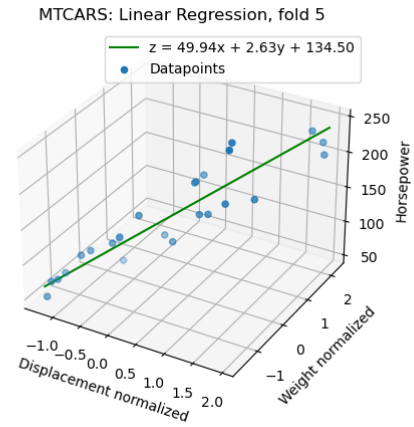
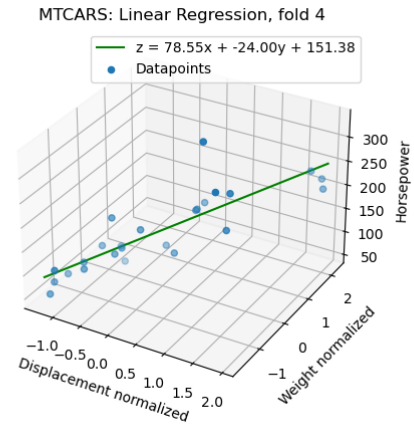
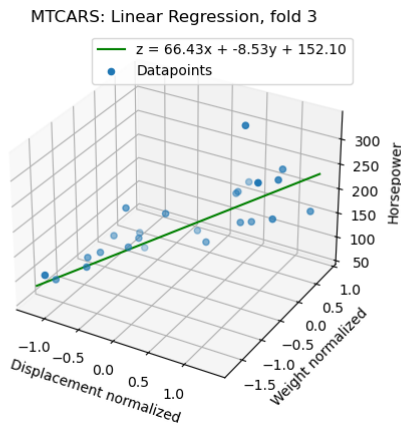
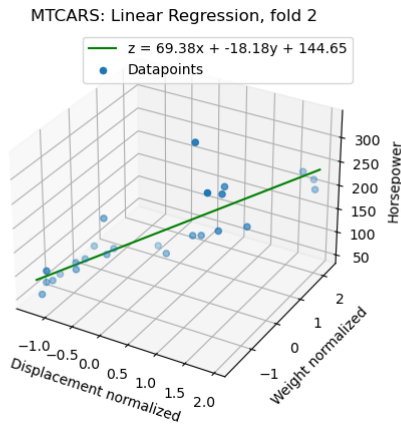
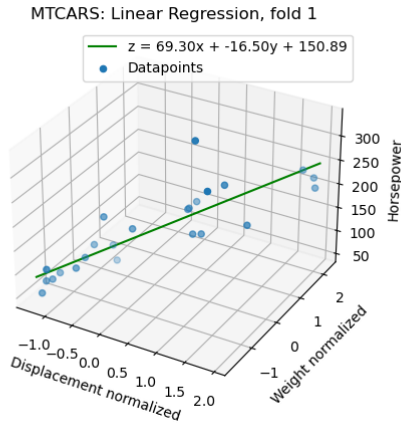
5.2 Reflexión de Mario

Considero que durante todo el desarrollo de la práctica pude obtener nuevos conocimientos: por un lado el uso de las librerías de Python sklearn, numpy, y pandas, y por otro en el ámbito de la estadística y los métodos numéricos aplicados a problemas reales como éste. La parte de la práctica que más disfruté fue implementar el algoritmo de gradiente descendente, ya que pude comprender la importancia de las librerías y cómo facilitan la creación de programas complejos, que en otros lenguajes sería un gran desafío.

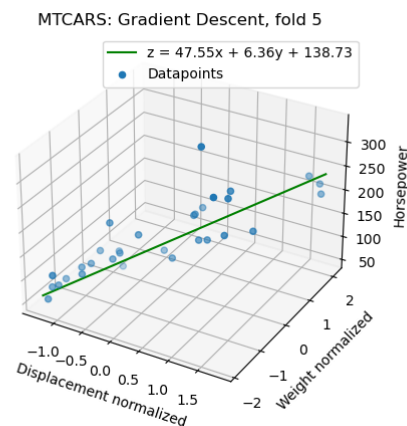
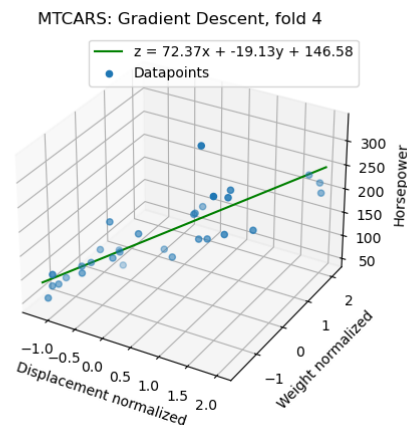
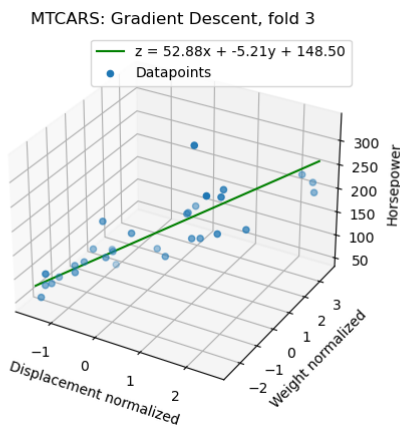
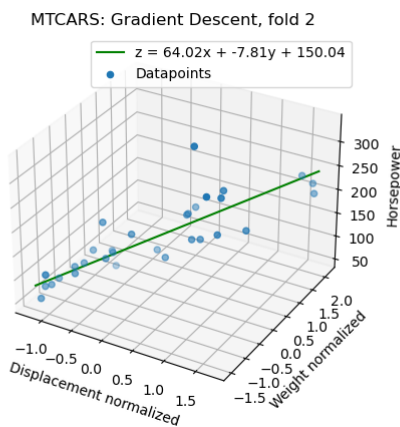
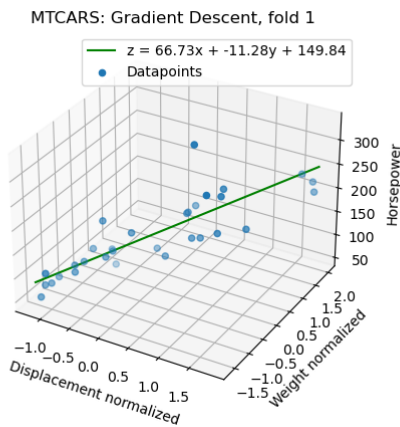
REFERENCES

- [1] Issam El Naqa and Martin J Murphy. 2015. What is machine learning? In *machine learning in radiation oncology*. Springer, 3–11.
- [2] S. Gunjal. 2020. *Multivariate Linear Regression From Scratch With Python*. Retrieved August 22, 2020 from https://satishgunjal.com/multivariate_lr/
- [3] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. 2011. Scikit-learn: Machine Learning in Python. *Journal of Machine Learning Research* 12 (2011), 2825–2830.

A GRÁFICAS DE LA METODOLOGÍA N-FOLD CROSS VALIDATION USANDO REGRESIÓN LINEAL EN EL DATASET MTCARS



B GRÁFICAS DE LA METODOLOGÍA N-FOLD CROSS VALIDATION USANDO GRADIENTE DESCENDENTE EN EL DATASET MTCARS



C CÓDIGO DE ANÁLISIS ESTADÍSTICO

```

1 # Dataset GENERO
2 genero = read.table("genero.txt", header=TRUE, sep=",")
3
4 paste("Resumen estadístico de GENERO")
5 summary(genero)
6
7 paste("Generando boxplot para GENERO")
8 png(file="generoBoxplot.png")
9 par(mfcol=c(1,2))
10 boxplot(genero$Height, ylab="inches", xlab="Height", main="Height distribution")
11 boxplot(genero$Weight, ylab="pounds", xlab="Weight", main="Weight distribution")
12 dev.off()
13
14 paste("Generando gráfica de dispersión para GENERO")
15 png(file="generoDispersion.png")
16 plot(x=genero$Weight, y=genero$Height, xlab="Weight (lbs)", ylab="Height (in)",
17      main="GENERO - Weight vs Height")
18 dev.off()
19
20 cat("\n\n\n")
21
22 # Dataset MTCARS
23 mtcars = read.table("mtcars.txt", header=TRUE)
24
25 paste("Resumen estadístico de MTCARS")
26 summary(mtcars)
27
28 paste("Generando boxplot para MTCARS")
29 png(file="mtcarsBoxplot.png")
30 par(mfcol=c(1,3))
31 boxplot(mtcars$disp, xlab="engine displacement",
32        main="engine displacement distribution")
33 boxplot(mtcars$wt, xlab="weight", main="weight distribution")
34 boxplot(mtcars$hp, xlab="horsepower", main="horsepower distribution")
35 dev.off()
36
37 paste("Generando gráficas de dispersión para MTCARS")
38 png(file="mtcarsDispersion.png")
39 pairs(~disp+hp+wt, data=mtcars, main="MTCARS - Dispersion plots")
40 dev.off()

```


D CÓDIGO DE REGRESIÓN LINEAL DEL DATASET GENERO

```
1 import numpy as np
2 import pandas as pd
3 from sklearn.linear_model import LinearRegression
4 from sklearn.model_selection import train_test_split
5 from sklearn.metrics import mean_squared_error, r2_score
6
7 from GradientDescent import normalize
8 from graphs import graphGENERO
9
10 print(" ~ Reading genero.txt and generating train and test sets")
11 generoData = pd.read_csv('genero.txt')
12 height = generoData.iloc[:, 1].values.reshape(-1, 1)
13 height, _, _ = normalize(height)
14 weight = generoData.iloc[:, 2].values.reshape(-1, 1)
15 heightTrain, heightTest, weightTrain, weightTest = train_test_split(
16     height, weight, test_size=0.2)
17
18 print(" ~ Creating linear regression model")
19 regressor = LinearRegression()
20 regressor.fit(heightTrain, weightTrain)
21
22 print(" ~ Testing linear regression model")
23 predictedWeight = regressor.predict(heightTest)
24 regressorMSE = mean_squared_error(weightTest, predictedWeight)
25 print(" → MSE:", regressorMSE)
26 regressorR2 = r2_score(weightTest, predictedWeight)
27 print(" → R2:", regressorR2)
28
29 theta = np.append(regressor.coef_.copy(), regressor.intercept_)
30 graphGENERO(height, weight, theta, "Linear Regression", "genero_linearReg.png")
```

E CÓDIGO DE REGRESIÓN LINEAL DEL DATASET MTCARS

```

1 import numpy as np
2 import pandas as pd
3 from sklearn.linear_model import LinearRegression
4 from sklearn.model_selection import KFold
5 from sklearn.metrics import mean_squared_error, r2_score
6
7 from GradientDescent import normalize
8 from graphs import graphMTCARS
9
10 print(" ~ Reading mtcars.txt")
11 data = pd.read_csv('mtcars.txt', sep=' ')
12 data["disp"], _, _ = normalize(data["disp"])
13 data["wt"], _, _ = normalize(data["wt"])
14
15 print(" ~ Creating the k-fold cross validation model")
16 kfold = KFold()
17 mses = []
18 r2s = []
19
20 print(" ~ Creating the regression model for every train/test division")
21 for trainIdxs, testIdxs in kfold.split(data):
22     # Fit model to data
23     train = data.iloc[trainIdxs]
24     trainX = train[["disp", "wt"]]
25     trainY = train[["hp"]]
26     regressor = LinearRegression()
27     regressor.fit(trainX, trainY)
28
29     # Graph model
30     idx = len(mses) + 1
31     theta = np.append(regressor.coef_.copy(), regressor.intercept_)
32     graphMTCARS(trainX["disp"], trainX["wt"], trainY["hp"],
33                 theta, "Linear Regression, fold {}".format(idx),
34                 "mtcars_linearReg_{}.png".format(idx))
35
36     # Predict data using model
37     test = data.iloc[testIdxs]
38     testX = test[["disp", "wt"]]
39     testY = test[["hp"]]
40     predictedY = regressor.predict(testX)
41
42     # Score predicted data
43     regressorMSE = mean_squared_error(testY, predictedY)
44     regressorR2 = r2_score(testY, predictedY)
45     print(" → Partial MSE:", regressorMSE, "\tPartial R2: ", regressorR2)
46     mses.append(regressorMSE)
47     r2s.append(regressorR2)
48
49 print(" ~ Calculating final scores")
50 print(" → Final MSE:", np.average(mses))
51 print(" → Final R2:", np.average(r2s))

```

F NUESTRA IMPLEMENTACIÓN DE GRADIENTE DESCENDENTE

```
1 from typing import Tuple
2 import numpy as np
3 import pandas as pd
4 from sklearn.metrics import mean_squared_error
5
6
7 def normalize(x: pd.DataFrame) -> Tuple[np.ndarray, float, float]:
8     mu = np.mean(x, axis=0)
9     sigma = np.std(x, axis=0, ddof=1)
10    x_norm = (x - mu) / sigma
11    return x_norm, mu, sigma
12
13
14 class GradientDescent:
15     def __init__(self, learning_rate: float = 0.1, max_iterations: int = 200,
16                  precision: float = 0.00001):
17         self.learning_rate = learning_rate
18         self.max_iterations = max_iterations
19         self.precision = precision
20         self.theta = None
21         self.mu = None
22         self.sigma = None
23
24     def fit(self, x: pd.DataFrame, y: pd.DataFrame):
25         x, self.mu, self.sigma = normalize(x)
26         x = np.hstack((x, np.ones((x.shape[0], 1))))
27         self.theta = np.zeros(x.shape[1])
28         prev_cost = -1
29         for _ in range(self.max_iterations):
30             predictions = x.dot(self.theta)
31             errors = np.subtract(predictions, y)
32             sum_delta = (self.learning_rate / x.shape[0]) * x.T.dot(errors)
33             self.theta -= sum_delta
34             cost = mean_squared_error(y, predictions)
35             if abs(cost - prev_cost) < self.precision:
36                 break
37             prev_cost = cost
38
39     def predict(self, x: pd.DataFrame):
40         if self.theta is None or self.mu is None or self.sigma is None:
41             raise Exception(
42                 "GradientDescent::predict() called before model was trained")
43
44         x = (x - self.mu) / self.sigma
45         x = np.hstack((x, np.ones((x.shape[0], 1))))
46         return x.dot(self.theta)
```

G CÓDIGO DE GRADIENTE DESCENDENTE DEL DATASET GENERO

```

1 import pandas as pd
2 from sklearn.model_selection import train_test_split
3 from sklearn.metrics import mean_squared_error, r2_score
4
5 from GradientDescent import GradientDescent
6 from graphs import graphGENERO
7
8 print(" ~ Reading genero.txt and generating train and test sets")
9 data = pd.read_csv('genero.txt')
10 height = data.iloc[:, 1].values.reshape(-1, 1)
11 weight = data.iloc[:, 2].values.reshape(-1, 1)
12 heightTrain, heightTest, weightTrain, weightTest = train_test_split(
13     height, weight, test_size=0.2)
14
15 print(" ~ Creating gradient descent model")
16 regressor = GradientDescent()
17 regressor.fit(heightTrain, weightTrain.ravel())
18
19 print(" ~ Testing gradient descent model")
20 predictedWeight = regressor.predict(heightTest)
21 regressorMSE = mean_squared_error(weightTest, predictedWeight)
22 print(" → MSE:", regressorMSE)
23 regressorR2 = r2_score(weightTest, predictedWeight)
24 print(" → R2:", regressorR2)
25
26 normHeight = (height - regressor.mu) / regressor.sigma
27 graphGENERO(normHeight, weight, regressor.theta,
28     "Gradient Descent", "genero_gradient.png")

```

H CÓDIGO DE GRADIENTE DESCENDENTE DEL DATASET MTCARS

```
1 import numpy as np
2 import pandas as pd
3 from sklearn.model_selection import KFold
4 from sklearn.metrics import mean_squared_error, r2_score
5
6 from GradientDescent import GradientDescent
7 from graphs import graphMTCARS
8
9 print(" ~ Reading mtcars.txt")
10 data = pd.read_csv('mtcars.txt', sep=' ')
11
12 print(" ~ Creating the k-fold cross validation model")
13 kfold = KFold()
14 mses = []
15 r2s = []
16
17 print(" ~ Creating the gradient descent model for every train/test division")
18 for trainIdxs, testIdxs in kfold.split(data):
19     # Fit model to data
20     train = data.iloc[trainIdxs]
21     trainX = train[["disp", "wt"]]
22     trainY = train[["hp"]]
23     regressor = GradientDescent()
24     regressor.fit(trainX, trainY.values.ravel())
25
26     # Graph model
27     idx = len(mses) + 1
28     normX = (data[["disp", "wt"]] - regressor.mu) / regressor.sigma
29     graphMTCARS(normX["disp"], normX["wt"], data["hp"], regressor.theta,
30                 "Gradient Descent, fold {}".format(idx),
31                 "mtcars_gradient_{}.png".format(idx))
32
33     # Predict data using model
34     test = data.iloc[testIdxs]
35     testX = test[["disp", "wt"]]
36     testY = test[["hp"]]
37     predictedY = regressor.predict(testX)
38
39     # Score predicted data
40     regressorMSE = mean_squared_error(testY, predictedY)
41     regressorR2 = r2_score(testY, predictedY)
42     print(" → Partial MSE:", regressorMSE, "\tPartial R2: ", regressorR2)
43     mses.append(regressorMSE)
44     r2s.append(regressorR2)
45
46 print(" ~ Calculating final MSE")
47 print(" → Final MSE:", np.average(mses))
48 print(" → Final R2:", np.average(r2s))
```

I CÓDIGO DE GENERACIÓN DE GRÁFICAS

```

1 import sys
2 import matplotlib.pyplot as plt
3 import numpy as np
4
5 shouldDisplay = "--display-graphs" in sys.argv
6 shouldSave = "--save-graphs" in sys.argv
7
8
9 def graphGENERO(height, weight, theta, title, filename):
10     fig, ax = plt.subplots()
11     ax.set_title("GENERO: " + title)
12     ax.set_xlabel("Height normalized")
13     ax.set_ylabel("Weight")
14
15     # Add scatter plot
16     ax.scatter(height, weight, label="Datapoints")
17
18     # Add line from theta parameters
19     points = 60
20     lineX = np.linspace(height.min(), height.max(), points)
21     constant = np.ones(points)
22     lineY = np.dot(np.vstack((lineX, constant)).T, theta)
23     equation = "y = {:.2f}x + {:.2f}".format(theta[0], theta[1])
24     ax.plot(lineX, lineY, 'g', label=equation)
25
26     ax.legend()
27     if shouldSave:
28         fig.savefig(filename)
29     if shouldDisplay:
30         plt.show()
31
32
33 def graphMTCARS(displacement, weight, horsepower, theta, title, filename):
34     fig = plt.figure()
35     ax = fig.add_subplot(projection='3d')
36     ax.set_title("MTCARS: " + title)
37     ax.set_xlabel("Displacement normalized")
38     ax.set_ylabel("Weight normalized")
39     ax.set_zlabel("Horsepower")
40
41     # Add scatter plot
42     ax.scatter(displacement, weight, horsepower, label="Datapoints")
43
44     # Add line from theta parameters
45     points = 60
46     lineX = np.linspace(displacement.min(), displacement.max(), points)
47     lineY = np.linspace(weight.min(), weight.max(), points)
48     constant = np.ones(points)
49     lineZ = np.dot(np.vstack((lineX, lineY, constant)).T, theta)
50     equation = "z = {:.2f}x + {:.2f}y + {:.2f}".format(
51         theta[0], theta[1], theta[2])
52     ax.plot(lineX, lineY, lineZ, 'g', label=equation)
53
54     ax.legend()
55     if shouldSave:
56         fig.savefig(filename)
57     if shouldDisplay:
58         plt.show()

```