

# APLICACIÓN: REGRESIÓN LINEAL

*Nexus-Probability*

## CURSO 3 (PROCESOS ESTOCÁSTICOS I)

### PARTE 1 / LECCIÓN 3

## 1. La esperanza condicional en un modelo de regresión lineal

La **esperanza condicional** es un concepto clave en probabilidad y estadística, que describe el valor esperado de una variable aleatoria  $Y$  dado un conjunto de variables  $X$ . En el contexto de un modelo de regresión lineal, la esperanza condicional tiene una interpretación central, ya que el modelo busca describir cómo  $Y$  depende sistemáticamente de  $X$ , aislando la variabilidad aleatoria.

### La esperanza condicional en regresión lineal simple

Un modelo de regresión lineal simple asume que la relación entre  $Y$  y  $X$  está dada por:

$$Y = \beta_0 + \beta_1 X + \epsilon,$$

donde:

- $\beta_0$  y  $\beta_1$  son los parámetros del modelo (intercepto y pendiente, respectivamente),
- $\epsilon$  es el término de error aleatorio, que satisface:

$$\mathbb{E}[\epsilon] = 0, \quad \text{y generalmente } \epsilon \sim \mathcal{N}(0, \sigma^2).$$

Dado este modelo, la esperanza condicional de  $Y$  dado  $X$  es:

$$\mathbb{E}[Y \mid X] = \mathbb{E}[\beta_0 + \beta_1 X + \epsilon \mid X].$$

Debido a la linealidad de la esperanza y al hecho de que  $\epsilon$  es independiente de  $X$  con  $\mathbb{E}[\epsilon] = 0$ , se obtiene:

$$\mathbb{E}[Y | X] = \beta_0 + \beta_1 X.$$

Esto muestra que la esperanza condicional es una función lineal de  $X$ , que describe la relación promedio entre  $Y$  y  $X$ .

## Propiedades probabilísticas del modelo

En un modelo de regresión lineal, la distribución condicional de  $Y$  dado  $X$  tiene las siguientes propiedades:

- La esperanza condicional  $\mathbb{E}[Y | X]$  es la **tendencia central** de la distribución de  $Y$  para un valor dado de  $X$ .
- La varianza condicional de  $Y$  dado  $X$  es constante e igual a  $\sigma^2$ , es decir:

$$\text{Var}(Y | X) = \mathbb{E}[(Y - \mathbb{E}[Y | X])^2 | X] = \sigma^2.$$

Esto significa que la dispersión de  $Y$  alrededor de  $\mathbb{E}[Y | X]$  no depende de  $X$ .

- La independencia entre  $X$  y el término de error  $\epsilon$  implica que el modelo captura toda la relación sistemática entre  $Y$  y  $X$  en  $\mathbb{E}[Y | X]$ .

## Regresión lineal múltiple

Para un modelo de regresión lineal múltiple con  $k$  variables independientes  $X_1, X_2, \dots, X_k$ , el modelo se escribe como:

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_k X_k + \epsilon.$$

En este caso, la esperanza condicional de  $Y$  dado el vector de variables independientes  $\mathbf{X} = (X_1, X_2, \dots, X_k)$  es:

$$\mathbb{E}[Y | \mathbf{X}] = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_k X_k.$$

Esto significa que:

- La relación esperada entre  $Y$  y las variables  $X_1, X_2, \dots, X_k$  es una función lineal de estas últimas.
- La distribución condicional de  $Y$  dado  $\mathbf{X}$  sigue siendo normal, con:

$$Y | \mathbf{X} \sim \mathcal{N}(\mathbb{E}[Y | \mathbf{X}], \sigma^2).$$

## Interpretación en términos de predicción

En el modelo de regresión lineal, la esperanza condicional  $\mathbb{E}[Y | X]$  o  $\mathbb{E}[Y | \mathbf{X}]$  tiene una interpretación como el valor esperado de  $Y$  para un valor dado de  $X$  (o  $\mathbf{X}$ ). Este valor representa la mejor predicción promedio de  $Y$  bajo el supuesto de mínimos cuadrados. La regresión lineal se basa en encontrar los coeficientes  $\beta_0, \beta_1, \dots, \beta_k$  que minimicen la suma de los errores al cuadrado:

$$\min_{\beta_0, \beta_1, \dots, \beta_k} \sum_{i=1}^n (Y_i - \mathbb{E}[Y | \mathbf{X}_i])^2.$$

La esperanza condicional  $\mathbb{E}[Y | \mathbf{X}]$  es el corazón del modelo de regresión lineal, ya que describe la relación sistemática entre  $Y$  y las variables independientes. Su base probabilística garantiza que el modelo no solo capture la tendencia promedio, sino que también explique cómo el ruido  $\epsilon$  se distribuye de manera independiente alrededor de esta tendencia. Este enfoque proporciona una herramienta poderosa para modelar y predecir relaciones en datos observados.

## 2. Aplicación.

El siguiente modelo fue implementado en Python utilizando una variedad de bibliotecas importantes para el tratamiento de datos, visualización y modelado. Para la manipulación de datos se emplearon pandas y numpy. Las visualizaciones se realizaron con matplotlib y seaborn. El preprocesado y modelado se llevó a cabo utilizando sklearn, que incluye módulos como LogisticRegression y train\_test\_split, además de statsmodels para análisis estadísticos más avanzados. Finalmente, se utilizó la biblioteca warnings para gestionar advertencias y asegurar una ejecución fluida del código.

El siguiente slot ejecuta todas las bibliotecas necesarias para este scrip.

```
1      # Tratamiento de datos
2      import pandas as pd
3      import numpy as np
4
5      #Graficos
6      import matplotlib.pyplot as plt
7      from matplotlib import style
8      import seaborn as sns
9
10     #Procesado y modelado
11     from scipy.stats import pearsonr
12     from sklearn.linear_model import LinearRegression
13     from sklearn.model_selection import train_test_split
14     from sklearn.metrics import r2_score
15     from sklearn.metrics import mean_squared_error
```

```

16     import statsmodels.api as sm
17     import statsmodels.formula.api as smf
18     from scipy import stats
19     import statsmodels as sms
20     from tqdm import tqdm
21     import time
22
23     # configuracion de matplotlib
24     plt.rcParams['image.cmap']="bwr"
25     plt.rcParams['figure.dpi']="100"
26     plt.rcParams['savefig.bbox']="tight"
27     style.use('ggplot') or plt.style.use('ggplot')
28
29     # configuracion de warnings
30     import warnings
31     warnings.filterwarnings('ignore')
32

```

## 2.1. Descripción de la base de datos.

**Jamboree** es una empresa educativa con sede en la India, especializada en la preparación y asesoramiento para estudiantes que buscan ingresar a universidades tanto en su país como en el extranjero. Ofrecen cursos intensivos de preparación para exámenes estandarizados como el SAT, ACT, GRE, GMAT y TOEFL, ayudando a mejorar los puntajes de los estudiantes. **Jamboree** se destaca por su enfoque personalizado y su amplia experiencia en la educación internacional.

A continuación se detalla el contexto y el significado de cada variable:

1. **GRE Score:** Puntaje obtenido en el examen GRE (Graduate Record Examination), utilizado para la admisión en programas de posgrado.
2. **TOEFL Score:** Puntaje obtenido en el examen TOEFL (Test of English as a Foreign Language), que evalúa la competencia en inglés.
3. **University Rating:** Calificación o clasificación de la universidad donde el estudiante ha solicitado admisión.
4. **SOP:** Statement of Purpose (Declaración de Propósito), un ensayo donde el estudiante explica sus motivaciones y metas académicas.
5. **LOR:** Letter of Recommendation (Carta de Recomendación), una carta escrita por un profesor o supervisor que evalúa las habilidades y aptitudes del estudiante.
6. **CGPA:** Cumulative Grade Point Average (Promedio de Calificaciones), el promedio acumulativo de las calificaciones obtenidas durante el grado anterior.

7. **Research:** Experiencia en investigación académica o científica del estudiante.
8. **Chance of Admit:** Probabilidad de admisión, una estimación de las posibilidades del estudiante de ser admitido en el programa deseado.

En este contexto específico:

- **Chance of Admit:** Probabilidad de admisión, una estimación de las posibilidades del estudiante de ser admitido en el programa deseado.
- **Variables predictoras:** Usaremos las variables mencionadas anteriormente en la lista, enumeradas del 1 al 7.

Al considerar todas estas variables como predictores, se obtiene una visión holística y completa del perfil académico, profesional y personal del estudiante, lo cual es crucial para hacer predicciones informadas sobre su probabilidad de admisión en un programa universitario específico.

**Jamboree** es reconocida por su compromiso en ayudar a los estudiantes a alcanzar sus metas académicas mediante una preparación rigurosa y personalizada para el proceso de admisión universitaria en instituciones de renombre a nivel mundial.

[Click para ver la Base de Datos de Jamboree](#)

Lectura de la base de datos con `pd.read_csv()`.

```
1 data = pd.read_csv("jamboree_dataset.csv") #Lectura de la base
    de datos
2 data.head() #Muestra las primeras 5 filas del DF.
```

	Serial No.	GRE Score	TOEFL Score	University Rating	SOP	LOR	CGPA	Research
Chance of Admit								
0	1	337	118	4	4.5	4.5	9.65	1
0.92								
1	2	324	107	4	4.0	4.5	8.87	1
0.76								
2	3	316	104	3	3.0	3.5	8.00	1
0.72								
3	4	322	110	3	3.5	2.5	8.67	1
0.80								
4	5	314	103	2	2.0	3.0	8.21	0
0.65								

Vemos que las columnas de **Serial No.** y de **Research** son categoricas, por tanto las desecharemos del análisis, usaremos `data.drop("", axis = 1)`

```
1 data = data.drop(['Serial No.', 'Research'], axis=1) #Elimina las
    columnas mencionadas.
2 data.head()
```

## 2.2. Análisis Descriptivo.

Usaremos `data.describe()` para generar un resumen estadístico de las columnas.

```
1 data.describe()  
2
```

### GRE Score

- **Media:** 316 puntos.
- **Desviación Estándar:** 11 puntos.
- **Rango:** Desde 290 puntos hasta 340 puntos.

### TOEFL Score

- **Media:** 107.192 puntos.
- **Desviación Estándar:** 6.081868 puntos.
- **Rango:** Desde 92 puntos hasta 120 puntos.

### University Rating

- **Media:** 3.114.
- **Desviación Estándar:** 1.143512.
- **Rango:** Desde 1 hasta 5.

### SOP

- **Media:** 3.374 puntos.
- **Desviación Estándar:** 0.991004 puntos.
- **Rango:** Desde 1 hasta 5.

### LOR

- **Media:** 3.484 puntos.
- **Desviación Estándar:** 0.92545 puntos.
- **Rango:** Desde 1 hasta 5.

## CGPA

- **Media:** 8.57644 puntos.
- **Desviación Estándar:** 0.604813 puntos.
- **Rango:** Desde 6.8 puntos hasta 9.92 puntos.

## Chance of Admit

- **Media:** 0.72174.
- **Desviación Estándar:** 0.14114.
- **Rango:** Desde 0.34 % hasta 0.97 %.

Ahora procedemos a la verificación de datos faltantes con `isna().sum()`

```
1 data.isna().sum() #Cuenta el numero de celdas vacias.
```

```
GRE Score      0
TOEFL Score    0
University Rating 0
SOP            0
LOR            0
CGPA           0
Chance of Admit 0
dtype: int64
```

En efecto, no se encontraron datos faltantes.

Ahora usaremos `.corr()`, calcula la correlación de Pearson entre las columnas de un DataFrame.

```
1 correlacion = data.corr() #Matriz de correlacion del DataFrame
2
3 #Creacion de mapa de calor de correlacion
4 sns.heatmap(correlacion, annot=True, cmap="YlGnBu")
5 plt.title('Mapa de calor de correlacion')
6 plt.show()
```

A continuación se muestra la salida del `s/ot` previo.

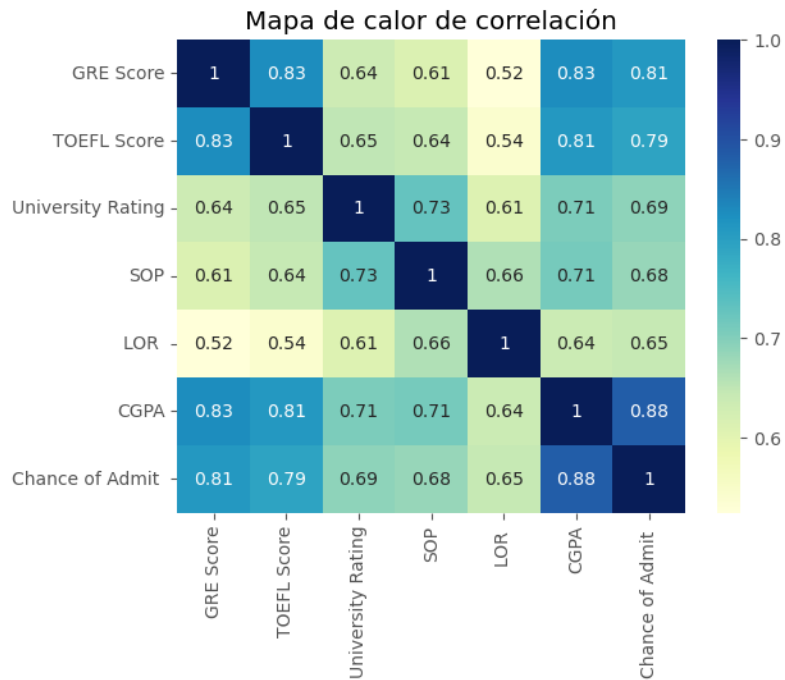


Figura 1: Mapa de Calor de Correlación de la Base de Datos.

Ahora usaremos `sns.pairplot()` para visualizar las relaciones entre múltiples variables numéricas de nuestro DataFrame.

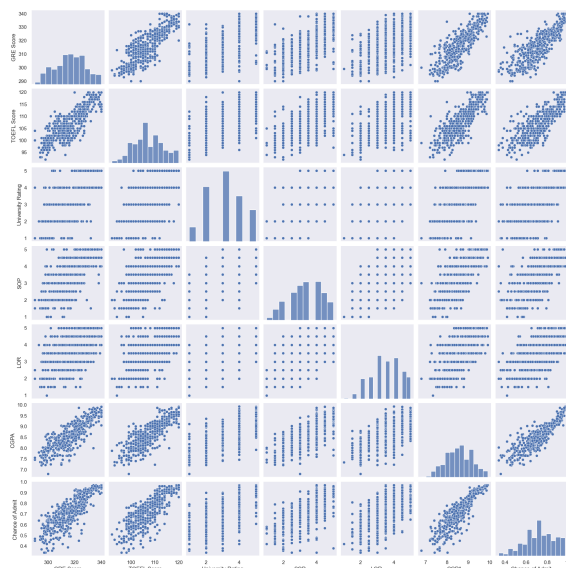


Figura 2: Pair-Plot de la Base de Datos.



El gráfico **pairplot** permite observar las relaciones por pares entre las distintas variables. En la diagonal principal se muestran los histogramas de cada variable, donde se puede notar que aquellas de carácter numérico continuo (como los scores) presentan mejores comportamientos. Las variables categóricas, por otro lado, presentan un nivel de visualización más complejo.

En los **gráficos de dispersión**, se observa que las variables que exhiben una mayor relación y un comportamiento más adecuado para predecir la variable **Chance of Admit** son: TOEFL Score, GRE Score y CGPA. Las variables categóricas como LOR, SOP y University Ranking, aunque parecen no tener un buen comportamiento para predecir Chance of Admit.<sup>a</sup> simple vista, muestran una correlación mayor a 0.6 en el diagrama de correlación. Por lo tanto, no podemos descartarlas de nuestro análisis.

Procedemos a la verificación de posibles **Outliers**, usaremos la función `sns.boxplot()`

```
1 colors = ['#00b8a9', '#f7b801', '#db2b30', '#6a4c93', '#67b7dc',
2         '#b7d63c', '#7f8fa6']
3
4 # Plot 1: GRE Score
5 sns.boxplot(data[["GRE Score"]], color=colors[0], ax=axes[0, 0])
6 axes[0, 0].set_title('Boxplot de GRE Score')
7
8 # Plot 2: TOEFL Score
9 sns.boxplot(data[["TOEFL Score"]], color=colors[1], ax=axes[0,
10         1])
11 axes[0, 1].set_title('Boxplot de TOEFL Score')
12
13 # Plot 3: University Rating
14 sns.boxplot(data[["University Rating"]], color=colors[2], ax=axes
15         [1, 0])
16 axes[1, 0].set_title('Boxplot de University Rating')
17
18 # Plot 4: SOP
19 sns.boxplot(data[["SOP"]], color=colors[3], ax=axes[1, 1])
20 axes[1, 1].set_title('Boxplot de SOP')
21
22 # Plot 5: LOR
23 sns.boxplot(data[["LOR "]], color=colors[4], ax=axes[2, 0])
24 axes[2, 0].set_title('Boxplot de LOR')
25
26 # Plot 6: CGPA
27 sns.boxplot(data[["CGPA"]], color=colors[5], ax=axes[2, 1])
28 axes[2, 1].set_title('Boxplot de CGPA')
29
30 # Plot 7: Chance of Admit
```

```

29 sns.boxplot(data[["Chance of Admit "]], color=colors[6], ax=axes
    [3, 0])
30 axes[3, 0].set_title('Boxplot de Chance of Admit')
31 fig.tight_layout()
32 plt.show()

```

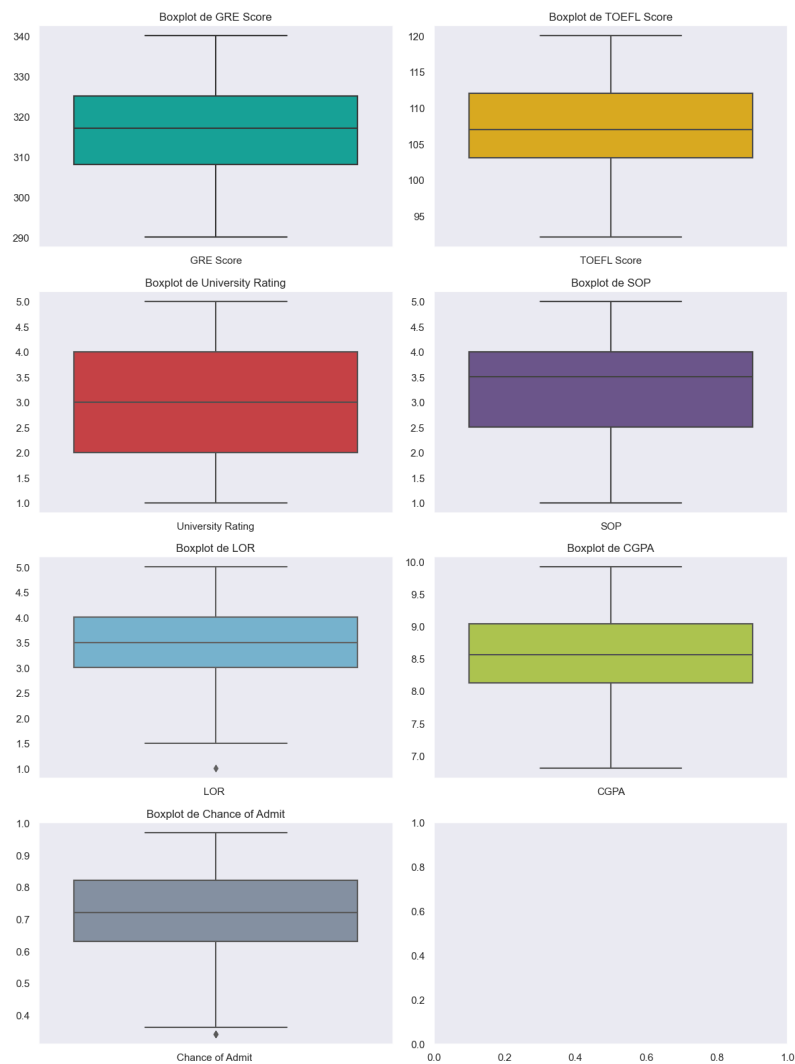


Figura 3: BoxPlot de las Variables Predictoras.

Basandonos en los boxplots anteriores notamos que las variables con presencia de posibles outliers son LOR y Chance of admit.

La variables LOR muestra un puntaje bastante bajo en la carta de recomendación de un alumno mientras que Chance of admite asigna una probabilidad por debajo del 4 % de acceder a la universidad desada.

Algunas posibles causas de un valor tan bajo en LOR pueden ser las siguientes:

- Desempeño Académico Excepcional o Deficiente: Un puntaje bajo en la carta de recomendación podría deberse a un desempeño académico muy deficiente del alumno, lo que puede llevar a que los recomendadores no proporcionen una calificación alta.
- Conflictos Personales: El alumno podría haber tenido conflictos personales con la persona que escribió la carta de recomendación, lo que puede influir negativamente en la evaluación.
- Errores de Evaluación: Puede haber un error o sesgo en la evaluación realizada por la persona que escribió la carta, como malentendidos o errores al interpretar el desempeño del alumno.
- Falta de Información Completa: El recomendador puede no tener información completa o precisa sobre el desempeño y las cualidades del alumno, lo que puede llevar a una evaluación baja.

## 2.3. Creación del Modelo

Recuerde que en `X_orig`, únicamente van almacenadas las variables predictoras.

```
1 X_orig = data.drop("Chance of Admit ",axis=1)
2 X_orig.head()
```

```
1 y_orig = data["Chance of Admit "]
2 y_orig.head()
```

```
0    0.92
1    0.76
2    0.72
3    0.80
4    0.65
```

Name: Chance of Admit , dtype: float64

```
1 # Dividir los datos originales X_orig (características) y y_orig
   (etiquetas) en conjuntos de entrenamiento y prueba
2 X_train, X_test, y_train, y_test = train_test_split(X_orig,
   y_orig, random_state=230624, train_size=0.80)
3
4 # Crear un nuevo DataFrame df combinando las características y
   etiquetas de entrenamiento
5 df = pd.concat([X_train, y_train], axis=1)
```

El siguiente slot, es la parte primordial del código, aquí se entrena todo el modelo, y como salida nos proporciona un resumen de toda la información de la sesión.

```

1 # Anadir una columna constante (intercepto) a X_train
2 X_train = sm.add_constant(X_train, prepend=True)
3
4 # Crear un modelo de regresion lineal ordinaria (OLS) utilizando
   Statsmodels
5 modelo1 = sm.OLS(y_train, X_train)
6
7 # Ajustar el modelo utilizando los datos de entrenamiento
8 modelo1 = modelo1.fit()
9
10 # Imprimir un resumen detallado del modelo ajustado
11 print(modelo1.summary())

```

OLS Regression Results						
=====						
Dep. Variable:	Chance of Admit		R-squared:		0.818	
Model:	OLS		Adj. R-squared:		0.816	
Method:	Least Squares		F-statistic:		354.0	
Date:	Sun, 23 Jun 2024		Prob (F-statistic):		2.87e-143	
Time:	19:35:52		Log-Likelihood:		558.78	
No. Observations:	400		AIC:		-1106.	
Df Residuals:	394		BIC:		-1082.	
Df Model:	5					
Covariance Type:	nonrobust					
=====						
	coef	std err	t	P> t	[0.025	0.975]
-----						
const	-1.3889	0.107	-12.933	0.000	-1.600	-1.178
GRE Score	0.0021	0.001	4.028	0.000	0.001	0.003
TOEFL Score	0.0028	0.001	3.024	0.003	0.001	0.005
University.R	0.0105	0.004	2.636	0.009	0.003	0.018
LOR	0.0181	0.004	4.096	0.000	0.009	0.027
CGPA	0.1205	0.011	11.309	0.000	0.100	0.141
=====						
Omnibus:	89.788	Durbin-Watson:		1.964		
Prob(Omnibus):	0.000	Jarque-Bera (JB):		213.566		
Skew:	-1.118	Prob(JB):		4.21e-47		
Kurtosis:	5.796	Cond. No.		1.19e+04		
=====						

```

1 from complete_analisis import *
2
3 coeficientes(modelo1)

```

```

Beta_0: -1.38976
Beta_1: 0.00214
Beta_2: 0.00285
Beta_3: 0.01060
Beta_4: -0.00044
Beta_5: 0.01817
Beta_6: 0.12069

```

Ahora, podemos construir el modelo de Regresión Lineal Múltiple, vea que:

$$\hat{y} = -\underbrace{1.38976}_{\beta_0} + \underbrace{0.00214}_{\beta_1}x_1 + \underbrace{0.00285}_{\beta_2}x_2 + \underbrace{0.01060}_{\beta_3}x_3 - \underbrace{0.00044}_{\beta_4}x_4 + \underbrace{0.01817}_{\beta_5}x_5 + \underbrace{0.12069}_{\beta_6}x_6$$

## Prueba Global

Ahora, vamos a evaluar si al menos uno de los predictores (**variables independientes**) tiene un efecto significativo sobre la variable de respuesta (**variable dependiente**).

$H_0$  : No hay efecto significativo de las variables predictoras en la variable de respuesta.

$H_1$  : Al menos una de las variables predictoras tiene un efecto significativo en la variable de respuesta.

Ahora en símbolos:

$$H_0 : \beta_1 = \beta_2 = \dots = \beta_6 = 0$$

$$H_1 : \beta_1 \neq \beta_2 \neq \dots \neq \beta_6 \neq 0$$

```

1 prueba_global(modelo1)

```

```

Estadística F : 294.26472705879445

```

```

P_value_f: 0.0000000000

```

```

Progreso: 100% 100/100 [00:01<00:00, 96.89it/s]H_0 ES rechazado.

```

```

Al menos una de las variables predictoras tiene un efecto
significativo en la variable de respuesta.

```

En efecto, vea que al menos una de las variables predictoras tiene un efecto significativo en la variable de respuesta.

$$P\text{-value} = 0.000 < 0.05 = \alpha \text{ (True)}$$

Por lo tanto,  $H_0$  es rechazado y podemos inferir que:

$$\beta_1 \neq \beta_2 \neq \dots \neq \beta_6 \neq 0$$

Con un nivel  $\alpha = 0.05$  de significancia.

## Pruebas Individuales.

$H_0$  : No hay efecto significativo de la variable  $\beta_i$  en la variable de respuesta.

$H_1$  : La variables predictora  $\beta_i$  tiene un efecto significativo en la variable de respuesta.

Ahora en simbolos:

$$H_0 : \beta_i = 0$$

$$H_1 : \beta_i \neq 0$$

```
1 prueba_ind(modelo1)
```

PRUEBA PARA BETA\_0

Estadística t para Beta\_0: -12.86207

P-value para Beta\_0: 0.00000

H\_0 ES rechazado.

Beta\_0 tiene un efecto significativo en la variable de respuesta.

-----  
PRUEBA PARA BETA\_1

Estadística t para Beta\_1: 4.01574

P-value para Beta\_1: 0.00007

H\_0 ES rechazado.

Beta\_1 tiene un efecto significativo en la variable de respuesta.

-----  
PRUEBA PARA BETA\_2

Estadística t para Beta\_2: 3.02080

P-value para Beta\_2: 0.00269

H\_0 ES rechazado.

Beta\_2 tiene un efecto significativo en la variable de respuesta.

-----  
PRUEBA PARA BETA\_3

Estadística t para Beta\_3: 2.44204

P-value para Beta\_3: 0.01504

H\_0 ES rechazado.

Beta\_3 tiene un efecto significativo en la variable de respuesta.

-----  
PRUEBA PARA BETA\_4

Estadística t para Beta\_4: -0.08248

P-value para Beta\_4: 0.93431

H\_0 NO es rechazado.

Beta\_4 NO tiene efecto significativo en la variable de respuesta.

#### PRUEBA PARA BETA\_5

Estadística t para Beta\_5: 3.91547

P-value para Beta\_5: 0.00011

H\_0 ES rechazado.

Beta\_5 tiene un efecto significativo en la variable de respuesta.

---

#### PRUEBA PARA BETA\_6

Estadística t para Beta\_6: 11.04250

P-value para Beta\_6: 0.00000

H\_0 ES rechazado.

Beta\_6 tiene un efecto significativo en la variable de respuesta.

---

En resumen, tenemos de las pruebas individuales:

$\beta_0 \rightarrow$  Contribuye al modelo

$\beta_1 \rightarrow$  Contribuye al modelo

$\beta_2 \rightarrow$  Contribuye al modelo

$\beta_3 \rightarrow$  Contribuye al modelo

$\beta_4 \rightarrow$  NO contribuye al modelo

$\beta_5 \rightarrow$  Contribuye al modelo

$\beta_6 \rightarrow$  Contribuye al modelo

#### Ecuación de la Mejora del Modelo.

Ahora, podemos construir el modelo de Regresión Lineal Múltiple, vea que:

$$\hat{y} = -\underbrace{1.38889}_{\beta_0} + \underbrace{0.00214}_{\beta_1}x_1 + \underbrace{0.00285}_{\beta_2}x_2 + \underbrace{0.01045}_{\beta_3}x_3 + \underbrace{0.01805}_{\beta_4}x_4 + \underbrace{0.12050}_{\beta_5}x_5$$

La esperanza condicional  $\mathbb{E}[Y \mid \mathbf{X}]$  es el corazón del modelo de regresión lineal, ya que describe la relación sistemática entre  $Y$  y las variables independientes. Su base probabilística garantiza que el modelo no solo capture la tendencia promedio, sino que también explique cómo el ruido  $\epsilon$  se distribuye de manera independiente alrededor de esta tendencia. Este enfoque proporciona una herramienta poderosa para modelar y predecir relaciones en datos observados.