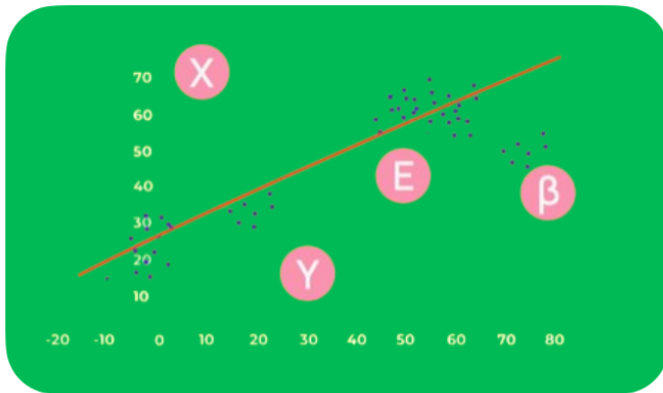




Regresión Lineal

Consideraciones para el análisis de datos en los cursos de laboratorio de Física



- ① El problema de la regresión
- ② La regresión lineal simple
- ③ Método de mínimos cuadrados
- ④ Coeficiente de regresión
- ⑤ Coeficiente de correlación lineal
- ⑥ Inferencias acerca de los parámetros
- ⑦ Condiciones y supuestos para modelo de regresión lineal
- ⑧ Un ejemplo en donde no se cumplen los supuestos

1. El problema de la regresión

- La regresión es una técnica estadística para investigar y modelar relaciones entre variables.
- Las relaciones estadísticas difieren de las funcionales porque no son perfectas; las observaciones no caen directamente sobre una curva.
- Se supone una relación entre una respuesta cuantitativa y y k predictores x_1, x_2, \dots, x_k de la forma general:

$$y = f(x) + \varepsilon$$

donde f es una función desconocida de x_1, \dots, x_k y ε es un término de error aleatorio independiente de x con media cero.

- La función f representa la información sistemática que x proporciona sobre y .
- El método paramétrico más utilizado asume que f es lineal en x :

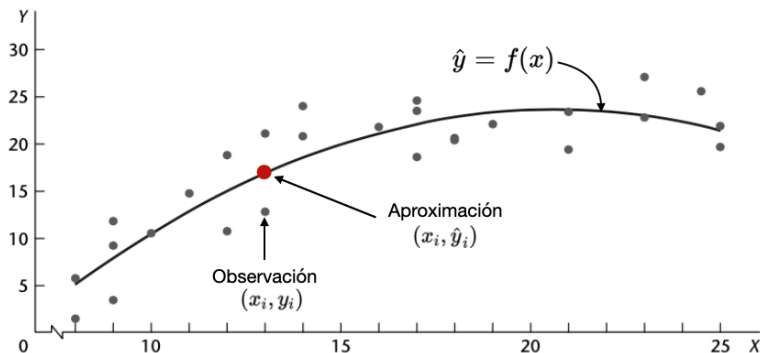
$$f(x) = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_k x_k$$

- Para ajustar el modelo lineal, se estiman los parámetros $\beta_0, \beta_1, \dots, \beta_k$ de manera que:

$$y \approx \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_k x_k$$

Con técnicas de regresión de una variable y (respuesta), sobre una variable x (regresor), se busca una función que sea una buena aproximación de una nube de puntos, mediante una curva del tipo:

$$\hat{y} = f(x)$$



2. La regresión lineal simple

- Modelo:

$$y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$$

donde y_i es la respuesta, x_i es el regresor, β_0 es la intersección, β_1 es la pendiente, y ε_i es el término de error aleatorio.

- Características del Modelo:

- y_i es una variable aleatoria compuesta por $\beta_0 + \beta_1 x_i$
- La intersección β_0 y la pendiente β_1 son constantes desconocidas y ε es un componente de error aleatorio.
- Se supone que los errores tienen media cero y varianza desconocida σ^2 .

- Normalmente se asume que los errores no están correlacionados.
- El modelo asume varianza constante para y_i :

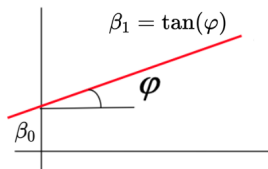
$$\sigma^2 \{y_i\} = \sigma^2$$

- Los términos de error ε_i no están correlacionados entre sí, lo que implica que las respuestas y_i tampoco lo están.

Un modelo con un único regresor x y que tiene una relación con la respuesta y en la forma de una línea recta es lo que se denomina un modelo de regresión lineal simple:

$$y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$$

En donde β_0 es la ordenada en el origen (el valor que toma y cuando x vale 0), β_1 es la pendiente de la recta (e indica cómo cambia y al incrementar x en una unidad) y ε una variable que incluye un conjunto grande de factores, cada uno de los cuales influye en la respuesta sólo en pequeña magnitud, a la que llamaremos error.



Por lo tanto, x e y son variables aleatorias, por lo que no se puede establecer una relación lineal exacta entre ellas.

- **Ejemplo:** Ley de enfriamiento de Steinhart-Hart.

Para describir con precisión la relación entre la resistencia R de un termistor y su temperatura T se utiliza la ecuación de Steinhart-Hart:

$$\frac{1}{T} = A + B \ln(R) + C[\ln(R)]^3$$

donde:

- T es la temperatura.
- R es la resistencia en ohmios.
- A, B, C son constantes específicas del termistor.

En un rango pequeño de temperaturas, la relación entre R y T puede aproximarse a:

$$R \approx R_0 + \alpha (T - T_0)$$

donde: R_0 es la resistencia a una temperatura de referencia T_0 y α es el coeficiente de temperatura, que indica la tasa de cambio de la resistencia con la temperatura.

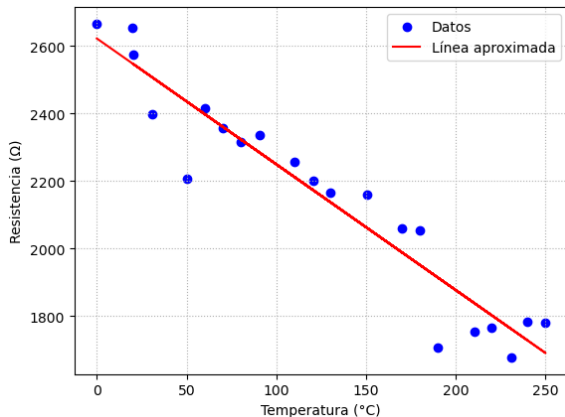
Experimento: Medición de la resistencia de un termistor NTC

Descripción del Experimento:

- **Objetivo:** Determinar la relación lineal aproximada entre la resistencia y la temperatura de un termistor NTC en un rango limitado de temperaturas.
- **Variables:** y , resistencia del material y x la temperatura.

Procedimiento:

- Calentar el termistor: coloca el termistor en un baño de agua caliente para llevarlo a una temperatura significativamente más alta que la ambiente.
- Medir la resistencia y temperatura: retirar el termistor del agua caliente y medir su resistencia a intervalos regulares mientras se enfría. Simultáneamente, medir la temperatura del termistor.
- Registrar los Datos: anotar las mediciones de resistencia R y temperatura T



3. Método de mínimos cuadrados

El método de los mínimos cuadrados permite estimar β_0 y β_1 de forma que la suma de los cuadrados de las diferencias entre las observaciones y_i y la recta sea un mínimo.

$$y_i = \beta_0 + \beta_1 x_i + \varepsilon_i, \quad i = 1, 2, \dots, n \quad (1)$$

La ecuación (1) es un modelo de regresión muestral, escrito en términos de los n pares de datos (y_i, x_i) , con $i = 1, 2, \dots, n$.

A los estimadores obtenidos por mínimos cuadrados β_0 y β_1 , los llamaremos b y m , respectivamente, y deben satisfacer una relación lineal de la forma:

$$\hat{y} = mx + b,$$

donde \hat{y} es la variable dependiente y x es la variable independiente, en nuestro caso la magnitud controlada por el experimentador. El método de mínimos cuadrados consiste en minimizar suma de los cuadrados de los errores:

$$\sum_{i=1}^n e_i^2 = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

Es decir, la suma de los cuadrados de las diferencias entre los valores reales observados (y_i) y los valores estimados (\hat{y}_i) .

Minimizar suma de los cuadrados de los errores se logra calculando las derivadas parciales de la suma con respecto a m y con respecto a b , e igualándolas a cero.

Con este método, las expresiones que se obtiene para b y m son las siguientes:

$$m = \frac{S_{xy}}{S_x^2}, \quad b = \bar{y} - m\bar{x},$$

En donde \bar{x} e \bar{y} denotan las medias muestrales de x y y (respectivamente),

$$\bar{x} = \frac{\sum_{i=1}^n x_i}{n}, \quad \bar{y} = \frac{\sum_{i=1}^n y_i}{n},$$

S_x^2 es la varianza muestral de x y S_{xy} es la covarianza muestral entre x y y . Estos parámetros se calculan como:

$$S_x^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n}, \quad S_y^2 = \frac{\sum_{i=1}^n (y_i - \bar{y})^2}{n}, \quad S_{xy} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{n}.$$

La cantidad m se denomina coeficiente de regresión de y sobre x , lo denotamos por $m_{y/x}$.

En nuestro ejemplo, los estadísticos descriptivos anteriores para las ($n = 21$) variables temperatura y tiempo del enfriamiento son los siguientes:

$$\bar{x} = 125,0119, \quad \bar{y} = 2156,81$$

$$S_x^2 = 12896,0685, \quad S_y^2 = 196582,6099$$

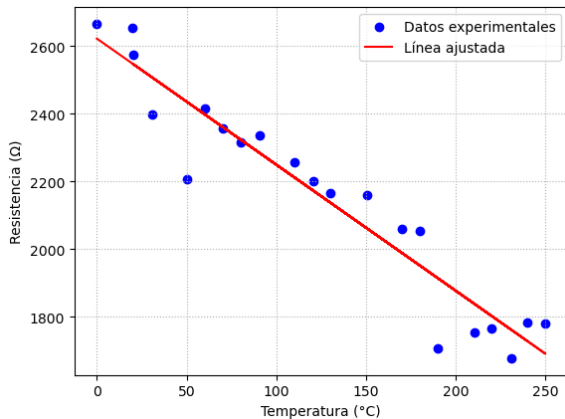
$$S_{xy} = -48017,6465$$

$$m = -3,7234, \quad b = 2622,2834$$

La recta de regresión ajustada es la siguiente:

$$\hat{y} = 2622,2834 - 3,7234x,$$

donde \hat{y} es la resistencia y x la temperatura.



Hay varias propiedades que se cumplen para los mínimos cuadrados:

- ① La suma de los valores observados y_i es igual a la suma de los valores ajustados \hat{y}_i , o bien

$$\sum_{i=1}^n y_i = \sum_{i=1}^n \hat{y}_i$$

- ② La suma de los residuos que contiene un intercepto β_0 es siempre cero

$$\sum_{i=1}^n (y_i - \hat{y}_i) = \sum_{i=1}^n e_i = 0$$

- ③ La recta siempre pasa por el centroide, el punto (\bar{x}, \bar{y}) de los datos.

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i, \quad \bar{y} = \frac{1}{n} \sum_{i=1}^n y_i = m\bar{x} + b.$$

- ④ La suma de los residuos ponderada por el valor de la variable regresora es cero

$$\sum_{i=1}^n x_i e_i = 0$$

- ⑤ La suma de los residuos ponderada por el valor ajustado siempre es cero

$$\sum_{i=1}^n \hat{y}_i e_i = 0$$

4. Coeficiente de regresión

El coeficiente de regresión nos da información sobre el comportamiento de la variable y frente a la variable x , de manera que:

- ① Si $m_{y/x} = 0$, para cualquier valor de x la variable y es constante.
- ② Si $m_{y/x} > 0$, esto nos indica que al aumentar el valor de x , también aumenta el valor de y .
- ③ Si $m_{y/x} < 0$, esto nos indica que al aumentar el valor de x , el valor de y disminuye.

En el ajuste de regresión lineal para la resistencia y la temperatura resultó en

$$\hat{y} = 2622,2834 - 3,7234x,$$

El coeficiente de regresión que se obtuvo fue $m_{y/x} = -3,7234 < 0$ y esto indica que al aumentar x disminuye y .

5. Coeficiente de correlación lineal

El coeficiente de correlación lineal entre x e y viene dado por:

$$r = \frac{S_{xy}}{S_x S_y}$$

y es una medida estadística que cuantifica la intensidad de la relación lineal entre dos variables. Su cuadrado se denomina coeficiente de determinación r^2 .

Propiedades del coeficiente de correlación:

- 1 No tiene dimensión, y siempre toma valores en $[-1,1]$.
- 2 Si las variables son independientes, entonces $r = 0$, el inverso no tiene por qué ser cierto.
- 3 Si existe una relación lineal exacta entre x e y , entonces $r = 1$ (relación directa) ó $r = -1$ (relación inversa).
- 4 Si $r > 0$, indica una relación directa entre las variables (si aumenta x , aumenta y).
- 5 Si $r < 0$, indica una relación inversa entre las variables (si aumenta x , disminuye y).

Para nuestro ejemplo el valor de r es

$$r = \frac{S_{xy}}{S_x S_y} = \frac{-48017,6465}{\sqrt{12896,0685} \sqrt{196582,6099}} = -0,9537$$

El menos indica que existe una relación inversa entre las variables. Además su valor es próximo a 1 indicando una dependencia lineal muy fuerte.

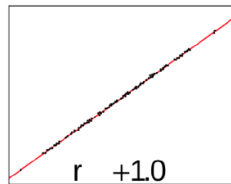
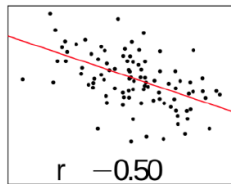
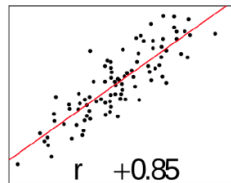
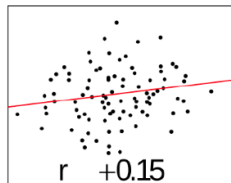
El coeficiente de determinación al cuadrado es $r^2 = 0,9095$.

La relación entre los coeficientes de regresión y de correlación:

$$m_{y/x} = r \frac{S_y}{S_x} = -3,7234,$$

$$m_{x/y} = r \frac{S_x}{S_y} = -0,2443$$

Los dos coeficientes de regresión y el coeficiente de correlación tienen el mismo signo: a medida que x aumenta y tiende a disminuir.



- Descomposición de la variabilidad:

- 1 Variabilidad Explicada (Sum of Squares for Regression, SSR): Representa la parte de la variabilidad de y que se explica por el modelo de regresión.

$$SSR = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2$$

- 2 Variabilidad Residual (Sum of Squares for Error, SSE): Representa la parte de la variabilidad de y que no se puede explicar por el modelo de regresión.

$$SSE = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

- 3 Variabilidad Total (Total Sum of Squares, SST): Representa la variabilidad total de la variable dependiente y respecto a su media.

$$SST = \sum_{i=1}^n (y_i - \bar{y})^2$$

- La variabilidad total se puede descomponer en la variabilidad explicada por el modelo y la variabilidad residual

$$SST = SSR + SSE \Rightarrow \sum_{i=1}^n (y_i - \bar{y})^2 = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 + \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

- Coeficiente de determinación r^2 : El coeficiente de determinación es una medida que cuantifica la proporción de la variabilidad total de y que es explicada por el modelo de regresión. Se calcula como:

$$r^2 = \frac{\sum (\hat{y}_i - \bar{y})^2}{\sum (y_i - \bar{y})^2} = \frac{\text{SSR}}{\text{SST}}$$

- Para el modelo de temperaturas y resistencias:

$$\text{SSR} = 1787904,8827, \text{SSE} = 177921,2163, \text{SST} = 1965826,099$$

$$r^2 = \frac{\text{SSR}}{\text{SST}} = \frac{1787904,8827}{1965826,099} = 0,9095$$

6. Inferencias acerca de los parámetros

Existen pruebas estadísticas para evaluar la significancia de los coeficientes en un modelo de regresión. Ayudan a determinar si las relaciones observadas entre las variables independientes (predictores) y la variable dependiente (respuesta) son estadísticamente significativas.

① Hipótesis en el Contraste de Regresión:

- Hipótesis Nula (H_0): Indica que el coeficiente de regresión no tiene un efecto significativo en la variable dependiente. Matemáticamente es $\beta_i = 0$, donde β_i es el coeficiente de regresión de una variable independiente específica.
- Hipótesis Alternativa (H_a) Sugiere que el coeficiente de regresión tiene un efecto significativo en la variable dependiente. Esto se expresa como $\beta_i \neq 0$.

② Estadístico t (t-Estadístico) Para cada coeficiente de regresión β_i , se calcula un t-estadístico para evaluar su significancia. Este se define como:

$$t_i = \frac{\hat{\beta}_i}{SE(\hat{\beta}_i)}$$

donde $\hat{\beta}_i$ es el estimador del coeficiente de regresión y $SE(\hat{\beta}_i)$ es el error estándar del coeficiente de regresión.

- ③ Valor p (p-Value): el valor p asociado con el t -estadístico se utiliza para tomar decisiones respecto a las hipótesis. Si el valor p es menor que un nivel de significancia predefinido (generalmente 0.05), se rechaza la hipótesis nula, indicando que el coeficiente es significativamente diferente de cero.
- ④ Contraste Global del Modelo: además de evaluar cada coeficiente individualmente, se puede realizar un contraste global del modelo para verificar si al menos una de las variables independientes tiene un efecto significativo. Esto se hace usando el estadístico F (F-statistic)

$$F = \frac{MSR}{MSE} = \frac{\frac{SSR}{k-1}}{\frac{SSE}{n-k}} = \frac{\frac{SSR}{k-1}}{\hat{\sigma}^2}$$

donde k es el número de coeficientes estimados en el modelo de regresión y n el número de observaciones. $MSE = \hat{\sigma}^2$ = Error Mean Square o Residual Mean Square.

Para el ejemplo en estudio podemos calcular los t-estadístico, primero calculamos el error estándar residual y los errores estándar de la pendiente Δm y el término independiente Δb de la regresión lineal.

$$SE(m)^2 = \Delta m = \frac{\hat{\sigma}^2}{\sum_{i=1}^n (x_i - \bar{x})^2}, \quad SE(b)^2 = \Delta b = \hat{\sigma}^2 \left[\frac{1}{n} + \frac{\bar{x}^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \right]$$

donde

$$\hat{\sigma}^2 = \frac{SSE}{n-2} = 9364,2745$$

$$\Delta m = \sqrt{\frac{\hat{\sigma}^2}{\sum_{i=1}^n (x_i - \bar{x})^2}} = 0,2695$$

$$\Delta b = \sqrt{\hat{\sigma}^2 \left[\frac{1}{n} + \frac{\bar{x}^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \right]} = 39,7582$$

$$t_m = \frac{m}{\Delta m} = -13,8177, \quad t_b = \frac{b}{\Delta b} = 65,9557$$

- El nivel de confianza refleja la probabilidad de que el intervalo de confianza contenga el valor verdadero del parámetro de la población. Se construye a través de un valor crítico $t_{\alpha/2, n-2}$ o t_{critical} que es un punto específico en la distribución (t -Student). Un nivel de confianza del 95 % corresponde a $\alpha = 0,05$, porque $1 - 0,95 = 0,05$. Dado que el nivel de confianza es bilateral, se divide en dos colas de la distribución t , lo que significa que cada cola tiene un área de $\alpha/2$.
- Fórmula del intervalo de confianza:

$$\hat{\beta}_j \pm t_{\alpha/2, n-2} \cdot \text{SE}(\hat{\beta}_j)$$

Aquí, $\hat{\beta}_j$ es el estimador del parámetro (intersección o pendiente), y $\text{SE}(\hat{\beta}_j)$ es el error estándar del estimador.

- Estimaciones de los intervalos de confianza para los parámetros:

En el contexto del ejemplo para tener un nivel de confianza del 95 %, $\alpha = 0,05$. Como los grados de libertad son 19, el valor crítico $t_{\alpha/2, \text{dof}}$ para $\alpha/2 = 0,025$ y 19 grados de libertad es de 2.093. (prueba t de Student)

Para $\hat{\beta}_1 = m = -3,7234$ y $\text{SE}(m) = 0,2695$ Para $\hat{\beta}_0 = b = 2622,2834$ y $\text{SE}(b) = 39,7582$

$$\text{CI} = \hat{\beta}_1 \pm t_{\alpha/2, n-2} \cdot \text{SE}(\hat{\beta}_1)$$

$$\text{CI} = -3,7234 \pm 2,093 \cdot 0,2695$$

$$\text{CI} = -3,7234 \pm 0,5640$$

$$\text{CI} = [-4,2874, -3,1594]$$

$$\text{CI} = \hat{\beta}_0 \pm t_{\alpha/2, n-2} \cdot \text{SE}(\hat{\beta}_0)$$

$$\text{CI} = 2622,2834 \pm 2,093 \cdot 39,7582$$

$$\text{CI} = 2622,2834 \pm 83,2140$$

$$\text{CI} = [2539,0694, 2705,4975]$$

Estos intervalos indican que con un 95 % de confianza, el verdadero valor de la pendiente m y el intercepto b se encuentra dentro de este rango.

- La significancia global del modelo de regresión

Para el cálculo de estadístico F (F-statistic)

$$F = \frac{MSR}{MSE} = \frac{\frac{SSR}{2 - 1}}{\frac{SSE}{21 - 2}}$$

Resulta

$$F = \frac{1787904,8827}{9364,2745} = 190,9283$$

Este valor indica que el modelo es estadísticamente significativo, lo que significa que la relación entre las variables independientes y la variable dependiente es poco probable que se deba al azar.

- statsmodels: es un módulo de Python que proporciona clases y funciones para la estimación de modelos estadísticos diferentes, así como para la realización de pruebas estadísticas, y la exploración de datos estadísticos. Para cada estimador se dispone de una extensa lista de estadísticas de resultados. Los resultados se comprueban con paquetes estadísticos existentes para garantizar que son correctos. La documentación en línea se encuentra en statsmodels.org.

OLS Regression Results						
=====						
Dep. Variable:	y	R-squared:	0.909			
Model:	OLS	Adj. R-squared:	0.905			
Method:	Least Squares	F-statistic:	190.9			
Date:	Thu, 30 May 2024	Prob (F-statistic):	2.31e-11			
Time:	19:33:52	Log-Likelihood:	-124.77			
No. Observations:	21	AIC:	253.5			
Df Residuals:	19	BIC:	255.6			
Df Model:	1					
Covariance Type:	nonrobust					
=====						
	coef	std err	t	P> t	[0.025	0.975]

const	2622.2834	39.758	65.956	0.000	2539.068	2705.498
x1	-3.7234	0.269	-13.818	0.000	-4.287	-3.159
=====						
Omnibus:	6.366	Durbin-Watson:	2.027			
Prob(Omnibus):	0.041	Jarque-Bera (JB):	4.490			
Skew:	-1.114	Prob(JB):	0.106			
Kurtosis:	3.411	Cond. No.	278.			
=====						

Condiciones y supuestos para modelo de regresión lineal

Hemos visto cómo aproximar el modelo de regresión lineal simple

$$y = \beta_0 + \beta_1 x + \varepsilon,$$

por la recta

$$\hat{y} = a + bx.$$

Para garantizar que una aproximación usando una regresión lineal simple es válida, deben cumplirse varias condiciones y supuestos. Estas condiciones aseguran que el modelo lineal es una buena representación de la relación entre las variables y que las inferencias realizadas a partir del modelo son fiables.

- ➊ **Linealidad:** La relación entre la variable dependiente y y la variable independiente x debe ser lineal. Esto significa que los puntos de datos deben seguir un patrón que se asemeja a una línea recta.
- ➋ **Independencia:** Los errores o residuos e_i deben ser independientes entre sí. Esto significa que el error de una observación no debe influir en el error de otra observación.

- ③ **Homoscedasticidad:** La varianza de los errores $e_i = (\hat{y}_i - y_i)$ debe ser constante en todos los niveles de la variable independiente x . Esto implica que la dispersión de los puntos alrededor de la línea de regresión debe ser aproximadamente la misma para todos los valores de x .
- ④ **Normalidad de los Errores:** Los errores e_i deben seguir una distribución normal con media cero y varianza constante. Esto es especialmente importante para la validación de las pruebas de hipótesis y los intervalos de confianza.
- Verificación de Condiciones
 - Linealidad: Diagramas de Dispersión (Scatter Plots).
 - Homoscedasticidad: Gráficos de Residuos (Residual Plots).
 - Independencia: Prueba de Durbin-Watson.
 - Normalidad: Histograma de Residuos o Gráfico Q-Q (Quantile-Quantile Plot).

- La prueba de Durbin-Watson

Es una prueba estadística utilizada para detectar la presencia de autocorrelación en los residuos (errores) de un modelo de regresión. La autocorrelación ocurre cuando los residuos no son independientes entre sí, lo que puede invalidar muchas de las inferencias que se pueden hacer a partir del modelo.

La ecuación para el estadístico DW es:

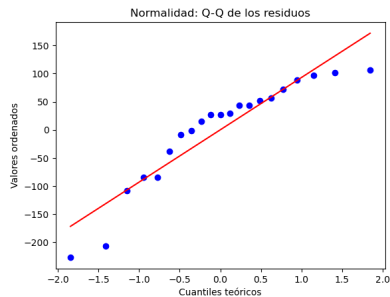
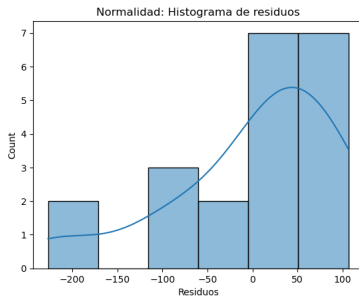
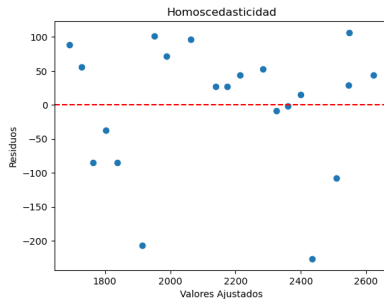
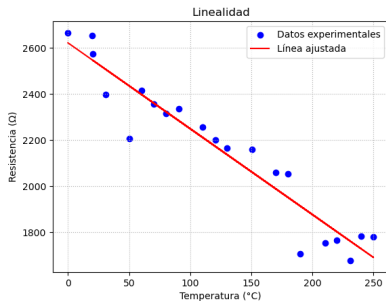
$$DW = \frac{\sum_{i=2}^n (e_i - e_{i-1})^2}{\sum_{i=1}^n e_i^2}$$

donde e_i son los residuos del modelo.

Interpretación del Estadístico de Durbin-Watson

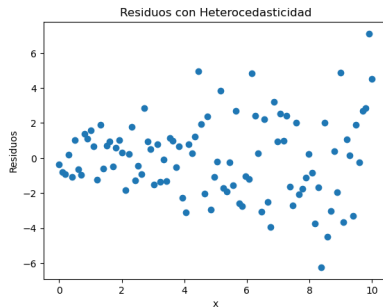
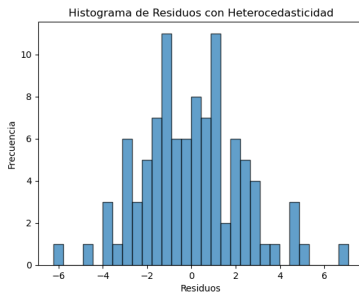
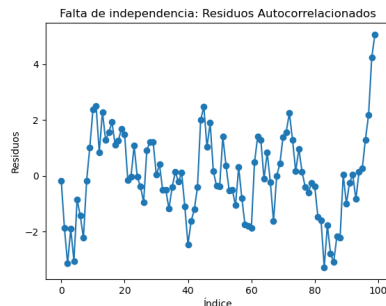
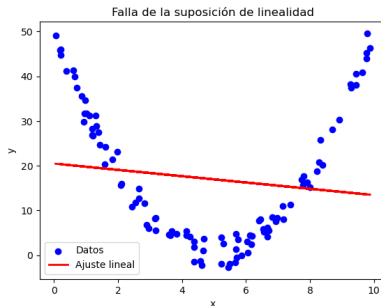
- $DW \approx 2$: No hay autocorrelación.
- $DW < 2$ Autocorrelación positiva.
- $DW > 2$: Autocorrelación negativa.

Para los datos considerados en el ejemplo: Medición de la resistencia de un termistor

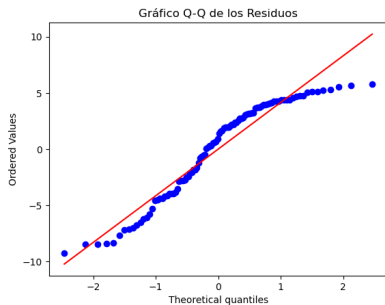
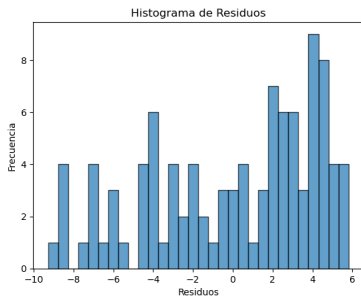
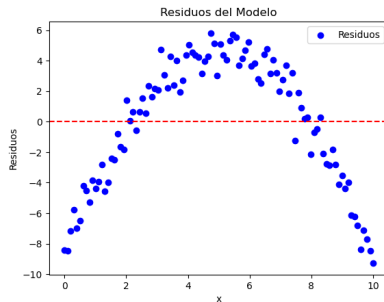
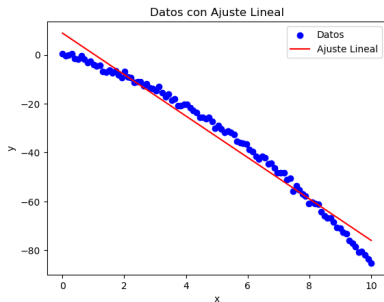


Estadístico de Durbin-Watson: 2.027

- Algunos casos en los que no se cumplen los supuestos



En el ejemplo mostrado pareciera que a simple vista el ajuste lineal es adecuado pero el resto de criterios falla



OLS Regression Results

```

=====
Dep. Variable:          y      R-squared:          0.972
Model:                  OLS    Adj. R-squared:      0.971
Method:                 Least Squares    F-statistic:      3364.
Date:                   Mon, 03 Jun 2024    Prob (F-statistic): 1.13e-77
Time:                   16:24:10    Log-Likelihood:    -285.95
No. Observations:      100    AIC:              575.9
Df Residuals:          98    BIC:              581.1
Df Model:               1
Covariance Type:       nonrobust
=====

```

	coef	std err	t	P> t	[0.025	0.975]
const	8.9013	0.847	10.511	0.000	7.221	10.582
x1	-8.4862	0.146	-58.002	0.000	-8.777	-8.196

```

=====
Omnibus:                15.781    Durbin-Watson:          0.099
Prob(Omnibus):           0.000    Jarque-Bera (JB):       8.325
Skew:                    -0.525    Prob(JB):               0.0156
Kurtosis:                 2.054    Cond. No.                11.7
=====

```