

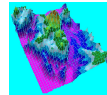
GEOESTADÍSTICA APLICADA

Tema: Análisis Exploratorio de Datos

Dr. Martín A. Díaz Viera,
Dr. Ricardo Casar González



UNIVERSIDAD NACIONAL
AUTÓNOMA DE MÉXICO
mdiazv@imp.mx



Contenido I

1 Análisis Exploratorio de Datos

- ¿Qué es el AED?
- Importancia del AED
- Etapas de un AED
- Herramientas del AED

2 Estadística univariada

- Función de Distribución de Probabilidad
- Percentiles o cuantiles de una distribución
- Valor esperado y momentos de una VA
- Distribuciones Normal y Lognormal

3 Estadística bivariada

- Función de Distribución de Probabilidad Bivariada
- Covarianza y semivarianza

Contenido II

- Coeficiente de correlación lineal
- Coeficientes de correlación de rango

4 Estadística multivariada

- Estadística multivariada
- Regresión lineal y mínimos cuadrados
- Análisis de regresión lineal

¿Qué es el AED?

- *Es un conjunto de técnicas estadísticas y gráficas que permiten establecer un buen entendimiento básico del comportamiento de los datos y de las relaciones existentes entre las variables que se estudian.*

¿Qué es el AED?

- El análisis exploratorio de datos (AED) es un paso previo e indispensable para la aplicación exitosa de cualquier método estadístico.
- En particular permite la detección de fallos en el diseño y toma de datos, el tratamiento y/o la evaluación de datos ausentes, la identificación de valores atípicos y la comprobación de los supuestos requeridos por parte de las técnicas geoestadísticas.

Etapas de un AED

- Realizar un examen gráfico de la naturaleza de las variables individuales y un análisis descriptivo numérico que permita cuantificar algunos aspectos gráficos de los datos.
- Realizar un examen gráfico de las relaciones entre las variables y un análisis descriptivo numérico que cuantifique el grado de interrelación existente entre ellas.
- Evaluar algunos supuestos básicos subyacentes a muchas técnicas estadísticas, por ejemplo, normalidad, linealidad y homocedasticidad.
- Identificar los posibles valores atípicos (outliers) y evaluar el impacto potencial que puedan ejercer en análisis estadísticos posteriores.
- Evaluar, el impacto potencial que pueden tener los datos ausentes (missing) sobre la representatividad de los datos analizados.

Herramientas del AED

- Estadística univariada
- Estadística multivariada
- Regresión lineal y mínimos cuadrados

Estadística univariada

- **Variable Aleatoria (V.A.)** Es una variable Z que puede tomar una serie de valores o realizaciones (z_i) cada una de las cuales tienen asociadas una probabilidad de ocurrencia (p_i).
- Ejemplo: Al lanzar un dado puede resultar $\{1, 2, 3, 4, 5 \text{ o } 6\}$ con una probabilidad de ocurrencia igual a $1/6$.
- Las probabilidades cumplen las condiciones:

$$a) p_i \geq 0 \quad \forall i \quad b) \sum_i p_i = 1 \quad (1)$$

Estadística univariada

- **Variable Aleatoria Discreta** cuando el número de ocurrencias es finito o contable, se conoce como variable aleatoria discreta.
- Ejemplo: tipos de facies en un yacimiento.
- **Variable Aleatoria Continua** si el número de ocurrencias posibles es infinito.
- Ejemplo: el valor de la porosidad de un medio se encuentra en el intervalo $[0, 100 \ %]$.

Función de Distribución de Probabilidad (FDP)

- La **FDP** caracteriza completamente a la **VA**
- Se define como:

$$F(z) = Pr \{Z \leq z\} \in [0, 1] \quad (2)$$

- Su gráfica es el histograma acumulativo

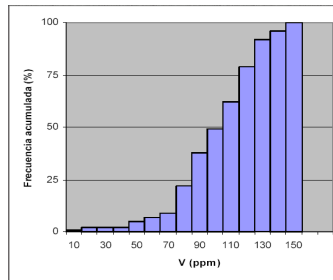


Figura 1: Histograma acumulativo.

Función de Densidad de Probabilidad (fdp)

- Se define como:

$$f(z) = \frac{dF(z)}{dz} \quad (3)$$

- Su gráfica es el histograma

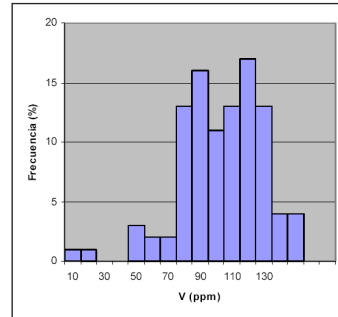


Figura 2: Histograma.

Percentiles o cuantiles de una distribución

- El percentil de una distribución $F(z)$ es el valor z_p de la **V.A.** que corresponde a un valor p de probabilidad acumulada, es decir:

$$F(z_p) = p \quad (4)$$

- Si existe la función inversa se puede expresar como:

$$z_p = F^{-1}(p) \quad (5)$$

Percentiles o cuantiles de una distribución

Algunos cuantiles de interés:

- Mediana $p = 0.5$ $M = F^{-1}(0.5)$
- Cuartiles
 - Primer cuartil o inferior $p = 0.25$ $z_{0.25} = F^{-1}(0.25)$
 - Tercer cuartil o superior $p = 0.75$ $z_{0.75} = F^{-1}(0.75)$
 - Rango o intervalo intercuartil (IR) $[z_{0.25}, z_{0.75}]$

Estadística univariada

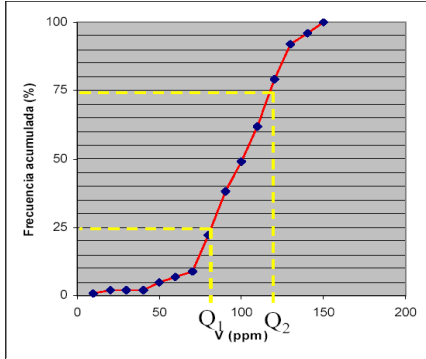


Figura 3: Ejemplo de cuartiles y rango intercuartil.

Valor esperado o esperanza matemática de una VA. I

- Es el valor más probable que puede tomar una VA. Se conoce también como valor medio o media.
- Se define como:

$$m = E[Z] = \int_{-\infty}^{+\infty} z dF(z) \quad (6)$$

- Se calcula como el promedio de todas las observaciones de la variable Z

$$m^* = \frac{1}{N} \sum_{i=1}^N z_i \quad (7)$$

Valor esperado o esperanza matemática de una VA. II

- Es muy sensible a los valores atípicos (outliers)

Estadística univariada

- **Momento de orden r de una FDP**

$$m_r = E[Z^r] = \int_{-\infty}^{+\infty} z^r dF(z) = \int_{-\infty}^{+\infty} z^r f(z) dz \quad (8)$$

- **Momento centrado de orden r de una FDP**

$$\mu_r = E[(Z - m)^r] = \int_{-\infty}^{+\infty} (z - m)^r dF(z) = \int_{-\infty}^{+\infty} (z - m)^r f(z) dz \quad (9)$$

Estadística univariada

- **Varianza de una VA (2do. momento centrado)**

- Se define como

$$\sigma^2 = \text{Var}[Z] = E[(Z - m)^2] \geq 0 \quad (10)$$

- Caracteriza la dispersión de la distribución respecto al valor medio
- Se estima

$$(\sigma^2)^* = \frac{1}{N-1} \sum_{i=1}^N (z_i - m)^2 \quad (11)$$

Distribución Normal o Gaussiana

- Esta distribución está completamente caracterizada por sus dos parámetros: media y varianza y se designa mediante

$$N(m, \sigma^2) \quad (12)$$

- La *fdp* normal o Gaussina está dada por

$$g(z) = \frac{1}{\sigma\sqrt{2\pi}} \exp \left[-\frac{1}{2} \left(\frac{z - m}{\sigma} \right)^2 \right] \quad (13)$$

Distribución Normal o Gaussiana

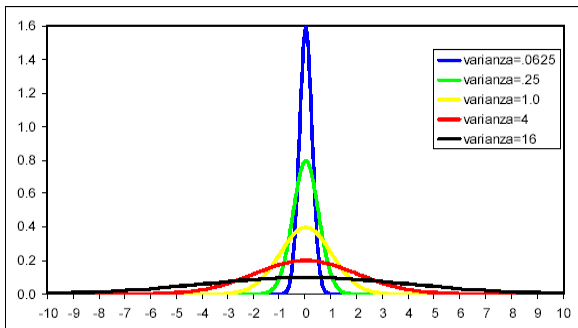


Figura 4: Ejemplos de distribuciones Gaussianas.

Distribución LogNormal

- Una VA positiva \mathbf{Y} se dice que tiene una distribución lognormal si su logaritmo $\ln(\mathbf{Y})$ está normalmente distribuido.

$$Y > 0 \rightarrow \log N(m, \sigma^2), \quad \text{si } X = \ln(Y) \rightarrow N(\alpha, \beta^2) \quad (14)$$

- Muchas distribuciones experimentales en Ciencias de la Tierra tienden a ser asimétricas y la mayoría de las variables toman valores no negativos.

Ejemplos de distribuciones Lognormales

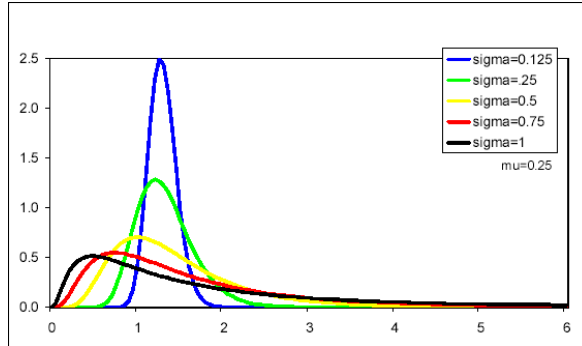


Figura 5: Ejemplos de distribuciones Lognormales.

Estadística univariada I

- Desviación Estándar

$$\sigma = \sqrt{\text{Var}[Z]} \quad (15)$$

- Coeficiente de variación (dispersión relativa)

$$CV = \frac{\sigma}{m} \quad (16)$$

- Coeficiente de simetría (medida de la simetría)

$$\alpha_1 = \frac{\mu_3}{\mu_2^{\frac{3}{2}}} \quad (17)$$

Estadística univariada II

- Coeficiente de curtosis (medida del achatamiento)

$$\alpha_2 = \frac{\mu_4}{\mu_2^2} - 3 \quad (18)$$

Estadística univariada

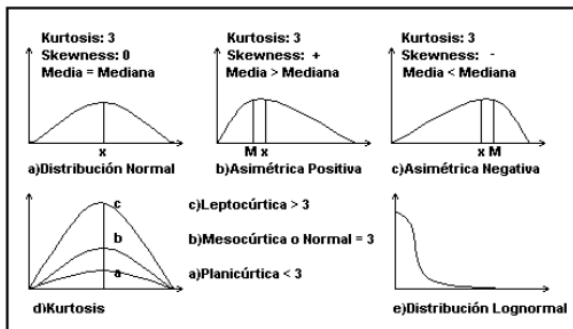


Figura 6: Simetría y curtosis de una distribución.

Estadística univariada

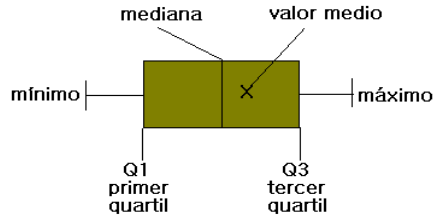


Figura 7: Gráfica de cajas sin valores atípicos.

Estadística univariada

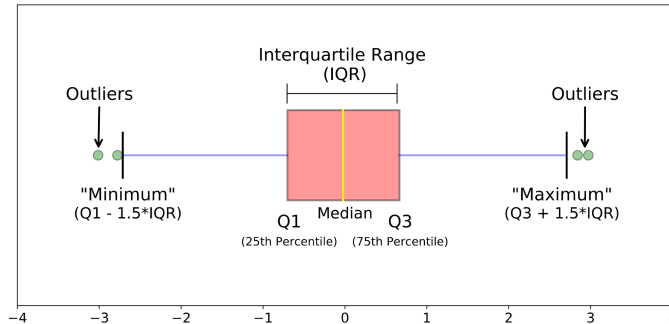


Figura 8: Gráfica de cajas con valores atípicos.

Figura 9: Histograma de la porosidad.

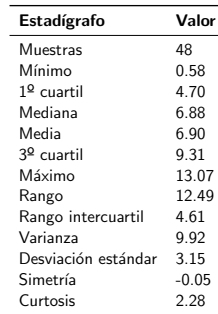


Tabla 1: Estadística básica.

Estadística univariada

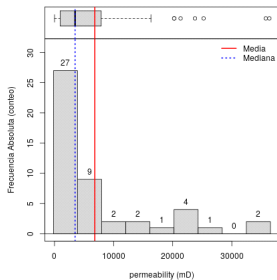
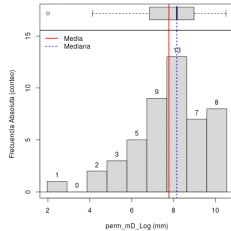


Figura 10: Histograma de la permeabilidad.

Estadígrafo	Valor
Muestras	48
Mínimo	7.40
1º cuartil	1,002.45
Mediana	3,482.20
Media	6,818.24
3º cuartil	7,743.62
Máximo	36,347.4
Rango	36,340
Rango intercuartil	6,741.17
Varianza	83,532,706.36
Desviación estándar	9,139.62
Simetría	1.83
Curtosis	5.62

Tabla 2: Estadística básica.

Estadística univariada



Estadígrafo	Valor
Muestras	48
Mínimo	2.00
1º cuartil	6.91
Mediana	8.15
Media	7.77
3º cuartil	8.95
Máximo	10.50
Rango	8.49
Rango intercuartil	2.04
Varianza	3.21
Desviación estándar	1.79
Simetría	-0.84
Curtosis	3.87

Figura 11: Transformación logarítmica de la Permeabilidad.

Tabla 3: Estadística básica.

Estadística univariada

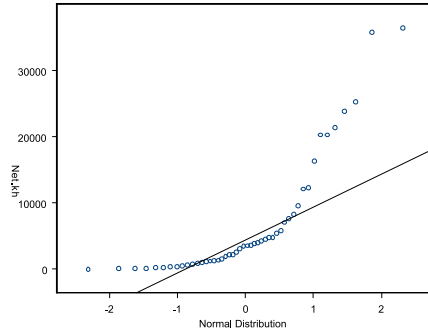


Figura 12: Q-Q plot de la permeabilidad antes de transformar.

Estadística univariada

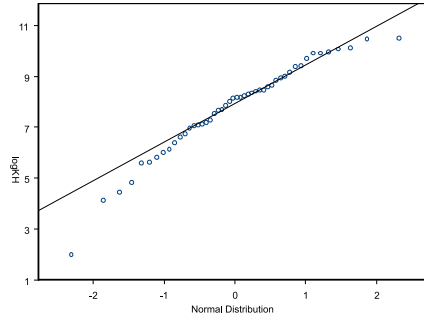


Figura 13: Q-Q plot de la permeabilidad después de transformar.

Estadística univariada

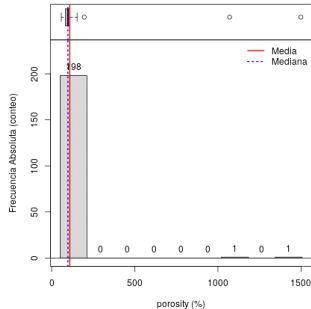


Figura 14: Con valores atípicos (outliers).

Estadígrafo	Valor
Muestras	200
Mínimo	58.2
1º cuartil	82.25
Mediana	97.85
Media	108.9925
3º cuartil	110.325
Máximo	1499
Rango	1440.8
Rango intercuartil	28.075
Varianza	14873.08823
Desviación estándar	121.95527
Simetría	9.92162
Curtosis	104.73871

Tabla 4: Estadística básica.

Estadística univariada

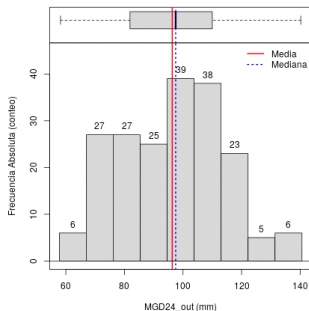


Figura 15: Sin valores atípicos (outliers).

Estadígrafo	Valor
Muestras	196
Mínimo	58.2
1º cuartil	82
Mediana	97.5
Media	96.3265
3º cuartil	110
Máximo	140.2
Rango	82
Rango intercuartil	28
Varianza	319.7503
Desviación estándar	17.8816
Simetría	0.0291
Curtosis	2.3889

Tabla 5: Estadística básica.

Estadística univariada

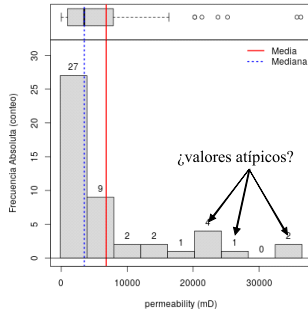


Figura 16: ¿Serán valores atípicos?.

Estadígrafo	Valor
Muestras	48
Mínimo	7.4
1º cuartil	1002.45
Mediana	3482.205
Media	6818.24521
3º cuartil	7743.625
Máximo	36347.4
Rango	36340
Rango intercuartil	6741.175
Varianza	83532706.36
Desviación estándar	9139.62288
Simetría	1.83579
Curtosis	5.62603

Tabla 6: Estadística básica.

Estadística univariada

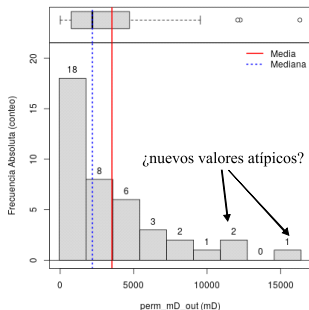


Figura 17: Después de eliminar los valores atípicos.

Estadígrafo	Valor
Muestras	41
Mínimo	7.4
1º cuartil	748
Mediana	2188.7
Media	3521.0285
3º cuartil	4720.5
Máximo	16315.9
Rango	16308.5
Rango intercuartil	3972.5
Varianza	14353741.71
Desviación estándar	3788.6332
Simetría	1.5704
Curtosis	5.1874

Tabla 7: Estadística básica.

Estadística bivariada

- Hasta el momento, sólo hemos considerado a las variables aleatorias por separado, sin que exista ninguna interrelación entre éstas.
- En muchos campos de aplicación y en particular, en las Ciencias de la Tierra, es frecuentemente más importante conocer el patrón de dependencia que relaciona a una variable aleatoria X (porosidad) con otra variable aleatoria Y (permeabilidad).
- Por lo que le dedicaremos especial atención al análisis conjunto de dos variables aleatorias, conocido como análisis bivariado.

Función de Distribución de Probabilidad Bivariada

- La distribución de probabilidad conjunta de un par de variables aleatorias **X** y **Y** se define como:

$$F_{XY}(x, y) = Pr \{X \leq x, Y \leq y\} \quad (19)$$

- En la práctica se estima mediante la proporción de pares de valores de **X** y **Y** que se encuentran por debajo del umbral x, y respectivamente.

Diagrama de Dispersión (Scattergram)

- El equivalente bivariado del histograma es el diagrama de dispersión o scattergram, donde cada par (x_i, y_i) es un punto.
- El grado de dependencia entre dos variables aleatorias **X** y **Y** puede ser caracterizado por el diagrama de dispersión alrededor de cualquier línea de regresión.

Covarianza

- Se define la covarianza de manera análoga a los momentos centrales univariados, como

$$\text{Cov}(X, Y) = \sigma_{XY} = E \{ (X - m_X)(Y - m_Y) \} \quad (20)$$

- Se estima como

$$\sigma_{XY}^* = \frac{1}{N} \sum_{i=1}^N (x_i - m_X)(y_i - m_Y) = \frac{1}{N} \sum_{i=1}^N x_i y_i - m_X m_Y \quad (21)$$

Semivarianza

- Se define como

$$\gamma_{XY} = \frac{1}{2} E[(X - Y)^2] \quad (22)$$

- Se interpreta como el momento de inercia del diagrama de dispersión con respecto a una línea con pendiente de 45° .
- Se estima como

$$\gamma_{XY}^* = \frac{1}{N} \sum_{i=1}^N [d_i]^2 = \frac{1}{2N} \sum_{i=1}^N [x_i - y_i]^2 \quad (23)$$

- Permite caracterizar la carencia de dependencia.

Semivarianza

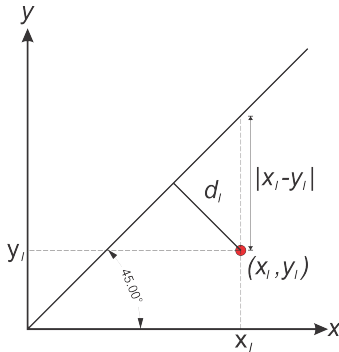


Figura 18: Semivarianza.

Mientras mayor sea el valor de la semivarianza más dispersos estarán los valores en el diagrama de dispersión y menor será la dependencia entre las dos variables aleatorias.

Observe que por el teorema de Pitágoras tenemos:

$$2[d_i]^2 = [x_i - y_i]^2 \quad (24)$$

Coeficiente de correlación lineal de Pearson

- Se define como:

$$r_{XY} = \frac{\sigma_{XY}}{\sigma_X \sigma_Y} = \frac{\text{Cov}\{X, Y\}}{\sqrt{\text{Var}\{X\} \text{Var}\{Y\}}} \in [-1, 1] \quad (25)$$

- Caracteriza el grado de dependencia lineal o correlación entre dos variables aleatorias.
- Por ejemplo si $Y = aX + b$, entonces se cumple que:

$$r_{XY} = \begin{cases} 1 & \text{para } a > 0 \\ -1 & \text{para } a < 0 \end{cases} \quad (26)$$

Coeficiente de correlación de rango: ρ de Spearman

- Se define como:

$$\rho = 1 - \frac{6 \sum D^2}{N(N^2 - 1)} \quad (27)$$

- Para calcular ρ , las parejas de datos X y Y se ordenan y son reemplazados por su respectivo orden donde D es la diferencia $X-Y$ entre los estadísticos de orden y N es el número de parejas de datos.
- Oscila entre -1 y $+1$, indicándonos asociaciones negativas o positivas respectivamente, cero, significa no correlación pero no independencia.
- Es menos sensible a los valores atípicos que el coeficiente de Pearson.

Coeficiente de correlación de rango: τ de Kendall

- Se define como:

$$\tau = \frac{(\text{número de pares concordantes}) - (\text{número de pares discordantes})}{\binom{n}{2}} \quad (28)$$

- Un par es concordante si el orden de ambos está de acuerdo de lo contrario se dice que son discordantes.
- Si X y Y son independientes, entonces esperaríamos que el coeficiente sea aproximadamente cero.
- Es menos sensible a los valores atípicos que el coeficiente de Pearson.

Antes de transformar

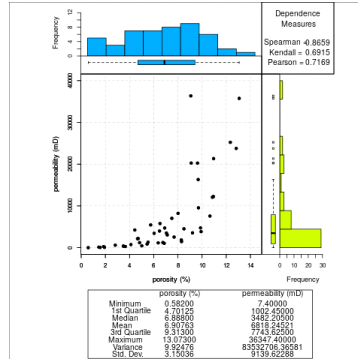


Figura 19: Coeficiente de correlación lineal = 0.71.

Después de transformar

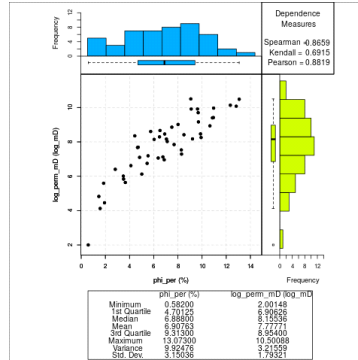


Figura 20: Coeficiente de correlación lineal = 0.88.

Estadística multivariada

Existen muchas técnicas multivariadas:

- Análisis de Regresión
- Análisis de Conglomerados
- Análisis de Componentes Principales
- Análisis Factorial
- Análisis Discriminante, etc

Regresión lineal y Mínimos cuadrados

- La **regresión** trata de establecer relaciones funcionales entre variables aleatorias.
- En particular la **regresión lineal** consiste en establecer una relación descrita mediante una recta.
- Los **modelos de regresión** nos permiten hacer predicciones o pronósticos a partir del modelo establecido.
- El método que se emplea para estimar los parámetros del modelo de regresión es el de los **Mínimos Cuadrados**

Regresión lineal I

Dados N valores de dos V.A. X y Y suponemos que:

- 1 X es una variable independiente
- 2 Y depende de X en forma lineal

Modelo lineal:

$$Y = \beta_0 + \beta_1 X \quad (29)$$

Donde

$$y_i = \beta_0 + \beta_1 x_i + e_i, \quad i = 1, \dots, N \quad (30)$$

- β_0, β_1 son los parámetros del modelo
- e_i son los errores o residuos del modelo

Mínimos Cuadrados Ordinarios (MCO)

- **Mínimos Cuadrados Ordinarios** consiste en hallar los parámetros del modelo de manera que la suma de los cuadrados de los residuos o errores sea mínima.

$$SCR = \sum_{i=1}^N e_i^2 = \sum_{i=1}^N [y_i - \hat{y}_i]^2 = \sum_{i=1}^N \left[y_i - (\hat{\beta}_0 + \hat{\beta}_1 x_i) \right]^2 \quad (31)$$

- Sistema de ecuaciones a resolver

$$\frac{\partial SCR}{\partial \beta_0} = 0, \frac{\partial SCR}{\partial \beta_1} = 0 \quad (32)$$

Coeficiente de determinación R^2

- El coeficiente de determinación se define como:

$$R^2 = 1 - \frac{SCR}{SCT} \quad (33)$$

donde $SCT = \sum_{i=1}^N [y_i - m_y^*]^2$ es la suma de cuadrados total (proporcional a la varianza de los datos)

- Para los modelos lineales

- 1 Mide el **grado de la bondad del ajuste**
- 2 Es igual al coeficiente de correlación lineal al cuadrado
- 3 Representa la proporción de varianza explicada por la regresión lineal.

Criterios de la bondad del ajuste

- Si $R^2 \approx 1$, el ajuste es bueno (Y se puede calcular de modo bastante aproximado a partir de X y viceversa).
- Si $R^2 \approx 0$, las variables X y Y no están relacionadas (linealmente al menos), por tanto no tiene sentido hacer un ajuste lineal.
- Sin embargo no es seguro que las dos variables no posean ninguna relación en el caso $r = 0$, ya que si bien el ajuste lineal puede no ser procedente, tal vez otro tipo de ajuste sí lo sea.

Regresión lineal

- Condiciones que deben cumplir los residuos

- 1 Valor esperado cero: $E\{e_i\} = 0$
- 2 Varianza constante: $Var\{e_i\} = \sigma_e^2$
- 3 No correlacionados: $Cov\{e_i, e_j\} = 0, \quad \forall i \neq j$
- 4 Distribución normal: $e \sim N(0, \sigma_e^2)$

Figura 21: Permeabilidad vs. porosidad antes de transformar.

Después de transformar

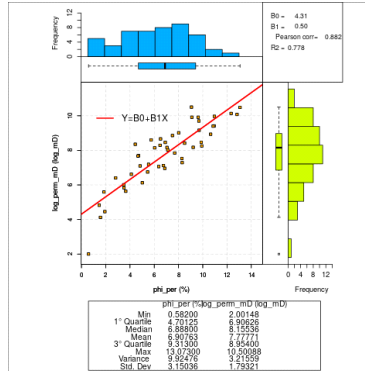


Figura 22: Permeabilidad vs. porosidad después de transformar.

Análisis de los residuos

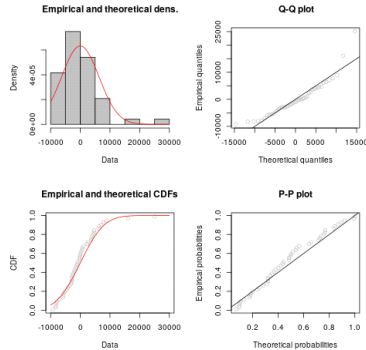


Figura 23: Residuos antes de transformar.

Análisis de los residuos

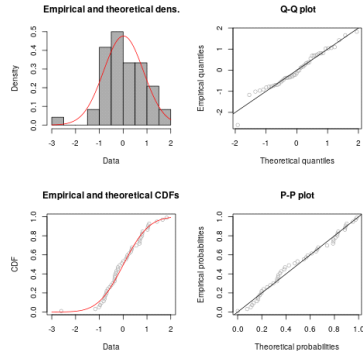


Figura 24: Residuos después de transformar.

Figure 1: Histogram of absolute residuals for the control group. The x-axis represents the absolute residuals of the log-transformed permittivity data (log mD). The y-axis represents the absolute frequency. The distribution is centered around zero, with a mean (Media) indicated by a red line and a median (Mediana) indicated by a blue dashed line. An inset box plot shows the distribution with whiskers and an outlier.

Estadígrafo	Valor
Muestras	48
Mínimo	-2.5995
1º cuartil	-0.5856
Mediana	-0.0955
Media	0.0139
3º cuartil	0.6961
Máximo	1.8249
Rango	4.4244
Rango intercuartil	1.2817
Varianza	0.7147
Desviación estándar	0.8454
Simetría	-0.1914
Curtosis	3.5273

A set of small navigation icons typically found in Beamer presentations, including symbols for back, forward, search, and other slide controls.

Análisis de los residuos

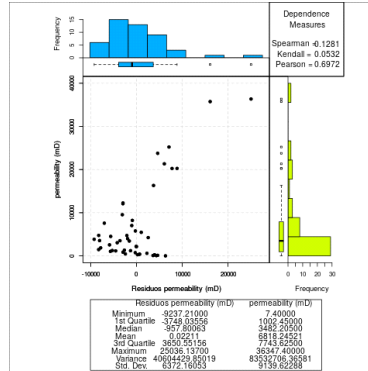


Figura 26: Correlación de la permeabilidad vs. los residuos antes de transformar.

Análisis de los residuos

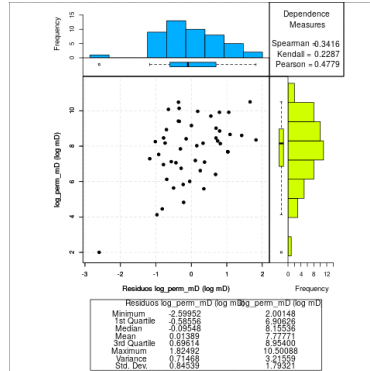


Figura 27: Correlación de la permeabilidad vs. los residuos después de transformar.

Agradecimiento especial

Al estudiante de doctorado M. en C. Daniel Vázquez Ramírez, por su desinteresado apoyo en la conversión de esta presentación del curso de Powerpoint a Latex con Beamer.

Siguiente tema: Funciones Aleatorias