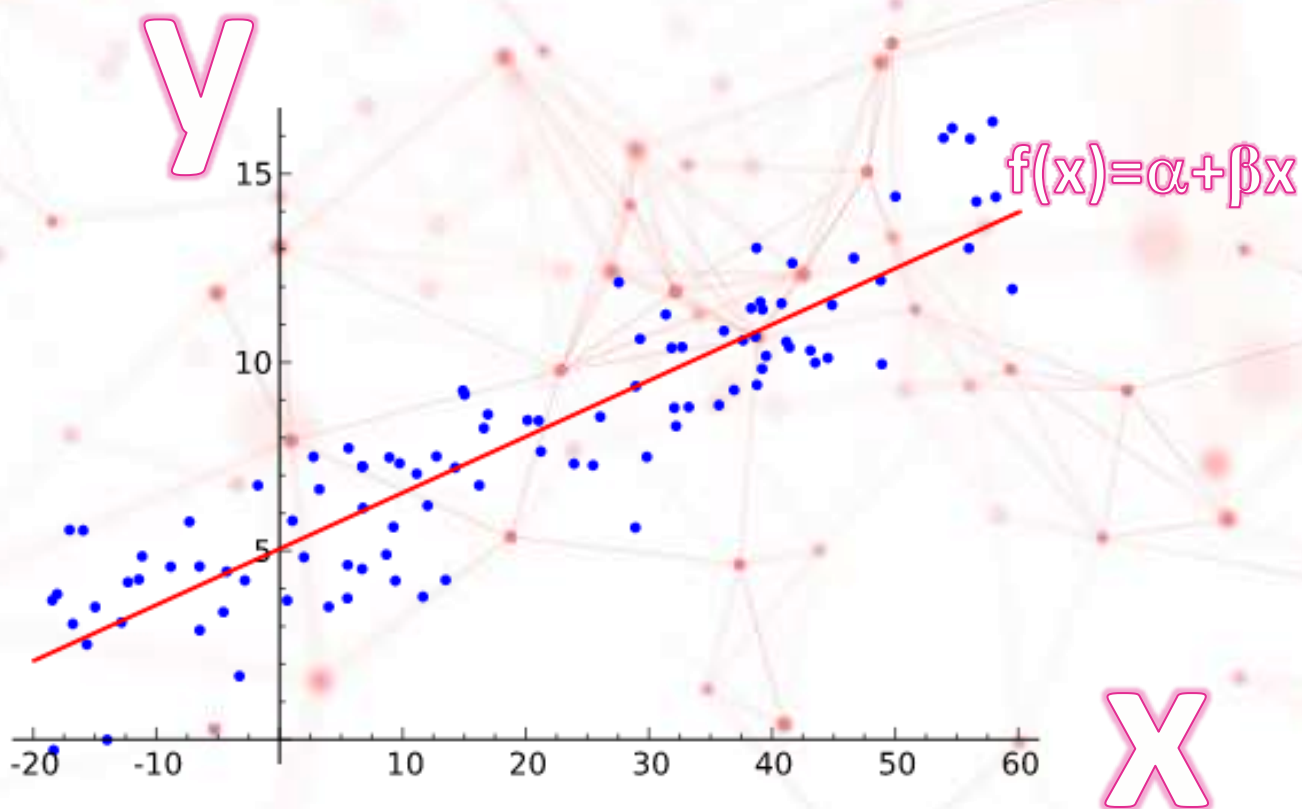


REGRESIÓN LINEAL SIMPLE

Aplicaciones con



Edgar E. CARPIO ■ Alcides RAMOS ■ Fredy H. VILLASANTE ■ Teresa P. ALVAREZ

Presentación

La Estadística ha jugado un papel primordial en el desarrollo de la investigación y la sociedad moderna, al proporcionar herramientas metodológicas generales para analizar la variabilidad, determinar relaciones entre variables, diseñar de forma óptima experimentos, realizar predicciones y la toma de decisiones en situaciones de incertidumbre. Una de las técnicas de este desarrollo es la regresión lineal que es una técnica de modelado estadístico que permite generar un modelo lineal en el que el valor de la variable dependiente o respuesta (Y continua) se determina a partir de un conjunto de variables independientes llamadas predictoras (X_1, X_2, X_3, \dots). Es una extensión de la regresión lineal simple, por lo que es fundamental comprender esta última. Los modelos de regresión múltiple pueden emplearse para predecir el valor de la variable dependiente o para evaluar la influencia que tienen los predictores sobre ella (esto último se debe analizar con cautela para no malinterpretar causa-efecto), entonces, puede ayudar a comprender y predecir el comportamiento de sistemas complejos o a analizar datos experimentales, financieros y biológicos y de otras áreas.

REGRESIÓN LINEAL SMPLE

Aplicaciones con R

Autores:

Edgar Eloy Carpio Vargas
Fredy Heric Villasante Saravia
Teresa Paola Alvarez Rozas
Alcides Ramos Calcina

Editado por:

Fredy Heric Villasante Saravia
Jr. Iquitos 181
Puno – Perú

Primera edición digital, marzo 2023

Hecho el depósito Legal en la Biblioteca Nacional del Perú

Registro N° 2023-02904

ISBN: 978-612-00-8501-1



Publicación electrónica en: <https://repositorio.unap.edu.pe/handle/20.500.14082/19667>

Contenido

1.	De los modelos	8
2.	Definición de modelo	8
3.	Proceso de formulación de un modelo.	8
	Primera etapa: especificación del modelo.....	8
	Segunda etapa: estimación del modelo	9
	Tercera etapa: evaluación o verificación.	11
	Cuarta etapa: utilización de los modelos.....	11
	Quinta etapa: modificación o adecuación del modelo	11
4.	Elementos que constituyen un modelo	12
	4.1 Ecuaciones o relaciones.	12
	4.2 Variables.....	13
	4.3 Constantes.....	15
5.	Clasificación de los modelos:.....	15
	a) De acuerdo a su construcción lógico empírico.	15
	b) De acuerdo con el dominio de la investigación:	16
	c) de acuerdo con sus fines o utilidad práctica.	17
6.	Tipos de modelos.....	18
	a) Modelos de series temporales.....	18
	b) Modelos de regresión uniecuacionales.	18
	c) Modelos de simulación multiecuacionales.	18
7.	Regresión y correlación	18
8.	Modelo lineal simple	19
9.	Relaciones entre variables.....	20
10.	Supuestos (hipótesis) fundamentales en los que se basa	22
	10.1 Hipótesis relativas a los perturbadores.....	22
	10.2 Hipótesis relativas a las variables.....	23
	10.3 Hipótesis relativas a los parámetros	24
11.	Inferencia.....	24
	11.1 Significancia e intervalo de confianza para β_0 y β_1	24
	11.2 Inferencia (Hipótesis) para β_1	25
	11.3 Inferencia (Hipótesis) para β_0	26
	11.4 Cálculo de la Varianza residual (varianza no explicada)	26
	11.5 Los intervalos de confianza en regresión	27
	11.6 Intervalos de confianza para β_0	27
	11.7 Intervalo de confianza para el parámetro β_1	28
12.	Residuos del modelo.....	29
13.	Análisis de varianza en R.L.S. (ANVA, ANOVA)	29

14.	Bondad de ajuste del modelo	33
15.	Condiciones para el uso de la regresión lineal.	35
16.	Predicción de valores.....	36
17.	Interpretación geométrica del método.....	40
18.	Ejemplo de aplicación manual.	42
	1) Diagrama de dispersión entre X e Y.	42
	2) Obtención de parámetros del modelo:	43
	3) Expresión de la línea de regresión.....	46
	4) Cálculo de las varianzas	47
	5) Prueba de hipótesis de los parámetros poblacionales	48
	6) Intervalo de confianza	49
	7) Tabla de Análisis de Varianza (ANVA).	49
	8) Coeficiente de Determinación.....	50
	9) Coeficiente de correlación.....	51
	10) prueba de supuestos.	51
19.	Ejemplo 1 en R.....	51
	Regresión lineal simple con variable independiente numérica.	51
	a) Observar las variables descriptivamente.....	52
	b) Representación gráfica de las observaciones.....	54
	c) Comportamiento del coeficiente de correlación.....	55
	d) Cálculo del modelo de regresión lineal simple	56
	e) Intervalos de confianza para los parámetros del modelo	58
	f) Representación gráfica del modelo.....	59
	g) Verificar los supuestos del modelo lineal	64
	h) Identificación de valores atípicos: <i>outliers</i> , <i>leverage</i> y <i>bservaciones</i> influyentes.71	
	i) Prediciendo nuevos valores.	75
20.	Ejemplo manual para mostrar los cálculos de las Bandas de Confianza:	76
21.	Ejemplo 2 en R.	81
	Regresión lineal simple con valor atípico.....	81
	a) Observar las variables descriptivamente.....	82
	b) Representación gráfica de las observaciones.....	83
	c) comportamiento del coeficiente de correlación.....	84
	d) Cálculo del modelo de regresión lineal simple	85
	e) Intervalos de confianza para los parámetros del modelo	87
	f) Representación gráfica del modelo.....	87
	g) Verificar condiciones para poder aceptar un modelo lineal.....	92
	h) Identificación de valores atípicos: <i>outliers</i> , <i>leverage</i> y <i>observaciones</i> influyentes. 100	
	i) Modelo final.....	109

22.	Regresión lineal con un predictor categórico de dos niveles.	125
	a) Diagrama de dispersión	127
	b) Estadísticos descriptivos para cada variable.	127
	c) Modelo lineal	130
	d) Residuales	135
23.	Apuntes varios (miscellaneous).....	139
	a) Origen del método de mínimos cuadrados y regresión	139
	b) Significado de modelo lineal.....	139
	c) Ventajas del método de mínimos cuadrados para estimar los coeficientes de un modelo lineal	140
	d) Ecuaciones de curvas de aproximación:	140
	e) Método de los mínimos cuadrados	141
	f) Propiedades de la regresión lineal	149
24.	Bibliografía	152

1. De los modelos

La regresión se considera como una ciencia empírica, ya que la formulación de sus modelos se inicia con la teoría y observación del comportamiento empírico de los sujetos que en ella intervienen.

2. Definición de modelo

Según José Luis San Pedro, modelo es una representación simplificada en símbolos matemáticos de cierto conjunto de relaciones.

Corrado Gini "un modelo es una representación simplificada de la forma en que ciertos fenómenos están constituidos y/o de manera que se desenvuelven".

Edmund Malinvaud "un modelo es una representación formal de ideas y conocimientos relativos a un fenómeno".

Enders A. Robinson "Un modelo es una abstracción simplificada e idealizada cuyo objetivo es representar en forma aproximada el comportamiento de un sistema"

André Regnier "Modelo es un objeto concreto un objeto abstracto cuya descripción es considerada como una descripción de dicho objeto concreto".

Características mínimas que debe reunir todo modelo:

- a) Que represente un fenómeno real
- b) Que la representación sea simplificada
- c) Que se haga en términos matemáticos

3. Proceso de formulación de un modelo.

Primera etapa: especificación del modelo

- a) Conocimiento de la Teoría.
- b) Construcción del modelo. expresar en términos matemáticos la teoría objeto de estudio.
- c) Análisis de las especificaciones del modelo:

- Enumeración de las variables consideradas relevantes, su clasificación en endógenas, exógenas etc.
 - Análisis del tipo de relación existente entre las variables lo que nos permitirá determinar si el modelo es interdependiente o recursivo.
 - Indicación del tipo de observación a utilizar, o sea si se trata de series de tiempo, o atemporales, a fin de establecer si el modelo será estático (histórico o ahistórico) o dinámico.
 - Especificación de las relaciones existentes (comportamiento, institucional, tecnológicas, de definición o de equilibrio móvil) con Indicación de las variables que entran en cada relación (ecuación) y de su forma funcional.
- d) Identificación de los parámetros estructurales. Este Análisis se realiza en los modelos multiecuacionales ya que no existe problema de identificación en los modelos uniecuacionales. Consiste en determinar cuál ecuación corresponde a los parámetros por estimar. Usualmente en forma matricial reducida.

Segunda etapa: estimación del modelo

- a) Recolección de información estadística de las variables incluidas en el modelo, Estos datos pueden ser de varios tipos:

Tipos de información (datos)

- Series de tiempo. son aquellos que se almacenan durante un determinado período (Ruiz, 2020)
- Series de corte transversal. Se realiza su recolección en algún momento del tiempo
- Datos combinados. son provenientes de datos de series de tiempo y de corte transversal (sección cruzada), (Gujarati, 2010).
- Datos de ingeniería.
- Regulaciones institucionales y otras legislaciones (ejemplo, tasas de impuestos)
- Datos contruidos: variables dummy (para factores cualitativos)

- b) Examen del grado de correlación entre las variables. hay que evitar que la variable explicada reciba múltiple influencia a través de una variable explicativa.

Si la colinealidad es alta, los resultados pueden ser seriamente distorsionados. Estamos en el problema de multicolinealidad (dependencia lineal).

- c) Estimación del modelo. (Estimación de parámetros estructurales) Se refiere a la selección de la técnica de estimación más apropiada dentro de las diversas técnicas.

Métodos Para Estimar Parámetros

Se pueden clasificar en dos grupos:

Técnicas de ecuación simple:

- Momentos
- Mínimos cuadrados ordinarios (MCO)
- Mínimos Cuadrados Indirectos (MCI)
- Máxima verosimilitud de información Limitada (MVIL)
- Varios Métodos de estimación mixta.
- De la mínima X^2
- Bayesianos, etc.

Técnicas de ecuaciones simultaneas:

- Mínimos Cuadrados Bietapicos (MCB)
- Mínimos Cuadrados Trietápicos (MCT)
- Máxima verosimilitud de información completa (MVIC)

- d) Examen de las propiedades de los estimadores obtenidos de cada técnica. Un buen estimador posee las siguientes propiedades:

- Insesgamiento
- Consistencia
- Eficiencia

- Suficiencia
- o alguna combinación de las propiedades anteriores.

Tercera etapa: evaluación o verificación.

a) verificación de hipótesis estadísticas.

Pruebas de primer orden

- Pruebas "t" de student para establecer la significación estadística de cada variable explicativa individualmente consideradas.
- Cálculo del coeficiente de determinación (R^2) para establecer como están explicando en su conjunto las variables independientes. y cálculo del coeficiente de correlación (r). La prueba F de confianza del modelo.

Pruebas de segundo orden

- Análisis de la varianza residual a fin de establecer si la variancia de la variable aleatoria es constante (homocedasticidad) o es variable (heteroscedasticidad).
- Prueba de hipótesis para establecer si existe o no violación del supuesto de autocorrelación en los perturbadores, basada en la prueba de Durbin-Watson o la razón de Von-Newman.

Cuarta etapa: utilización de los modelos

a) Uso de los modelos (Cybertesis, 2020)

- Proyectar la variable dependiente (modelos predictivos)
- Tomar decisiones de políticas existentes y probar consistencia de planes de desarrollo (modelos de decisión)
- Describir una realidad (modelos descriptivos)
- Establecer causas relevantes que han originado una determinada realidad (modelos explicativos).

Quinta etapa: modificación o adecuación del modelo

Debe revisarse frecuentemente y modificarse según los cambios de la realidad.

4. Elementos que constituyen un modelo

4.1 Ecuaciones o relaciones.

Un modelo que queda expresado mediante una ecuación se denomina modelo uniecuacional, y si un modelo queda especificado por varias ecuaciones se denomina modelos multiecuacionales.

Según Morocho, (2015) su contenido empírico las ecuaciones de un modelo se clasifican en:

- **Ecuaciones de comportamiento.** explica el modo de actuar de los sujetos (conducta) (individuo, familia, empresa estado).
- **Ecuaciones Institucionales o legales.** reflejan los efectos que producen en un modelo, la existencia de leyes. Expresa una relación entre una variable económica y una de tipo político o jurídico. Ej. impuesto, ecuación de la masa monetaria en función de la tasa de interés.
- **Ecuaciones tecnológicas.** Explican los modelos de producción incorporados a la actividad económica. Ej. producción (modelo Cobb-Duglas)
- **Ecuaciones de definición e identidad.** son relaciones que se verifican, ya sea por su construcción lógica o por definición contable. ejem: demanda.
- **Ecuación de equilibrio móvil.** Son aquellas igualdades que resultan de una condición impuesta o postulado introducido. Entidades de tipo contable.

Solo las tres primeras clases son resultado de axiomas empíricos comprobables. De la observación empírica se obtendrá: (1) variables relevantes (2) permanencia o regularidad de las variables (3) sus relaciones de causalidad.

4.2 Variables.

Formado por un universo en donde podemos distinguir, cosas, entes, personas, etc. los que tienen determinadas rasgos característicos o propiedades.

Clasificación de las variables

a. En los modelos estructurales.

- Variables endógenas
- Variables predeterminadas
- Variables aleatorias o estocásticas
- Variables expectativas

Variables endógenas. Son aquellos cuyos valores estimados, van a ser determinados por las soluciones particulares del sistema de ecuaciones. Son las llamadas variables *dependientes* en el análisis matemático.

Variable predeterminada. Son el conjunto de variables exógenas y variables con rezagos. Son aquellas cuyos valores no se obtienen por la solución del modelo, sino que provienen de afuera del mismo. Ellas contribuyen a explicar el comportamiento de las variables endógenas sin ser explicadas por el modelo mismo. Comprende dos categorías (Cybertesis, 2020)

(1) **Variable exógena.** que incluye variables propias de otros sistemas distintos del económico. (Sistema Físico, político, social, biológico) ejemplo: lluvia, tasa de interés.

(2) **variable endógena con retardo (variables retardadas, históricas).** estas variables intervienen como variables explicativas. Expresan un valor pasado Ejemplo, ingreso anterior.

Variables aleatorias o estocásticas (perturbaciones). (Morocho, 2015)son variables "no observables" que caracterizan a los modelos estocásticos o probabilísticos, por oposición a los modelos deterministas. Estas variables cumplen con recoger el conjunto de

"causas" que no se encuentran explícitamente incorporadas en un modelo; tales como:

- (1) **Omisión de variables explicativas.** intervienen las variables más relevantes
- (2) **Errores de especificación.** Suponiendo incluidas las variables relevantes, la variable aleatoria recoge los efectos especificación incorrecta sobre la ley matemática.
- (3) **Errores de medida sobre las variables endógenas.** Se considera que dichos errores son aleatorios y se los incorpora en la variable estocástica de cada ecuación del modelo, Se supone que las variables exógenas están medidas sin error "modelos con errores en las variables".

Variables expectativas (esperadas). Son variables no observables cuya introducción exige al enunciado de un post lado adicional en el que se especifica su comportamiento en función de variables observables. Ejemplo, precio normal esperado, ingreso normal esperado (Econometricos, 2010).

b. En los modelos de decisión.

Estos modelos nos sirven en la programación del crecimiento y desarrollo de un sector, región o nación. Ejemplo, tasa de empleo, nivel de empleo, tasa mínima de analfabetismo producto nacional, etc. Las variables en estos modelos se clasifican en:

$$\begin{array}{l} \text{Endogenas} \left\{ \begin{array}{l} \text{objetivas} \\ \text{no objetivas} \end{array} \right. \\ \text{exogenas} \left\{ \begin{array}{l} \text{controlables} \left\{ \begin{array}{l} \text{instrumentales} \\ \text{no instrumentales} \end{array} \right. \\ \text{no controlables} \end{array} \right. \end{array}$$

Variables endógenas objetivas. A estas variables se les fija un nivel por alcanzar o un comportamiento temporal. ejemplo: ingreso nacional, distribución del ingreso, nivel de ocupación, volumen de la demanda total etc.

VARIABLES ENDÓGENAS CONTROLABLES. Son variables exógenas sobre las cuales pueden actuar directamente los sujetos de las decisiones. ejemplo: Impuestos.

4.3 Constantes.

Son valores numéricos que intervienen en un modelo. Las constantes tienen un triple significado: matemático, estadístico y económico. Cuando la constante se obtiene de una población se denomina parámetro (β_0, β_1 , etc.) y si se obtiene de una muestra se denomina estimador y se representa por $\hat{\beta}_0, \hat{\beta}_1$, etc.

5. Clasificación de los modelos:

a) De acuerdo a su construcción lógico empírico (Tareas, 2013).

Lineales. Son ecuaciones algebraicas de primer grado

No lineales. Son aquellas en las cuales por lo menos una de sus ecuaciones no es algebraica entera de primer grado en sus variables.

Deterministas. Supone la existencia de variables que exactamente satisfacen las ecuaciones que conforman el modelo. Estas no son aleatorias ni estocásticas.

Estocásticas o aleatorias. Son aquellas que incorporan las variables aleatorias o residuales, el caso límite de estos modelos lo constituyen los modelos deterministas cuando la probabilidad de presentación de valores estimados de variables endógenas condicionadas por los valores de las variables predeterminadas.

Completos. Cuando el número de variables endógenas del modelo es igual número de ecuaciones, entonces puede determinarse trabajando con el álgebra matricial y si admite una "única solución" se dice que el modelo es completo. Para que un modelo sea completo es condición necesaria que " el número de ecuaciones sea igual al número de variables endógenas (incógnitas)."

Condición necesaria y suficiente es que la matriz de coeficientes sea no singular.

Incompletos. Son aquellos que tienen infinitas soluciones.

Condición necesaria: que el número de ecuaciones sea menor que el número de variables endógenas.

Condición necesaria y suficiente la matriz de coeficiente sea singular.

b) De acuerdo con el dominio de la investigación (Cybertesis, 2020):

Uniecuacionales. Son los que tratan de describir y/o explicar una actividad mediante una relación funcional entre la variable endógena o explicada por el modelo y las variables predeterminadas o explicativas.

Multiecuacionales. Son los tratan de describir y/o explicar las interrelaciones de una o varias actividades, mediante ecuaciones estructurales.

Macroeconómicos. Se caracterizan por utilizar variables agregadas, como el consumo nacional, ingreso nacional.

Microeconómicos. emplean variables directamente observables, es decir variables tipo individual, ejemplo demanda, oferta etc.

Dinámicos. Cuando su comportamiento temporal está determinado por ecuaciones funcionales cuyas variables están referidas a distintos momentos del tiempo en una forma esencial. variables con retardo (Martinez, 1992).

Estáticos. Son aquellas cuyas ecuaciones no contienen variables referidas al mismo momento del tiempo (estático-históricas) o sin referencia temporal alguna (Estático ahistóricas). Estos modelos se obtienen usando datos de corte transversal. A Los modelos que incluyen variable tiempo se les denomina MODELOS DINAMICOS, las cuales se estiman mediante series cronológicas.

Interdependientes. son aquellos modelos que no incluyen subsistemas completamente contenidos (independientes). Cuando la matriz no es triangular.

Recursivos o de cadenas causales. Son aquellos en los que la matriz de coeficientes de las variables endógenas es triangular. En estos modelos existe una causalidad en cadena inter-ecuacional.

Particionables. Se caracterizan por tener una matriz de coeficientes que puede transformarse en una matriz triangular o diagonal por bloques. Entonces cada subsistema es completamente independiente.

Abiertos. Son aquellos modelos en los que intervienen el comercio exterior, generalmente macroeconómicos (Cybertesis, 2020).

Cerrados. no interviene el comercio exterior.

c) de acuerdo con sus fines o utilidad práctica.

Descriptivos. Solamente buscan la realidad, responden a una pregunta, ¿QUE HA PASADO?, existen dos tipos:

De tendencia. que relacionan las variables endógenas con el tiempo, la cual no es explicativa en sí misma. Ejemplo, los conjuntos de indicadores económicos, referidos a un problema específico

Explicativos. No se limitan a describir solamente una realidad, sino tratan de descubrir las causas relevantes, responden a la pregunta, ¿POR QUE HA PASADO?

Predictivos. responden a la pregunta del tipo, ¿QUE PASARÁ?

De decisión. se dividen en:

Los deterministas. Utilizan técnicas de solución bastante avanzadas.

Los estocásticos. donde la incertidumbre juega un papel primordial.

Terminología y notación

- | | |
|------------------------|--------------------------------|
| - Variable dependiente | variable independiente |
| - Variable explicada | Variable explicativa |
| - Variable predicha | predictor |
| - Variable regresada | regresor |
| - Variable respuesta | Variable de control o estímulo |
| - Variable endógena | Variable exógena. |
| - Efecto | Causa |

6. Tipos de modelos

a) Modelos de series temporales.

En esta clase de modelos se supone que no se sabe nada sobre las relaciones causales que en realidad afectan a la variable que se trata de predecir. En cambio, examinamos el comportamiento de una serie temporal en el pasado para inferir cuál será su comportamiento en el futuro. El método de las series temporales para la obtención de una predicción puede implicar la utilización de un modelo determinístico simple, como la extrapolación lineal; o la utilización de un modelo estocástico complejo para la predicción adaptativa.

b) Modelos de regresión uniecuacionales.

En esta clase de modelos la variable objeto de estudio se explica mediante una única función (lineal o no lineal) de otras variables explicativas. La ecuación será en otros casos dependiente del tiempo, de forma que será posible predecir la respuesta de la variable.

c) Modelos de simulación multiecuacionales.

En esta clase de modelos, la variable objeto de estudio puede ser una función de varias variables explicativas, pero dichas variables se relacionan entre sí, así como la variable de estudio en un conjunto de ecuaciones.

7. Regresión y correlación

El análisis de la regresión está estrechamente ligado con el análisis de la correlación, pero conceptualmente las dos definiciones son diferentes. En el análisis de la correlación, el objetivo fundamental es la medición de la fuerza o grado de asociación lineal o covariabilidad entre dos variables (aquí se dice que las dos variables son aleatorias), mientras que los métodos de regresión se usan para determinar la mejor relación funcional entre variables, con este análisis de regresión también se intenta predecir o estimar el valor promedio de una variable en base a otros valores fijos de otras variables (grado de dependencia de y en términos de x).

En el análisis de regresión existe una asimetría en la manera como se manejan las variables dependientes y explicativas. Mientras que en el análisis de correlación se manejan las dos variables simétricamente (no existe distinción entre las variables dependientes y explicativas).

8. Modelo lineal simple

(Calcina, 2019) El modelo lineal simple (mls) es un modelo uniecuacional por que consta de una relación lineal entre dos variables, donde la variable Y es la variable dependiente o explicada por el modelo, y la X es la variable independiente o explicativa, siendo su representación básica una línea recta.

Función Lineal. $Y = \beta_0 + \beta_1 X$

β_0 : ordenada en el origen (valor de Y cuando X = 0)

β_1 : pendiente (cambio de Y al aumentar X en 1)

Modelo de regresión lineal simple: Una relación lineal se expresa como: $E(Y/X) = \beta_0 + \beta_1 X$, donde la Y observada puede diferir aleatoriamente de $E(Y/X)$ y esta diferencia también es denotada por u_i .

$$Y = \beta_0 + \beta_1 X_i + \mu_i \quad \text{o} \quad Y = E(Y/X) + \mu_i \quad \text{y} \quad \mu_i = Y - E(Y/X)$$

Modelo matemático o teórico. es la expresión matemática de una determinada teoría. $Y = \beta_0 + \beta_1 X_i$, donde: β_0 y β_1 son parámetros a encontrar.

Modelo estadístico. Es la expresión matemática de una determinada teoría que incluye en su expresión el término de perturbación o error.

$$Y = \beta_0 + \beta_1 X_i + \mu_i$$

En este modelo no se consideran los errores de observación, pero si se considera el error en la especificación del modelo.

$\mu_i =$ **término de perturbación o error.** es una variable aleatoria (estocástica) con propiedades probabilísticas bien definidas. Es aquí donde están involucradas todos los errores que no han sido consideradas en el modelo, que pueden provenir de (Trujillo, 2017):

- De variables explicativas relevantes que no se consideraron en el modelo.

- De errores de especificación en la relación de correspondencia entre las variables (se considera una relación lineal, cuando en realidad es una función no lineal) (Morocho, 2015)
- De errores de medida sobre las variables endógenas, considerándose dichos errores como aleatorios y se les incorpora en la variable estocástica de cada ecuación de un modelo.

Las relaciones exactas como la indicada no se dan a menudo en la práctica, por cuanto si muestreamos un par de valores (X, Y) no se distribuirán exactamente a lo largo de la recta, ello se debe a que el modelo es teórico y no es más que una simplificación de la realidad. En estadística al realizar tipos de muestreo, expresar un modelo etc. se cometen ciertos errores o perturbadores aleatorios que, para tapar esta brecha, entre la realidad y la práctica, en nuestro modelo introducimos la "u". Entonces el modelo quedara de la siguiente manera:

Si el estudio se está realizando con muestras el modelo quedara especificado como:

$$y = \hat{\beta}_0 + \hat{\beta}_1 X_i + e_i \quad (\text{modelo estadístico muestral})$$

9. Relaciones entre variables

- a) Una variable X puede influir en otra variable Y, esto es $X \rightarrow Y$.

Ejemplo: La edad influye sobre la actividad mental del niño, la harina influye en el volumen del pan, la lluvia influye en la cosecha, el peso de un animal vivo influye en el peso de la carcasa, la temperatura influye en la intensidad del ataque de los insectos, etc.

- b) Dos variables pueden estar influenciadas entre sí; esto es $X \leftrightarrow Y$.

Ejemplo: Precio y producción de un artículo, peso y volumen del pan, peso y altura de los individuos, nubosidad y horas del sol, uso de tabaco y afecciones cardiacas, etc.

- c) Dos variables sin estar influenciadas, pueden estar relacionadas entre sí (concomitantes), por estar ambas influenciadas por una tercera variable.

Ejemplo: El peso de los hermanos y el peso de las hermanas, el peso del pan y el precio de las papas (influencia del aumento de costo de vida), las notas

de química y de bioquímica relacionadas por la afición de los alumnos a los cursos de ciencias. (correlación).

En resumen, **la regresión lineal** simple consiste en generar un modelo de regresión (ecuación de una recta) que permita explicar la relación lineal que existe entre dos variables. A la variable dependiente o respuesta se le identifica como Y y a la variable predictora o independiente como X (Amat, 2016).

El modelo de regresión lineal simple se describe de acuerdo a la ecuación:

$$Y = \beta_0 + \beta_1 X_1 + u_i$$

Modelo muestral:

$$y = \hat{\beta}_0 + \hat{\beta}_1 X_i + e_i$$

siendo $\hat{\beta}_0$ la ordenada en el origen, $\hat{\beta}_1$ la pendiente y e_i el error aleatorio. Este último representa la diferencia entre el valor ajustado por la recta y el valor real. Recoge el efecto de todas aquellas variables que influyen en Y pero que no se incluyen en el modelo como predictores. Al error aleatorio también se le conoce como residuo (Amat, 2016).

En la gran mayoría de casos, los valores β_0 y β_1 poblacionales son desconocidos, por lo que, a partir de una muestra, se obtienen sus estimaciones $\hat{\beta}_0$ y $\hat{\beta}_1$. Estas estimaciones se conocen como coeficientes de regresión o *least square coefficient estimates*, ya que toman aquellos valores que minimizan la suma de cuadrados residuales, dando lugar a la recta que pasa más cerca de todos los puntos. (Existen alternativas al método de mínimos cuadrados para obtener las estimaciones de los coeficientes).

(existen diferentes ecuaciones para calcular $\hat{\beta}_0$ y $\hat{\beta}_1$, veamos algunas.

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

$$\hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \bar{X}$$

$\hat{\beta}_0$ es el valor esperado de la variable Y cuando $X = 0$, es decir, la intersección de la recta con el eje Y . Es un dato necesario para generar la recta, pero, en ocasiones, no tiene interpretación práctica (situaciones en las que X no puede adquirir el valor 0).

Una recta de regresión puede emplearse para diferentes propósitos y dependiendo de ellos es necesario satisfacer distintas condiciones. En caso de querer medir la relación lineal entre dos variables, la recta de regresión lo va a indicar de forma directa (ya que calcula la correlación). Sin embargo, en caso de querer predecir el valor de una variable en función de la otra, no solo se necesita calcular la recta, sino que además hay que asegurar que el modelo sea bueno (Amat, 2016).

10. Supuestos (hipótesis) fundamentales en los que se basa

La estimación de parámetros requiere de un conjunto de hipótesis que deben de cumplirse en un modelo de regresión lineal simple.

10.1 Hipótesis relativas a los perturbadores

Los perturbadores constituyen un conjunto de variables individualmente poco relevantes, y si nos basamos en el supuesto de que estas variables son independientes entre sí, la perturbación tendrá en virtud del teorema del límite central una distribución aproximadamente normal; además, es lógico suponer que estas variables irrelevantes actúan en dirección positiva o negativa y por lo tanto su media será cero; finalmente parece verosímil que la distribución muestral de las perturbaciones tendrá una dispersión constante (varianza constante) e independiente del valor de X (homocedasticidad) de lo expresado se desprende lo siguiente:

- a) El valor esperado del término de error es igual a cero; es decir:

$$E(\mu_i) = 0: i = 1, 2, \dots, n$$

El hecho de aceptar esta hipótesis significa que las perturbaciones van a tener valores positivos y negativos.

- b) Todas las perturbaciones aleatorias tienen varianza constante, es decir $E(\mu_i^2) = \sigma^2; 1, 2, \dots, n$

Esta propiedad se conoce como **homocedasticidad**, o sea que la variable aleatoria error tiene varianza finita, constante e independiente de X_i .

- c) Las perturbaciones (μ_i) son independientes entre sí (incorrelacionadas en sentido estadístico), es decir la correlación de los errores correspondientes a observaciones distintas es igual a cero: $E(\mu_i\mu_j) = cov(\mu_i\mu_j) = 0$ para $i \neq j$

Esta propiedad se conoce como no **autocorrelación** (no están autocorrelacionadas), Los valores que asume μ_i para cada i son completamente independientes de todos los valores precedentes.

- d) El termino de error sigue una **distribución normal** con media cero y varianza σ^2 . $\mu_i = N(0, \sigma^2)$. Si e_i se distribuye normalmente, entonces Y_i también se distribuyen normalmente, esto se deduce puesto que Y es una combinación lineal de los errores los cuales son todos normales.
- e) Asunción adicional. μ_i, μ_j no son solamente correlacionadas sino necesariamente independientes.

10.2 Hipótesis relativas a las variables.

El modelo en estudio es de dos variables. La estimación suele hacerse suponiendo que X es una variable fija, es decir no aleatoria e independiente del muestreo; En cuanto a Y es una variable aleatoria que puede descomponerse en dos partes:

- a) La variable X no es aleatoria, es decir es fija, los valores vienen fijados.
- b) La relación entre X e Y es una relación lineal.
- c) La variable endógena Y es evidentemente una variable aleatoria, cuyos parámetros serán:

$$Y = \hat{\beta}_0 + \hat{\beta}_1 x + e_i$$

$$E(Y) = E(\hat{\beta}_0 + \hat{\beta}_1 x + e_i) \quad \text{tomando valor esperado}$$

$$E(Y) = E(\hat{\beta}_0) + E(\hat{\beta}_1 X_i) + E(e_i) \quad \text{empleando propiedades y conociendo que } E(e_i) = 0$$

$$E(Y) = E(\hat{\beta}_0) + \hat{\beta}_1 E(X_i) \quad \text{Xi constante}$$

$$E(Y) = \hat{\beta}_0 + \hat{\beta}_1 X_i \quad \text{es la media}$$

Por otro lado, por definición de varianza tenemos:

$$\sigma_y^2 = E[Y_i - E[Y_i]]^2 = E[(\hat{\beta}_0 + \hat{\beta}_1 X_i + e_i)]^2$$

$$E[(\hat{\beta}_0 + \hat{\beta}_1 X_i + e_i - \hat{\beta}_0 - \hat{\beta}_1 X_i)]^2$$

$$E[(e_i)]^2 = \sigma^2 \quad \text{es la varianza.}$$

d) Las variables Y e X se obtienen sin errores de observación.

10.3 Hipótesis relativas a los parámetros

El valor de los parámetros especifica una determinada estructura entre las variables, que se supone constante para todos los elementos de la población y de la muestra.

Los parámetros estructurales son constantes para todas las unidades de la muestra y no existe ninguna restricción para ellos.

11. Inferencia

11.1 Significancia e intervalo de confianza para β_0 y β_1

En la mayoría de casos, aunque el estudio de regresión se aplica a una muestra, el objetivo último es obtener un modelo lineal que explique la relación entre las dos variables en toda la población. Esto significa que el modelo generado es una estimación de la relación poblacional a partir de la relación que se observa en la muestra y, por lo tanto, está sujeta a variaciones. Para cada uno de los parámetros de la ecuación de regresión lineal simple (β_0 y β_1) se puede calcular su significancia (p-value) y su intervalo de confianza. El test estadístico más empleado es el t-test (existen alternativas no paramétricas) (Amat, 2016).

11.2 Inferencia (Hipótesis) para β_1 .

- 1) Hipótesis estadística.

H_0 : No hay relación lineal entre ambas variables por lo que la pendiente del modelo lineal es cero. Y es independiente de X, no hay coherencia, la pendiente es cero, X no influye en Y. $\beta_1 = 0$.

H_a : Hay relación lineal entre ambas variables por lo que la pendiente del modelo lineal es distinta de cero. X influye en Y, existe relación lineal, Y es dependiente de X. $\beta_1 \neq 0$.

- 2) nivel de significancia = α

- 3) prueba estadística:

Cálculo del estadístico y del p-value:

$$z = \frac{\hat{\beta}_1 - \beta}{S_{\hat{\beta}_1}}; \text{ muestra grande}$$

$$t = \frac{\hat{\beta}_1 - \beta}{S_{\hat{\beta}_1}}; \text{ muestra pequeña}$$

Varianza del estimador:

$$S_{\hat{\beta}_1}^2 = \frac{\hat{\beta}_1 - \beta \sqrt{\sum x_i^2}}{S_e} = \frac{S_e^2}{\sum x_i^2} = S_e^2 \left[\frac{1}{n} + \frac{\bar{x}^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \right]$$

Desviación estándar del estimador:

$$S_{\hat{\beta}_1} = \sqrt{S_{\hat{\beta}_1}^2}$$

- 4) región crítica, decisión:

Utilizando valor puntual, se rechaza H_0 si:

Muestra grande $Z > Z_{\alpha/2}$ y muestra pequeña: $t > t_{(n-2), \alpha/2}$

Utilizando probabilidades: $p < \alpha$

11.3 Inferencia (Hipótesis) para β_0 .

1) Hipótesis estadística.

$$H_0: \beta_0 = 0$$

$$H_a: \beta_0 \neq 0$$

2) nivel de significancia = α

3) prueba estadística:

$$Z = \frac{\hat{\beta}_0 - 0}{S_{\hat{\beta}_0}} ; \text{muestra grande}$$

$$t = \frac{\hat{\beta}_0 - 0}{S_{\hat{\beta}_0}} ; \text{muestra pequeña}$$

Varianza del estimador:

$$S_{\hat{\beta}_0}^2 = \frac{S_e^2 \sum X_i^2}{n \sum x_i^2} = S_e^2 \left[\frac{1}{n} + \frac{\bar{X}^2}{\sum x_i^2} \right]$$

$$S_{\hat{\beta}_0} = \sqrt{S_{\hat{\beta}_0}^2}$$

4) región crítica, decisión:

Utilizando valor puntual, se rechaza H_0 si:

Muestra grande $Z > Z_{\alpha/2}$ y muestra pequeña: $t > t_{(n-2), \alpha/2}$

Utilizando probabilidades: $p < \alpha$

11.4 Cálculo de la Varianza residual (varianza no explicada)

La varianza residual σ^2 es desconocida, siendo su estimador insesgado, entonces:

$$MSE = S_e^2 = \frac{\sum (Y_i - \hat{\beta}_0 X_i - \hat{\beta}_1)^2}{n - 2}$$

$$S_e^2 = \frac{\sum e_i^2}{n - 2} = \frac{\sum Y_i^2 - \hat{\beta}_0 \sum Y_i - \hat{\beta}_1 \sum X_i Y_i}{n - 2} = \frac{(n - 1)(S_y^2 - \hat{\beta}_1^2 S_x^2)}{n - 2} = \frac{\sum y_i^2 - \hat{\beta}_1^2 \sum x_i^2}{n - 2}$$

La varianza del error σ^2 se estima a partir del Residual Standard Error (RSE), que puede entenderse como la diferencia promedio que se desvía la

variable respuesta de la verdadera línea de regresión. En el caso de regresión lineal simple, RSE equivale a (Amat, 2016):

$$RSE = \sqrt{\frac{1}{n-2} RSS} = \sqrt{\frac{1}{n-2} \sum_{i=1}^n (y_i - \hat{y}_i)^2}$$

Grados de libertad (df) = número observaciones - 2 = número observaciones - número predictores - 1, suma del cuadrado de cada residuo (RSS).

11.5 Los intervalos de confianza en regresión

Los intervalos de confianza se utilizan para evaluar la precisión de las estimaciones de los parámetros de la regresión y para hacer inferencias sobre la relación entre las variables independientes y la variable dependiente. Los intervalos de confianza son un rango de valores dentro del cual se espera que se encuentre el verdadero valor del parámetro con un cierto nivel de confianza.

En la regresión lineal simple, los intervalos de confianza se pueden calcular para los coeficientes de la regresión (intercepto y pendiente) y para la predicción de valores individuales de la variable dependiente.

En resumen, los intervalos de confianza en regresión son una herramienta útil para evaluar la precisión de las estimaciones de los parámetros de la regresión y para hacer inferencias sobre la relación entre las variables independientes y la variable dependiente.

11.6 Intervalos de confianza para β_0

$$p(-z_{\alpha/2} < Z < z_{\alpha/2}) = 1 - \alpha$$

$$p\left(-z_{\alpha/2} < \frac{\hat{\beta}_0 - \beta_0 \sqrt{n \sum x_i^2}}{\sigma_e \sqrt{\sum X_i^2}} < z_{\alpha/2}\right) = 1 - \alpha$$

$$p\left(\hat{\beta}_0 - z_{\alpha/2} \frac{\sigma_e \sqrt{\sum X_i^2}}{\sqrt{n \sum x_i^2}} < \beta_0 < \hat{\beta}_0 + z_{\alpha/2} \frac{\sigma_e \sqrt{\sum X_i^2}}{\sqrt{n \sum x_i^2}}\right) = 1 - \alpha$$

Muestra grande:

$$IC = p(\hat{\beta}_0 - z_{\alpha/2}S_{\hat{\beta}_0} < \beta_0 < \hat{\beta}_0 + z_{\alpha/2}S_{\hat{\beta}_0}) = 1 - \alpha$$

Muestra pequeña

$$IC = p(\hat{\beta}_0 - t_{(n-2),\alpha/2}S_{\hat{\beta}_0} < \beta_0 < \hat{\beta}_0 + t_{(n-2),\alpha/2}S_{\hat{\beta}_0}) = 1 - \alpha$$

11.7 Intervalo de confianza para el parámetro β_1

Muestra grande:

$$p(\hat{\beta}_1 - z_{\alpha/2}S_{\hat{\beta}_1} < \beta_1 < \hat{\beta}_1 + z_{\alpha/2}S_{\hat{\beta}_1}) = 1 - \alpha$$

Muestra pequeña:

$$p(\hat{\beta}_1 - t_{(n-2),\alpha/2}S_{\hat{\beta}_1} < \beta_1 < \hat{\beta}_1 + t_{(n-2),\alpha/2}S_{\hat{\beta}_1}) = 1 - \alpha$$

Cuanto menor es el número de observaciones n , menor la capacidad para calcular el error estándar del modelo. Como consecuencia, la exactitud de los coeficientes de regresión estimados se reduce. Esto tiene importancia sobre todo en la regresión múltiple (Amat, 2016).

En R, cuando se genera el modelo de regresión lineal, se devuelve junto con el valor de la pendiente y la ordenada en el origen el valor del estadístico t obtenido para cada uno y los p-value correspondientes. Esto permite saber, además de la estimación de β_0 y β_1 , si son significativamente distintos de 0.

También es posible que se quiera determinar una región de confianza para la estimación simultánea de los parámetros, sabiendo que:

$$Q = \frac{[n(\hat{\beta}_0 - \beta_0)^2 - 2n\bar{X}(\hat{\beta}_0 - \beta_0)(\hat{\beta}_1 - \beta_1) + (\hat{\beta}_1 - \beta_1)^2 \sum X^2]}{\sigma_e^2} \rightarrow X_2^2$$

$$\frac{(n-2)S_e^2}{\sigma_e^2} \rightarrow X_{(n-2)}^2$$

Se puede ver que:

$$F = \frac{[n(\hat{\beta}_0 - \beta_0)^2 - 2n\bar{X}(\hat{\beta}_0 - \beta_0)(\hat{\beta}_1 - \beta_1) + (\hat{\beta}_1 - \beta_1)^2 \sum X^2]}{2S_e^2} = F_{2,(n-2)}$$

12. Residuos del modelo

El residuo de una estimación se define como la diferencia entre el valor observado y el valor esperado acorde al modelo. A la hora de sumar el conjunto de residuos hay dos posibilidades:

- La suma del valor absoluto de cada residuo.
- La suma del cuadrado de cada residuo (RSS). Esta es la aproximación más empleada (mínimos cuadrados) ya que magnifica las desviaciones más extremas. En R, cuando se genera un modelo los residuos también se calculan automáticamente y se almacenan dentro del modelo.

Cuanto mayor es la sumatoria del cuadrado de los residuos, menor es la precisión con la que el modelo puede predecir el valor de la variable dependiente a partir de la variable predictora. Los residuos son muy importantes, puesto que, en ellos, se basan las diferentes medidas de la bondad de ajuste del modelo (Amat, 2016).

13. Análisis de varianza en R.L.S. (ANVA, ANOVA)

El análisis de varianza (ANOVA) en regresión se utiliza para determinar si la regresión lineal simple es estadísticamente significativa y para evaluar la importancia relativa de las variables independientes en la explicación de la variabilidad de la variable dependiente. El ANOVA en regresión se basa en la comparación de la varianza explicada por la regresión con la varianza no explicada (residual)

El ANOVA en regresión se divide en dos partes: el análisis de varianza global y el análisis de varianza individual. El análisis de varianza global se utiliza para determinar si la regresión en su conjunto es estadísticamente significativa. El análisis de varianza individual se utiliza para determinar la importancia relativa de cada variable independiente en la explicación de la variabilidad de la variable dependiente.

El análisis de varianza global se realiza mediante la prueba F de Fisher. El estadístico F se calcula dividiendo la varianza explicada por la regresión entre la varianza no explicada (residual). Si el valor de F es mayor que el valor crítico

correspondiente para un nivel de significancia determinado, se rechaza la hipótesis nula de que la regresión no es significativa.

Es importante tener en cuenta que el ANOVA en regresión asume que los residuos siguen una distribución normal con media cero y varianza constante, y que los residuos son independientes. Si estos supuestos no se cumplen, los resultados del ANOVA pueden ser sesgados y poco confiables. Por lo tanto, es importante realizar una evaluación cuidadosa de los supuestos de la regresión antes de realizar cualquier análisis.

Tomemos la desviación $Y_i - \bar{Y}$ (variación total), sumando y restando \hat{Y}_i para la media de Y_i .

$$Y_i - \bar{Y} = \underbrace{\hat{Y}_i - \bar{Y}} + \underbrace{Y_i - \hat{Y}_i}$$

Variación explicada + variación no explicada

Elevando al cuadrado y aplicando sumatorias

$$\sum (Y_i - \bar{Y})^2 = \sum (\hat{Y}_i - \bar{Y})^2 + \sum (Y_i - \hat{Y}_i)^2$$

$$\sum y_i^2 = \sum \hat{Y}_i^2 + \sum e_i^2$$

$$SCT = SCR + SCE$$

(Canavos, 1995) Si todas las observaciones se encuentran sobre la recta estimada, el valor de todos los residuos es cero y SCE=0. Entre más grande es el valor de SCE mayor es la contribución de la componente de error a la variación de las observaciones, o mayor es la incertidumbre. La SCR representa la variación de la observación que es atribuible al efecto lineal de X sobre Y. Si la pendiente es cero, entonces SCR= 0. de esta forma entre más grande es la proporción de SCR con respecto a SCT, mayor será la cantidad de la variación en las observaciones que puede explicarse mediante el termino lineal BX. Las SC son asociados con números llamados grados de libertad (g.l.). Este número indica la fracción de información independiente que existe en n números dependientes y_1, \dots, y_n .

Para la **SCT** existen $(n-1)$ g.l. ya que se pierde uno por causa de la restricción lineal $\sum(Y_i - \bar{Y}) = 0$ entre las observaciones Y_i .

Para la **SCE** existen $(n-2)$ g.l. Ya que se pierden dos g.l. a causa de dos restricciones lineales dadas (se estima dos parámetros).

Para la **SCR** debido a que es aditivo, entonces, g.l.SCR= g.l.SCT - g.l.SCE, esto es; $(n-1) - (n-2) = 1$.

A la ecuación fundamental dividimos sobre n, se obtiene varianzas poblacionales sesgadas, esto es:

$$\frac{\sum y_i^2}{n} = \frac{\sum \hat{Y}_i^2}{n} + \frac{\sum e_i^2}{n}$$

$$\sigma_y^2 = \sigma_{\hat{y}}^2 + \sigma_e^2$$

Para corregir se usa los g.l. y obtenemos varianzas insesgadas.

$$\sigma_y^2 = S_y^2 = \frac{\sum y_i^2}{n-1} \quad \sigma_{\hat{y}}^2 = S_{\hat{y}}^2 = \frac{\sum \hat{Y}_i^2}{1} \quad \sigma_e^2 = S_e^2 = \frac{\sum e_i^2}{n-2}$$

Con todos estos datos podemos consolidar un cuadro ANVA

Tabla 1.

ANVA para el modelo lineal simple

Fuente de variación	Grados de libertad	Suma de cuadrados (SC)	Cuadrados medios (CM)	F
Regresión	1	SCR	CMR	CMR/CME
Error	n-2	SCE	CME	-
Total	n-1	SCT	-	-

Donde:

Suma de cuadrados de la regresion: $SCR = \sum \hat{y}_i^2 = \hat{\beta}_1^2 \sum x_i^2 = (n)\hat{\beta}_1^2 S_x^2$

Suma de cuadrados del error: $SCE = \sum e_i^2 = \sum Y_i^2 - \hat{\beta}_0 \sum Y_i - \hat{\beta}_1 \sum X_i Y_i$

Suma de cuadrados del total: $SCT = \sum y_i^2 = \sum Y_i^2 - \frac{(\sum Y_i)^2}{n} = \sum Y_i^2 - n\bar{Y}^2$

Cuadrado medio de la regresion: $CMR = SCR/1$

Cuadrado medio del error: $CME = SCE/(n - 2) = S_e^2$

Tabla 2.

ANVA más general corregida por la media

Fuente de variación	Grados de libertad	Suma de cuadrados (SC)	Cuadrados medios (CM)	F
Regresión	1	SCR	CMR	CMR/CME
Error	n-2	SCE	CME	-
Total	n-1	SCT	-	-
debido a b*	1	SCC	CMC=SCC/1	CMC/CME
Total	n	SCTC		

$$SCC = (\sum y_i)^2/n = n\bar{Y}^2$$

$$SCTC = \sum Y_i^2$$

Prueba de hipótesis de F para la significación de la regresión.

1) Hipótesis estadística:

Ho: $\beta_0 = \beta_1 = 0$; no existe relación lineal entre X e Y

Ha: $\beta_0 \neq \beta_1 \neq 0$; existe relación lineal entre X e Y

2) nivel de significancia = α

3) prueba estadística, estadístico de contraste

$$F = \frac{CMR}{CME}$$

4) región crítica, decisión:

Se rechaza Ho si:

Muestra grande $F > F_{(1,n-2),\alpha/2}$

14. Bondad de ajuste del modelo

Una vez que se ha ajustado un modelo es necesario verificar su eficiencia, ya que aun siendo la línea que mejor se ajusta a las observaciones de entre todas las posibles, el modelo puede ser malo. Las medidas más utilizadas para medir la calidad del ajuste son: error estándar de los residuos, el test F y el coeficiente de determinación R^2 .

Error estándar de los residuos (Residual Standard Error, RSE): Mide la desviación promedio de cualquier punto estimado por el modelo respecto de la verdadera recta de regresión poblacional. Tiene las mismas unidades que la variable dependiente Y . Una forma de saber si el valor del RSE es grande consiste en dividirlo entre el valor medio de la variable respuesta, obteniendo así un % de la desviación.

$$RSE = \sqrt{\frac{1}{n - p - 1} RSS}$$

En modelos lineales simples, dado que hay un único predictor

$$(n - p - 1) = (n - 2)$$

Coeficiente de determinación R^2 :

Describe la proporción de variabilidad observada en la variable dependiente Y explicada por el modelo y relativa a la variabilidad total. Su valor está acotado entre 0 y 1. Al ser adimensional presenta la ventaja frente al RSE de ser más fácil de interpretar.

$$R^2 = \frac{\text{Suma de cuadrados totales} - \text{Suma de cuadrados residuales}}{\text{Suma de cuadrados totales}} =$$

$$1 - \frac{SCE}{SCT} =$$

$$1 - \frac{\sum(\hat{y}_i - y_i)^2}{\sum(y_i - \bar{y})^2}$$

En los modelos de regresión lineal simple el valor de R^2 se corresponde con el cuadrado del coeficiente de correlación de Pearson (r) entre X e Y , no siendo así en regresión múltiple. Existe una modificación de R^2 conocida como R^2 -ajustado que se emplea principalmente en los modelos de regresión

múltiple. Introduce una penalización cuantos más predictores se incorporan al modelo. En los modelos lineales simples no se emplea.

Test F: El test F es un test de hipótesis que considera como hipótesis nula que todos los coeficientes de correlación estimados son cero, frente a la hipótesis alternativa de que al menos uno de ellos no lo es. Se emplea en modelos de regresión múltiple para saber si al menos alguno de los predictores introducidos en el modelo contribuye de forma significativa. En modelos lineales simples, dado que solo hay un predictor, el p-value del test F es igual al p-value del t-test del predictor (Amat, 2016).

Desarrollo:

Los residuos pueden contribuir a proporcionar una medida útil de hasta qué punto la recta de regresión estimada se ajusta a los datos o que porcentaje de la variación es explicada por el modelo.

Tomemos la variación total de Y con respecto a su media.

$$(Y_i - \bar{Y}) = (\hat{Y}_i - \bar{Y}) + (Y_i - \hat{Y}_i) \quad \text{Como desviaciones de puede representar:}$$

$$y_i = \hat{y}_i + e_i \quad \text{Elevando al cuadrado y aplicando sumatorias}$$

$$\sum y_i^2 = \sum (\hat{y}_i + e_i)^2$$

$$\sum y_i^2 = \sum \hat{y}_i^2 + 2 \sum \hat{y}_i e_i + \sum e_i^2 \quad \text{Como } \sum \hat{y}_i e_i = 0 \quad \text{queda}$$

$$\sum y_i^2 = \sum \hat{y}_i^2 + \sum e_i^2$$

$$SCT = SCR + SCE$$

$$R^2 = \frac{SCR}{SCT} = 1 - \frac{SCE}{SCT}$$

$$R^2 = \left(\frac{S_{XY}}{S_X S_Y} \right)^2$$

Coefficiente de no determinación: $1 - R^2$

Coefficiente de alejamiento: $\sqrt{1 - R^2}$

Ejemplo: $R^2 = 90\%$. La regresión mínimo cuadrática explica el 90% de la varianza de Y. La proporción de variación total alrededor de la media explicada por la regresión es 90%. La reducción de la suma de cuadrados del error total fue del 90%. Entre más cercano se encuentre R^2 a 100%, mayor será el grado de asociación lineal entre X e Y.

15. Condiciones para el uso de la regresión lineal.

La regresión lineal simple es una técnica estadística que se utiliza para examinar la relación entre dos variables continuas. Sin embargo, la validez de los resultados de la regresión depende de ciertos supuestos. Estos supuestos son los siguientes:

1. **Linealidad:** (Amat, 2016) La relación entre ambas variables debe ser lineal, lo que significa que los cambios en una variable deben estar asociados con cambios proporcionales en la otra variable. Si la relación es no lineal, la regresión lineal simple no es apropiada.

Para comprobarlo se puede recurrir a:

- Graficar ambas variables a la vez (scatterplot o diagrama de dispersión), superponiendo la recta del modelo generado por regresión lineal.
 - Calcular los residuos para cada observación acorde al modelo generado y graficarlos (scatterplot). Deben distribuirse de forma aleatoria en torno al valor 0.
2. **Distribución Normal de los residuos:** Los errores de la regresión deben seguir una distribución normal, con media 0. Si los errores no siguen una distribución normal, los intervalos de confianza y las pruebas de hipótesis pueden no ser precisos. La normalidad se puede comprobar con un histograma, con la distribución de cuantiles (`qqnorm()` + `qqline()`) o con un test de hipótesis de normalidad. Los valores extremos suelen ser una causa frecuente por la que se viola la condición de normalidad (Amat, 2016).
 3. **Varianza de residuos constante (homocedasticidad):** La varianza de los errores debe ser constante en todo el rango de valores de la variable independiente (X). Si la varianza no es constante, se dice que hay

heterocedasticidad. La heterocedasticidad puede afectar la precisión de los coeficientes de la regresión y los intervalos de confianza

Se puede comprobar mediante gráficos (scatterplot) de los residuos de cada observación (formas cónicas son un claro indicio de falta de homocedasticidad) o mediante contraste de hipótesis mediante el test de Breusch-Pagan (Amat, 2016).

4. **Independencia, Autocorrelación:** Las observaciones deben ser independientes entre sí. Esto significa que el valor de la variable independiente no debe estar relacionado con el valor de la variable dependiente en ninguna otra observación. La violación de este supuesto puede conducir a resultados sesgados y poco confiables.
5. **Valores atípicos y de alta influencia:** Los valores atípicos pueden tener un efecto significativo en los resultados de la regresión (pueden generar una falsa correlación que realmente no existe, u ocultar una existente.). Si hay valores atípicos en los datos, deben ser identificados y evaluados para determinar si deben ser excluidos de la regresión.

Es importante tener en cuenta que la violación de estos supuestos no siempre invalida los resultados de la regresión lineal simple. Sin embargo, la validez de los resultados se verá afectada y puede ser necesario utilizar técnicas más avanzadas para analizar los datos. Por lo tanto, es importante evaluar cuidadosamente los supuestos de la regresión lineal simple antes de realizar cualquier análisis.

16. Predicción de valores

(Amat, 2016) Una vez generado un modelo que se pueda considerar válido, es posible predecir el valor de la variable dependiente Y para nuevos valores de la variable predictora X . Es importante tener en cuenta que las predicciones deben, a priori, limitarse al rango de valores dentro del que se encuentran las observaciones con las que se ha generado el modelo. Esto es importante puesto que solo en esta región se tiene certeza de que se cumplen las condiciones para que el modelo sea válido. Para calcular las predicciones se emplea la ecuación generada por regresión.

Dado que el modelo generado se ha obtenido a partir de una muestra y por lo tanto las estimaciones de los coeficientes de regresión tienen un error asociado, también lo tienen los valores de las predicciones. Existen dos formas de medir la incertidumbre asociada con una predicción:

- **Intervalos de confianza (predicción puntual):** Responden a la pregunta ¿Cuál es el intervalo de confianza del valor promedio de la variable respuesta Y para un determinado valor del predictor X ? (Amat, 2016)

$$\hat{y} \pm t_{\alpha/2, n-2} \sqrt{MSE \left(\frac{1}{n} + \frac{(x_k - \bar{x})^2}{\sum (x_i - \bar{x})^2} \right)}$$

Obtención del intervalo.

a) error de predicción: (fuente de error) se muestra por la ecuación.

$$\mu_i = Y_p - \hat{Y}_i$$

b) valor medio esperado de μ_i :

$$\begin{aligned} E(\mu_i) &= E(Y_p - \hat{Y}_i) \\ E(\mu_i) &= E(Y_p - \hat{\beta}_0 - \hat{\beta}_1 X_i) \\ E(\mu_i) &= E(Y_p) - E(\hat{\beta}_0) - X_i E(\hat{\beta}_1) \\ E(\mu_i) &= \hat{\beta}_0 + \hat{\beta}_1 X_i - \hat{\beta}_0 - \hat{\beta}_1 X_i = 0 \end{aligned}$$

c) Varianzas de \hat{Y}_p :

Se sabe que $\hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \bar{X}$

Tomemos:

$$\begin{aligned} V(\hat{Y}_i) &= V(\hat{\beta}_0 + \hat{\beta}_1 X_p) \\ V(\hat{Y}_i) &= V(\bar{Y} - \hat{\beta}_1 \bar{X} + \hat{\beta}_1 X_p) \\ V(\hat{Y}_i) &= V(\bar{Y}) + (X_p - \bar{X})^2 V(\hat{\beta}_1) \\ V(\hat{Y}_i) &= \frac{\sigma_e^2}{n} + (X_p - \bar{X})^2 \frac{\sigma_e^2}{\sum (X - \bar{X})^2} \\ V(\hat{Y}_i) &= \sigma_e^2 \left[\frac{1}{n} + \frac{(X_p - \bar{X})^2}{\sum (X - \bar{X})^2} \right] \end{aligned}$$

Estimando σ_e^2 por S_e^2

$$V(\hat{Y}_i) = S_{\hat{Y}_p}^2 = S_e^2 \left[\frac{1}{n} + \frac{(X_p - \bar{X})^2}{\sum (X - \bar{X})^2} \right]$$

La desviación estándar es:

$$S_{\hat{Y}_p} = S_e \sqrt{\frac{1}{n} + \frac{(X_p - \bar{X})^2}{\sum (X - \bar{X})^2}} = S_e \sqrt{\frac{1}{n} + \frac{(X_p - \bar{X})^2}{(n-1)S_X^2}}$$

con:

$$t = \frac{\hat{Y}_p - Y_p}{S_{\hat{Y}_p}} \rightarrow t_{(n-2)}$$

Intervalo de confianza: para $E(y/x)$

$$\left[\hat{Y}_p - t_{\alpha/2(n-2)} S_{\hat{Y}_p} < E(Y/X) < \hat{Y}_p + t_{\alpha/2(n-2)} S_{\hat{Y}_p} \right] = 1 - \alpha$$

- **Intervalos de predicción:** Responden a la pregunta ¿Dentro de que intervalo se espera que esté el valor de la variable respuesta Y para un determinado valor del predictor X ?

$$\hat{y} \pm t_{\alpha/2, n-2} \sqrt{MSE \left(1 + \frac{1}{n} + \frac{(x_k - \bar{x})^2}{(x_i - \bar{x})^2} \right)}$$

Desarrollo:

a) **valor medio ordenado:**

$$E(\hat{Y}_p - Y_p) = E(\hat{\beta}_0 + \hat{\beta}_1 X_p - (\alpha + \beta X_p + \mu_p))$$

$$E(\hat{Y}_p - Y_p) = E(\hat{\beta}_0 + \hat{\beta}_1 X_p - \alpha - \beta X_p - \mu_p)$$

$$E(\hat{Y}_p - Y_p) = E(\hat{\beta}_0) + X_p E(\hat{\beta}_1) - E(\alpha) + X_p E(\beta) + E(\mu_p)$$

$$E(\hat{Y}_p - Y_p) = \alpha + X_p \beta - \alpha - X_p \beta - 0$$

$$E(\hat{Y}_p - Y_p) = 0$$

b) **Varianza estimada** de una E individual predicha para una X dada:

$$\sigma_{(\hat{Y}_p - Y_p)}^2 = V(\hat{\beta}_0 + \hat{\beta}_1 X_p - e_p)$$

$$\sigma_{(\hat{Y}_p - Y_p)}^2 = V(\bar{Y} - \hat{\beta}_1 \bar{X} + \hat{\beta}_1 X_p - e_p)$$

$$\sigma_{(\hat{Y}_p - Y_p)}^2 = V(\bar{Y}) + (X_p - \bar{X})^2 V(\hat{\beta}_1) + V(e_p)$$

$$\sigma_{(\hat{Y}_p - Y_p)}^2 = \frac{\sigma_e^2}{n} + (X_p - \bar{X})^2 \frac{\sigma_e^2}{\sum(X_i - \bar{X})^2} + \sigma_e^2$$

$$\sigma_{(\hat{Y}_p - Y_p)}^2 = \sigma_e^2 \left[1 + \frac{1}{n} + \frac{(X_p - \bar{X})^2}{\sum(X_i - \bar{X})^2} \right] \quad \text{entonces:}$$

$$S_{(\hat{Y}_p - Y_p)}^2 = S_e^2 \left[1 + \frac{1}{n} + \frac{(X_p - \bar{X})^2}{\sum(X_i - \bar{X})^2} \right]$$

Intervalo de predicción

$$I.C = \left[\hat{Y}_p - t_{\alpha/2, (n-2)} S_{(\hat{Y}_p - Y_p)} < Y_o < \hat{Y}_p + t_{\alpha/2, (n-2)} S_{(\hat{Y}_p - Y_p)} \right] = 1 - \alpha$$

Si bien ambas parecen similares, la diferencia se encuentra en que los intervalos de confianza se aplican al valor promedio que se espera de Y para un determinado valor de X , mientras que los intervalos de predicción no se aplican al promedio. Por esta razón los segundos siempre son más amplios que los primeros. En R se puede emplear la función `predict()` que recibe como argumento el modelo calculado, un dataframe con los nuevos valores del predictor X y el tipo de intervalo (confidence o prediction).

Una característica que deriva de la forma en que se calcula el margen de error en los intervalos de confianza y predicción, es que el intervalo se ensancha a medida que los valores de X se aproximan a los extremos el rango observado (Amat, 2016).

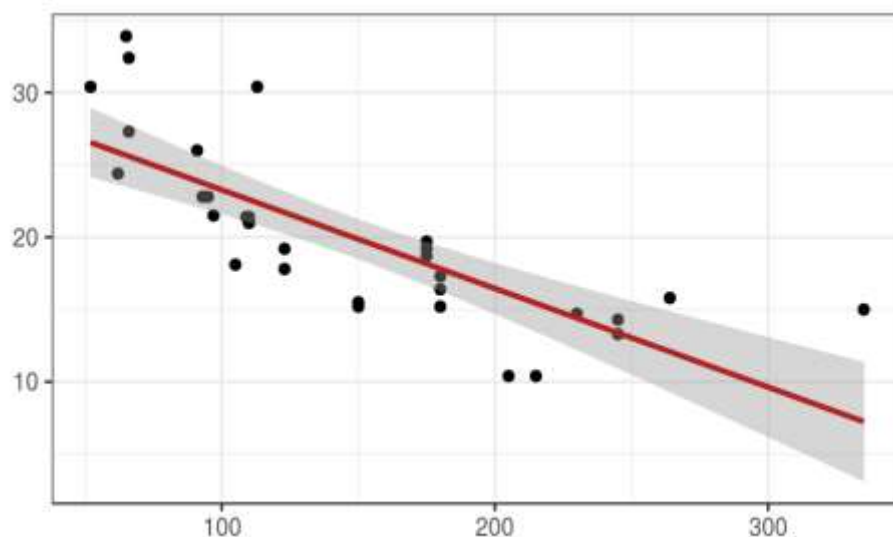


Figura 1. Bandas de confianza para la línea de regresión

¿Por qué ocurre esto? Prestando atención a la ecuación del error estándar del intervalo de confianza, el numerador contiene el término $(x_k - \bar{x})^2$ (lo mismo ocurre para el intervalo de predicción) (Amat, 2016).

$$\sqrt{MSE \left(\frac{1}{n} + \frac{(x_k - \bar{x})^2}{(x_i - \bar{x})^2} \right)}$$

Este término se corresponde con la diferencia al cuadrado entre el valor x_k para el que se hace la predicción y la media \bar{x} de los valores observados del predictor X . Cuanto más se aleje x_k de \bar{x} mayor es el numerador y por lo tanto el error estándar.

17. Interpretación geométrica del método

En regresión lineal tenemos un conjunto de observaciones pareadas (X, Y) , si graficamos podemos obtener en el plano cartesiano diversos diagramas de dispersión, luego de determinado la forma se obtiene una función matemática que mejor ajusta a los datos obteniendo los parámetros para luego interpretar, predecir, etc.

La Técnica estadística que consiste en determinar la curva que mejor asocie o ajuste los datos se llama **técnica o método de ajuste de curvas**.

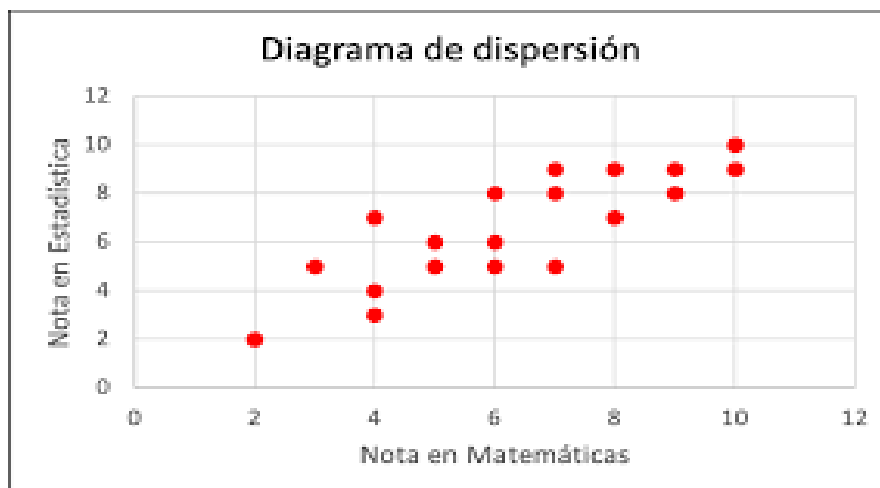


Figura 2. Diagrama de dispersión

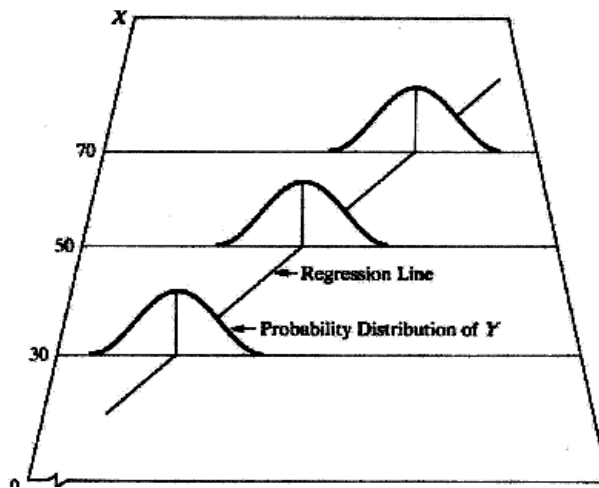


Figura 3. Distribución de probabilidad en regresión

Al diagrama de dispersión también se le conoce como diagrama de puntos, ploteo, nube de puntos, diagramas de Esparcimiento. Note que para un valor de X dado puede existir un rango de observaciones de Y y viceversa.

El diagrama de puntos permite dar una primera observación de la tendencia de los puntos (lineal, exponencial, etc.), en este caso podemos ver que podría ser lineal sin ningún inconveniente, asimismo las líneas presentadas inicialmente son trazadas al ojo (posible solución) pero la línea verdadera será trazada mediante la ecuación lineal $\hat{Y} = \hat{\beta}_0 + \hat{\beta}_1 X$ con el menor error posible entre \hat{Y} e Y.

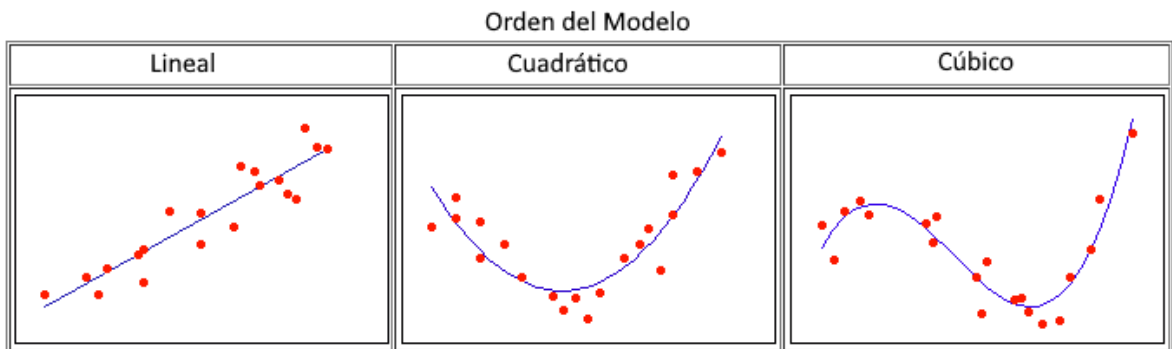


Figura 4. Ajuste de curvas.

18. Ejemplo de aplicación manual.

Ejemplo. De una determinada empresa se conocen los siguientes datos, referidos al volumen de ventas (en millones de pesetas) y al gasto en publicidad (en miles de pesetas).

Determinamos la variable dependiente e independiente considerando los lineamientos teóricos de la regresión, en este caso teniendo en cuenta causa efecto. Los gastos en publicidad influyen en el volumen de ventas.

Y_i : volumen de ventas

X_i : gastos en publicidad

La tabla siguiente muestra los cálculos manuales necesarios para realizar análisis de regresión lineal simple y entender su procedimiento.

$$y_i = (Y_i - \bar{Y}) \quad x_i = (X_i - \bar{X}) \quad e_i = (Y_i - \hat{Y})$$

Tabla 3.

Cálculos necesarios para obtener un modelo de regresión lineal simple.

dato	Datos		Observaciones			Desviaciones					estimados	errores	
	Y_i	X_i	Y_i^2	X_i^2	$X_i Y_i$	y_i	x_i	$x_i y_i$	y_i^2	x_i^2	\hat{Y}	e_i	e_i^2
1	10	16	100	256	160	-12	-28.0	328.5	138	784	12.804812	-2.804812	7.866969
2	15	32	225	1024	480	-7	-12.0	80.8	45	144	17.906825	-2.906825	8.449629
3	20	48	400	2304	960	-2	4.0	-6.9	3	16	23.008837	-3.008837	9.053102
4	22	56	484	3136	1232	0	12.0	3.2	0	144	25.559844	-3.559844	12.67249
5	30	64	900	4096	1920	8	20.0	165.3	68	400	28.11085	1.8891498	3.568887
6	32	80	1024	6400	2560	10	36.0	369.6	105	1296	33.212863	-1.212863	1.471037
7	12	20	144	400	240	-10	-24.0	233.6	95	576	14.080315	-2.080315	4.32771
8	16	22	256	484	352	-6	-22.0	126.1	33	484	14.718067	1.2819334	1.643353
9	23	50	529	2500	1150	1	6.0	7.6	2	36	23.646589	-0.646589	0.418077
10	29	52	841	2704	1508	7	8.0	58.1	53	64	24.284341	4.7156594	22.23744
11	31	75	961	5625	2325	9	31.0	287.3	86	961	31.618484	-0.618484	0.382522
12	14	35	196	1225	490	-8	-9.0	69.6	60	81	18.863452	-4.863452	23.65317
13	17	20	289	400	340	-5	-24.0	113.6	22	576	14.080315	2.919685	8.52456
14	27	55	729	3025	1485	5	11.0	57.9	28	121	25.240968	1.759032	3.094194
15	28	35	784	1225	980	6	-9.0	-56.4	39	81	18.863452	9.136548	83.47651
sumas	326	660	7862	34804	16182	0.00	0.0	1838	776.933	5764		-1.3E-05	190.8396

1) Diagrama de dispersión entre X e Y.

Es importante representar los datos para observar si siguen una tendencia lineal o no lineal.

El diagrama de dispersión se obtiene graficando los valores de Y y X.

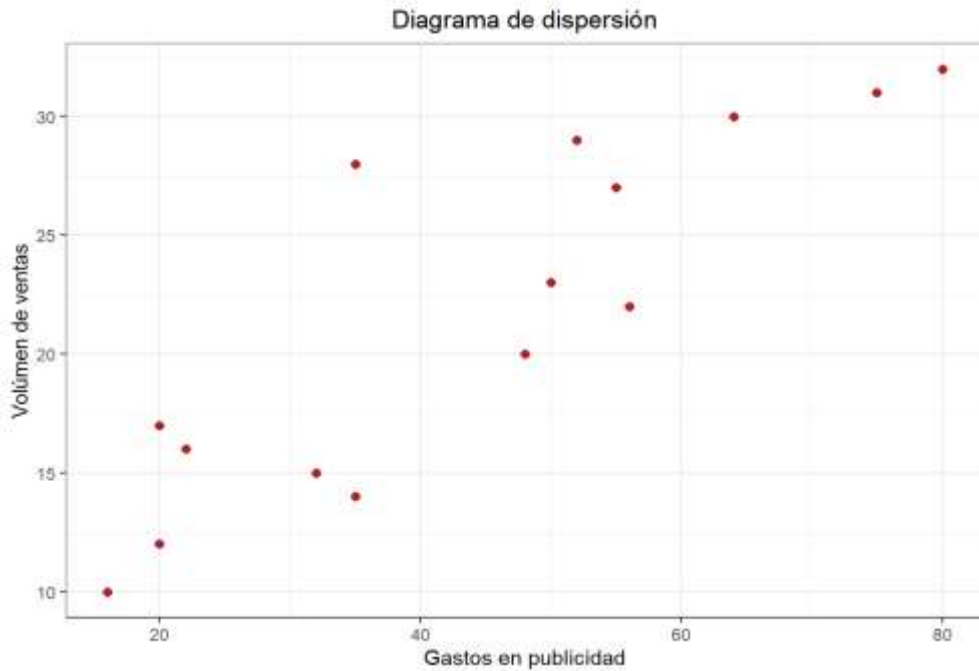


Figura 5. Diagrama de dispersión

La figura muestra que los datos pueden ajustarse a una regresión lineal.

2) Obtención de parámetros del modelo:

Se usarán diferentes fórmulas y procesos para determinar los valores de parámetros (estimadores) y cálculos necesarios:

Cálculo de estimadores $\hat{\beta}_0$ y $\hat{\beta}_1$ con datos observados

Con ecuaciones normales

$$\sum_{i=1}^n Y_i = n\hat{\beta}_0 + \hat{\beta}_1 \sum_{i=1}^n X_i$$

$$\sum_{i=1}^n X_i Y_i = \hat{\beta}_0 \sum_{i=1}^n X_i + \hat{\beta}_1 \sum_{i=1}^n X_i^2$$

$$\left. \begin{aligned} 326 &= 15\hat{\beta}_0 + \hat{\beta}_1(660) \\ 16182 &= \hat{\beta}_0(660) + \hat{\beta}_1(34804) \end{aligned} \right\}^{-660}_{15}$$

$$\begin{aligned} -215160 &= -9900\hat{\beta}_0 - 435600\hat{\beta}_1 \\ 242730 &= 9900\hat{\beta}_0 + 522060\hat{\beta}_1 \\ \hline 27570 &= 86460\hat{\beta}_1 \end{aligned}$$

despejando queda: $\hat{\beta}_1 = 27570/86460 = 0.3189$ (pendiente positiva)

Tomando una de las ecuaciones y reemplazando el estimador obtenido:

$$\begin{aligned} 326 &= 15\hat{\beta}_0 + 660\hat{\beta}_1 \\ 326 &= 15\hat{\beta}_0 + 660(0.3189) \\ -15\hat{\beta}_0 &= -326 + 210.45802 \\ \hat{\beta}_0 &= 7.702799 \end{aligned}$$

El modelo obtenido es:

$$\hat{Y} = \hat{\beta}_0 + \hat{\beta}_1 X_i = 7.702799 + 0.3189 X_i$$

Empleando formulas:

$$\hat{\beta}_1 = \frac{n \sum X_i Y_i - (\sum X_i)(\sum Y_i)}{n \sum X_i^2 - (\sum X_i)^2} = \frac{15(16182) - (660)(326)}{15(34804) - (660)^2} = \frac{27570}{86460} = 0.3188758$$

Mas formulas:

$$\hat{\beta}_1 = \frac{\sum X_i Y_i - [(\sum X_i)(\sum Y_i)]/n}{\sum X_i^2 - (\sum X_i)^2/n} = \frac{(16182) - [(660)(326)]/15}{(34804) - (660)^2/15} = \frac{27570}{86460} = 0.3188758$$

Mediante la regla de Cramer

$$\begin{aligned} \hat{\beta}_1 &= \frac{\begin{vmatrix} n & \sum Y_i \\ \sum X_i & \sum X_i Y_i \end{vmatrix}}{\begin{vmatrix} n & \sum X_i \\ \sum X_i & \sum X_i^2 \end{vmatrix}} = \frac{\begin{vmatrix} 15 & 326 \\ 660 & 16182 \end{vmatrix}}{\begin{vmatrix} 15 & 660 \\ 660 & 34804 \end{vmatrix}} = \frac{15(16182) - (660)(326)}{15(34804) - (660)^2} = \frac{27570}{86460} \\ &= 0.3188758 \end{aligned}$$

$$\begin{aligned} \hat{\beta}_0 &= \frac{(\sum Y_i)(\sum X_i^2) - (\sum X_i)(\sum X_i Y_i)}{n \sum X_i^2 - (\sum X_i)^2} = \frac{(326)(34804) - (660)(16182)}{15(34804) - (660)^2} = \frac{665984}{86460} \\ &= 7.702799 \end{aligned}$$

$$\hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \bar{X} = 326 - (0.3189 * 44) = 7.702799$$

Donde,

$$\begin{aligned} \bar{X} &= \frac{\sum X_i}{n} = \frac{55}{15} = 44 \\ \bar{Y} &= \frac{\sum Y_i}{n} = \frac{27}{15} = 326 \end{aligned}$$

Cálculo de estimadores $\hat{\beta}_0$ y $\hat{\beta}_1$ Utilizando desviaciones

Podemos usar los resultados de la tabla o las siguientes transformaciones, note la diferencia en las fórmulas (se usa minúsculas).

$$\sum x_i^2 = \sum X_i^2 - \frac{(\sum X_i)^2}{n} = 64 - \frac{18^2}{6} = 5764$$

$$\sum x_i y_i = \sum X_i Y_i - \frac{(\sum X_i)(\sum Y_i)}{n} = 106 - \frac{18(36)}{6} = 1838$$

$$\hat{\beta}_1 = \frac{\sum x_i y_i}{\sum x_i^2} = \frac{1838}{5764} = 0.3188758$$

$$\hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \bar{X} = 326 - (0.3189 * 44) = 7.702799$$

Otras fórmulas:

$$d_{xy} = \sum X_i Y_i - n \bar{X} \bar{Y} = (16182) - (15 * 21.733 * 44) = 1838$$

$$S_{xx} = \sum X_i^2 - n \bar{X}^2 = 34804 - (15 * (44)^2) = 5764$$

$$S_{yy} = \sum Y_i^2 - n \bar{Y}^2 = 7862 - (15 * (21.733)^2) = 777.15$$

con estos datos:

$$\hat{\beta}_1 = \frac{d_{xy}}{S_{xx}} = \frac{1838}{5764} = 0.3188758$$

Usando Covarianza y varianzas (**como datos muestrales**)

Nota: Al hallar la varianza tenga cuidado, existe diferencias entre la poblacional y la muestral, los resultados difieren.

$$S_{\bar{X}}^2 = \frac{S_{xx}}{n-1} = \frac{5764}{14} = 411.7143$$

$$S_{\bar{Y}}^2 = \frac{S_{yy}}{n-1} = \frac{777.15}{14} = 55.51076$$

$$S_{XY} = \frac{d_{xy}}{n-1} = \frac{1838}{14} = 131.2857$$

$$\hat{\beta}_1 = \frac{S_{XY}}{S_{\bar{X}}^2} = \frac{131.29}{411.71} = 0.3188791$$

3) Expresión de la línea de regresión.

Se puede encontrar valores de \hat{Y}_i (recta estimada por regresión) para cada valor de X, aunque con solo 2 puntos se puede trazar la recta. Veamos:

El modelo obtenido es: $\hat{Y} = \hat{\beta}_0 + \hat{\beta}_1 X_i = 7.702799 + 0.3188758X_i$.

Reemplazando valores de cada X, puede realizar el ajuste gráficamente.

$$\hat{Y}_1 = 7.702799 + 0.3189(16) = 12.80481$$

$$\hat{Y}_2 = 7.702799 + 0.3189(32) = 17.90682$$

$$\hat{Y}_3 = 7.702799 + 0.3189(48) = 23.00884$$

$$\hat{Y}_4 = 7.702799 + 0.3189(56) = 25.55984$$

$$\hat{Y}_5 = 7.702799 + 0.3189(64) = 28.11085$$

.

.

$$\hat{Y}_{15} = 7.702799 + 0.3189(35) = 18.86345$$

La línea de regresión se obtiene graficando los valores de \hat{Y} y X

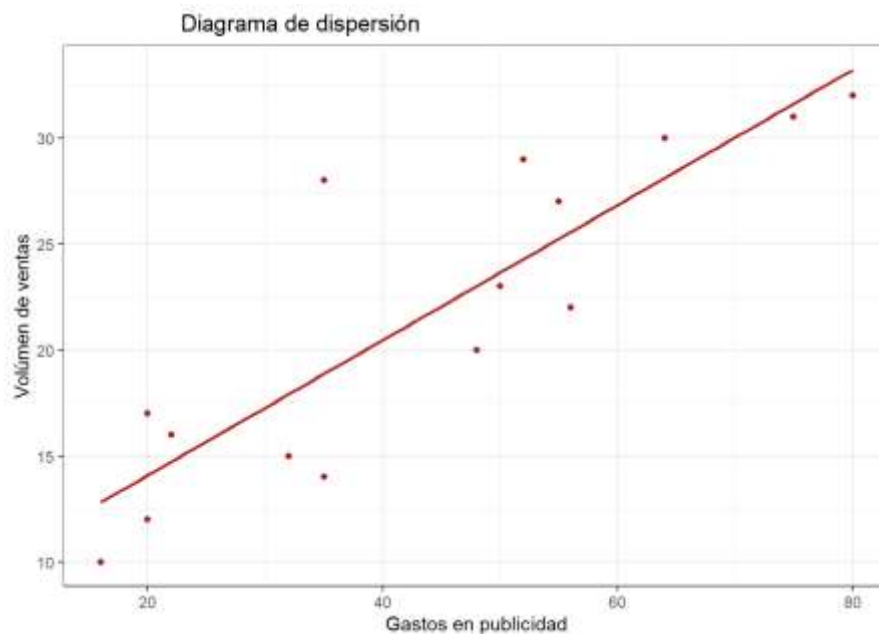


Figura 6. Representación de la línea de regresión.

4) Cálculo de las varianzas

a) Varianza residual (Varianza no explicada)

Con datos observados:

$$S_e^2 = \frac{\sum Y_i^2 - \hat{\beta}_0 \sum Y_i - \hat{\beta}_1 \sum X_i Y_i}{n - 2}$$

$$= \frac{7862 - (7.702799 * 326) - (0.3188758 * 16182)}{15 - 2} = \frac{190.8393}{13}$$

$$= 14.67995$$

$$S_e^2 = \frac{(n - 1)(S_Y^2 - \hat{\beta}_1^2 S_X^2)}{n - 2} = \frac{14(55.51076 - (0.3188758)^2(411.7143))}{n - 2}$$

$$= \frac{191.0569}{13} = 14.69668$$

Error residual= error residual estándar

$$S_e = \sqrt{14.67995} = 3.831442$$

En función de residuos y desviaciones:

$$S_e^2 = \frac{\sum e^2}{n - 2} = \frac{190.8396}{13} = 14.67997$$

$$S_e^2 = \frac{\sum y_i^2 - \hat{\beta}_1^2 \sum x_i^2}{n - 2} = \frac{776.933 - (0.3188758^2 * 5764)}{13} = 14.67994$$

b) Cálculo de las varianzas de los estimadores:

$$Var(\hat{\beta}_1) = S_{\hat{\beta}_1}^2 = \frac{S_e^2}{\sum x_i^2} = \frac{14.67995}{5764} = 0.0025468$$

$$Var(\hat{\beta}_1) = \frac{S_e^2}{S_{xx}} = \frac{S_e}{\sqrt{(n)S_X^2}}$$

$$Var(\hat{\beta}_0) = S_{\hat{\beta}_0}^2 = \frac{S_e^2 \sum X_i^2}{n \sum x_i^2} = \frac{14.67995(34804)}{15(5764)} = 5.909334$$

$$Var(\hat{\beta}_0) = S_{\hat{\beta}_0}^2 = S_e^2 \left[\frac{1}{n} + \frac{\bar{X}^2}{\sum x_i^2} \right] = 14.67995 \left[\frac{1}{15} + \frac{44^2}{5764} \right] = 5.909334$$

5) Prueba de hipótesis de los parámetros poblacionales

Prueba de coherencia de los estimadores

a) Prueba de hipótesis para el parámetro β_1 .

1) Hipótesis estadística:

Ho: $\beta_1 = 0$; no existe relación lineal, Y es independiente de X, no hay coherencia, la pendiente es cero, X no influye en Y

Ha: $\beta_1 \neq 0$; la pendiente no es cero, X influye en Y, existe relación lineal, Y es dependiente de X.

2) Nivel de significancia = $\alpha = 0.05$

3) Prueba estadística, estadístico de contraste: (n pequeño)

$$t = \frac{\hat{\beta}_1 - \beta}{S_{\hat{\beta}_1}} = \frac{0.3188758 - 0}{\sqrt{0.0025468}} = 6.318648$$

4) Región crítica, decisión:

Muestra pequeña: $t_t = t_{(n-2),\alpha/2} = t_{13,0.025} = 2.160$

Entonces: $|t| = 6.378 > |t_t| = 2.160$. Se rechaza Ho. La prueba es significativa al 5%. La pendiente $\hat{\beta}_1$ es diferente de cero, X influye en Y, existe relación lineal, Y es dependiente de X.

b) Prueba de hipótesis para el parámetro $\hat{\beta}_0$:

1) Hipótesis estadística:

Ho: $\beta_0 = 0$

Ha: $\beta_0 \neq 0$

2) Nivel de significancia = $\alpha = 0.05$

3) Prueba estadística, estadístico de contraste: (n pequeño)

$$t = \frac{\hat{\beta}_0 - \beta_0}{S_{\hat{\beta}_0}} = \frac{7.702799 - 0}{\sqrt{5.909334}} = 3.1686$$

4) Región crítica, decisión:

Muestra pequeña: $t_t = t_{(n-2),\alpha/2} = t_{13,0.025} = 2.160$

$|t| = 3.1686 > |t_t| = 2.160$. Se rechaza Ho.

$\hat{\beta}_0$ proviene de una población con parámetro β_0 . Debe estar en el modelo.

6) Intervalo de confianza

a) Intervalo de confianza para el parámetro β_0 (muestra pequeña)

$$p(\hat{\beta}_0 - t_{(n-2),\alpha/2}S_{\hat{\beta}_0} < \beta_0 < \hat{\beta}_0 + t_{(n-2),\alpha/2}S_{\hat{\beta}_0}) = 1 - \alpha$$

$$p(7.702799 - 2.160(\sqrt{5.9065}) < \beta_0 < 7.702799 + 2.160(\sqrt{5.9065}))$$

$$= 1 - 0.05$$

$$p(2.45 < \beta_0 < 12.95) = 0.95$$

Si 100 muestras del mismo tamaño son escogidas, podemos esperar que 95 de ellas contengan el valor del parámetro y 5 caigan fuera.

b) Intervalo de confianza para el parámetro β_1

$$p\left(\hat{\beta}_1 - t_{(n-2),\frac{\alpha}{2}}S_{\hat{\beta}_1} < \beta_1 < \hat{\beta}_1 + t_{(n-2),\frac{\alpha}{2}}S_{\hat{\beta}_1}\right) = 1 - \alpha$$

$$p(0.3189 - 2.160(\sqrt{0.0025}) < \beta_1 < 0.3189 + 2.160(\sqrt{0.0025}))$$

$$= 1 - 0.05$$

$$p(0.2109 < \beta_1 < 0.4269) = 0.95$$

7) Tabla de Análisis de Varianza (ANVA).

Prueba de confiabilidad del modelo, prueba F para la significación de la regresión.

Numero de pares de datos: $n = 15$

Grados de libertad de la regresión: $GLR = 1$

Grados de libertad del error: $GLE = n-2 = 13$

Grados de libertad del total: $GLT = n-1 = 14$

Sumas de cuadrados:

$$SCR = \sum \hat{y}_i^2 = \hat{\beta}_1^2 \sum x_i^2 = (0.3188758)^2(5764) = 586.0938$$

$$SCR = \hat{\beta}_1 d_{xy} = (0.3188758)(1838) = 586.0937$$

$$SCT = \sum y_i^2 = \sum Y_i^2 - \frac{(\sum Y_i)^2}{n} = 7862 - \frac{326^2}{15} = 776.93$$

$$SCT = \sum Y_i^2 - n\bar{Y}^2 = 7862 - 15(21.733)^2 = 776.93$$

$$SCE = \sum Y_i^2 - \hat{\beta}_0 \sum Y_i - \hat{\beta}_1 \sum X_i Y_i =$$

$$= 7862 - 7.702799(326) - (0.3188758)(16182) = 190.8393$$

$$SCE = SCT - SCR = 776.93 - 586.0938 = 190.8362$$

Tabla 4.

Análisis de varianza

Fuente de variación	Grados de libertad	Suma de cuadrados (SC)	Cuadrados medios (CM)	F
Regresión	1	586.0938	586.18272	39.932
Error	13	190.8362	14.67938	
Total	14	776.93		

$$F_t = F_{(1,13),0.95} = 4.667$$

Entonces: $F(39.932) > F_t(4.667)$, Se rechaza la H_0 al 5% de significancia, Utilizando probabilidades $p(0.00002665) < \alpha(0.05)$. Se rechaza la H_0 . Hay relación lineal entre X e Y, el modelo en su conjunto es significativo.

El modelo puede ser tomado para realizar predicciones.

8) Coeficiente de Determinación

$$R^2 = \frac{SCR}{SCT} = \frac{586.0938}{776.93} = 0.7543 \quad \circ$$

$$1 - \frac{SCE}{SCT} = 1 - \frac{190.8362}{776.93} = 0.7543 \quad \circ$$

$$R^2 = \left(\frac{S_{XY}}{S_X S_Y} \right)^2$$

9) Coeficiente de correlación.

$$r = \frac{n \sum X_i Y_i - (\sum X_i)(\sum Y_i)}{\sqrt{n \sum X_i^2 - (\sum X_i)^2} \sqrt{n \sum Y_i^2 - (\sum Y_i)^2}}$$

$$= \frac{15(16182) - (660 * 326)}{\sqrt{15(34804) - (660)^2} \sqrt{15(7862) - (326)^2}} = 0.8685$$

$$r = \frac{\sum x_i y_i}{\sqrt{\sum x_i^2} \sqrt{\sum y_i^2}} = \frac{1838}{\sqrt{776.933} * \sqrt{5764}} = 0.8685$$

$$r = \frac{S_{XY}}{S_X S_Y} = \frac{131.286}{\sqrt{411.7143} * \sqrt{55.1076}} = 0.8715$$

10) prueba de supuestos.

Se realizará en software R.

19. Ejemplo 1 en R.**Regresión lineal simple con variable independiente numérica.**

Con los datos del ejemplo anterior.

```
# regresión lineal simple
# Y: volumen de ventas
# X: gastos en publicidad
Y <- c(10, 15, 20, 22, 30, 32, 12, 16, 23, 29, 31, 14, 17, 27, 28)
X <- c(16, 32, 48, 56, 64, 80, 20, 22, 50, 52, 75, 35, 20, 55, 35)
```

Creando un *data.frame* y observando los primeros 6 datos

```
datos <- data.frame(X, Y)
head(datos)
```

```
##      X  Y
## 1 16 10
## 2 32 15
## 3 48 20
```

```
## 4 56 22
## 5 64 30
## 6 80 32
```

La estructura (*str*) sirve para observar la estructura de ingreso de los datos (datos numéricos, o cualitativos)

```
str(datos)
```

```
## 'data.frame': 15 obs. of 2 variables:
## $ X: num 16 32 48 56 64 80 20 22 50 52 ...
## $ Y: num 10 15 20 22 30 32 12 16 23 29 ...
```

Observamos que la data set está conformada por 15 observaciones y 2 variables:

- X, como variable numérica

Y, como variable numérica.

El volumen de ventas (Y) estará en función al gasto en publicidad (X), dicho de otra manera, el gasto en publicidad influye el volumen de ventas de la empresa.

a) Observar las variables descriptivamente.

```
# representación grafica
# histograma para las variables
require(ggplot2)
par(mfrow = c(1, 2)) # dos figuras en una fila
hist(datos$X, breaks = 10, main = "", xlab = "Gastos en publicidad",
      border = "darkred")
hist(datos$Y, breaks = 10, main = "", xlab = "Volumen de ventas",
      border = "blue")
```

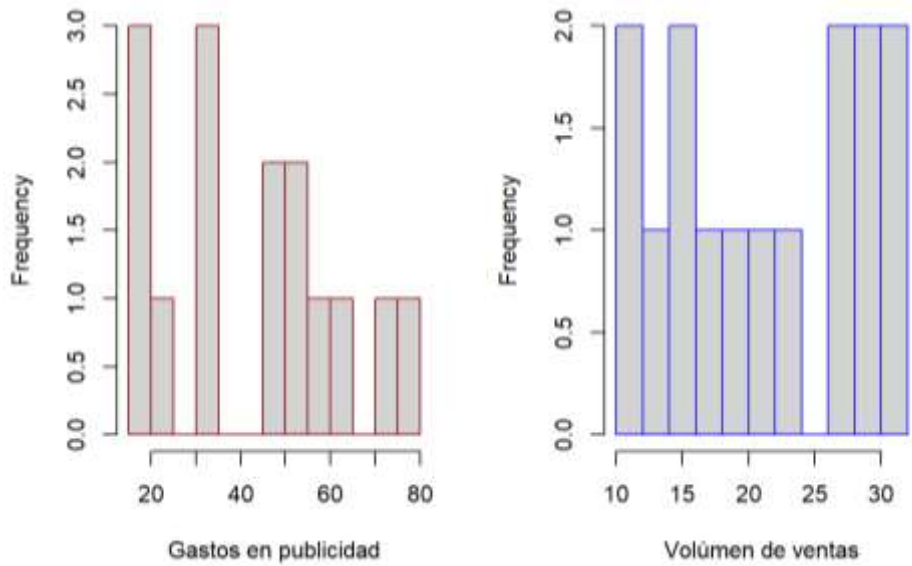


Figura 7. Histograma de frecuencias para el número de bateos y numero de corridas

La figura muestra una tendencia a valores bajos en gastos en publicidad y valores altos en volumen de venta, al parecer sin tendencia a una distribución normal.

```
# estadísticos descriptivos
summary(datos)
```

```
##           X           Y
## Min.      :16.0    Min.      :10.00
## 1st Qu.:27.0    1st Qu.:15.50
## Median :48.0    Median :22.00
## Mean   :44.0    Mean   :21.73
## 3rd Qu.:55.5    3rd Qu.:28.50
## Max.   :80.0    Max.   :32.00
```

El promedio de volumen de ventas (Y) es 21.73 millones de pesetas con ventas mínimas de 10 millones de pesetas y máximo de ventas de 32 millones de pesetas. El promedio de gastos en publicidad (X) es de 44 mil pesetas, con valores mínimo y máximo de 16 y 80 mil pesetas

respectivamente. La mediana del volumen de ventas es 22 millones de pesetas y la mediana del gasto en publicidad es 80 mil pesetas.

b) Representación gráfica de las observaciones

Es necesario graficar un diagrama de dispersión, para ver la tendencia de los datos.

```
# diagrama de dispersión
require(ggplot2)
ggplot(data = datos, mapping = aes(x = X, y = Y)) +
  geom_point(color = "firebrick", size = 2) +
  (labs(title = "Diagrama de dispersión", x="Gastos en publicidad",
        y= "Volúmen de ventas"))+
  theme_bw() +
  theme(plot.title = element_text(hjust =0.5))
```

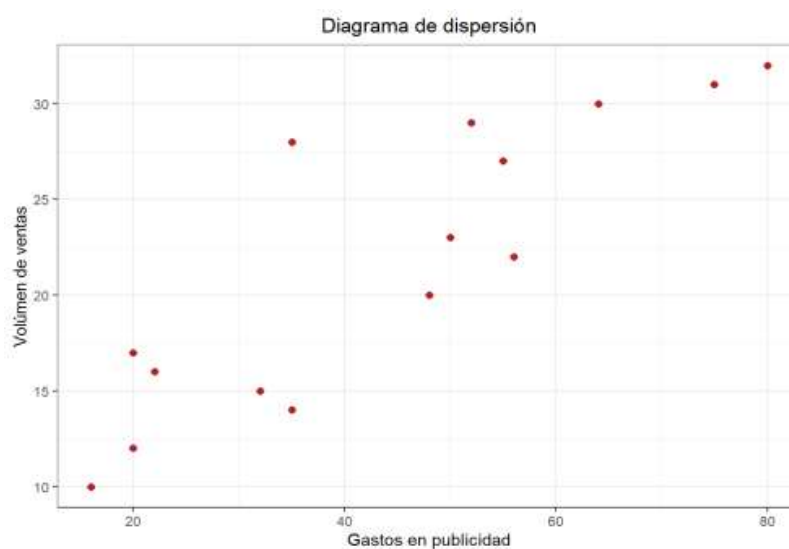


Figura 8. Diagrama de dispersión

La figura muestra un comportamiento (tendencia) positivo, podría interpretarse como, a más gastos en publicidad se incrementa el volumen de ventas. No se observa claramente la presencia de datos atípicos (outliers) o valores influyentes.

c) comportamiento del coeficiente de correlación

otro indicador de una posible relación lineal es el coeficiente de correlación, por lo que es recomendable calcular el coeficiente.

Si se observa en el diagrama de dispersión y un valor bajo del coeficiente, no tiene sentido seguir adelante generando un modelo lineal, sería mejor buscar alternativas no lineales.

Estadístico de correlación de Pearson (para datos cuantitativos)

```
# correlacion de Pearson
cor.test(x = datos$X, y = datos$Y, method = "pearson")
```

```
## Pearson's product-moment correlation
##
## data:  datos$X and datos$Y
## t = 6.3186, df = 13, p-value = 2.665e-05
## alternative hypothesis: true correlation is not equal to 0
## 95 percent confidence interval:
##  0.6418579 0.9556266
## sample estimates:
##          cor
## 0.8685436
```

Resumen:

Coeficiente de correlación	Test de significancia	Intervalo de confianza de r	p-value
$r = 0.8685$	$t = 6.3186$	0.6419 - 0.9556	0.00002665

1) hipótesis estadística:

Ho: la correlación es igual a 0

Ha: la correlación no es igual a 0

2) nivel de significancia: $\alpha = 0.05$

3) coeficiente de correlación $r = 0.868$

test de significancia: $t = 6.3186$

p-value: $p = 0.00002665$

intervalo de confianza: $IC = (0.6419 - 0.9556)$

- 4) decisión: El test de correlación muestra una relación lineal significativa $p(0.00002665) < \alpha(0.05)$, de intensidad considerable ($r=0.868$). Tiene sentido intentar generar un modelo de regresión lineal que permita predecir el volumen de ventas. Se espera que el coeficiente de correlación oscile entre 0.64 y 0.95 utilizando el IC al 95%.

d) Cálculo del modelo de regresión lineal simple

La regresión es una herramienta estadística que se utiliza para analizar la relación entre dos variables. Es un modelo que se ajusta a los datos de la muestra y se utiliza para predecir los valores de una variable (la variable dependiente) a partir de los valores de otra variable (la variable independiente).

$$y = a + bx$$

donde y es la variable dependiente, x es la variable independiente, a es el intercepto (el valor de y cuando x es igual a cero) y b es la pendiente de la recta.

La pendiente de la recta indica el cambio en y por unidad de cambio en x . Es decir, si la pendiente es positiva, significa que a medida que x aumenta, y también aumenta. Si la pendiente es negativa, significa que a medida que x aumenta, y disminuye. La magnitud de la pendiente indica la fuerza de la relación entre las dos variables.

Para obtener los estimadores de la ecuación, utilizamos el siguiente script

```
# Calculo del modelo de regresión lineal simple
modelo_lineal <- lm(Y ~ X, data=datos) # selección de variables y data
summary(modelo_lineal)
```



```
## Call:
## lm(formula = Y ~ X, data = datos)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -4.8635 -2.8558 -0.6466  1.8241  9.1365
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  7.70280    2.43091   3.169  0.0074 **
## X            0.31888    0.05047   6.319 2.66e-05 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 3.831 on 13 degrees of freedom
## Multiple R-squared:  0.7544, Adjusted R-squared:  0.7355
## F-statistic: 39.92 on 1 and 13 DF, p-value: 2.665e-05
```

Tanto el intercepto como la variable X (gastos en publicidad) son estadísticamente significativas al 5%, en ambos casos [$p(0.0074) < \alpha(0.05)$ y $p(0.0000266) < \alpha(0.05)$], con tendencia positiva. El modelo tiene un coeficiente de determinación de 75.44% (el modelo calculado explica el 75.44% de la variabilidad presente en la variable respuesta (volumen de ventas) mediante la variable independiente (gastos en publicidad). La ecuación de regresión se expresa como:

$$Y = 7.70280 + 0.31888(X)$$

Por cada unidad (mil pesetas) que se invierte en gastos de publicidad, el volumen de ventas se incrementa en promedio 0.3188 millones de pesetas.

El modelo en su conjunto es también significativo (ANVA) con $F = 39.92$ y $p(0.000026665) < \alpha(0.05)$. el error residual (RSS) = 3.831.

Cálculo del error cuadrático medio RMSE:

```
# Cálculo del error cuadrático medio RMSE:
sqrt(mean(modelo_lineal$residuals^2))
```

```
## [1] 3.566881
```

Un modelo completo de análisis de varianza se puede obtener con el siguiente script.

```
anova(modelo_lineal)
```

```
## Analysis of Variance Table
##
## Response: Y
##           Df Sum Sq Mean Sq F value    Pr(>F)
## X           1  586.09   586.09  39.925 2.665e-05 ***
## Residuals  13  190.84    14.68
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

Hasta aquí se puede concluir que el modelo es bueno y existe relación entre las variables.

e) Intervalos de confianza para los parámetros del modelo

```
# Intervalos de confianza para los parámetros del modelo
confint(modelo_lineal)
```

```
##           2.5 %      97.5 %
## (Intercept) 2.4511286 12.9544694
## X           0.2098502  0.4279014
```

Los intervalos de confianza contienen a los coeficientes obtenidos. Por cada unidad que se incrementa los gastos en publicidad, el volumen de ventas aumenta en promedio entre 0.2098 y 0.4279 unidades.

f) Representación gráfica del modelo

La recta de regresión es una herramienta estadística que se utiliza para analizar la relación entre dos variables. Es una línea recta que se ajusta a los datos de la muestra y se utiliza para predecir los valores de una variable (la variable dependiente) a partir de los valores de otra variable (la variable independiente).

La recta de regresión se calcula utilizando la técnica de mínimos cuadrados, que implica encontrar la línea recta que minimiza la distancia entre los puntos de datos y la línea recta. La recta de regresión se puede expresar en la forma de una ecuación lineal:

$$y = a + bx$$

```
# Representación gráfica del modelo
ggplot(data = datos, mapping = aes(x = X, y = Y)) +
  geom_point(size=3) +
  (labs(title = "Diagrama de dispersión", x="Gastos en publicidad",
        y= "Volumen de ventas"))+
  geom_smooth(method = "lm", se = FALSE, color = "firebrick") +
  theme_bw() +
  theme(plot.title = element_text(hjust = 0.2))
```

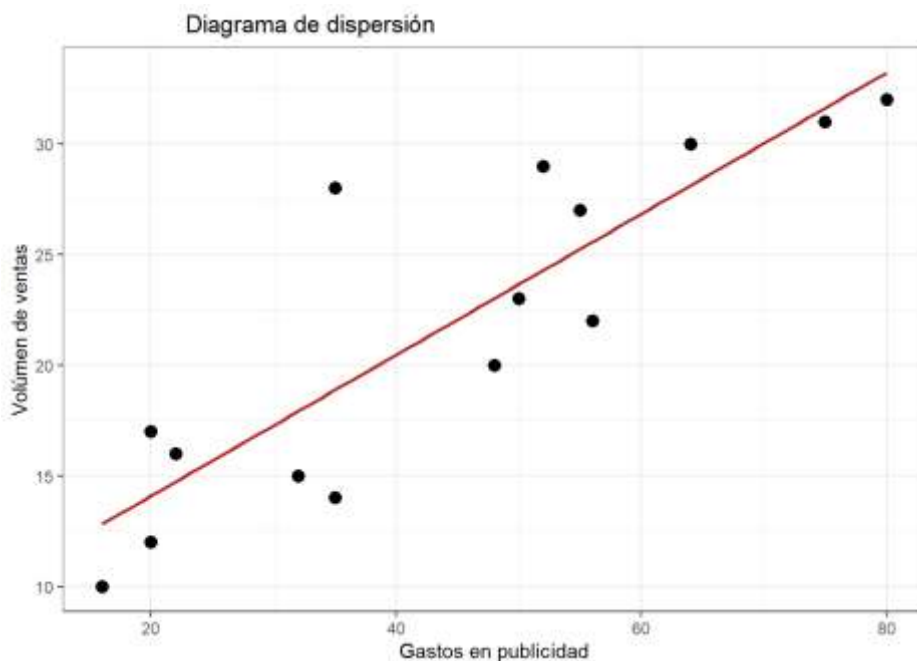


Figura 9. Representación de la línea de regresión

```
names(modelo_lineal)
```

```
## [1] "coefficients" "residuals" "effects" "rank"  
## [5] "fitted.values" "assign" "qr" "df.residual"  
## [9] "xlevels" "call" "terms" "model"
```

Con la opción *names* muestra los resultados guardados por el modelo (ejemplo: `modelo_lineal$coefficients`) se muestra solo los coeficientes del modelo.

Además de la línea de mínimos cuadrados es recomendable incluir los límites superior e inferior del intervalo de confianza. Esto permite identificar la región en la que, según el modelo generado y para un determinado nivel de confianza, se encuentra el valor promedio de la variable dependiente.

Para poder representar el intervalo de confianza a lo largo de todo el modelo se recurre a la función `predict()` para predecir valores que abarquen todo el eje *X*. Se añaden al gráfico líneas formadas por los límites superiores e inferiores calculados para cada predicción (Amat, 2016).

Obtenemos los valores predichos para cada valor de *X*.

```
# predecir  
nuevos_datos <- data.frame(X= seq(min(X),max(X)))  
predict_value <- predict(modelo_lineal)  
head(predict_value) # muestra 6 valores predichos
```

```
##      1      2      3      4      5      6  
## 12.80481 17.90682 23.00884 25.55984 28.11085 33.21286
```

Graficando las bandas de confianza:

Intervalos de confianza de la respuesta media

```
# solo una banda
par(mfrow = c(1, 1))
puntos <- seq(from = min(datos$X),
              to = max(datos$X), length.out = 100)
limites_intervalo <- predict(object = modelo_lineal,
                             newdata = data.frame(X = puntos),
                             interval = "confidence", level = 0.95)
head(limites_intervalo, 3)
```

```
##          fit          lwr          upr
## 1 12.80481  9.078325 16.53130
## 2 13.01095  9.341982 16.67992
## 3 13.21710  9.605180 16.82901
```

```
plot(datos$X, datos$Y, col = "firebrick", pch = 19,
      ylab = "Gastos en publicidad", xlab = "Volúmen de ventas",
      main = "Volúmen de ventas ~ Gastos en publicidad")
abline(modelo_lineal, col = 1)
lines(x = puntos, y = limites_intervalo[, 2], type = "l", col = 2, lty =
3)
lines(x = puntos, y = limites_intervalo[, 3], type = "l", col = 3, lty =
3)
```

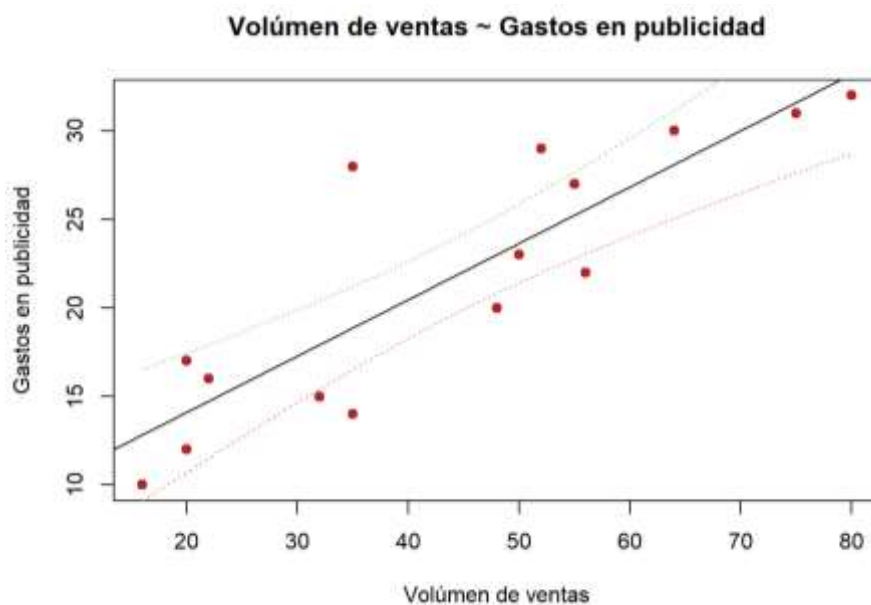


Figura 10. Bandas de confianza para el ajuste de la regresión

La función `geom_smooth()` del paquete `ggplot2` genera la regresión y su intervalo de forma directa.

```
# con geom plot
ggplot(data = datos, mapping = aes(x = X, y = Y)) +
  geom_point(color = "firebrick", size = 2) +
  labs(title = "Diagrama de dispersion", x = "numero de bateos") +
  geom_smooth(method = "lm", se = TRUE, color = "black") +
  theme_bw() +
  theme(plot.title = element_text(hjust = 0.5))
```

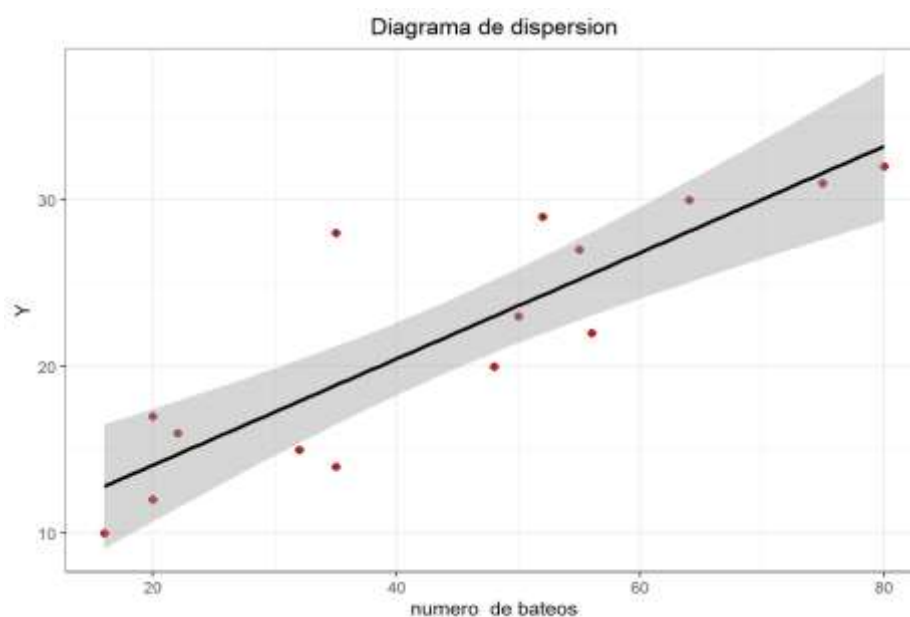


Figura 11. *Bandas de confianza para la respuesta media*

Intervalo de confianza para la respuesta media y para la predicción:

```
# dos bandas de confianza
par(mfrow = c(1, 1))
# Grafico dispersion y recta
plot(datos$X, datos$Y, col = "firebrick", pch = 19,
      ylab = "Y", xlab = "X",
      main = "")
abline(modelo_lineal, col = 1)

# Intervalos de confianza de la respuesta media:
# valores medios
```

```

ic <- predict(modelo_lineal, nuevos_datos, interval = 'confidence')
lines(nuevos_datos$X, ic[, 2], lty = 2)
lines(nuevos_datos$X, ic[, 3], lty = 2)

# Intervalos de predicción
# para cualquier valor
ic <- predict(modelo_lineal, nuevos_datos, interval = 'prediction')
lines(nuevos_datos$X, ic[, 2], lty = 2, col = 'red')
lines(nuevos_datos$X, ic[, 3], lty = 2, col = 'red')

```

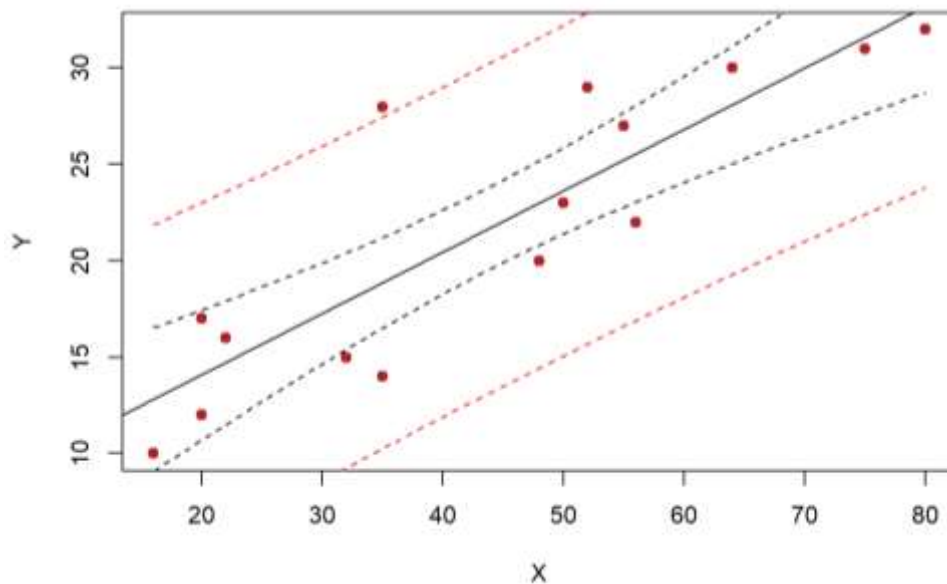


Figura 12. Bandas de confianza para la respuesta media y la predicción

Observamos un punto fuera de la banda de confianza de la predicción, posible punto outlier.

El siguiente script muestra los datos reales, los datos predichos y los errores calculados.

```

datos$prediccion <- modelo_lineal$fitted.values # valores ajustados
datos$residuos <- modelo_lineal$residuals # residuales
head(datos)

```

```
##      X  Y prediccion  residuos
## 1 16 10   12.80481 -2.804811
## 2 32 15   17.90682 -2.906824
## 3 48 20   23.00884 -3.008836
## 4 56 22   25.55984 -3.559843
## 5 64 30   28.11085  1.889151
## 6 80 32   33.21286 -1.212861
```

g) Verificar los supuestos del modelo lineal

Linealidad. Se calculan los residuos para cada observación y se representan gráficamente (scatterplot). Si las observaciones siguen la línea del modelo, los residuos se deben distribuir aleatoriamente entorno al valor 0 (Amat, 2016).

```
# linealidad
ggplot(data = datos, aes(x = prediccion, y = residuos)) +
  geom_point(aes( color = residuos)) +
  scale_color_gradient2(low = "blue3", mid = "grey", high = "red") +
  geom_hline(yintercept = 0) +
  geom_segment(aes(xend = prediccion, yend = 0), alpha = 0.2) +
  labs(title = "Distribucion de los residuos", x = "prediccion modelo",
       y = "residuo") +
  theme_bw() +
  theme(plot.title = element_text(hjust = 0.5), legend.position = "none")
```

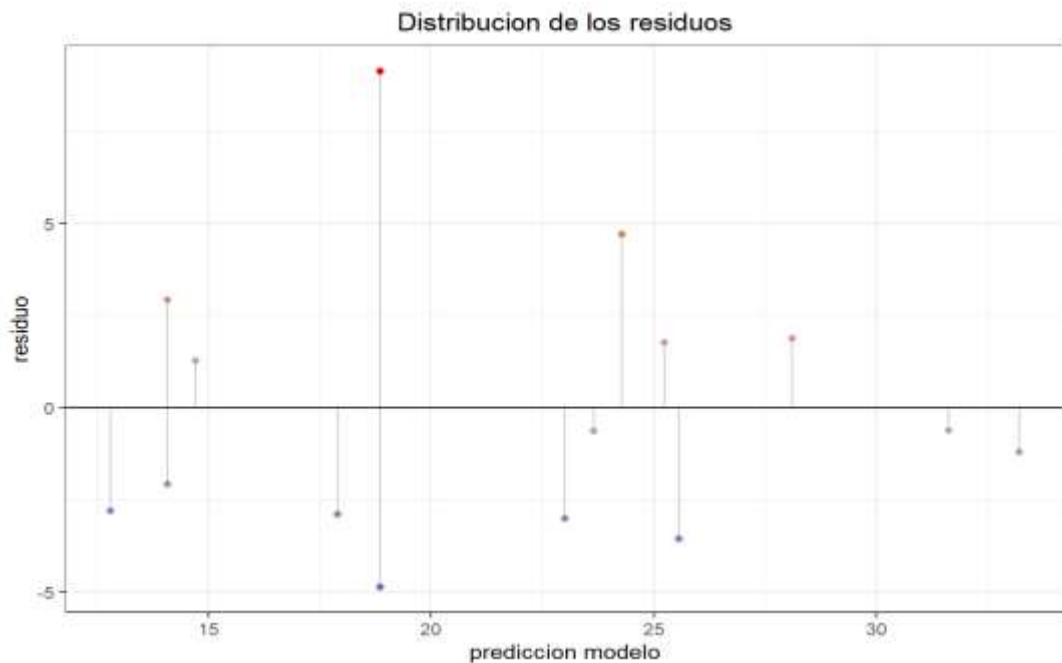



Figura 13. Distribución de los residuos de modelo.

Los residuos se distribuyen de forma aleatoria entorno al 0 por lo que se acepta la linealidad.

Normalidad de los residuos:

Los residuos se deben distribuir de forma normal con media 0. Para comprobarlo se recurre a histogramas, a los cuantiles normales o a un test de contraste de normalidad.

Histograma de residuos.

```
# Distribucion normal de los residuos:
ggplot(data = datos, aes(x = residuos)) +
  geom_histogram(aes(y = ..density..)) +
  labs(title = "Histograma de los residuos") +
  theme_bw() +
  theme(plot.title = element_text(hjust = 0.5))
```

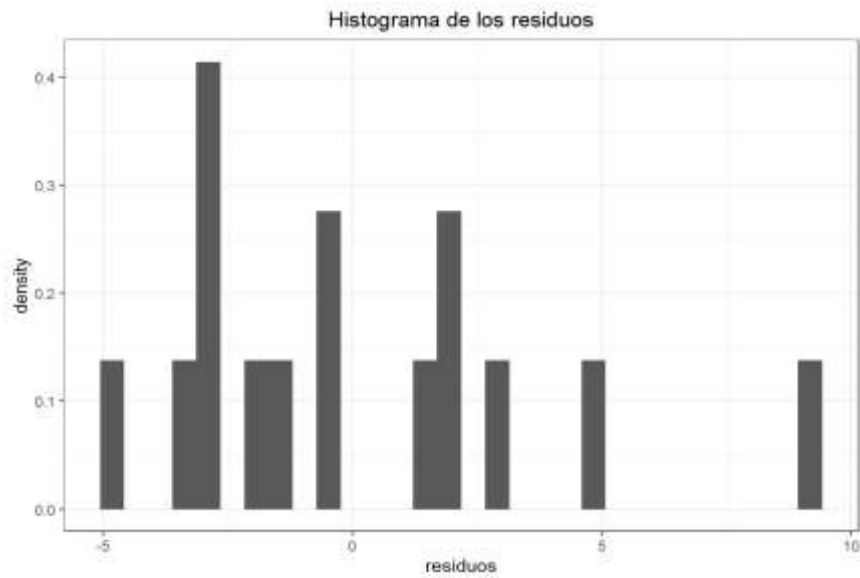


Figura 14. Histograma de densidad de los residuos

No se evidencia una clara distribución normal de los residuos.

Gráfico de cuantiles.

```
# grafico de cuantiles  
qqnorm(modelo_lineal$residuals)  
qqline(modelo_lineal$residuals)
```

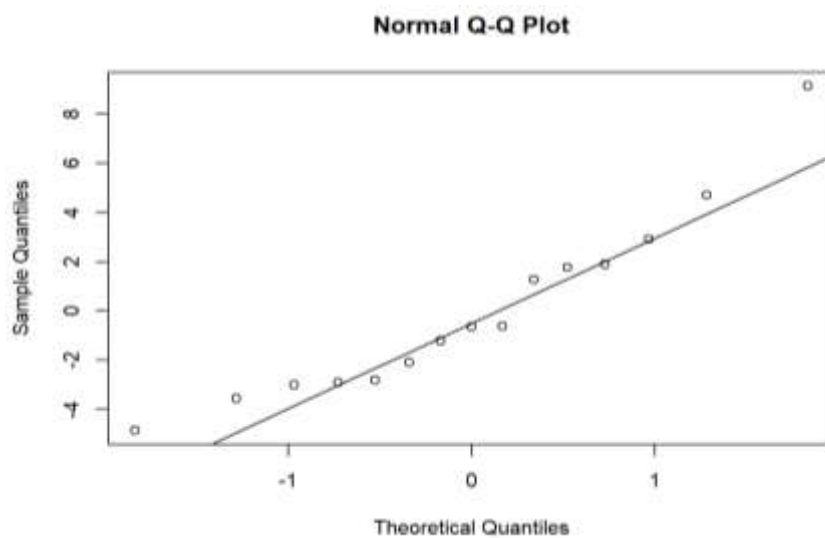


Figura 15. Gráfico de cuantiles de los residuos

Los puntos no se muestran alrededor de la línea, algunos se muestran alejados, posible no normalidad en los errores.

Finalmente realizamos una prueba de hipótesis de normalidad, por la cantidad de datos el estadístico más adecuado sería el test de Shapiro.

```
# test de normalidad
shapiro.test(modelo_lineal$residuals)
```

```
## Shapiro-Wilk normality test
##
## data:  modelo_lineal$residuals
## W = 0.92214, p-value = 0.2077
```

1) Prueba de hipótesis

Ho: existe normalidad en los residuos

Ha: no existe normalidad en los residuos

2) nivel de significancia 0.05

3) Prueba estadística. Shapiro y Wilks

W = 0.92214 con p=0.2077

4) decisión: $p(0.2077) > \alpha(0.05)$, se acepta la Ho de normalidad de los residuos.

Podemos obtener también el test de kolmogorov (muestras grandes), que, a diferencia del test de Shapiro, muestra normalidad de los residuos $p(0.7405) > \alpha(0.05)$.

```
# Kolmogorov test
ks.test(modelo_lineal$residuals, "pnorm",
        mean = mean(modelo_lineal$residuals),
        sd = sd(modelo_lineal$residuals))
```

```
## One-sample Kolmogorov-Smirnov test
##
## data: modelo_lineal$residuals
## D = 0.16652, p-value = 0.7405
## alternative hypothesis: two-sided
```

Varianza constante de los residuos (Homocedasticidad):

La variabilidad de los residuos debe de ser constante a lo largo del eje X.

```
# Varianza constante de los residuos (Homocedasticidad):
ggplot(data = datos, aes(x = prediccion, y = residuos)) + geom_point(aes(
color = residuos)) +
  scale_color_gradient2(low = "blue3", mid = "grey", high = "red") +
  geom_segment(aes(xend = prediccion, yend = 0), alpha = 0.2) +
  geom_smooth(se = FALSE, color = "firebrick") +
  labs(title = "Distribucion de los residuos", x = "prediccion modelo", y =
"residuo") +
  geom_hline(yintercept = 0) +
  theme_bw() +
  theme(plot.title = element_text(hjust = 0.5), legend.position = "none")
```

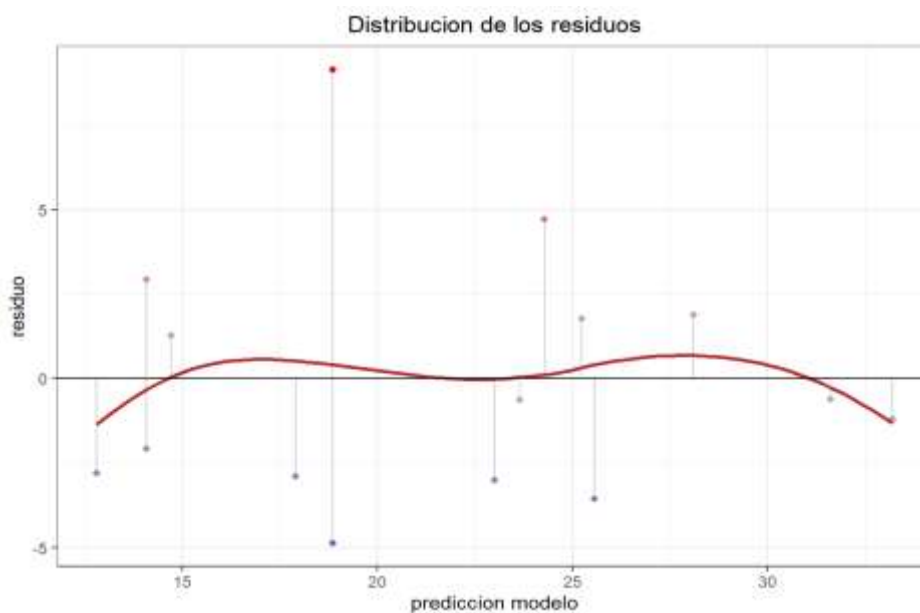


Figura 16. Distribución de los residuos.

Los residuos se comportan de manera aleatoria, no existe un patrón, existe homocedasticidad. Verifiquemos con un test de hipótesis.

Prueba de hipótesis de Breush Pagan

```
# Test de Breush-Pagan
```

```
library(lmtest)
```

```
bptest(modelo_lineal)
```

```
## studentized Breusch-Pagan test  
##  
## data: modelo_lineal  
## BP = 0.5092, df = 1, p-value = 0.4755
```

1) Prueba de hipótesis

Ho: existe homocedasticidad.

Ha: falta de homocedasticidad

2) nivel de significancia 0.05

3) Prueba estadística. Breuch Pagan

BP = 0.5092 con p=0.4755

4) decisión: $p(0.4755) > \alpha(0.05)$, se acepta la Ho de existencia de homocedasticidad.

(Amat, 2016) Ni la representación gráfica ni el contraste de hipótesis muestran evidencias que haga sospechar falta de homocedasticidad.

Autocorrelación de residuos:

Cuando se trabaja con intervalos de tiempo, es muy importante comprobar que no existe autocorrelación de los residuos, es decir que son independientes. Esto puede hacerse detectando visualmente patrones en la distribución de los residuos cuando se ordenan según han registrado o con el test de Durbin-Watson `dwt()` del paquete Car (Amat, 2016).

```
# Autocorrelacion de residuos:  
ggplot(data = datos, aes(x = seq_along(residuos), y = residuos)) +  
  geom_point(aes(color = residuos)) +  
  scale_color_gradient2(low = "blue3", mid = "grey", high = "red") +  
  geom_line(size = 0.3) +  
  labs(title = "Distribucion de los residuos", x = "index", y =  
  "residuo")+  
  geom_hline(yintercept = 0) +  
  theme(plot.title = element_text(hjust = 0.5), legend.position = "none")
```

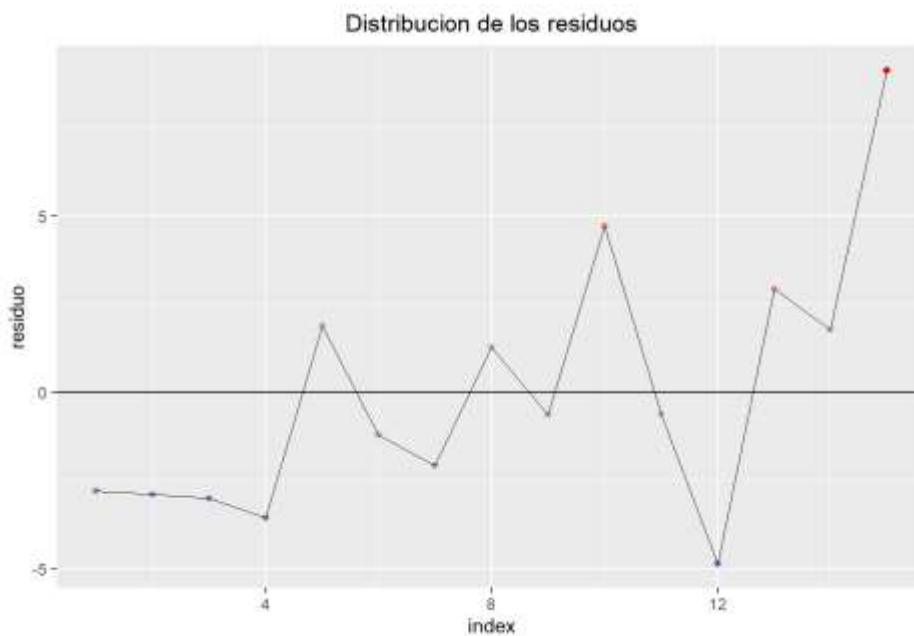


Figura 17. Distribución de los residuos.

En este caso, la representación de los residuos no muestra ninguna tendencia.

Test de Durbin Watson

```
#test de Durwin Watson  
library(lmtest)  
dwtest(modelo_lineal)
```

```
## Durbin-Watson test
##
## data: modelo_lineal
## DW = 1.2942, p-value = 0.06938
## alternative hypothesis: true autocorrelation is greater than 0
```

1) Prueba de hipótesis

Ho: no existe autocorrelación.

Ha: existe autocorrelación

2) nivel de significancia 0.05

3) Prueba estadística. Durwin watson

DW = 1.2942 con $p=0.06938$

4) decisión: $p(0.06938) > \alpha(0.05)$, se acepta la Ho de no existencia de autocorrelación.

h) Identificación de valores atípicos: *outliers*, *leverage* y *bservaciones influyentes* (Amat, 2016).

- **Outlier u observación atípica:** Observaciones que no se ajustan bien al modelo. El valor real se aleja mucho del valor predicho, por lo que su residuo es excesivamente grande. En una representación bidimensional se corresponde con desviaciones en el eje Y.
- **Observación influyente:** Observación que influye sustancialmente en el modelo, su exclusión afecta al ajuste. No todos los *outliers* tienen por qué ser influyentes.
- **Observación con alto leverage:** Observación con un valor extremo para alguno de los predictores. En una representación bidimensional se corresponde con desviaciones en el eje X. Son potencialmente puntos influyentes.

Independientemente de que el modelo se haya podido aceptar, siempre es conveniente identificar si hay algún posible *outlier*, *observación con alto leverage* u observación altamente influyente, puesto que podría estar

condicionando en gran medida el modelo. La eliminación de este tipo de observaciones debe de analizarse con detalle y dependiendo de la finalidad del modelo. Si el fin es predictivo, un modelo sin estas observaciones puede lograr mayor precisión en la mayoría de casos. Sin embargo, es muy importante prestar atención a estos valores ya que, de no ser errores de medida, pueden ser los casos más interesantes. El modo adecuado a proceder cuando se sospecha de algún posible valor atípico o influyente es calcular el modelo de regresión incluyendo y excluyendo dicho valor.

Test de Benferroni para detectar Outliers.

```
# Prueba de Benferroni para detectar outliers
library(car)
outlierTest(modelo_lineal, cutoff=Inf, n.max=4)
```

```
##      rstudent unadjusted p-value Bonferroni p
## 15  3.3004752          0.0063352    0.095028
## 12 -1.3674927          0.1965400         NA
## 10  1.3173787          0.2123100         NA
##  4 -0.9728475          0.3498300         NA
```

El posible valor atípico ubicado en la posición 15, no es significativo como outlier y mucho menos los otros valores (12,10,4).

Veamos ahora con otros procedimientos.

```
# Identificación de valores atípicos: outliers, leverage y observaciones
influyentes
library(ggrepel)
library(dplyr)
```

```
datos$studentized_residual <- rstudent(modelo_lineal)
ggplot(data = datos, aes(x = prediccion, y = abs(studentized_residual))) +
  geom_hline(yintercept = 3, color = "grey", linetype = "dashed") +
  # Se identifican en rojo residuos studentized absolutos > 3
  geom_point(aes(color = ifelse(abs(studentized_residual) > 3, "red",
"black"))) +
  scale_color_identity() +
```



```
#se muestra el equipo al que pertenece la observacion atipica,
geom_text_repel(data = filter(datos, abs(studentized_residual) > 3),
                aes(label = "" )) +
labs(title = "Distribucion de los residuos studentized", x = "prediccion
modelo") +
theme_bw()+
theme(plot.title = element_text(hjust = 0.5), legend.position = "none")
```

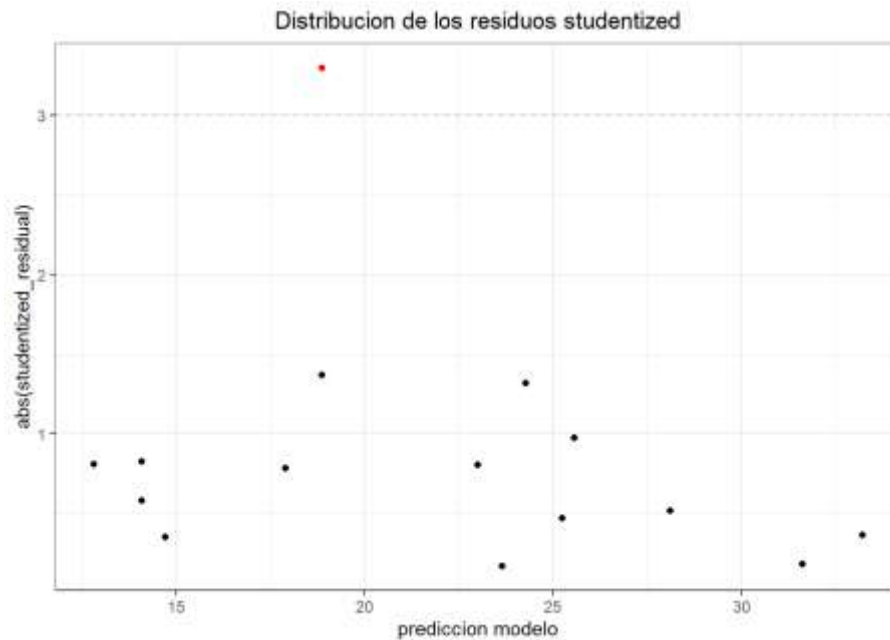


Figura 18. Distribución de los residuos estudiantizados para observar valores outliers

Se observa un punto por encima de 3. ¿Cuál es el valor reconocido como outliers?

```
datos %>% filter(abs(studentized_residual) > 3)
```

```
##      X Y prediccion residuos studentized_residual
## 1 35 28   18.86345 9.136549                3.300475
```

```
which(abs(datos$studentized_residual) > 3)
```

```
## [1] 15
```

El estudio de los residuos *studentized* identifica al dato 15 con valores (35 y 28), como una posible observación atípica.

El hecho de que un valor sea atípico o con alto grado de *leverage* no implica que sea influyente en el conjunto del modelo. Sin embargo, si un valor es influyente, suele ser o atípico o de alto *leverage*. En R se dispone de la función `outlierTest()` del paquete `car` y de las funciones `influence.measures()`, `influencePlot()` y `hatvalues()` para identificar las observaciones más influyentes en el modelo.

```
library(car)
summary(influence.measures(model = modelo_lineal))
```

```
## Potentially influential observations of
## lm(formula = Y ~ X, data = datos) :
##
##      dfb.1_ dfb.X dffit cov.r   cook.d hat
## 6   0.14  -0.20 -0.23  1.62_*  0.03  0.29
## 11  0.05  -0.08 -0.10  1.52_*  0.01  0.23
## 15  0.73  -0.41  0.98  0.35_*  0.27  0.08
```

No se detectan valores influyentes ni atípicos significativos que puedan modificar el modelo, pero podemos identificar gráficamente cuales eran esos posibles valores.

```
influencePlot(model = modelo_lineal)
```

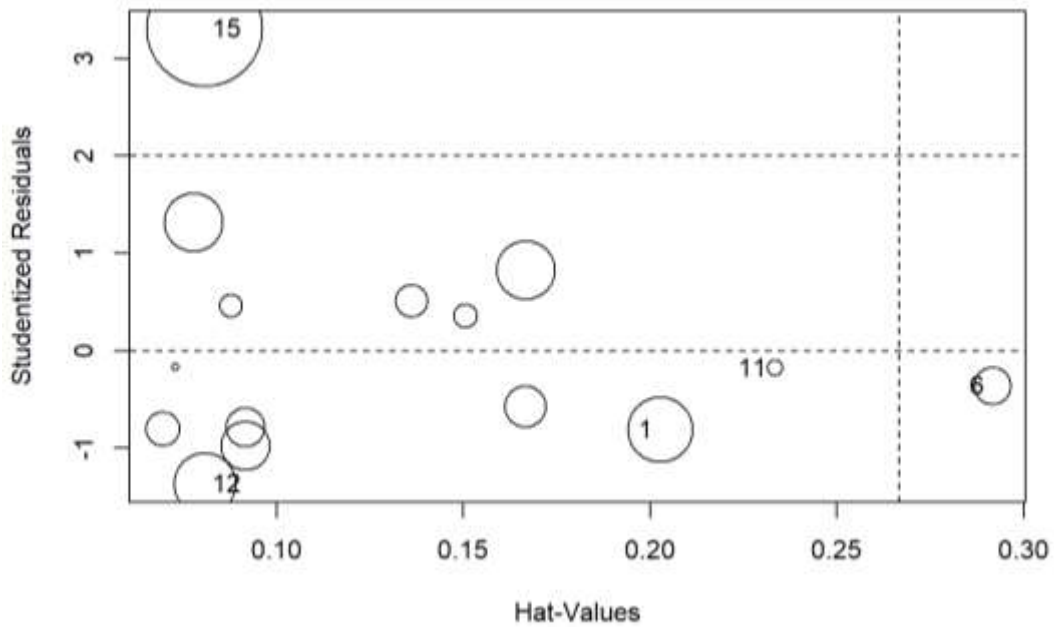


Figura 19. Posibles valores influyentes

##	StudRes	Hat	CookD
## 1	-0.8088578	0.20268332	0.085429608
## 6	-0.3633091	0.29151053	0.029097425
## 11	-0.1773640	0.23339116	0.005174119
## 12	-1.3674927	0.08071941	0.076951250
## 15	3.3004752	0.08071941	0.271575682

i) Prediciendo nuevos valores.

El modelo obtenido cumple con los supuestos y no presenta valores atípicos o influyentes:

$$Y = 7.70280 + 0.31888(X)$$

Tanto el intercepto como la variable X (gastos en publicidad) son estadísticamente significativas al 5%, en ambos casos [$p(0.0074) < \alpha(0.05)$ y $p(0.0000266) < \alpha(0.05)$], con tendencia positiva. El modelo tiene un coeficiente de determinación de 75.44% (el modelo calculado explica el 75.44% de la variabilidad presente en la variable respuesta (volumen de ventas) mediante la variable independiente (gastos en publicidad) y un ANVA significativo.

¿Cuánto será volumen de ventas, si se invierte 63 mil pesetas en publicidad?

```
# prediciendo nuevos valores, cuando X = 63
predict_value <- predict(modelo_lineal, data.frame(X= c(63)))
predict_value
```

```
##          1
## 27.79197
```

Si se invierten 63 mil pesetas en publicidad, el volumen de ventas será de 27.79 millones de pesetas.

Por cada unidad (mil pesetas) que se invierte en gastos de publicidad, el volumen de ventas se incrementa en promedio 0.3188 millones de pesetas.

20. Ejemplo manual para mostrar los cálculos de las Bandas de Confianza:

En la producción de herramientas la deformación del acero a cierta temperatura puede afectar la dureza del acero, para investigar esta relación se ha tomado la siguiente muestra.

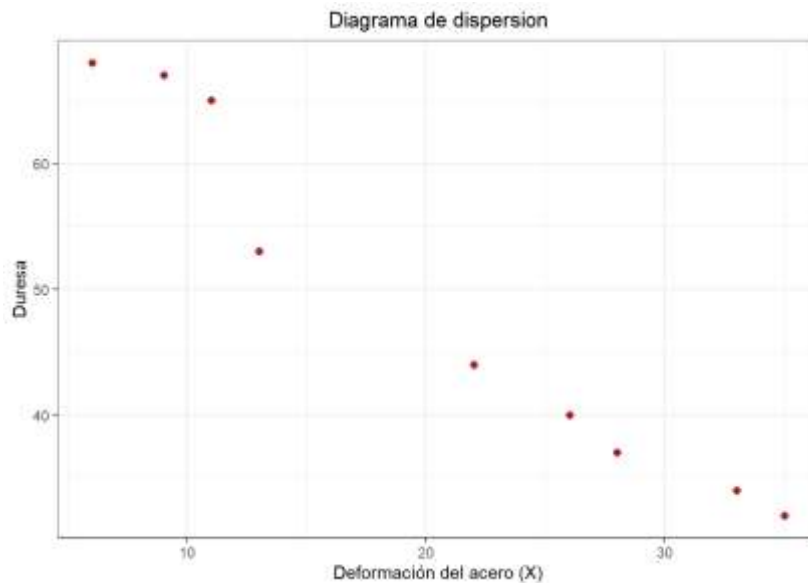
Deformación: mm ² (X)	6	9	11	13	22	26	28	33	35
Dureza en Kg/mm ² (Y)	68	67	65	53	44	40	37	34	32

Las sumas necesarias son:

$$\sum X_i = 183 \quad \sum Y_i = 440 \quad \sum X_i^2 = 4665 \quad \sum Y_i^2 = 23232 \quad \sum X_i Y_i = 7701$$

$$\bar{X} = 20.33 \quad \bar{Y} = 48.89$$

1) Diagrama de dispersión

**Figura 20.** Diagrama de dispersión

2) Obteniendo parámetros

Transformado los datos en desviaciones:

$$\sum x_i^2 = \sum X_i^2 - \frac{(\sum X_i)^2}{n} = 4665 - \frac{183^2}{9} = 944$$

$$\sum x_i y_i = \sum X_i Y_i - \frac{(\sum X_i)(\sum Y_i)}{n} = 7701 - \frac{183(440)}{9} = -1245.67$$

Calculando los parámetros

$$\hat{\beta}_1 = \frac{\sum x_i y_i}{\sum x_i^2} = \frac{-1245.67}{944} = -1.32$$

$$\hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \bar{X} = 48.89 - (-1.32)(20.33) = 75.73$$

La recta de regresión mínimo cuadrática es:

$$\hat{Y} = 75.73 - 1.32X_i$$

3) recta de regresión lineal simple

Usando puntos extremos para ajustar la recta:

$$\hat{Y}_1 = 75.73 - 1.32(6) = 67.81$$

$$\hat{Y}_9 = 75.73 - 1.32(35) = 29.53$$

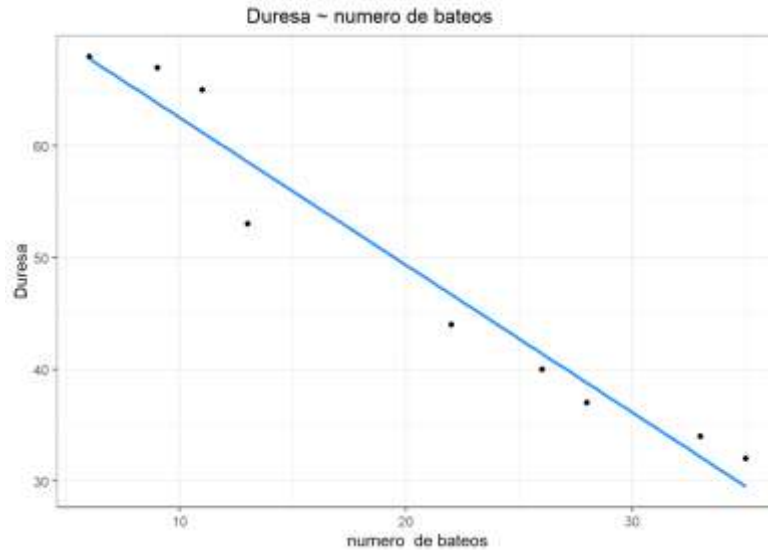


Figura 21. Modelo de rls

4) intervalo de confianza

Intervalo de confianza para el parámetro β

$$p(\hat{\beta}_1 - t_{(n-2),\alpha/2}S_b < \beta_1 < \hat{\beta}_1 + t_{(n-2),\alpha/2}S_b) = 1 - \alpha$$

$$p(-1.32 - 2.365(0.1073) < \beta_1 < -1.32 + 2.365(0.1073)) = 1 - 0.05$$

$$p(-1.57 < \beta_1 < -1.07) = 0.95$$

Donde:

$$S_e^2 = \frac{\sum Y_i^2 - \hat{\beta}_0 \sum Y_i - \hat{\beta}_1 \sum X_i Y_i}{n - 2} = \frac{23232 - (75.73)(440) - (-1.32)(7701)}{9 - 2} = 10.87$$

$$S_e = \sqrt{10.87} = 3.297$$

$$\text{Var}(\hat{\beta}_1) = S_{\hat{\beta}_1}^2 = \frac{S_e^2}{\sum x_i^2} = \frac{10.87}{944} = 0.0115$$

$$\text{entonces: } S_b^2 = \sqrt{0.0115} = 0.1073$$

Por otra fórmula:

$$S_b = \frac{S_e}{\sqrt{\sum x_i^2}} = \frac{3.297}{\sqrt{944}} = 0.1073$$

$$t_t = t_{(n-2),\alpha/2} = t_{7,0.025} = 2.365$$

Intervalo de confianza para el parámetro $Y = a + bX$ pronosticación ($E(Y/X)$)

$$p \left[\hat{Y}_p - t_{\alpha/2(n-2)} S_{\hat{Y}_p} < E(Y/X) < \hat{Y}_p + t_{\alpha/2(n-2)} S_{\hat{Y}_p} \right] = 1 - \alpha$$

$$t_t = t_{(n-2),\alpha/2} = t_{7,0.025} = 2.365$$

$$S_{\hat{Y}_p} = S_e \sqrt{\frac{1}{n} + \frac{(X_p - \bar{X})^2}{\sum x_i^2}} = 3.297 \sqrt{\frac{1}{9} + \frac{(X_p - 20.33)^2}{944}} = A_o$$

$$p \left[\hat{Y}_p - 2.365 A_o < E(Y/X) < \hat{Y}_p + 2.365 A_o \right] = 1 - \alpha$$

Para $X_p=6$:

$$p \left[67.81 \pm 2.365 \left((3.297) \sqrt{\frac{1}{9} + \frac{(X_p - 20.33)^2}{944}} \right) < E(Y/X) < \right] = 0.95$$

$$p[63.34 < E(Y/X) < 72.28] = 0.95$$

Zona de confianza para $E(Y/X)$

Cuanto más se desvía X del promedio, los valores tabulares son mayores y por lo tanto los intervalos son más amplios.

Intervalo de confianza para el parámetro Y predicción

$$p \left[\hat{Y}_p \pm t_{\alpha/2(n-2)} S_{\hat{Y}_p - Y_p} < E(Y/X) < \right] = 1 - \alpha$$

$$t_t = t_{(n-2),\alpha/2} = t_{7,0.025} = 2.365$$

$$S_{\hat{Y}_p - Y_p} = S_e \sqrt{1 + \frac{1}{n} + \frac{(X_p - \bar{X})^2}{\sum x_i^2}} = 3.297 \sqrt{1 + \frac{1}{9} + \frac{(X_p - 20.33)^2}{944}} = A_1$$

Reemplazando los valores de X, en la tabla

Se forma la siguiente tabla

Tabla 5.

X_i	\hat{Y}_p	$2.365A_o$	$\hat{Y}_p - 2.365A_o$	$\hat{Y}_p + 2.365A_o$	$\hat{Y}_p - 2.365A_1$	$\hat{Y}_p + 2.365A_1$
6	67.81	4.47	63.34	72.28	58.82	76.80
9	63.85	3.88	59.97	67.73	55.14	72.56
11	61.21	3.52	57.69	64.73	52.66	69.76
13	58.57	3.20	55.37	61.77	50.14	67.00
20.33	48.89	2.60	46.29	51.49	40.67	57.11
22	46.69	2.63	44.06	49.32	38.46	54.92
26	41.41	2.97	38.44	44.38	33.07	49.75
28	38.77	3.25	35.52	42.02	30.32	47.22
33	32.17	4.13	28.04	36.30	23.34	41.00
35	29.53	4.54	24.99	37.07	20.51	38.55

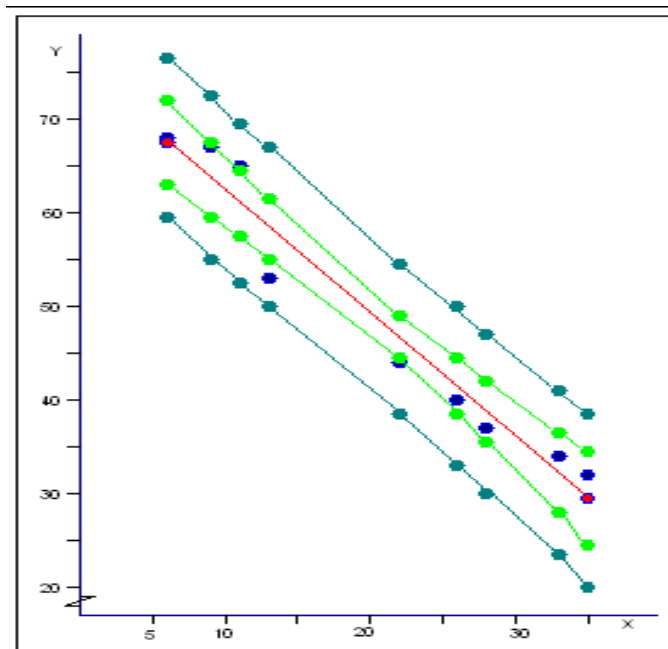


Figura 22. Zona de confianza para Y_o

Los límites pueden usarse para estimar el intervalo correspondiente a un X dado; ejemplo: estimar el valor de Y para X=18

$$\hat{Y} = 75.73 - 1.32(18) = 51.97$$

Halando el intervalo de predicción:

$$p[51.97 \pm 2.365(3.48) < Y_o <] = 0.95$$

$$p[43.74 , 60.20] = 0.05.$$

21. Ejemplo 2 en R.

22. Regresión lineal simple con valor atípico.

Un analista de deportes quiere saber si existe una relación entre el número de bateos que realiza un equipo de béisbol y el número de corridas que consigue. En caso de existir y de establecer un modelo, podría predecir el resultado del partido (Amat, 2016).

Ingreso de datos (también pueden ser importados de una base de datos)

```
equipos <- c("Texas", "Boston", "Detroit", "Kansas", "St.", "New_S.",
            "Ne_w_Y.", "Milwaukee", "Colorado", "Houston", "Baltimore",
            "Los_An.", "Chica go", "Cincinnati", "Los_P.", "Philadelphia",
            "Chicago", "Cleveland", "Ari zona", "Toronto", "Minnesota",
            "Florida", "Pittsburgh", "Oakland", "Tampa", "Atlanta",
            "Washington", "San.F", "San.I", "Seattle")
numero_bateos <- c(5659, 5710, 5563, 5672, 5532, 5600, 5518, 5447, 5544, 5598,
                  5585, 5436, 5549, 5612, 5513, 5579, 5502, 5509, 5421, 5559,
                  5487, 5508, 5421, 5452, 5436, 5528, 5441, 5486, 5417, 5421)
corridas <- c(855, 875, 787, 730, 762, 718, 967, 721, 735, 615, 708, 644, 654, 735,
             667, 713, 654, 704, 731, 743, 619, 625, 610, 645, 707, 641, 624, 570,
             593, 556)
```

Creando un data frame y observando una parte de los datos

```
datos <- data.frame(equipos, numero_bateos, corridas)
head(datos)
```

```
##   equipos numero_bateos corridas
## 1   Texas           5659      855
## 2   Boston           5710      875
## 3   Detroit          5563      787
## 4   Kansas           5672      730
## 5     St.            5532      762
## 6   New_S.          5600      718
```

La estructura (str) sirve para observar en R la estructura de ingreso de la data (numéricas, o cualitativas)

```
str(datos)
```

```
## 'data.frame':   30 obs. of  3 variables:
## $ equipos      : chr  "Texas" "Boston" "Detroit" "Kansas" ...
## $ numero_bateos: num  5659 5710 5563 5672 5532 ...
## $ corridas     : num  855 875 787 730 762 718 967 721 735 615 ...
```

Podemos ver que la data set está conformada por 30 observaciones y 3 variables:

- equipos está reconocido como variable cualitativa(chr)
- numero_bateos, como variable numérica
- runs, como variable numérica.

El número de corridas (Y) estará en función al número de bateos (X), dicho de otra manera, en número de bateos influye en el número de corridas.

a) Observar las variables descriptivamente.

```
# histograma para las variables
require(ggplot2)
par(mfrow = c(1, 2))
hist(datos$numero_bateos, breaks = 10, main = "", xlab = "número de bateos", ylab="Frecuencia", border = "darkred")
```

```
hist(datos$corridas, breaks = 10, main = "", xlab = "corridas", ylab="Frecuencia", border = "blue")
```

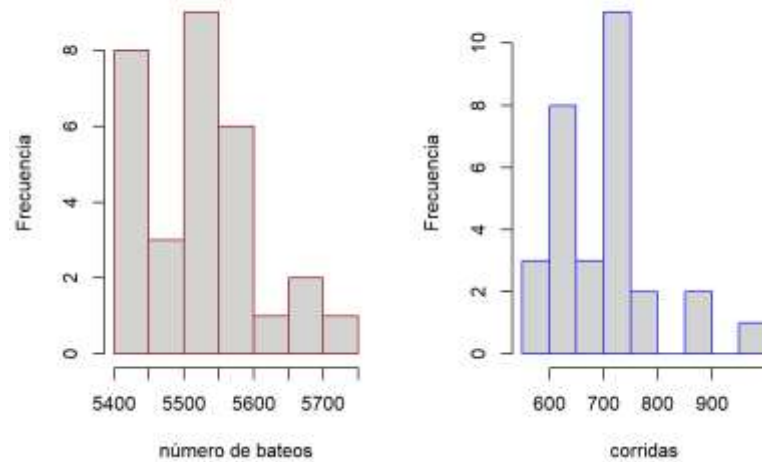


Figura 23. Histograma de frecuencias para el número de bateos y numero de corridas.

```
summary(datos)
```

```
##      equipos      numero_bateos      corridas
## Length:30      Min.   :5417      Min.   :556.0
## Class :character 1st Qu.:5448      1st Qu.:629.0
## Mode  :character Median :5516      Median :705.5
##                      Mean  :5524      Mean   :696.9
##                      3rd Qu.:5575      3rd Qu.:734.0
##                      Max.   :5710      Max.   :967.0
```

El promedio del número de bateos es 5524 con un mínimo de 5417 y máximo número de bateos de 5710. El promedio de corridas es 696 aproximadamente con un mínimo de corridas de 556 y máximo 967.

b) Representación gráfica de las observaciones

El primer paso antes de generar un modelo de regresión es representar los datos para poder intuir si existe una relación y cuantificar dicha relación mediante un coeficiente de correlación. Si en este paso no se detecta la posible relación lineal, no tiene sentido seguir adelante generando un modelo lineal (se tendrían que probar otros modelos no lineales).

Usando la librería ggplot2

```
ggplot(data = datos, mapping = aes(x = numero_bateos, y= corridas)) +
  geom_point(color = "firebrick", size = 3) +
  (labs(title = "Diagrama de dispersión", x = "número de bateos")) +
  theme_bw() +
  theme(plot.title = element_text(hjust =0.5))
```

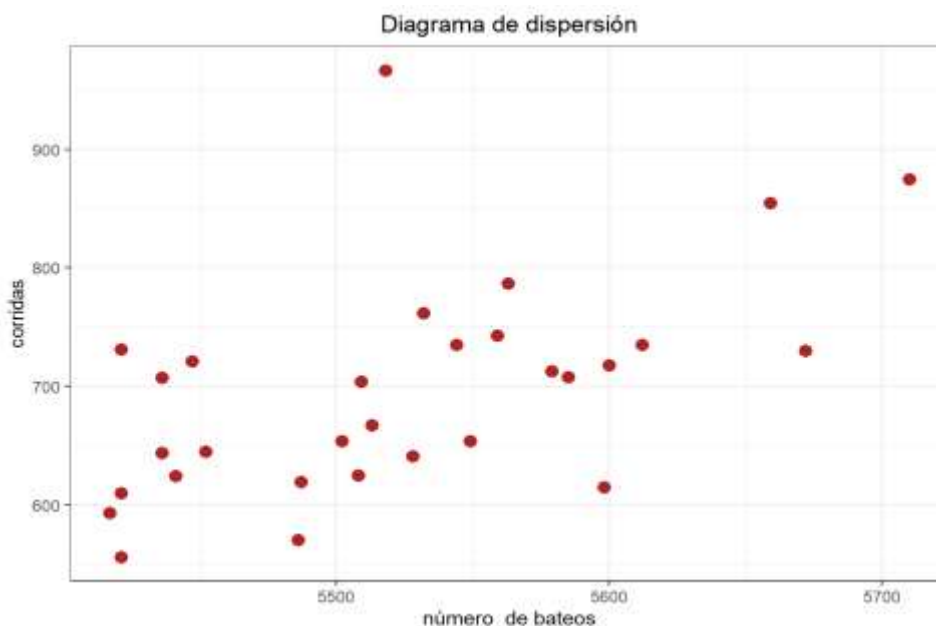


Figura 24. Diagrama de dispersión

La figura muestra un comportamiento (tendencia) positivo, a más bateos más corridas, así mismo, muestra posible outlier y puntos influyentes.

c) comportamiento del coeficiente de correlación

Estadístico de correlación de Pearson (para datos cuantitativos)

```
cor.test(x = datos$numero_bateos, y = datos$corridas, method =
"pearson")
```

```
## Pearson's product-moment correlation
##
## data:  datos$numero_bateos and datos$corridas
## t = 3.477, df = 28, p-value = 0.001673
## alternative hypothesis: true correlation is not equal to 0
## 95 percent confidence interval:
##  0.2354692 0.7592163
## sample estimates:
##      cor
## 0.5491526
```

Resumen:

Coeficiente de correlación	Test de significancia	Intervalo de confianza de r	p-value
$r = 0.549$	$t = 3.477$	0.2354692 - 0.7592163	0.001673

El test de correlación muestra una relación lineal significativa $p(0.001673) < \alpha(0.05)$, de intensidad considerable ($r = 0.549$). Tiene sentido intentar generar un modelo de regresión lineal que permita predecir el número de corridas en función del número de bateos.

d) Cálculo del modelo de regresión lineal simple

```
# Cálculo del modelo de regresión lineal simple
modelo_lineal <- lm(corridas ~ numero_bateos, data=datos)
summary(modelo_lineal)
```

```
## Call:
## lm(formula = corridas ~ numero_bateos, data = datos)
##
## Residuals:
```

```
##      Min      1Q  Median      3Q      Max
## -128.69 -50.54 -19.96   51.10  273.52
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept) -2769.4894   997.0462  -2.778  0.00966 **
## numero_bateos    0.6276    0.1805    3.477  0.00167 **
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 77.63 on 28 degrees of freedom
## Multiple R-squared:  0.3016, Adjusted R-squared:  0.2766
## F-statistic: 12.09 on 1 and 28 DF,  p-value: 0.001673
```

El modelo de regresión lineal simple es:

$$Y = -2769.4894 + 0.6276(\text{numero_bateos})$$

Por cada unidad que se incrementa el número de bateos, el número de corridas aumenta en promedio 0.6276 unidades.

Realizando la inferencia de los parámetros con la prueba t, ambos son significativos. $p < \alpha$, es decir, que tienen importancia en el modelo individualmente.

El coeficiente de determinación R^2 indica que el modelo calculado explica el 30.16 de la variabilidad presente en la variable respuesta (corridas) mediante la variable independiente (*número de bateos*). Podría indicarse un ajuste no muy bueno.

La prueba de Análisis de varianza observa al modelo en su conjunto. El *p-value* obtenido en el test $p(0.001673) < \alpha(0.05)$ determina que es significativamente superior la varianza explicada por el modelo en comparación a la varianza total. Es el parámetro que determina que el modelo en su conjunto es significativo.

Observando en detalle en ANVA

```
anova(modelo_lineal)
```

```
## Analysis of Variance Table
##
## Response: corridas
##           Df Sum Sq Mean Sq F value    Pr(>F)
## numero_bateos  1  72867   72867   12.09 0.001673 **
## Residuals    28 168761    6027
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

Hasta aquí se puede concluir que el modelo es bueno y existe relación entre las variables (excepto que el ajuste R^2 es bajo).

e) Intervalos de confianza para los parámetros del modelo

```
# Intervalos de confianza para los parámetros del modelo
confint(modelo_lineal, level=0.95)
```

```
##           2.5 %           97.5 %
## (Intercept) -4811.8458919 -727.1328213
## numero_bateos  0.2578569  0.9972975
```

Los intervalos de confianza contienen a los coeficientes obtenidos.

Por cada unidad que se incrementa el número de bateos, el número de corridas aumenta en promedio entre 0.2579 y 0.9973 unidades.

f) Representación gráfica del modelo

```
# Representación gráfica del modelo
ggplot(data = datos, mapping = aes(x = numero_bateos, y= corridas)) +
  geom_point(size=3) +
```

```
labs(title = "corridas~numero de bateos", x="numero de bateos") +
geom_smooth(method = "lm", se = FALSE, color = "red") +
theme_bw() +
theme(plot.title = element_text(hjust = 0.4))
```

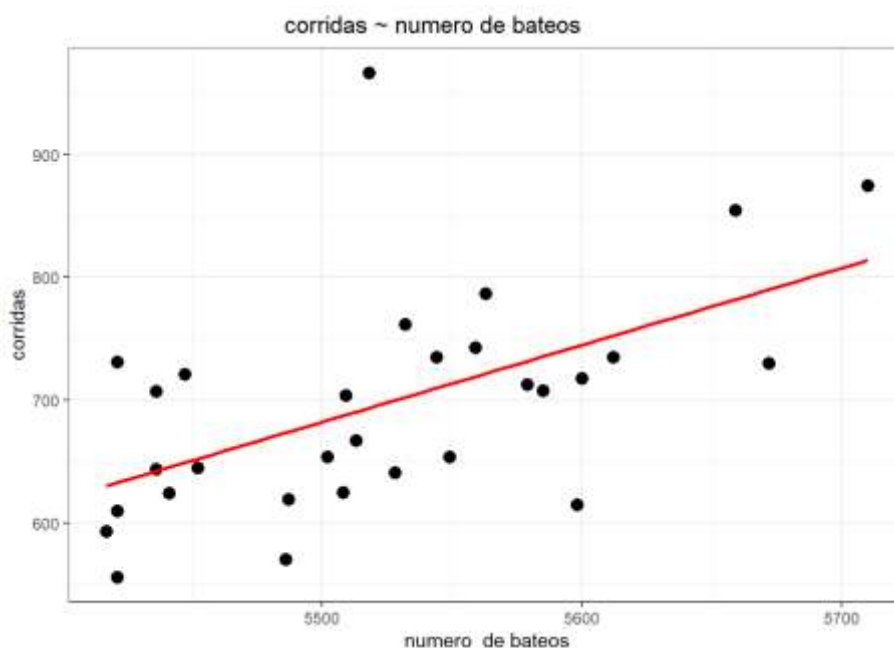


Figura 25. Representación de la línea de regresión

```
names(modelo_lineal)
```

```
## [1] "coefficients" "residuals" "effects" "rank"
## [5] "fitted.values" "assign" "qr" "df.residual"
## [9] "xlevels" "call" "terms" "model"
```

Con la opción `names` muestra los resultados guardados por el modelo (ejemplo: `modelo_lineal$coefficients`) se muestra solo los coeficientes del modelo.

Además de la línea de mínimos cuadrados es recomendable incluir los límites superior e inferior del intervalo de confianza. Esto permite identificar la región en la que, según el modelo generado y para un determinado nivel de confianza, se encuentra el valor promedio de la variable dependiente.

Para poder representar el intervalo de confianza a lo largo de todo el modelo se recurre a la función `predict()` para predecir valores que abarquen todo el eje X . Se añaden al gráfico líneas formadas por los límites superiores e inferiores calculados para cada predicción (Amat, 2016).

```
# predecir valores para y con valores de X originales
nuevos_datos <- data.frame(numero_bateos=
  seq(min(numero_bateos),max(numero_bateos)))
predict_value <- predict(modelo_lineal)
head(predict_value)
```

```
##           1           2           3           4           5           6
## 781.9700 813.9765 721.7226 790.1285 702.2677 744.9430
```

Calculando el error medio cuadrático (RMSE)

```
# Cálculo del error cuadrático medio RMSE:
sqrt(mean(modelo_lineal$residuals^2))
```

```
## [1] 75.00233
```

Bandas de confianza para valores medios.

```
# gráfico de bandas de confianza (una banda). valores medios
par(mfrow = c(1, 1))
puntos <- seq(from = min(datos$numero_bateos),
              to = max(datos$numero_bateos), length.out = 100)
limites_intervalo <- predict(object = modelo_lineal,
                             newdata = data.frame( numero_bateos =
puntos),
                             interval = "confidence", level = 0.95)
head(limites_intervalo, 3)
```

```
##           fit           lwr           upr
## 1 630.0964 581.1740 679.0188
## 2 631.9537 583.9076 679.9998
```

```
## 3 633.8111 586.6322 680.9900
```

```
plot(datos$numero_bateos, datos$corridas, col = "firebrick", pch =
19,
      ylab = "corridas", xlab = "número de bateos",
      main = "Banda de confianza para valores medios")
abline(modelo_lineal, col = 1)
lines(x = puntos, y = limites_intervalo[, 2], type = "l", col = 2,
      lty = 3)
lines(x = puntos, y = limites_intervalo[, 3], type = "l", col = 3,
      lty = 3)
```

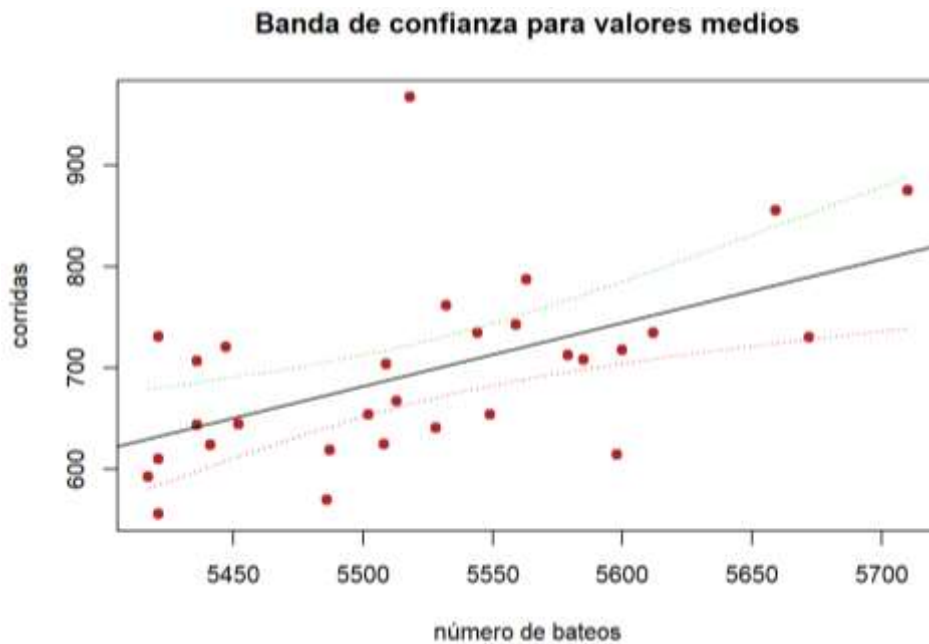


Figura 26. *Bandas de confianza para el ajuste de la regresión.*

La función `geom_smooth()` del paquete `ggplot2` genera la regresión y su intervalo de forma directa.

```
# con ggplot2
ggplot(data = datos, mapping = aes(x = numero_bateos, y = corridas)
) +
  geom_point(color = "firebrick", size = 2) +
  labs(title = "Diagrama de dispersion", x = "numero de bateos") +
  geom_smooth(method = "lm", se = TRUE, color = "black") +
```

```
theme_bw() +
theme(plot.title = element_text(hjust = 0.5))
```

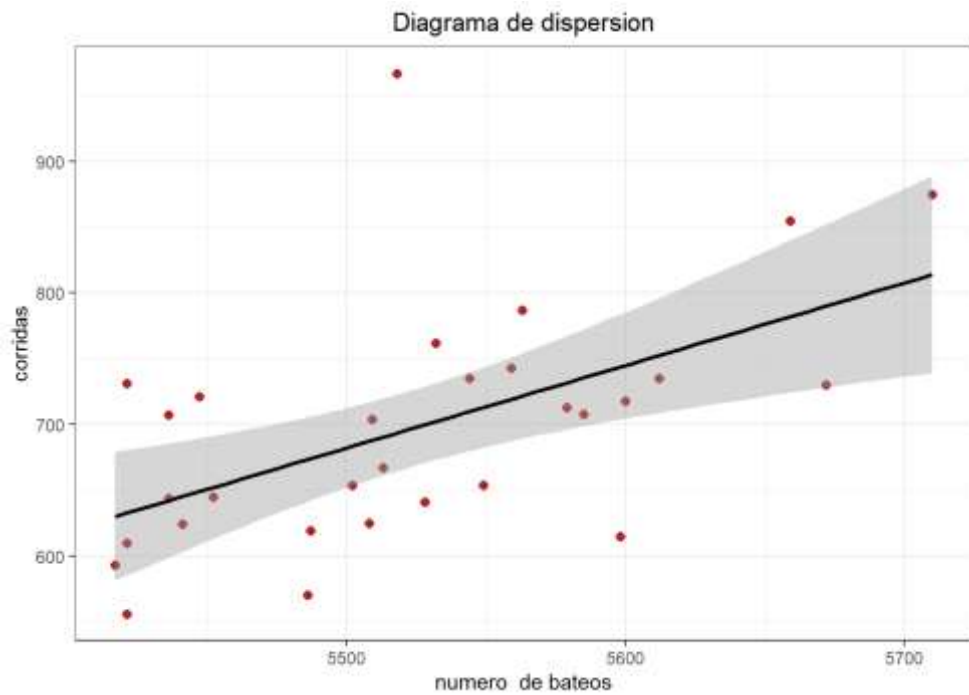


Figura 27. Bandas de confianza para la media y la predicción.

```
# Por defecto incluye la region de 95% de confianza

# dos bandas de confianza
# Grafico dispersion y recta
plot(datos$numero_bateos, datos$corridas, col = "firebrick", pch = 19,
      ylab = "corridas", xlab = "numero de bateos",
      main = "Bandas de confianza para valores medios y predicción")
abline(modelo_lineal, col = 1)

# Intervalos de confianza de la respuesta media:
# valores medios

ic <- predict(modelo_lineal, nuevos_datos, interval = 'confidence')
lines(nuevos_datos$numero_bateos, ic[, 2], lty = 2)
lines(nuevos_datos$numero_bateos, ic[, 3], lty = 2)

# Intervalos de predicción
# para cualquier valor
ic <- predict(modelo_lineal, nuevos_datos, interval = 'prediction')
```

```
lines(nuevos_datos$numero_bateos, ic[, 2], lty = 2, col = 'red')
lines(nuevos_datos$numero_bateos, ic[, 3], lty = 2, col = 'red')
```

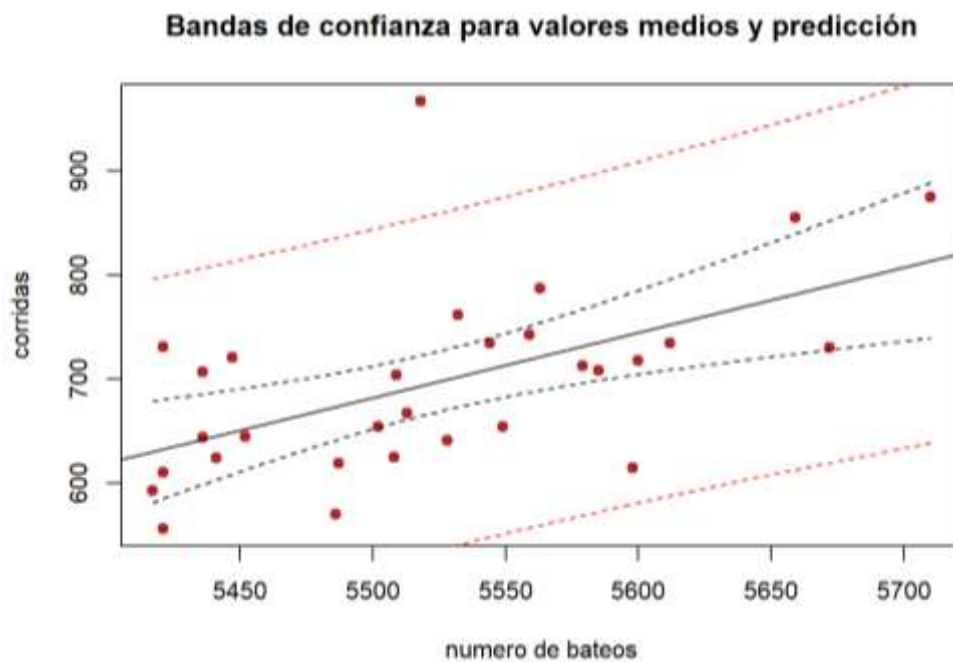


Figura 28. El punto fuera de la banda de confianza de la predicción implica un valor atípico.

Observando los valores reales, predichos y los residuales

```
# observando datos, valores predichos y residuales
datos$prediccion <- modelo_lineal$fitted.values
datos$residuos <- modelo_lineal$residuals
head(datos)
```

```
## equipos numero_bateos corridas prediccion residuos
## 1 Texas 5659 855 781.9700 73.02996
## 2 Boston 5710 875 813.9765 61.02352
## 3 Detroit 5563 787 721.7226 65.27737
## 4 Kansas 5672 730 790.1285 -60.12855
## 5 St. 5532 762 702.2677 59.73226
## 6 New_S. 5600 718 744.9430 -26.94299
```

g) Verificar condiciones para poder aceptar un modelo lineal

Relación lineal entre variable dependiente e independiente:

Se calculan los residuos para cada observación y se representan (scatterplot). Si las observaciones siguen la línea del modelo, los residuos se deben distribuir aleatoriamente entorno al valor 0.

```
# gráfico de residuales
ggplot(data = datos, aes(x = prediccion, y = residuos)) +
  geom_point(aes( color = residuos)) +
  scale_color_gradient2(low = "blue3", mid = "grey", high = "red") +
  geom_hline(yintercept = 0) +
  geom_segment(aes(xend = prediccion, yend = 0), alpha = 0.2) +
  labs(title = "Distribucion de los residuos", x = "prediccion modelo",
       y = "residuo") +
  theme_bw() +
  theme(plot.title = element_text(hjust = 0.5), legend.position = "none")
```

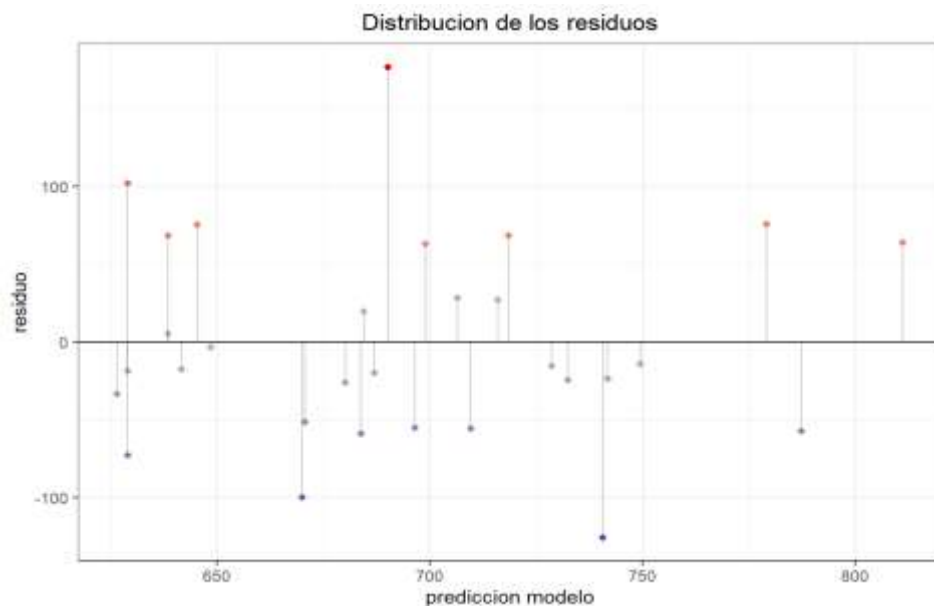


Figura 29. Distribución de residuos.

Los residuos se distribuyen de forma aleatoria entorno al 0 por lo que se acepta la linealidad.

Distribución normal de los residuos:

Los residuos se deben distribuir de forma normal con media 0. Para comprobarlo se recurre a histogramas, a los cuantiles normales o a un test de contraste de normalidad.

Histograma de residuos.

```
par(mfrow = c(1, 1))  
# Distribución normal de los residuos:  
ggplot(data = datos, aes(x = residuos)) +  
  geom_histogram(aes(y = ..density..)) +  
  labs(title = "Histograma de los residuos") +  
  theme_bw() +  
  theme(plot.title = element_text(hjust = 0.5))
```

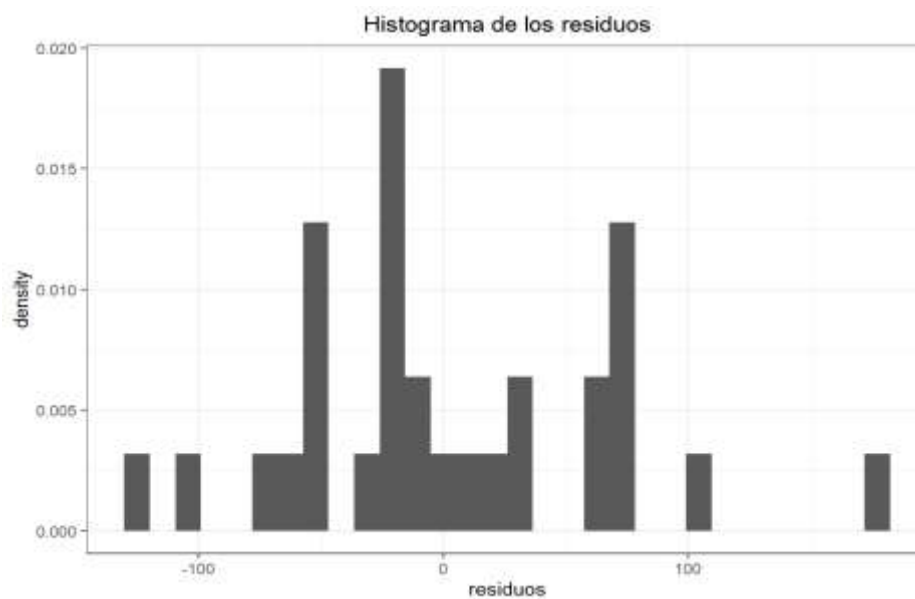


Figura 30. Histograma de densidad de los residuos.

No se evidencia una clara distribución normal de los residuos.

Gráfico de cuantiles.

```
# grafico de cuantiles  
qqnorm(modelo_lineal$residuals, main = "Gráfico de cuantiles QQ plot",  
        col = "darkred")  
qqline(modelo_lineal$residuals)
```

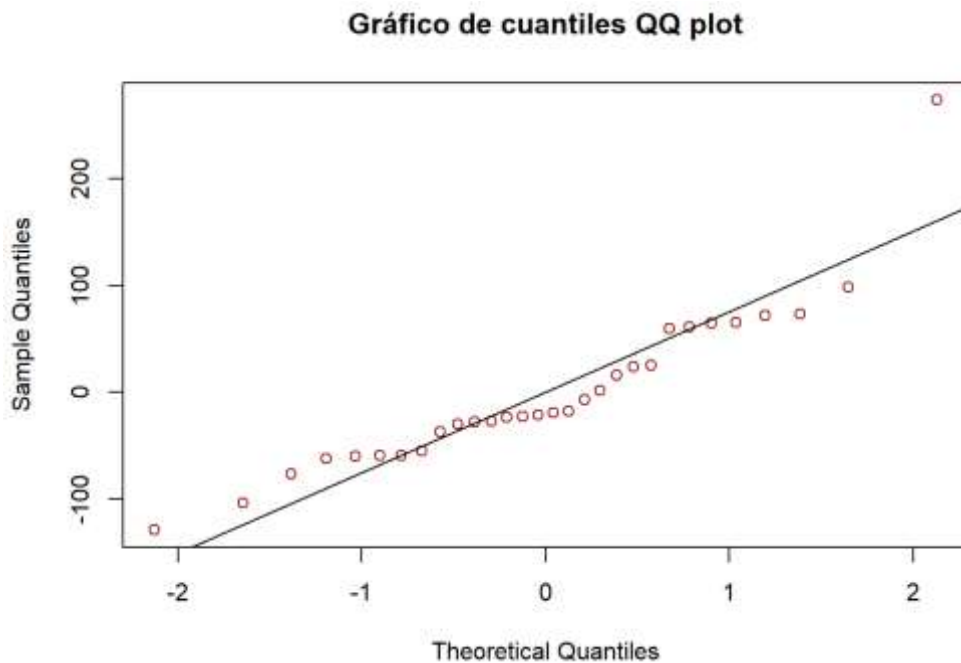


Figura 31. Gráfico de cuantiles de los residuos.

Un punto se aleja de la línea, posible no normalidad en los errores, finalmente realizamos una prueba de hipótesis de normalidad, por la cantidad de datos el estadístico más adecuado sería el test de Shapiro.

Prueba de hipótesis

```
# Test de normalidad
#shapiro.test
shapiro.test(modelo_lineal$residuals)
```

```
## Shapiro-Wilk normality test
##
## data:  modelo_lineal$residuals
## W = 0.88535, p-value = 0.003751
```

1) Prueba de hipótesis

H_0 : existe normalidad en los residuos

H_a : no existe normalidad en los residuos

2) nivel de significancia 0.05

3) Prueba estadística. Shapiro y Wilks

$W = 0.88535$ con $p=0.003751$

4) decisión: $p(0.003751) < \alpha(0.05)$, se rechaza la H_0 , los residuos no siguen una distribución normal, a diferencia de test de Kolmogorov.

Test de Kolmogorov- Smirnov

```
# Kolmogorov test
ks.test(modelo_lineal$residuals, "pnorm",
        mean = mean(modelo_lineal$residuals),
        sd = sd(modelo_lineal$residuals))
```

```
## One-sample Kolmogorov-Smirnov test
##
## data:  modelo_lineal$residuals
## D = 0.15726, p-value = 0.4062
## alternative hypothesis: two-sided
```

Varianza constante de los residuos (Homocedasticidad):

La variabilidad de los residuos debe de ser constante a lo largo del eje X.

```
# Varianza constante de los residuos (Homocedasticidad):

ggplot(data = datos, aes(x = prediccion, y = residuos)) + geom_point(aes(
color = residuos)) +
  scale_color_gradient2(low = "blue3", mid = "grey", high = "red") +
  geom_segment(aes(xend = prediccion, yend = 0), alpha = 0.2) +
  geom_smooth(se = FALSE, color = "firebrick") +
  labs(title = "Distribucion de los residuos", x = "prediccion modelo", y
= "residuo") +
  geom_hline(yintercept = 0) +
  theme_bw() +
  theme(plot.title = element_text(hjust = 0.5), legend.position = "none")
```

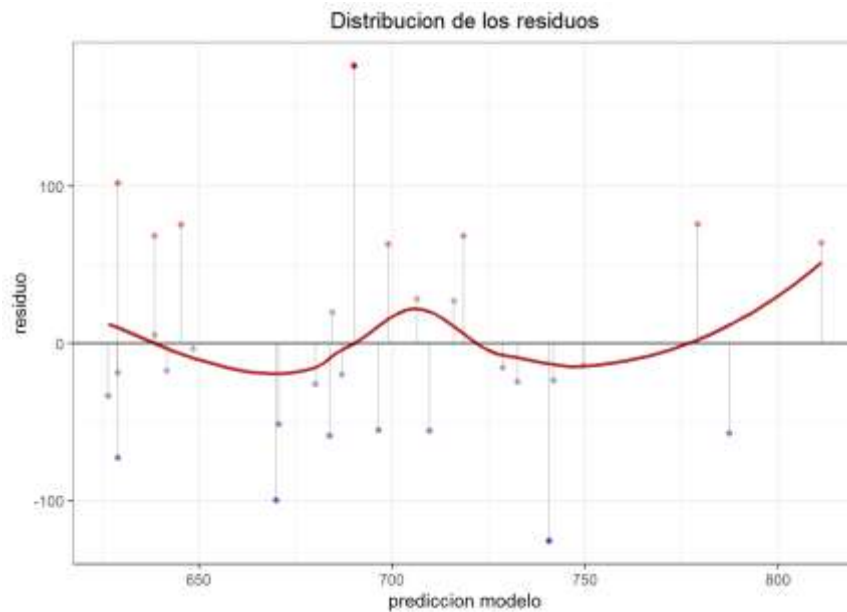



Figura 32. Distribución de los residuos con valores predichos.

Los residuos se comportan de manera aleatoria, no existe un patrón, existe homocedasticidad. Verifiquemos con un test de hipótesis.

Prueba de hipótesis de Pagan

```
# Test de Breush-Pagan
library(lmtest)
bptest(modelo_lineal)
```

```
## studentized Breusch-Pagan test
##
## data: modelo_lineal
## BP = 0.00011221, df = 1, p-value = 0.9915
```

Ni la representación gráfica ni el contraste de hipótesis muestran evidencias que haga sospechar falta de homocedasticidad.

1) Prueba de hipótesis

Ho: existe homocedasticidad.

Ha: falta de homocedasticidad

2) nivel de significancia 0.05

3) Prueba estadística. Pagan

$$W = 0.0001122 \text{ con } p=0.09915$$

5) decisión: $p(0.9915) > \alpha(0.05)$, se acepta la H_0 de existencia de homocedasticidad.

Autocorrelación de residuos:

Cuando se trabaja con intervalos de tiempo, es muy importante comprobar que no existe autocorrelación de los residuos, es decir que son independientes. Esto puede hacerse detectando visualmente patrones en la distribución de los residuos cuando se ordenan según han registrado o con el test de Durbin-Watson `dwt()` del paquete `Car` (Amat, 2016).

```
# Autocorrelacion de residuos:
ggplot(data = datos, aes(x = seq_along(residuos), y = residuos)) +
  geom_point(aes(color = residuos)) +
  scale_color_gradient2(low = "blue3", mid = "grey", high = "red") +
  geom_line(size = 0.3) +
  labs(title = "Distribucion de los residuos", x = "index", y =
"residuo")+
  geom_hline(yintercept = 0) +
  theme(plot.title = element_text(hjust = 0.5), legend.position =
"none")
```

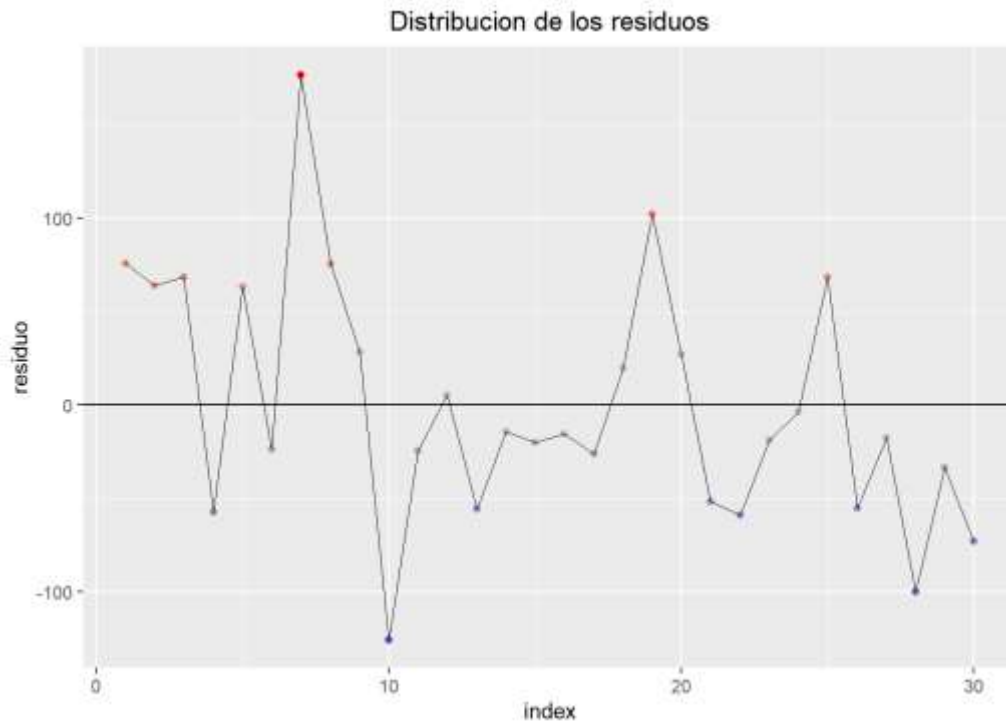


Figura 33. *Distribución de los residuos*

En este caso, la representación de los residuos no muestra ninguna tendencia.

Test de Durbin Watson

```
library(lmtest)
dwtest(modelo_lineal)
```

```
## Durbin-Watson test
##
## data: modelo_lineal
## DW = 1.59, p-value = 0.1045
## alternative hypothesis: true autocorrelation is greater than 0
```

1) Prueba de hipótesis

H_0 : no existe autocorrelación.

H_a : existe autocorrelación

2) nivel de significancia 0.05

3) Prueba estadística. Durwin watson

$$DW = 1.59 \text{ con } p=0.1045$$

4) decisión: $p(0.1045) > \alpha(0.05)$, se acepta la H_0 de no existencia de autocorrelación.

h) Identificación de valores atípicos: *outliers*, *leverage* y observaciones influyentes (Amat, 2016).

Outlier u observación atípica: Observaciones que no se ajustan bien al modelo. El valor real se aleja mucho del valor predicho, por lo que su residuo es excesivamente grande. En una representación bidimensional se corresponde con desviaciones en el eje Y . es una observación que es numéricamente distante del resto de los datos. **Observación influyente:** Observación que influye sustancialmente en el modelo, su exclusión afecta al ajuste. No todos los *outliers* tienen por qué ser influyentes. Punto que tiene impacto en las estimativas del modelo.

- **Observación con alto leverage:** Observación con un valor extremo para alguno de los predictores. En una representación bidimensional se corresponde con desviaciones en el eje X . Son potencialmente puntos influyentes.

Independientemente de que el modelo se haya podido aceptar, siempre es conveniente identificar si hay algún posible *outlier*, *observación con alto leverage* u observación altamente influyente, puesto que podría estar condicionando en gran medida el modelo. La eliminación de este tipo de observaciones debe de analizarse con detalle y dependiendo de la finalidad del modelo. Si el fin es predictivo, un modelo sin estas observaciones puede lograr mayor precisión en la mayoría de casos. Sin embargo, es muy importante prestar atención a estos valores ya que, de no ser errores de medida, pueden ser los casos más interesantes. El modo adecuado a proceder cuando se sospecha de algún posible valor atípico o influyente es calcular el modelo de regresión incluyendo y excluyendo dicho valor.

Prueba de Bonferroni para detectar outliers

```
# Prueba de Bonferroni para detectar outliers
library(car)
outlierTest(modelo_lineal, cutoff=Inf, n.max=4)
```

```
##      rstudent unadjusted p-value Bonferroni p
## 7  4.782991      5.4575e-05  0.0016372
## 10 -1.777555      8.6744e-02      NA
## 28 -1.381908      1.7833e-01      NA
## 19  1.347912      1.8889e-01      NA
```

Libro Fox presenta los detalles de la prueba

En la salida de arriba vemos las cuatro observaciones ($n.max=4$) que tienen los mayores valores de residual estudentizado. La observación ubicada en la línea 7 es la única con un valor-p muy pequeño y por lo tanto hay evidencias para considerar esa observación como un posible outlier (Hernandez, 2023).

Distancia de Cook

```
# Distancia de Cook, detección de valores influyentes
cutoff <- 4 / (26-2-2) # Cota
plot(modelo_lineal, which=4, cook.levels=cutoff, las=1)
abline(h=cutoff, lty="dashed", col="dodgerblue2")
```

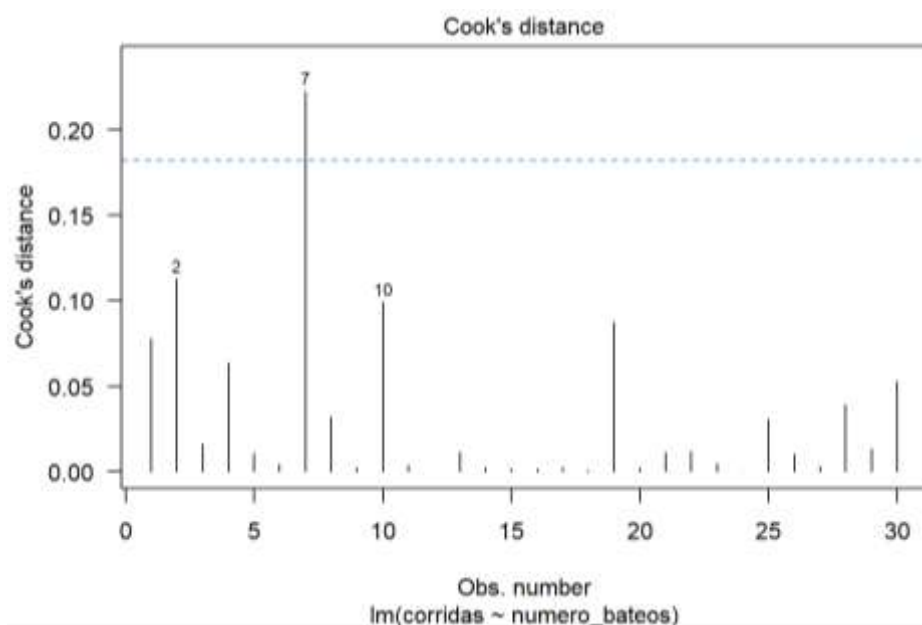


Figura 34. Valores influyentes.

De esta figura es claro que la observación 7 tiene Di por encima de la cota y se consideran observaciones influyentes. Podemos observar valores influyentes como el 2 y 10 influyentes, pero no significativos.

Utilizando otros indicadores.

```
# Identificación de valores atípicos: outliers, leverage y observaciones
influyentes
library(ggrepel)
library(dplyr)
datos$studentized_residual <- rstudent(modelo_lineal)
ggplot(data = datos, aes(x = prediccion, y = abs(studentized_residual))) +
  geom_hline(yintercept = 3, color = "grey", linetype = "dashed") +
  # Se identifican en rojo residuos studentized absolutos > 3
  geom_point(aes(color = ifelse(abs(studentized_residual) > 3, "red",
"black")))) +
  scale_color_identity() +
  #se muestra el equipo al que pertenece la observacion atipica,
  geom_text_repel(data = filter(datos, abs(studentized_residual) > 3),
  aes(label = equipos)) +
  labs(title = "Distribucion de los residuos studentized", x = "prediccion
modelo") +
  theme_bw()+
  theme(plot.title = element_text(hjust = 0.5), legend.position = "none")
```

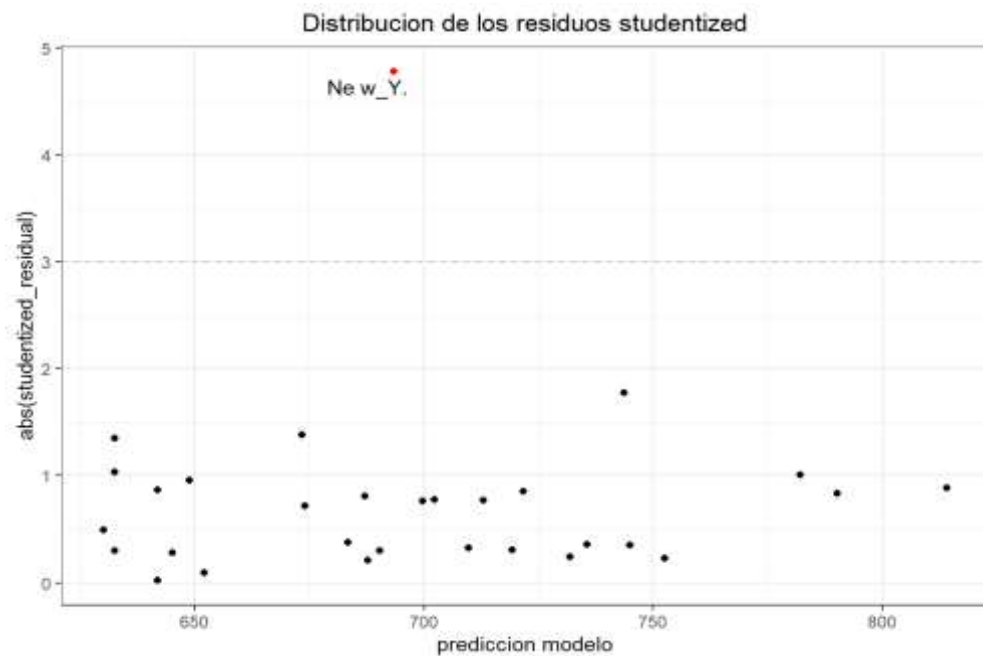


Figura 35. Distribución de los residuos estudiantizados para observar valores outliers.

Cuales son o cual es el o los valores reconocidos como outliers

```
datos %>% filter(abs(studentized_residual) > 3)
```

```
##   equipos numero_bateos corridas prediccion residuos
studentized_residual
## 1 New_Y.           5518      967   693.4817 273.5183
4.782991
```

```
which(abs(datos$studentized_residual) > 3)
```

```
## [1] 7
```

El estudio de los residuos *studentized* identifica al equipo de New_Y. como una posible observación atípica. Esta observación ocupa la posición 7 en la tabla de datos.

El hecho de que un valor sea atípico o con alto grado de *leverage* no implica que sea influyente en el conjunto del modelo. Sin embargo, si un valor es influyente,

suele ser o atípico o de alto *leverage*. En R se dispone de la función `outlierTest()` del paquete `car` y de las funciones `influence.measures()`, `influencePlot()` y `hatvalues()` para identificar las observaciones más influyentes en el modelo (Amat, 2016).

```
library(car)
summary(influence.measures(model = modelo_lineal))
```

```
## Potentially influential observations of
## lm(formula = corridas ~ numero_bateos, data = datos) :
##
##   dfb.1_ dfb.nmr_ dffit   cov.r   cook.d hat
## 2 -0.43   0.44    0.47  1.30_*  0.11  0.22_*
## 7  0.07  -0.06    0.89_*  0.33_*  0.22  0.03
```

Se detectan como influyente la observación que ocupa la segunda posición y es significativa.

```
influencePlot(model = modelo_lineal)
```

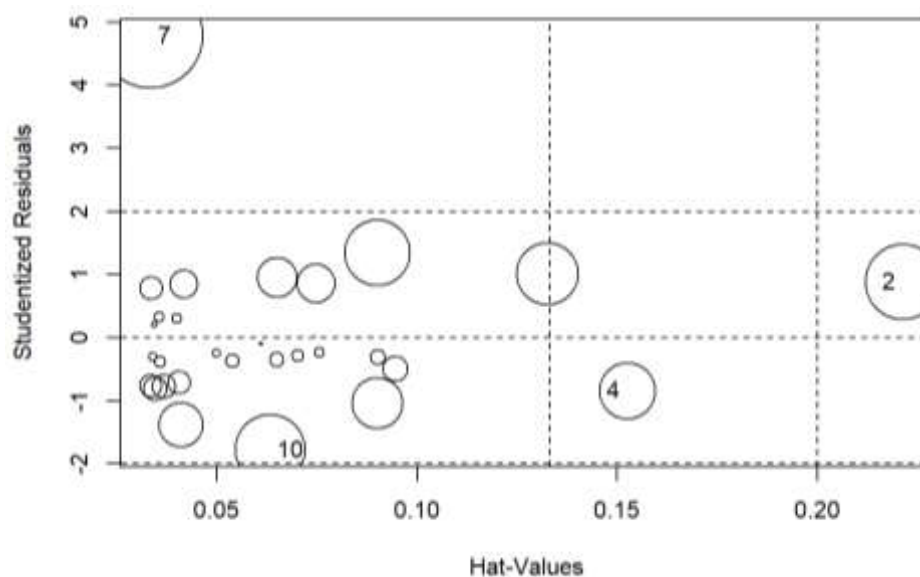


Figura 36. Datos influyentes


```
##      StudRes      Hat      CookD
## 2  0.8873817 0.22133381 0.11277084
## 4 -0.8368049 0.15252728 0.06369637
## 7  4.7829913 0.03349684 0.22254988
## 10 -1.7775546 0.06333282 0.09917231
```

Las funciones `influence.measures()` e `influencePlot()` detectan la observación 7 como atípica pero no significativamente influyente. Sí detectan como influyente la observación que ocupa la segunda posición. Para evaluar hasta qué punto condiciona el modelo, se recalcula la recta de mínimos cuadrados excluyendo esta observación.

Vista del modelo excluyendo la observación atípica. (7)

```
# exluyendo el valor atípico
ggplot(data = datos, mapping = aes(x = numero_bateos, y = corridas)) +
  geom_point(color = "grey50", size = 2) +
  geom_smooth(method = "lm", se = FALSE, color = "black") +
  # se resalta el valor excluido
  geom_point(data = datos[7, ], color = "red", size = 2) +
  # se anade la nueva recta de minimos cuadrados
  geom_smooth(data = datos[-7, ], method="lm", se =FALSE, color = "blue")
+
  labs(title = "Diagrama de dispersion", x = "numero de bateos") +
  theme_bw() +
  theme(plot.title = element_text(hjust = 0.5), legend.position = "none")
```

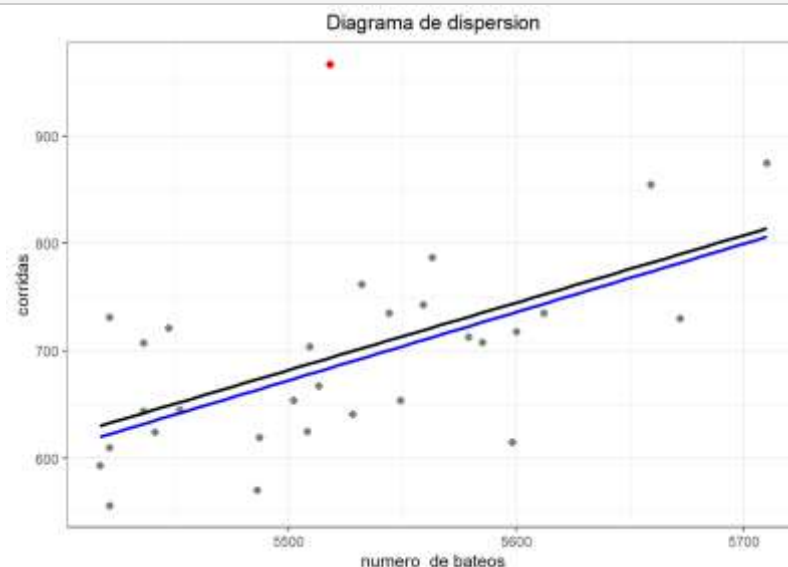


Figura 36. Modelo final

Observamos un cambio en la línea de regresión, veamos más detalladamente los parámetros. Obteniendo el modelo con todos los datos (modelo_inicial) y el modelo excluyendo el dato 7 (modelo_final).

```
modelo_inicial <- lm(formula = corridas ~ numero_bateos, data = datos)
summary(modelo_inicial)
```

```
## Call:
## lm(formula = corridas ~ numero_bateos, data = datos)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -128.69  -50.54  -19.96   51.10  273.52
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  -2769.4894   997.0462  -2.778  0.00966 **
## numero_bateos    0.6276    0.1805   3.477  0.00167 **
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 77.63 on 28 degrees of freedom
## Multiple R-squared:  0.3016, Adjusted R-squared:  0.2766
## F-statistic: 12.09 on 1 and 28 DF,  p-value: 0.001673
```

```
modelo_final <- lm(formula = corridas ~ numero_bateos, data = datos[-7,
])
summary(modelo_final)
```

```
## Call:
## lm(formula = corridas ~ numero_bateos, data = datos[-7, ])
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -119.88  -45.29  -11.03   34.46  108.69
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
```

```
## (Intercept)    -2825.3914    747.1319   -3.782 0.000786 ***
## numero_bateos    0.6360     0.1352    4.702 6.78e-05 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 58.17 on 27 degrees of freedom
## Multiple R-squared:  0.4503, Adjusted R-squared:  0.4299
## F-statistic: 22.11 on 1 and 27 DF,  p-value: 6.776e-05
```

En ambos, los estimadores $\hat{\beta}_0$ y $\hat{\beta}_1$ son significativos, la diferencia está en el coeficiente de determinación. En el modelo inicial se tiene un $R^2 = 30.16\%$ y en el modelo final $R^2 = 45.03\%$, mejorando el ajuste de la varianza. Hasta aquí el mejor modelo es el que excluye el valor atípico.

Vista del modelo excluyendo la observación influyente

En líneas arriba se conoce que el punto influyente significativo es el punto 2.

```
ggplot(data = datos, mapping = aes(x = numero_bateos, y = corridas)) +
  geom_point(color = "grey50", size = 2) +
  geom_smooth(method = "lm", se = FALSE, color = "black") +
  # se resalta el valor excluido
  geom_point(data = datos[2, ], color = "red", size = 2) +
  # se anade la nueva recta de minimos cuadrados
  geom_smooth(data = datos[-2, ], method="lm", se =FALSE, color = "blue") +
  labs(title = "Diagrama de dispersion", x = "numero de bateos") +
  theme_bw() +
  theme(plot.title = element_text(hjust = 0.5), legend.position = "none")
```

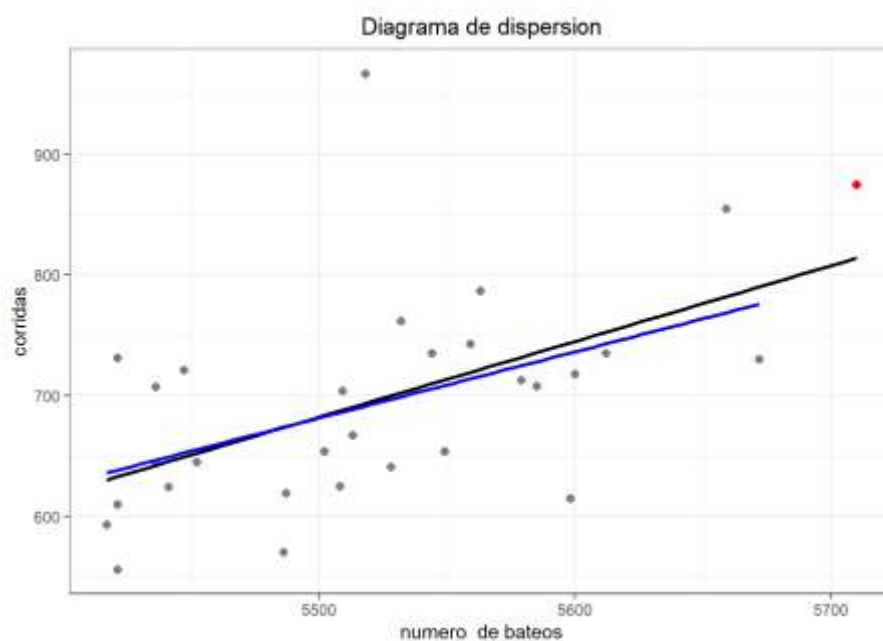


Figura 37. Modelo sin observación influyente.

La eliminación del valor identificado como influyente apenas cambia la recta de mínimos cuadrados. Para conocer con exactitud el resultado de excluir la observación se comparan los coeficientes del modelo inicial.

```
modelo_sin_influ <- lm(formula = corridas ~ numero_bateos, data = datos[-2, ])
summary(modelo_sin_influ)
```

```
## Call:
## lm(formula = corridas ~ numero_bateos, data = datos[-2, ])
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -120.19  -45.78  -18.29   29.43  275.70
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  -2335.7462   1113.8315  -2.097   0.0455 *
## numero_bateos    0.5486     0.2019   2.717   0.0113 *
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
```

```
## Residual standard error: 77.93 on 27 degrees of freedom
## Multiple R-squared:  0.2148, Adjusted R-squared:  0.1857
## F-statistic: 7.385 on 1 and 27 DF,  p-value: 0.01134
```

En ambos, los estimadores $\hat{\beta}_0$ y $\hat{\beta}_1$ son significativos, la diferencia está en el coeficiente de determinación. En el modelo inicial se tiene un $R^2 = 30.16\%$ y en el modelo sin observación influyente, $R^2 = 21.48\%$, bajando el ajuste de la varianza. Hasta aquí el mejor modelo es el que excluye el valor atípico.

i) Modelo final.

Análisis del modelo final excluyendo el dato atípico 7.

Datos completos.

```
#### MODELO FINAL #####
# Datos originales.
equipos <- c("Texas", "Boston", "Detroit", "Kansas", "St.", "New_S.",
            "Ne w_Y.", "Milwaukee", "Colorado", "Houston", "Baltimore",
            "Los_An.", "Chica go", "Cincinnati", "Los_P.", "Philadelphia",
            "Chicago", "Cleveland", "Ari zona", "Toronto", "Minnesota",
            "Florida", "Pittsburgh", "Oakland", "Tampa ", "Atlanta", "Washington",
            "San.F", "San.I", "Seattle")
numero_bateos <- c(5659, 5710, 5563, 5672, 5532, 5600, 5518, 5447, 5544, 5598,
                  5585, 5436, 5549, 5612, 5513, 5579, 5502, 5509, 5421, 5559,
                  5487, 5508, 5421, 5452, 5436, 5528, 5441, 5486, 5417, 5421)
corridas <- c(855, 875, 787, 730, 762, 718, 967, 721, 735, 615, 708, 644, 654, 735,
             667, 713, 654, 704, 731, 743, 619, 625, 610, 645, 707, 641, 624, 570,
             593, 556)
datos <- data.frame(equipos, numero_bateos, corridas)
head(datos)
```

```
#### MODELO FINAL
# eliminando el dato 7 y guardando los datos en un nuevo archivo
datos1 <- datos[-c(7), ]
head(datos1)
```

```
## equipos numero_bateos corridas
## 1 Texas 5659 855
## 2 Boston 5710 875
## 3 Detroit 5563 787
## 4 Kansas 5672 730
## 5 St. 5532 762
## 6 New_S. 5600 718
```

```
str(datos1)
```

```
## 'data.frame': 29 obs. of 3 variables:
## $ equipos : chr "Texas" "Boston" "Detroit" "Kansas" ...
## $ numero_bateos: num 5659 5710 5563 5672 5532 ...
## $ corridas : num 855 875 787 730 762 718 721 735 615 708 ...
```

Quedan 29 observaciones y 3 variables.

```
# histograma para las variables
library(ggplot2)
par(mfrow = c(1, 2))
hist(datos1$numero_bateos, breaks = 10, main = "", xlab = "número de
bateos",
      ylab="Frecuencia", border = "darkred")
hist(datos1$corridas, breaks = 10, main = "", xlab =
"corridas",ylab="Frecuencia",
      border = "blue")
```

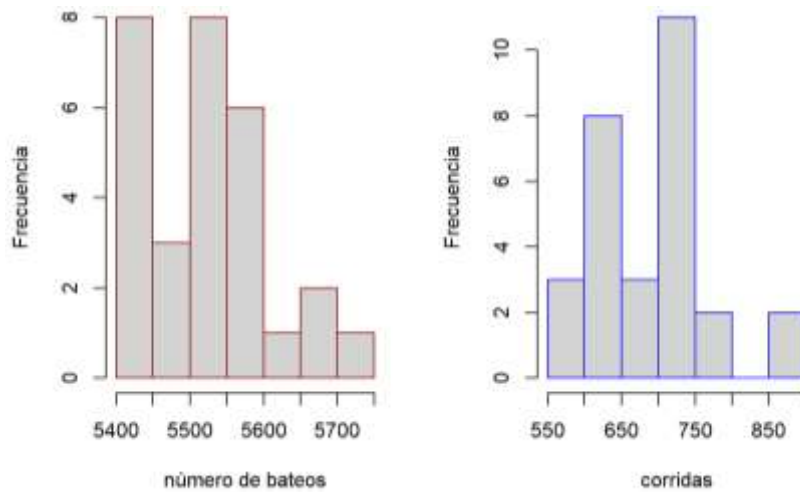


Figura 38. Histogramas

Resumen de estadísticos descriptivos.

```
summary(datos1)
```

```
##   equipos      numero_bateos      corridas
## Length:29      Min.   :5417      Min.   :556.0
## Class :character 1st Qu.:5447      1st Qu.:625.0
## Mode  :character Median :5513      Median :704.0
##                               Mean  :5524      Mean  :687.6
##                               3rd Qu.:5579      3rd Qu.:731.0
##                               Max.   :5710      Max.   :875.0
```

```
# representación gráfica
require(ggplot2)
ggplot(data = datos1, mapping = aes(x = numero_bateos, y = corridas)) +
  geom_point(color = "firebrick", size = 3) +
  (labs(title = "Diagrama de dispersión", x = "número de bateos")) +
  theme_bw() +
  theme(plot.title = element_text(hjust =0.5))
```

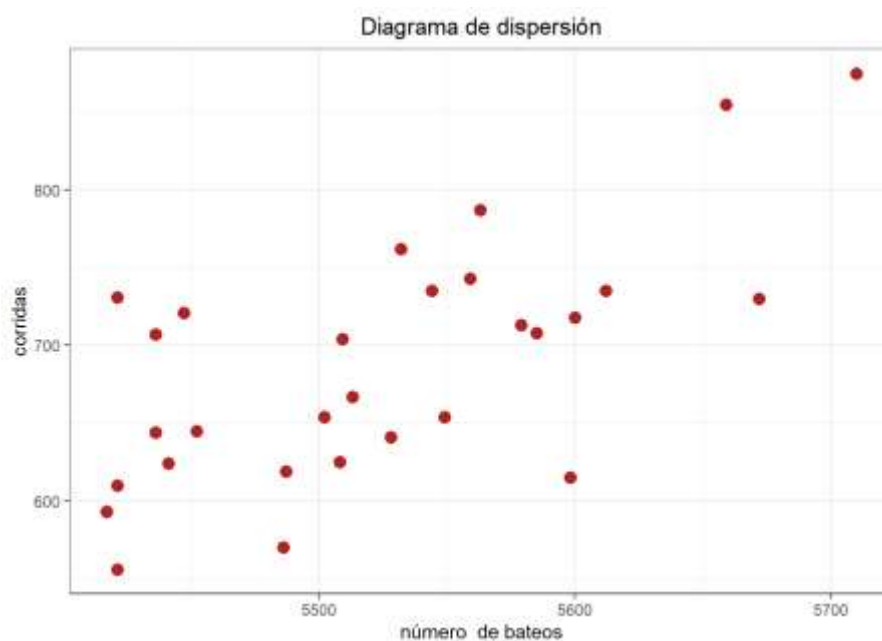


Figura 39. Diagrama de dispersión

```
cor.test(x = datos1$numero_bateos, y = datos1$corridas, method =
"pearson")
```

```
## Pearson's product-moment correlation
##
## data: datos1$numero_bateos and datos1$corridas
## t = 4.7025, df = 27, p-value = 6.776e-05
## alternative hypothesis: true correlation is not equal to 0
## 95 percent confidence interval:
## 0.4038110 0.8327235
## sample estimates:
## cor
## 0.6710081
```

El coeficiente de correlación se ha incrementado a 0.67 y es significativa con $p(0.0000677) < \alpha(0.05)$, los datos pueden ser ajustados a una regresión lineal.

```
# Cálculo del modelo de regresión lineal simple
modelo_final<- lm(corridas ~ numero_bateos, data=datos1)
summary(modelo_final)
```



```
## Call:
## lm(formula = corridas ~ numero_bateos, data = datos1)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -119.88  -45.29  -11.03   34.46  108.69
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  -2825.3914    747.1319  -3.782 0.000786 ***
## numero_bateos    0.6360     0.1352   4.702 6.78e-05 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 58.17 on 27 degrees of freedom
## Multiple R-squared:  0.4503, Adjusted R-squared:  0.4299
## F-statistic: 22.11 on 1 and 27 DF,  p-value: 6.776e-05
```

El intercepto y la variable X (numero de bateos) resultan ser significativos al 5% con un coeficiente de determinación de 45.03% y un ANVA significativo.

```
# Intervalos de confianza para los parámetros del modelo
confint(modelo_final, level=0.95)
```

```
##              2.5 %      97.5 %
## (Intercept) -4358.3793538 -1292.4034930
## numero_bateos    0.3584894    0.9134908
```

```
# Representación grafica del modelo
ggplot(data = datos1, mapping = aes(x = numero_bateos, y = corridas)) +
  geom_point(size=3) +
  labs(title = "corridas ~ numero de bateos", x = "numero de bateos") +
  geom_smooth(method = "lm", se = FALSE, color = "red") +
  theme_bw() +
  theme(plot.title = element_text(hjust = 0.4))
```

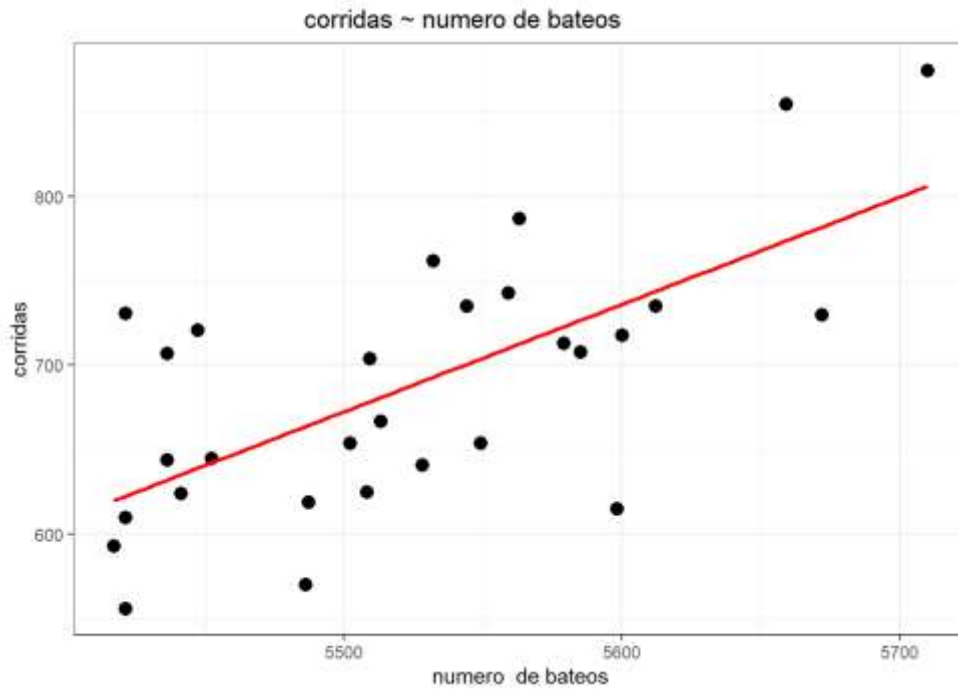


Figura 40. Modelo de regresión

```
# predecir
nuevos_datos <- data.frame(numero_bateos=
seq(min(numero_bateos),max(numero_bateos)))
predict_value <- predict(modelo_final)
head(predict_value)
```

```
##          1          2          3          4          5          6
## 773.6767 806.1122 712.6217 781.9446 692.9060 736.1533
```

```
# Cálculo del error cuadrático medio RMSE:
sqrt(mean(modelo_final$residuals^2))
```

```
## [1] 56.12652
```

```
# una banda
par(mfrow = c(1, 1))

puntos <- seq(from = min(datos1$numero_bateos),
              to = max(datos1$numero_bateos), length.out = 100)
limites_intervalo <- predict(object = modelo_final,
                             newdata = data.frame( numero_bateos = puntos
),
                             interval = "confidence", level = 0.95)
head(limites_intervalo, 3)
```

```
##          fit          lwr          upr
## 1 619.7671 582.7842 656.7501
## 2 621.6494 585.3206 657.9782
## 3 623.5317 587.8501 659.2132
```

```
plot(datos1$numero_bateos, datos1$corridas, col = "firebrick", pch = 19,
      ylab = "corridas", xlab = "numero de bateos",
      main = "corridas ~ numero de bateos")
abline(modelo_final, col = 1)
lines(x = puntos, y = limites_intervalo[, 2], type = "l", col = 2, lty =
3)
lines(x = puntos, y = limites_intervalo[, 3], type = "l", col = 3, lty =
3)
```

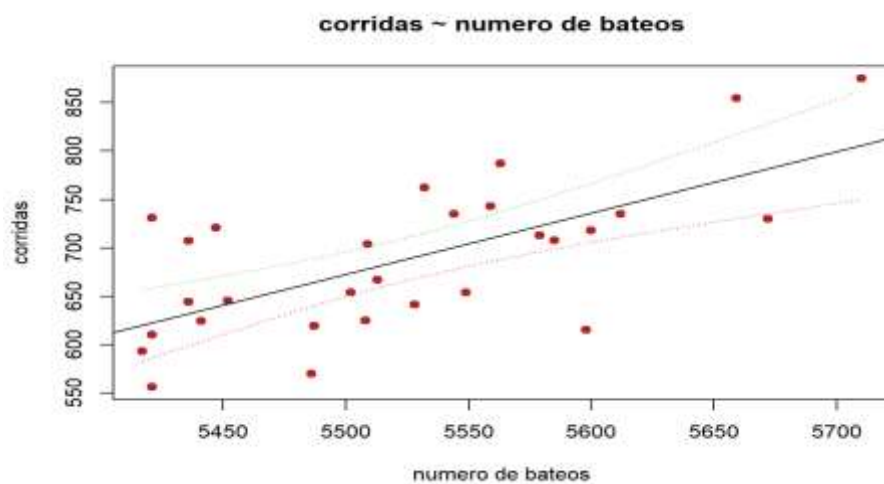


Figura 41. Modelo de regresión

```

# Por defecto incluye la región de 95% de confianza
# dos bandas
# Grafico dispersion y recta
plot(datos1$numero_bateos, datos1$corridas, col = "firebrick", pch = 19,
      ylab = "corridas", xlab = "numero de bateos",
      main = "corridas ~ numero de bateos")
abline(modelo_final, col = 1)

# Intervalos de confianza de la respuesta media:
# valores medios

ic <- predict(modelo_final, nuevos_datos, interval = 'confidence')
lines(nuevos_datos$numero_bateos, ic[, 2], lty = 2)
lines(nuevos_datos$numero_bateos, ic[, 3], lty = 2)

# Intervalos de prediccion
# para cualquier valor
ic <- predict(modelo_final, nuevos_datos, interval = 'prediction')
lines(nuevos_datos$numero_bateos, ic[, 2], lty = 2, col = 'red')
lines(nuevos_datos$numero_bateos, ic[, 3], lty = 2, col = 'red')

```

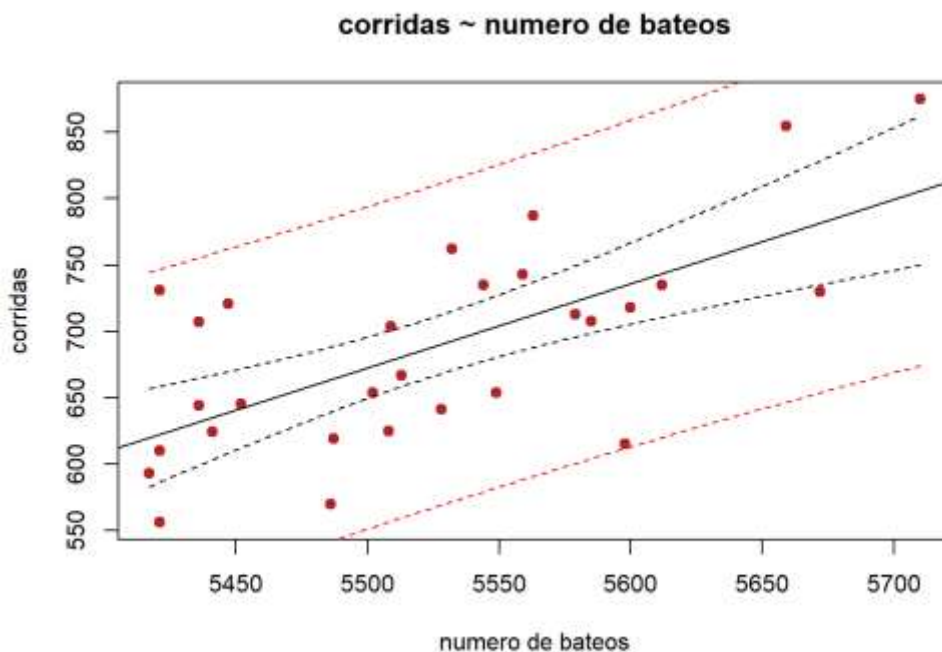


Figura 41. Modelo con intervalos de confianza

```
# observando datos, valores predichos y residuales
datos1$prediccion <- modelo_final$fitted.values
datos1$residuos <- modelo_final$residuals
head(datos1)
```

```
##   equipos numero_bateos corridas prediccion  residuos
## 1   Texas           5659      855   773.6767  81.32327
## 2   Boston           5710      875   806.1122  68.88777
## 3 Detroit           5563      787   712.6217  74.37832
## 4   Kansas           5672      730   781.9446 -51.94461
## 5     St.           5532      762   692.9060  69.09401
## 6 New_S.           5600      718   736.1533 -18.15332
```

```
# gráfico de residuales
ggplot(data = datos1, aes(x = prediccion, y = residuos)) +
  geom_point(aes( color = residuos)) +
  scale_color_gradient2(low = "blue3", mid = "grey", high = "red") +
  geom_hline(yintercept = 0) +
  geom_segment(aes(xend = prediccion, yend = 0), alpha = 0.2) +
  labs(title = "Distribucion de los residuos", x = "prediccion modelo",
       y = "residuo") +
  theme_bw() +
  theme(plot.title = element_text(hjust = 0.5), legend.position = "none")
```

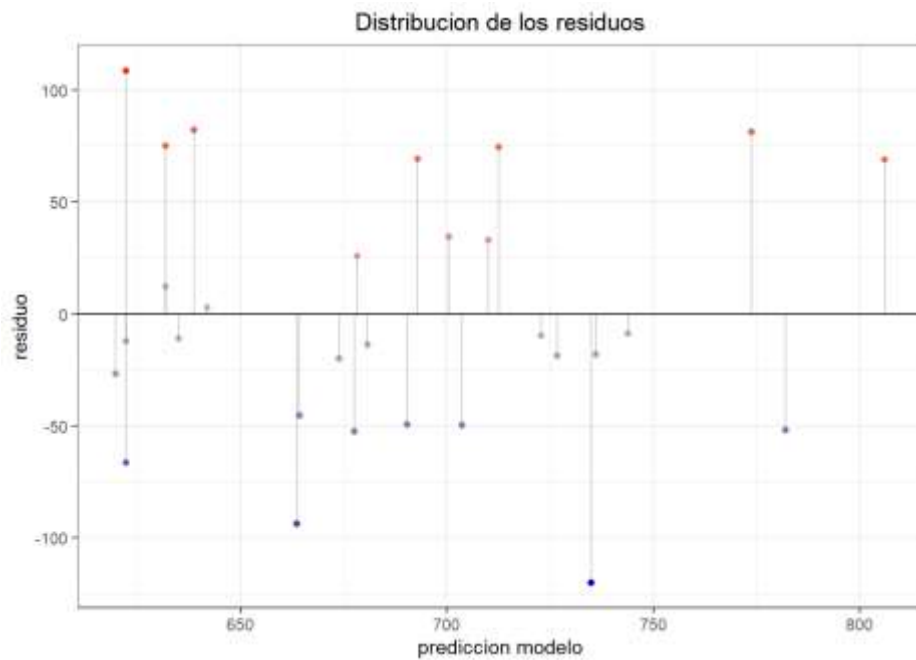


Figura 42. Gráfico de residuales

```
par(mfrow = c(1, 1))  
# Distribución normal de los residuos:  
ggplot(data = datos1, aes(x = residuos)) +  
  geom_histogram(aes(y = ..density..)) +  
  labs(title = "Histograma de los residuos") +  
  theme_bw() +  
  theme(plot.title = element_text(hjust = 0.5))
```

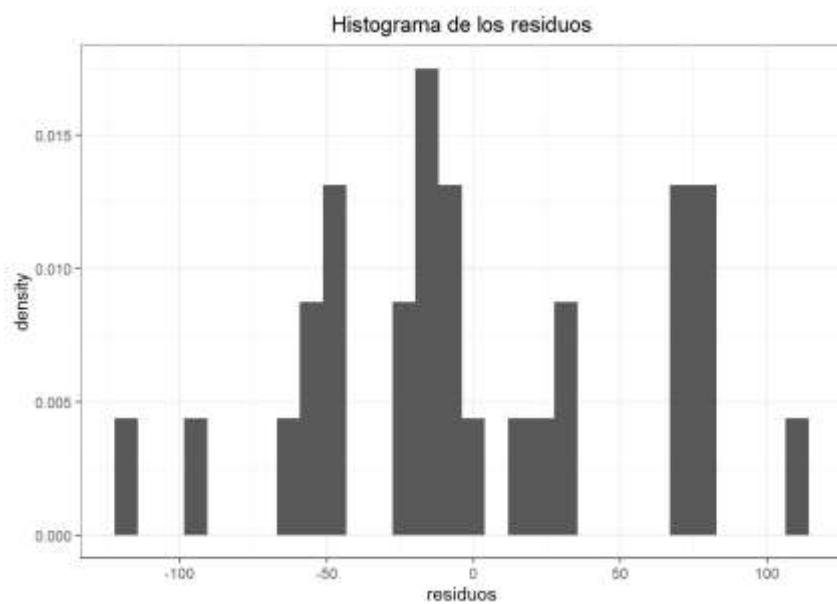


Figura 43. Histograma de los residuales

```
# gráfico de cuantiles
qqnorm(modelo_final$residuals, main = "Gráfico de cuantiles QQ plot",
       col = "darkred")
qqline(modelo_final$residuals)
```

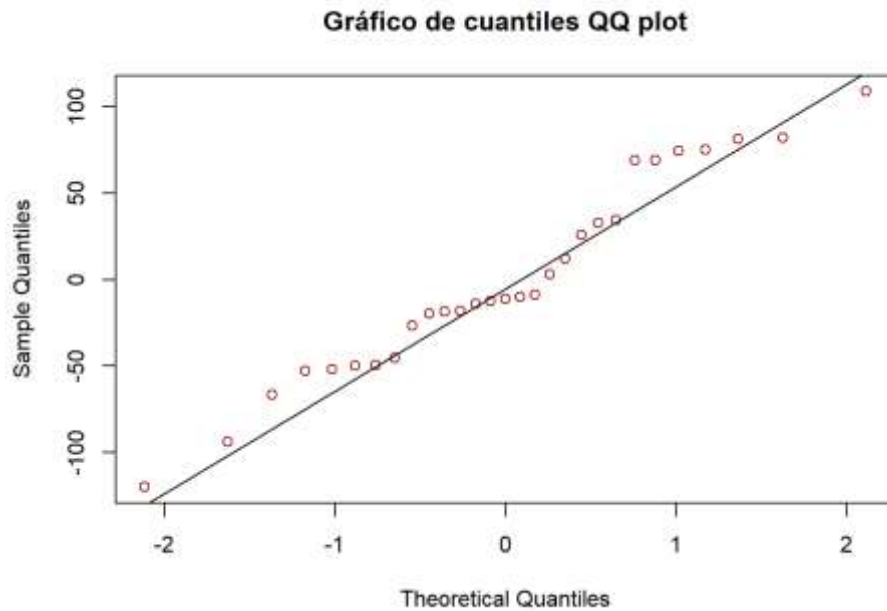


Gráfico 44. Gráfico de cuantiles QQ.

```
# Test de normalidad
#shapiro.test
shapiro.test(modelo_final$residuals)
```

```
## Shapiro-Wilk normality test
##
## data: modelo_final$residuals
## W = 0.96269, p-value = 0.3822
```

```
# Kolmogorov test
ks.test(modelo_final$residuals, "pnorm",
       mean = mean(modelo_final$residuals),
       sd = sd(modelo_final$residuals))
```

```
## One-sample Kolmogorov-Smirnov test
##
## data: modelo_final$residuals
## D = 0.14732, p-value = 0.5083
## alternative hypothesis: two-sided
```

```
# Varianza constante de los residuos (Homocedasticidad):
ggplot(data = datos1, aes(x = prediccion, y = residuos)) + geom_point(aes(
  color = residuos)) +
  scale_color_gradient2(low = "blue3", mid = "grey", high = "red") +
  geom_segment(aes(xend = prediccion, yend = 0), alpha = 0.2) +
  geom_smooth(se = FALSE, color = "firebrick") +
  labs(title = "Distribucion de los residuos", x = "prediccion modelo", y =
  "residuo") +
  geom_hline(yintercept = 0) +
  theme_bw() +
  theme(plot.title = element_text(hjust = 0.5), legend.position = "none")
```

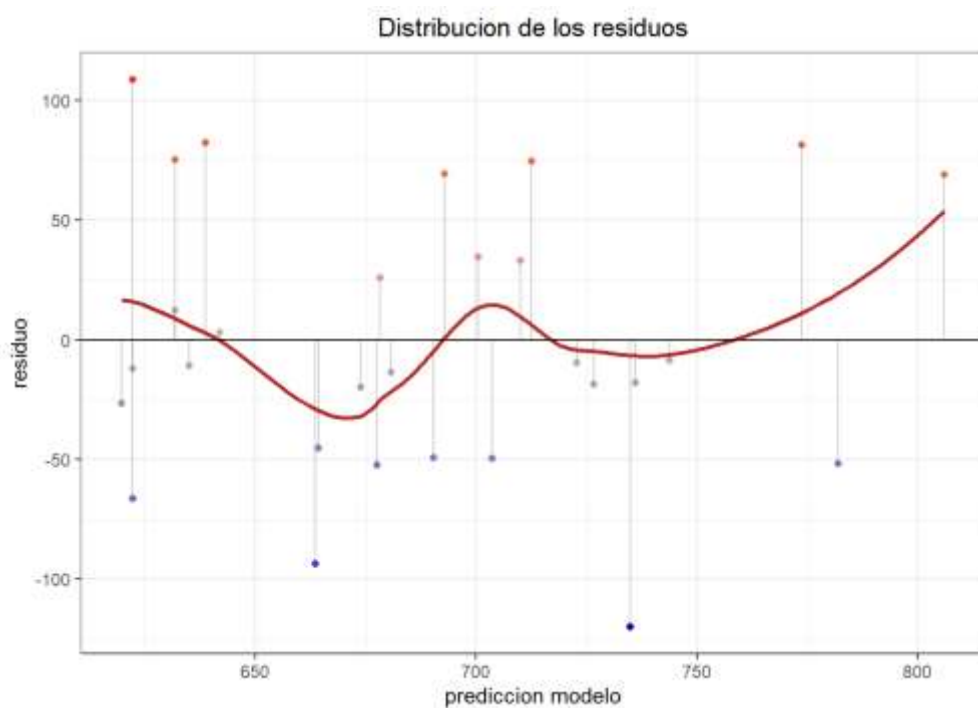


Figura 45. Gráfico de residuos

```
# Test de Breush-Pagan
```



```
library(lmtest)
bptest(modelo_final)
```

```
## studentized Breusch-Pagan test
##
## data: modelo_final
## BP = 0.076886, df = 1, p-value = 0.7816
```

```
# Autocorrelacion de residuos:
ggplot(data = datos1, aes(x = seq_along(residuos), y = residuos)) + geom_
point(aes(color = residuos)) +
  scale_color_gradient2(low = "blue3", mid = "grey", high = "red") + geom_
_line(size = 0.3) +
  labs(title = "Distribucion de los residuos", x = "index", y = "residuo"
)+
  geom_hline(yintercept = 0) +
  theme(plot.title = element_text(hjust = 0.5), legend.position = "none")
```

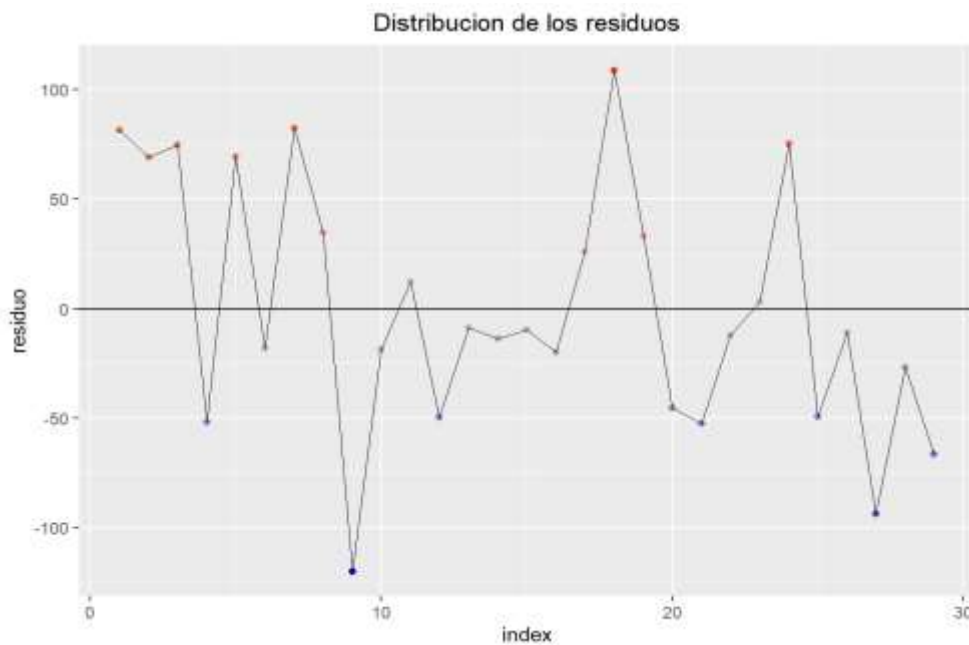


Figura 45. Gráfico de distribución de residuos.

```
#test de Durwin Watson
```

```
library(lmtest)
dwtest(modelo_final)
```

```
## Durbin-Watson test
##
## data: modelo_final
## DW = 1.633, p-value = 0.1327
## alternative hypothesis: true autocorrelation is greater than 0
```

```
# Prueba de Bonferroni para detectar outliers
library(car)
outlierTest(modelo_final, cutoff=Inf, n.max=4)
```

```
##      rstudent unadjusted p-value Bonferroni p
## 10 -2.292364          0.030226          0.87657
## 19  2.077211          0.047800             NA
## 28 -1.701835          0.100720             NA
## 1  1.539515          0.135760             NA
```

El punto 10 que es el nuevo outlier no es significativo.

```
# Distancia de Cook, detección de valores influyentes
cutoff <- 4 / (26-2-2) # Cota
plot(modelo_final, which=4, cook.levels=cutoff, las=1)
abline(h=cutoff, lty="dashed", col="dodgerblue2")
```

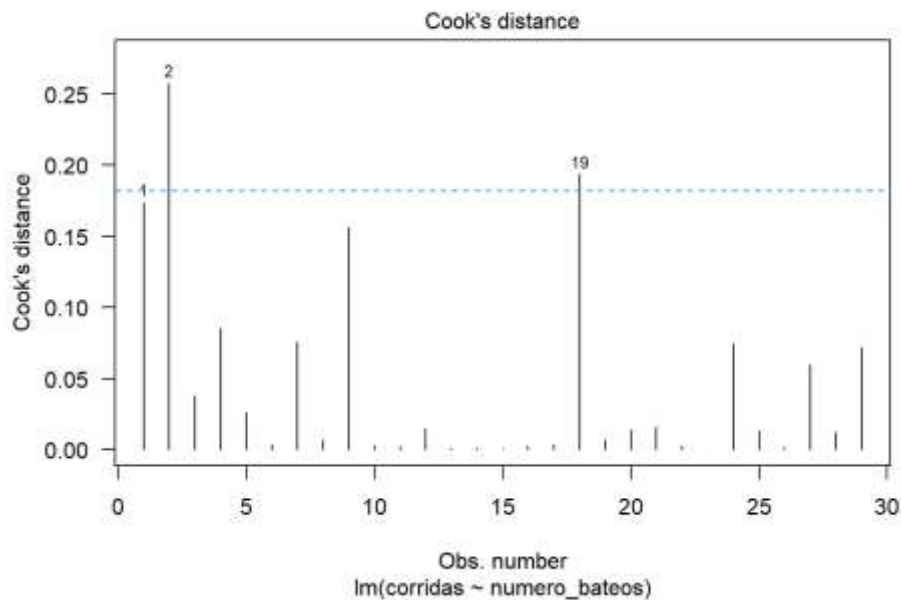


Figura 46. Gráfico de distancia de Cook.

La figura, muestra valores por encima de la línea como valores influyentes, veamos su significancia:

```
# Identificación de valores atípicos: outliers, leverage y observaciones
# influyentes

library(ggrepel)
library(dplyr)
datos1$studentized_residual <- rstudent(modelo_final)
ggplot(data = datos1, aes(x = prediccion, y = abs(studentized_residual)))
+
  geom_hline(yintercept = 3, color = "grey", linetype = "dashed") +
  # Se identifican en rojo residuos studentized absolutos > 3
  geom_point(aes(color = ifelse(abs(studentized_residual) > 3, "red", "black"))) +
  scale_color_identity() +
  # se muestra el equipo al que pertenece la observación atípica,
  geom_text_repel(data = filter(datos1, abs(studentized_residual) > 3),
    aes(label = equipos)) +
  labs(title = "Distribución de los residuos studentized", x = "predicción
modelo") +
  theme_bw() +
  theme(plot.title = element_text(hjust = 0.5), legend.position = "none")
```

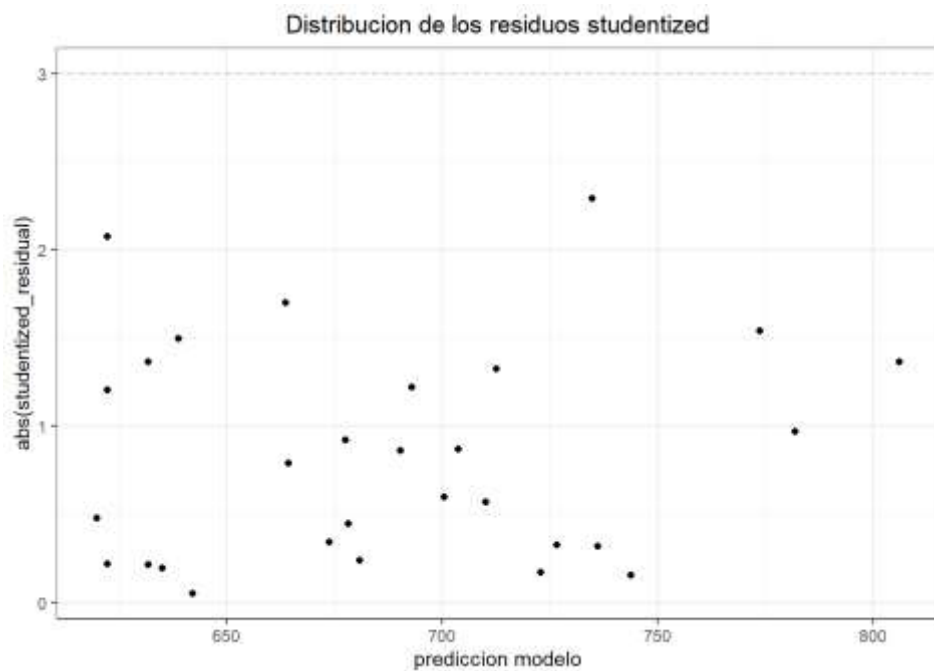


Figura 47. Gráfico de residuos estuzizado.

La figura, no muestra valores mayores a 3, no hay valores atípicos significativos

```
datos1 %>% filter(abs(studentized_residual) > 3)
```

```
## [1] equipos          numero_bateos      corridas  
## [4] prediccion        residuos           studentized_residual  
## <0 rows> (or 0-length row.names)
```

```
which(abs(datos1$studentized_residual) > 3)
```

```
## integer(0)
```

No se observa valores atípicos no influyentes

```
# prediciendo nuevos valores, cuando X = 5459 y 5455
prediciendo<- predict(modelo_final, data.frame(numero_bateos= c(5459,5455
)))
prediciendo
```

```
##          1          2
## 646.4787  643.9347
```

Conclusión

Dado que se satisfacen todas las condiciones para considerar válido un modelo de regresión lineal por mínimos cuadrados y que el *p-value* indica que el ajuste es significativo, se puede aceptar el modelo lineal.

23. Regresión lineal con un predictor categórico de dos niveles.

El set de datos `sexab` del paquete `faraway` contiene los resultados de un estudio en el que se investigó las secuelas que padecen mujeres adultas debido a abusos sufridos durante la infancia. En una clínica médica se midió el nivel de estrés post-traumático (*ptsd*) y nivel de abuso físico sufrido (*cpa*), ambos en escalas estandarizadas, en 45 mujeres que fueron víctimas en su infancia (*csa*). Las mismas mediciones se registraron para 31 mujeres que no sufrieron ningún tipo de abuso (Amat, 2016).

```
library(faraway)
library(ggplot2)
data(sexab)
head(sexab)
```

```
##      cpa      ptsd      csa
## 1 2.04786  9.71365 Abused
## 2 0.83895  6.16933 Abused
## 3 -0.24139 15.15926 Abused
## 4 -1.11461 11.31277 Abused
## 5 2.01468  9.95384 Abused
## 6 6.71131  9.83884 Abused
```

```
str(sexab)
```

```
## 'data.frame': 76 obs. of 3 variables:
## $ cpa : num 2.048 0.839 -0.241 -1.115 2.015 ...
## $ ptsd: num 9.71 6.17 15.16 11.31 9.95 ...
## $ csa : Factor w/ 2 levels "Abused","NotAbused": 1 1 1 1 1 1 1 1 1 ...
```

La data contiene 76 observaciones y 3 variables.

Resumen de estadísticos descriptivos

```
summary(sexab)
```

```
##           cpa           ptsd           csa
## Min.   : -3.1204  Min.   : -3.349   Abused  :45
## 1st Qu.:  0.8321  1st Qu.:  6.173   NotAbused:31
## Median :  2.0707  Median :  8.909
## Mean   :  2.3547  Mean   :  8.986
## 3rd Qu.:  3.7387  3rd Qu.: 12.240
## Max.   :  8.6469  Max.   : 18.993
```

a) Diagrama de dispersión

```
# representación grafica
require(ggplot2)
ggplot(data = sexab, mapping = aes(x = csa, y =ptsd )) +
  geom_point(color = "firebrick", size = 2) +
  (labs(title = "Diagrama de dispersion", x = "Deformación del acero
(X)")) +
  theme_bw() +
  theme(plot.title = element_text(hjust =0.5))
```

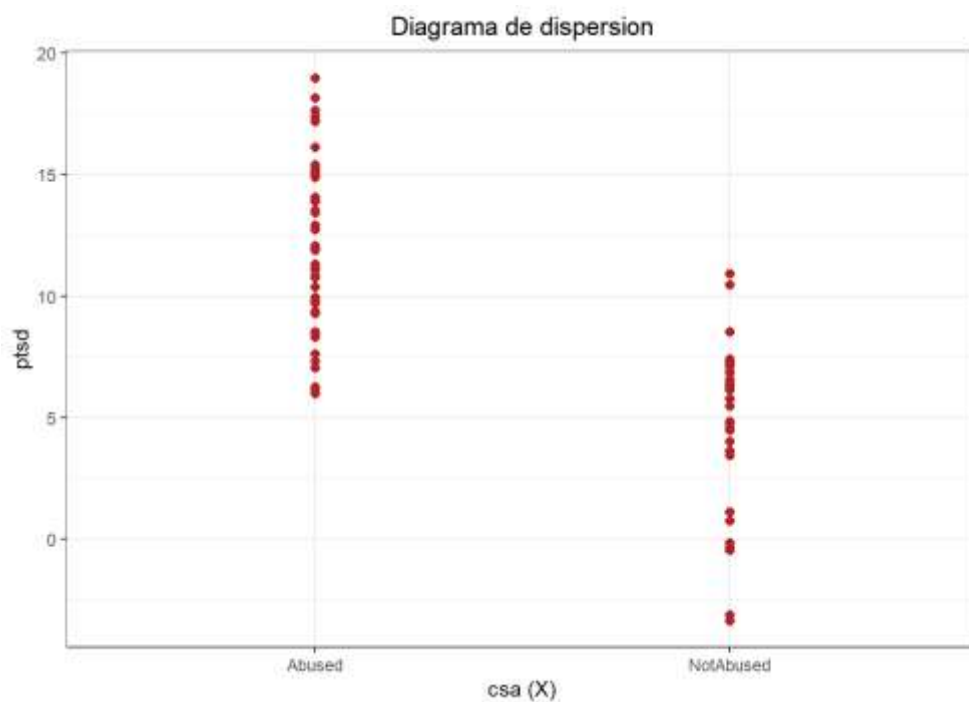


Figura 48. Diagrama de dispersión.

b) Estadísticos descriptivos para cada variable.

```
by(data = sexab, INDICES = sexab$csa, summary)
```

```
## sexab$csa: Abused
##      cpa      ptsd      csa
## Min.   :-1.115  Min.   : 5.985  Abused   :45
## 1st Qu.: 1.415  1st Qu. : 9.374  NotAbused: 0
## Median : 2.627  Median :11.313
```

```
## Mean : 3.075 Mean :11.941
## 3rd Qu.: 4.317 3rd Qu.:14.901
## Max. : 8.647 Max. :18.993
## -----
## sexab$csa: NotAbused
##      cpa      ptsd      csa
## Min. :-3.1204 Min. :-3.349 Abused :0
## 1st Qu.: -0.2299 1st Qu.: 3.544 NotAbused:31
## Median : 1.3216 Median : 5.794
## Mean : 1.3088 Mean : 4.696
## 3rd Qu.: 2.8309 3rd Qu.: 6.838
## Max. : 5.0497 Max. :10.914
```

```
ggplot(data = sexab, aes(x = csa, y = ptsd, colour = csa)) + geom_boxplot () +
  theme_bw() + theme(legend.position = "none")
```

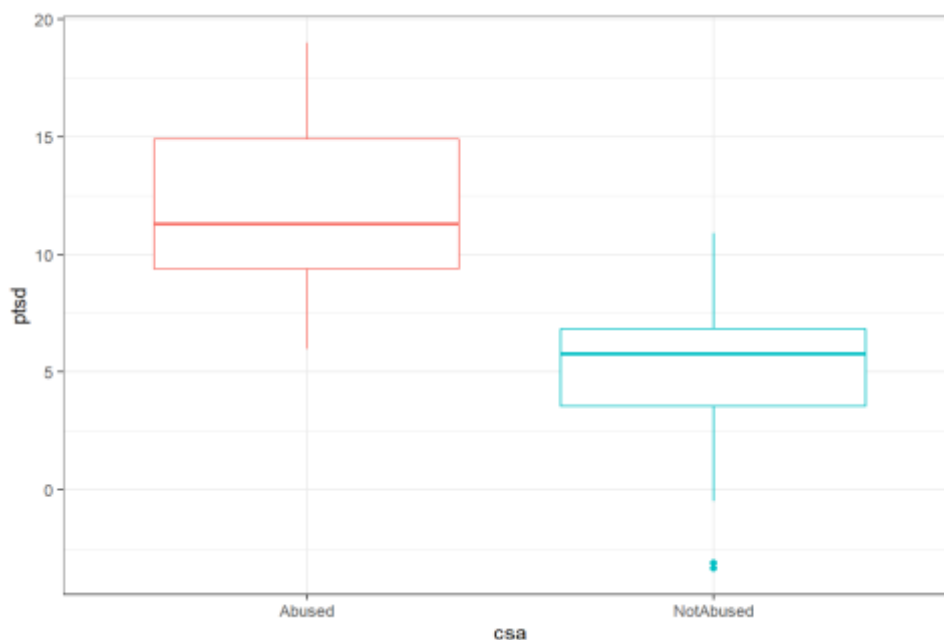


Figura 49. Gráfico BoxPlot

Se observa que las víctimas de abusos tienen niveles más altos de estrés postraumático en comparación a las mujeres que no han sufrido abusos (Amat, 2016).

Una forma de comparar si está diferencia es significativa, es mediante un t-test.

```
# Calculo de la varianza de cada grupo para determinar si son similares
```

```
aggregate(ptsd ~ csa, data = sexab, FUN = var)
```

```
##           csa      ptsd
## 1   Abused  11.83464
## 2 NotAbused 12.38859
```

```
t.test(formula = ptsd ~ csa, data = sexab, var.equal = TRUE)
```

```
## Two Sample t-test
##
## data: ptsd by csa
## t = 8.9387, df = 74, p-value = 2.172e-13
## alternative hypothesis: true difference in means is not equal to 0
## 95 percent confidence interval:
##  5.630165 8.860273
## sample estimates:
##  mean in group Abused mean in group NotAbused
##           11.941093           4.695874
```

El contraste de hipótesis muestra una clara significancia en la diferencia de medias ($p\text{-value} = 2.172e-13$).

Este mismo contraste puede realizarse desde la perspectiva de un modelo lineal incluyendo la variable cualitativa como predictor. Para hacerlo, cada uno de los niveles del predictor se tiene que convertir en una variable dummy cuyo valor puede ser 0 o 1.

$$ptsd = \beta_0 + \beta_1 dummy_{abused} + \beta_2 dummy_{NotAbused}$$

$$dummy_i = \begin{cases} 0 & \text{la observacion no pertenece al nivel } i \\ 1 & \text{la observacion pertenece al nivel } i \end{cases}$$

Para cada observación, únicamente una de las variables dummy toma el valor 1, por ejemplo, si una mujer sí ha sufrido abusos infantiles, el modelo queda:

$$ptsd = \beta_0 + \beta_1 dummy_{abused}$$

```
sexab$abused <- ifelse(test = sexab$csa == "Abused", yes = 1, no = 0)
sexab$not_abused <- ifelse(test = sexab$csa == "NotAbused", yes = 1, no = 0)
rbind(head(sexab, 3), tail(sexab, 3))
```

```
##      cpa   ptsd   csa abused   not_abused
## 1  2.04786  9.71365     Abused     1         0
## 2  0.83895  6.16933     Abused     1         0
## 3 -0.24139 15.15926     Abused     1         0
## 74 -1.85753 -0.46996   NotAbused    0         1
## 75  2.85253  6.84304   NotAbused    0         1
## 76  0.81138  7.12918   NotAbused    0         1
```

Se puede observar que la información de las dos variables dummy es redundante, al haber solo dos niveles, y dado que solo uno de ellos puede tomar el valor 1, conociendo uno se conoce el otro. Para evitar que aparezca el problema de la singularidad al ajustar el modelo, una de las dos variables se excluye del modelo y se considera como el nivel de referencia (Amat, 2016).

c) Modelo lineal

```
modelo <- lm(ptsd ~ abused, data = sexab)
summary(modelo)
```

```
## Call:
## lm(formula = ptsd ~ abused, data = sexab)
##
## Residuals:
##   Min     1Q   Median     3Q    Max
## -8.0451 -2.3123  0.0951  2.1645  7.0514
##
```

```
## Coefficients:
##           Estimate Std. Error  t value Pr(>|t|)
## (Intercept)  4.6959    0.6237   7.529  1.00e-10 ***
## abused       7.2452    0.8105   8.939  2.17e-13 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 3.473 on 74 degrees of freedom
## Multiple R-squared:  0.5192, Adjusted R-squared:  0.5127
## F-statistic: 79.9 on 1 and 74 DF, p-value: 2.172e-13
```

El modelo presenta el intercepto y la variable independiente abused. El modelo se escribe como:

$$y = 4.6959 + 7.2452(1) = 11.9411$$

No abused:

$$y = 4.6959 + 7.2452(0) = 4.6959$$

```
# media de ptsd en mujeres víctimas de abusos
mean(sexab[sexab$cscsa == "Abused", "ptsd"])
```

```
## [1] 11.94109
```

```
# media de ptsd en mujeres no víctimas de abusos
mean(sexab[sexab$cscsa == "NotAbused", "ptsd"])
```

```
## [1] 4.695874
```

```
mean(sexab[sexab$cscsa == "Abused", "ptsd"]) - mean(sexab[sexab$cscsa == "Not Abused", "ptsd"])
```

```
## [1] 7.245219
```

Independientemente del nivel que se escoja como referencia, el resultado es el mismo. Lo único que cambia es el valor de la intersección, ya que cambia el nivel de referencia, y el signo de la pendiente.

Podemos obtener también modelos utilizando las variables dummy para cada categoría.

```
modelo <- lm(ptsd ~ abused, data = sexab)
summary(modelo)
```

```
## Call:
## lm(formula = ptsd ~ abused, data = sexab)
##
## Residuals:
##   Min     1Q   Median     3Q    Max
## -8.0451 -2.3123  0.0951  2.1645  7.0514
##
## Coefficients:
##              Estimate Std. Error  t value Pr(>|t|)
## (Intercept)    4.6959     0.6237   7.529 1.00e-10 ***
## abused         7.2452     0.8105   8.939 2.17e-13 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 3.473 on 74 degrees of freedom
## Multiple R-squared:  0.5192, Adjusted R-squared:  0.5127
## F-statistic: 79.9 on 1 and 74 DF, p-value: 2.172e-13
```

El *p-value* obtenido para el predictor *abused* es el mismo que el obtenido en el *t-test* empleado anteriormente para comparar las medias. El valor estimado de la

intersección (4.6959), se corresponde con valor promedio de la variable respuesta en el nivel de referencia. La pendiente estimada (7.2452) se interpreta como el valor promedio de la influencia que tiene el predictor sobre la variable respuesta en comparación al nivel de referencia. En este caso, las mujeres que han sufrido abusos durante su infancia tienen en promedio 7.2452 unidades más de *ptsd* que las que no los han sufrido. Esta cantidad se corresponde con la diferencia de medias de ambos niveles (Amat, 2016).

```
modelo <- lm(ptsd ~ not_abused, data = sexab)
summary(modelo)
```

```
## Call:
## lm(formula = ptsd ~ not_abused, data = sexab)
##
## Residuals:
##   Min     1Q   Median     3Q    Max
## -8.0451 -2.3123  0.0951  2.1645  7.0514
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   11.9411    0.5177   23.067 < 2e-16 ***
## not_abused    -7.2452    0.8105  -8.939  2.17e-13 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 3.473 on 74 degrees of freedom
## Multiple R-squared:  0.5192, Adjusted R-squared:  0.5127
## F-statistic: 79.9 on 1 and 74 DF, p-value: 2.172e-13
```

En R la función `lm()` identifica automáticamente si un predictor es de tipo cualitativo y escoge como nivel de referencia el primero en base al orden alfabético.

(Amat, 2016) Si se desea especificar cuál debe ser el nivel de referencia empleado por `lm()`, se puede recurrir a la función `relevel()`

Es importante tener en cuenta que los *p-values* obtenidos por un *t-test* y por un modelo lineal que contenga un predictor cualitativo con dos niveles serán los mismos siempre y cuando el *t-test* no incluya una corrección de varianzas desiguales.

Representación grafica del modelo

```
ggplot(data = sexab, mapping = aes(x = csa, y = ptsd)) +  
  geom_point( size = 3) +  
  labs(title = "modelo", x = "ptsd") +  
  geom_smooth(method = "lm", se = FALSE, color = "red") +  
  theme_bw() +  
  theme(plot.title = element_text(hjust = 0.5))
```

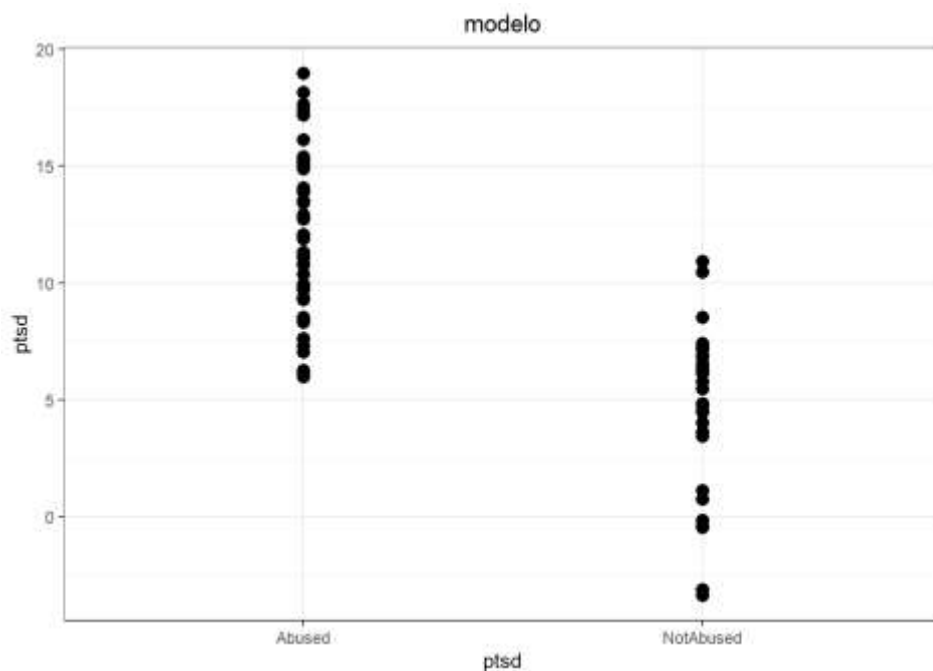


Figura 50. Gráfico del modelo

d) Residuales

```
# predecir
nuevos_datos <- data.frame(csa= seq(min(0),max(1)))
predict_value <- predict(modelo)
head(predict_value)
```

```
##          1          2          3          4          5          6
## 11.94109 11.94109 11.94109 11.94109 11.94109 11.94109
```

```
qqnorm(modelo$residuals)
qqline(modelo$residuals)
```

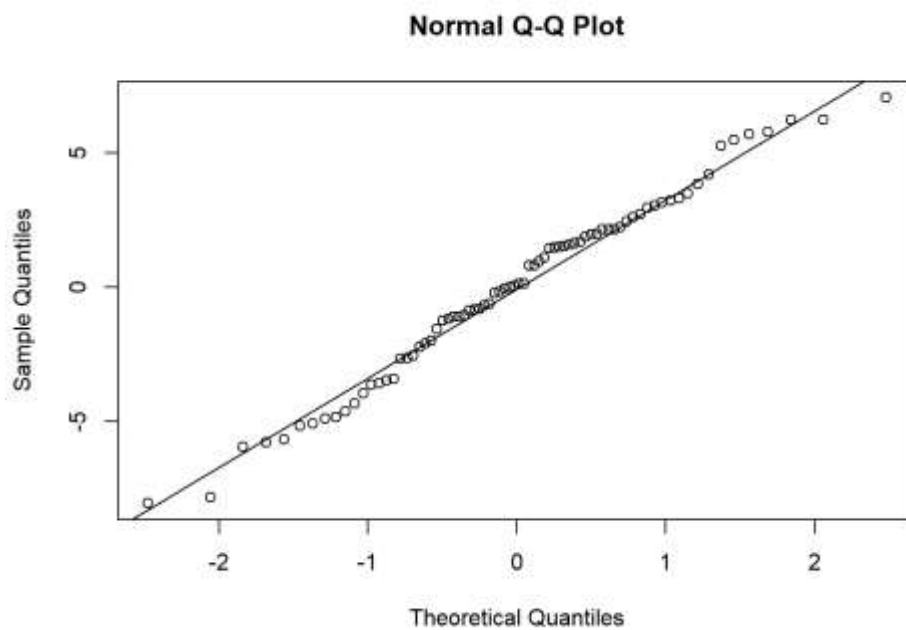


Figura 51. Gráfico Q-Q PLOT

```
shapiro.test(modelo$residuals)
```

```
##  
## Shapiro-Wilk normality test  
##  
## data:  modelo$residuals  
## W = 0.98418, p-value = 0.465
```

```
# Varianza constante de los residuos (Homocedasticidad):  
# Test de Breush-Pagan  
library(lmtest)
```

```
bptest(modelo)
```

```
## studentized Breusch-Pagan test  
##  
## data:  modelo  
## BP = 0.015186, df = 1, p-value = 0.9019
```

```
# Autocorrelacion de residuos:  
#test de Durwin Watson  
library(lmtest)  
dwtest(modelo)
```

```
##  
## Durbin-Watson test  
##  
## data:  modelo  
## DW = 1.8303, p-value = 0.1948  
## alternative hypothesis: true autocorrelation is greater than 0
```


Ejemplo con varias categorías de respuesta en variable cualitativa:

El siguiente ejemplo lo extraemos de Camacho, C. (2020), esta prueba es equivalente a la prueba de análisis de la varianza donde se estudia el efecto de una variable cualitativa de varias categorías con otra cuantitativa. Como se sabe, para aplicar el modelo de regresión lineal han de respetarse los supuestos del modelo: linealidad, normalidad y homocedasticidad. Los dos últimos son los mismos que los supuestos del análisis de la varianza.

La solución consiste en generar tantas variables independientes como categorías haya en el factor, y a continuación codificar cada una de estas variables con “ceros” y “unos” según la categoría a la que pertenezca los distintos sujetos.

Ejemplo.

Supongamos que aplicamos tres métodos de enseñanza (A, B y C) sobre tres grupos de sujetos, generaríamos tres variables: X_1 , X_2 y X_3 . Los sujetos que pertenecen al grupo A serían codificados como 1 (presencia en X_1) en la variable X_1 y 0 en las restantes (ausencia en X_2 y X_3). Así:

$$\begin{array}{ccc} X_1 & X_2 & X_3 \\ 1 & 0 & 0 \end{array}$$

Los sujetos que pertenecen al grupo B, tendrían la siguiente codificación:

$$\begin{array}{ccc} X_1 & X_2 & X_3 \\ 0 & 1 & 0 \end{array}$$

Y los sujetos pertenecientes al grupo C:

$$\begin{array}{ccc} X_1 & X_2 & X_3 \\ 0 & 0 & 1 \end{array}$$

Obsérvese que no es necesaria la variable X_3 . Con las dos primeras variables codificadas siempre estamos al tanto del grupo al que pertenecen los distintos sujetos. Si explícitamente están en X_1 o X_2 , no hay problemas, y si no están en ninguna de ellas, entonces se entiende que están en X_3 . Matemáticamente es conveniente hacerlo así, porque si no estaremos introduciendo una variable (cualquiera de ellas) que queda explicada por las otras, con lo que nos encontraremos con un problema de colinealidad, con matrices singulares y sin posible solución. Por tanto, generaremos dos variables con la siguiente codificación:

	X1	X2
Grupo A	1	0
Grupo B	0	1
Grupo C	0	0

Por otro lado, el hecho de plantear el análisis de la varianza como un problema de regresión múltiple permite salvar el supuesto de linealidad. De nuevo, cada una de las variables independientes sólo tiene dos posibles valores sobre los cuales establecer una recta. Ahora la ecuación de regresión corresponde geoméricamente con un plano y aunque las tres medias no estén alineadas en una recta (una dimensión) sí lo están en un plano (dos dimensiones)

Ejemplo. Supongamos que tenemos tres grupos de sujetos de estudiantes de matemáticas a los que hemos aplicado tres métodos de enseñanza distintos: A, B y C. Los resultados en esta materia son los siguientes:

	A	B	C
	6	5	7
	7	6	6
	6	5	6
	5	5	7
	4	4	8
	5	5	8
	5	5	7
	5	6	6

Si aplicáramos sobre estos datos un análisis de la varianza, configuraríamos la matriz de datos de la siguiente manera:

```

Método A:
Ŷ = 6.875 - 1.250X1 - 1.75X2 = 6.875 - 1.25*1 - 1.75*0 = 5.625

Método B:
Ŷ = 6.875 - 1.250X1 - 1.75X2 = 6.875 - 1.25*0 - 1.75*1 = 5.125

Método C:
Ŷ = 6.875 - 1.250X1 - 1.75X2 = 6.875 - 1.25*0 - 1.75*0 = 6.875

que son las medias de los grupos A, B y C respectivamente

```

Figura 52. Resultados de los grupos.

24. Apuntes varios (misceláneos)

a) Origen del método de mínimos cuadrados y regresión

Linear Models with R, by Julian J. Faraway

El método de mínimos cuadrados fue publicado en 1805 por Adrien Marie Legendre. El término *regresión* proviene de la publicación que hizo Francis Galton en 1885 llamada *Regression to mediocrity*. En ella, Galton emplea el método de mínimos cuadrados para demostrar que los hijos de parejas altas tienden a ser también altos, pero no tanto como sus padres y que los hijos de parejas de poca estatura tienden a ser bajos, pero no tanto como sus padres.

b) Significado de modelo lineal

Linear Models with R, by Julian J. Faraway

En los modelos lineales los parámetros se incorporan en la ecuación de forma lineal, sin embargo, los predictores no tienen por qué ser lineales. La siguiente ecuación muestra un modelo lineal en el que el predictor X_2 no es lineal:

$$y = \beta_0 + \beta_1 X_1 + \beta_2 \log(X_2) + \epsilon$$

En contraposición, el siguiente no es un modelo lineal:

$$y = \beta_0 + \beta_1 X_1^{\beta_2} + \epsilon$$

En ocasiones, algunas relaciones no lineales pueden transformarse de forma que se pueden expresar de manera lineal:

$$y = \beta_0 X_1^{\beta_1} \epsilon$$

$$\log(y) = \log(\beta_0) + \beta_1 \log(X_1) + \log(\epsilon)$$

c) Ventajas del método de mínimos cuadrados para estimar los coeficientes de un modelo lineal

Linear Models with R, by Julian J. Faraway

Si bien existen alternativas al método de mínimos cuadrados para obtener la estimación de los coeficientes de correlación $\hat{\beta}_i$ de un modelo lineal, este presenta una serie de ventajas:

- Tiene una interpretación geométrica, lo que facilita su comprensión,
- Si los errores son independientes y se distribuyen de forma normal, su solución equivale a la estimación de máxima verosimilitud (*likelihood*).
- Los $\hat{\beta}_i$ son estimaciones insesgadas.

d) Ecuaciones de curvas de aproximación:

Si el modelo es no lineal, entonces puede aproximarse a una de las siguientes curvas de aproximación y sus ecuaciones.

$y = a + bX$	lineal
$y = ab^x$ o $\text{Log } Y = \log a + (\log b)X = a + bX$	exponencial
$y = aX^b$ o $\log y = \log a + b \log X$	potencial, geométrica
$y = a + b(1/x)$	hiperbólico
$y = a + bX + CX^2$	parabólica o curva cuadrática
$Y = a + bX + CX^2 + dX^3$	Cúbica
$y = ab^x + g$	exponencial modificada
$y = aX^b + g$	Geométrica modificada
$\log y = \log p + b^x \log q = ab^x + g$	Gompertz

$y = 1/(ab^x + g)$ o $1/y = ab^x + g$ logística.

Etc.

e) Método de los mínimos cuadrados

El método de los mínimos cuadrados es equivalente al de máxima probabilidad, es un método que proporciona resultados aceptables y tiene la ventaja de normalidad.

El problema que se plantea es cómo calcular las cantidades a y b (estimadores desconocidos) a partir de un conjunto de n observaciones de forma que se minimice el error. Para esto sigamos los siguientes pasos.

Al tomar una muestra “ n ” tendremos n pares de observaciones para X, Y , sobre las que definimos $\hat{Y} = a + bX$ medimos el error que se comete al aproximar Y mediante \hat{Y} calculando la suma de las diferencias entre los valores reales y los aproximados al cuadrado (para que sean positivas y no se compensen los errores) (Marin, 2020):

Representación gráfica del concepto de mínimos cuadrados

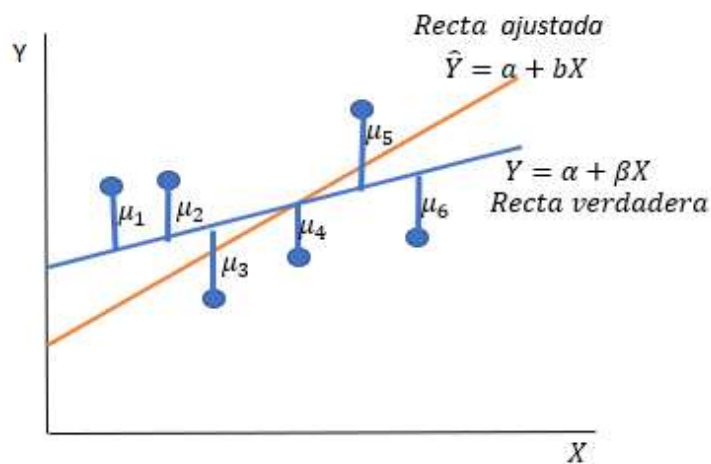


Figura 53. Representación gráfica del concepto de mínimos cuadrados

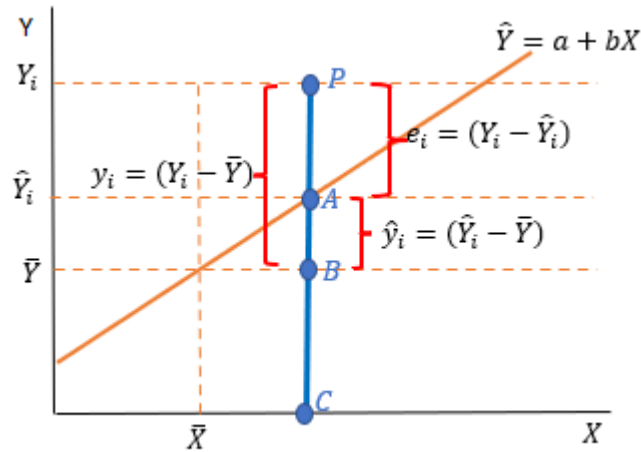


Figura 54. Representación gráfica mediante desviaciones

Donde:

$$e_i = (Y_i - \hat{Y}_i) = \text{variación no explicada por el modelo}$$

$$\hat{y}_i = (\hat{Y}_i - \bar{Y}) = \text{variación explicada por el modelo}$$

$$y_i = (Y_i - \bar{Y}) = \text{variación total}$$

Considere el punto P y trace una perpendicular al eje X, dicha perpendicular cortará a la línea estimada en A, a la línea de \bar{Y} en B y al eje X en C. entonces $\overline{OC} = x_i$; $\overline{PC} = Y_i$; $\overline{AC} = \hat{Y}_i$

Una vez que tenemos definido el error de aproximación las cantidades que lo minimizan se calculan derivando con respecto a ambas e igualando a cero (*procedimiento de los mínimos cuadrados*) (Baron, 1998):

Entonces el procedimiento mínimo cuadrático se basa en el hecho de minimizar la suma de cuadrados de los residuos (errores).

$$Y_i - \hat{Y}_i = e_i \text{ Elevando al cuadrado y aplicando sumatorias}$$

Iniciemos con la varianza no explicada

$$\sum_{i=1}^n e_i^2 = \sum_{i=1}^n (Y_i - \hat{Y}_i)^2 \quad \text{Además, } \hat{Y}_i = a + bX_i$$

$$\sum_{i=1}^n e_i^2 = \sum_{i=1}^n (Y_i - a - bX_i)^2$$

Para minimizar la suma de cuadrados de los errores derivar parcialmente respecto a cada estimador.

$$\frac{\partial \sum_{i=1}^n e_i^2}{\partial a} = 2 \sum_{i=1}^n (Y_i - a - bX_i)(-1) = -2 \sum_{i=1}^n (Y_i - a - bX_i)$$

$$\frac{\partial \sum_{i=1}^n e_i^2}{\partial b} = 2 \sum_{i=1}^n (Y_i - a - bX_i)(-X_i) = -2 \sum_{i=1}^n (Y_i - a - bX_i)(X_i)$$

Igualando ambas ecuaciones a cero para obtener un mínimo

$$-2 \sum_{i=1}^n (Y_i - a - bX_i) = 0$$

$$-2 \sum_{i=1}^n (Y_i - a - bX_i)(X_i) = 0$$

Dividiendo entre -2 y aplicando sumatorias

$$\sum_{i=1}^n (Y_i - a - bX_i) = 0$$

$$\sum_{i=1}^n (X_i Y_i - aX_i - bX_i^2) = 0 \quad ;$$

$$\sum_{i=1}^n Y_i - na - b \sum_{i=1}^n X_i = 0$$

$$\sum_{i=1}^n X_i Y_i - a \sum_{i=1}^n X_i - b \sum_{i=1}^n X_i^2 = 0$$

Despejando en función de la variable dependiente

$$\sum_{i=1}^n Y_i = na + b \sum_{i=1}^n X_i$$

$$\sum_{i=1}^n X_i Y_i = a \sum_{i=1}^n X_i + b \sum_{i=1}^n X_i^2$$

A esta fórmula se le denomina Ecuaciones normales,

Por ecuaciones simultáneas, multiplicando al numerador por $-\sum_{i=1}^n X_i$ y al denominador por n

$$\left. \begin{array}{l} \sum_{i=1}^n Y_i = na + b \sum_{i=1}^n X_i \\ \sum_{i=1}^n X_i Y_i = a \sum_{i=1}^n X_i + b \sum_{i=1}^n X_i^2 \end{array} \right\} \begin{array}{l} - \sum_{i=1}^n X_i \\ n \end{array}$$

$$\left. \begin{aligned} -\left(\sum X_i\right)\left(\sum Y_i\right) &= -na - \sum X_i - b\left(\sum X_i\right)^2 \\ n \sum X_i Y_i &= na \sum X_i + nb \sum X_i^2 \end{aligned} \right\}$$

$$\frac{-\left(\sum_{i=1}^n X_i\right)\left(\sum_{i=1}^n Y_i\right) = -b\left(\sum_{i=1}^n X_i\right)^2}{n \sum_{i=1}^n X_i Y_i = nb \sum_{i=1}^n X_i^2}$$

$$n \sum_{i=1}^n X_i Y_i - \left(\sum_{i=1}^n X_i\right)\left(\sum_{i=1}^n Y_i\right) = nb \sum_{i=1}^n X_i^2 - b\left(\sum_{i=1}^n X_i\right)^2$$

$$n \sum_{i=1}^n X_i Y_i - \left(\sum_{i=1}^n X_i\right)\left(\sum_{i=1}^n Y_i\right) = b\left(n \sum_{i=1}^n X_i^2 - \left(\sum_{i=1}^n X_i\right)^2\right)$$

Estas fórmulas se denominan POR OBSERVACIONES (*coeficiente de regresión de Y sobre X*)

$$b = \frac{n \sum_{i=1}^n X_i Y_i - \left(\sum_{i=1}^n X_i\right)\left(\sum_{i=1}^n Y_i\right)}{n \sum_{i=1}^n X_i^2 - \left(\sum_{i=1}^n X_i\right)^2}$$

$$a = \frac{\left(\sum_{i=1}^n Y_i\right)\left(\sum_{i=1}^n X_i^2\right) - \left(\sum_{i=1}^n X_i\right)\left(\sum_{i=1}^n X_i Y_i\right)}{n \sum_{i=1}^n X_i^2 - \left(\sum_{i=1}^n X_i\right)^2}$$

En *b* dividamos entre *n*:

$$b = \frac{\sum_{i=1}^n X_i Y_i - \left[\frac{\left(\sum_{i=1}^n X_i \cdot \sum_{i=1}^n Y_i\right)}{n}\right]}{\sum_{i=1}^n X_i^2 - \left[\frac{\sum_{i=1}^n X_i^2}{n}\right]}$$

$$b_1 = \frac{\sum XY - \bar{y} \sum X}{\sum X^2 - \bar{x} \sum X}$$

$$a = \frac{\sum Y - b \sum X}{N}$$

La desventaja del uso de observaciones aparece cuando los valores de las variables son bastante grandes, en este caso es factible usar **desviaciones**.

CALCULO DE LOS ESTIMADORES MEDIANTE DESVIACIONES

Tomemos la primera ecuación normal: $\sum_{i=1}^n Y_i = na + b \sum_{i=1}^n X_i$ y dividamos entre n:

$$\frac{\sum_{i=1}^n Y_i}{n} = \frac{na + b \sum_{i=1}^n X_i}{n}$$

queda $\bar{Y} = a + b\bar{X}$, en función del estimador a:

$$a = \bar{Y} - b\bar{X} \quad \text{y} \quad b = (\bar{Y} - a)/\bar{X}$$

Esta ecuación indica una propiedad de la línea de regresión o de los estimadores mínimos cuadráticos que dice “la línea estimada pasa por un punto cuyas coordenadas son (\bar{X}, \bar{Y}) ”.

Ahora, restemos al modelo $\hat{Y} = a + bX$ la ecuación $\bar{Y} = a + b\bar{X}$

$$\begin{array}{r} \hat{Y} = a + bX - \\ \bar{Y} = a + b\bar{X} \\ \hline \hat{Y} - \bar{Y} = b(X - \bar{X}) \\ \underbrace{\hat{y}} \qquad \underbrace{x} \end{array}$$

Queda: $\hat{y}_i = bx_i$ Modelo con **desviaciones**:

Recordemos:

$$y_i = Y_i - \bar{Y}; Y_i = y_i + \bar{Y} \dots \dots \dots (a)$$

$$\hat{y}_i = \hat{Y} - \bar{Y}; \hat{Y} = \hat{y}_i + \bar{Y} \dots \dots \dots (b)$$

$$e_i = Y_i - \hat{Y}_i; \dots \dots \dots (c)$$

En la ecuación (c) reemplacemos (a) y (b):

$$e_i = (y_i + \bar{Y}) - (\hat{y}_i + \bar{Y}) ; \text{Además se conoce que } \hat{y}_i = bx_i$$

$$e_i = (y_i + \bar{Y}) - (bx_i + \bar{Y})$$

$$e_i = y_i + \bar{Y} - bx_i - \bar{Y}$$

$$e_i = y_i - bx_i$$

Aplicando el principio de los mínimos cuadrados:

$$\sum e_i^2 = \sum (y_i - bx_i)^2$$

$$\frac{\partial \sum e_i^2}{\partial b} = 2 \sum (y_i - bx_i)(-x_i)$$

$$\frac{\partial \sum e_i^2}{\partial b} = -2 \sum (y_i - bx_i)(x_i)$$

Dividiendo entre -2, desarrollando, aplicando sumatorias e igualando a cero

$$= \sum (x_i y_i - bx_i^2)$$

$$\sum x_i y_i - b \sum x_i^2 = 0 \quad \text{despejando } b; \quad b = \frac{\sum x_i y_i}{\sum x_i^2}$$

Que por Observaciones puede escribirse como:

$$b = \frac{\sum (X_i - \bar{X})(Y_i - \bar{Y})}{\sum (X_i - \bar{X})^2}$$

Otra forma de obtener los datos como desviaciones es mediante las siguientes formulas:

$$\sum x_i y_i = \sum X_i Y_i - \frac{(\sum X_i)(\sum Y_i)}{n}$$

$$\sum x_i^2 = \sum X_i^2 - \frac{(\sum X_i)^2}{n}$$

Hallando la covarianza (S_{xy}):

Tomemos $b = \frac{\sum (X_i - \bar{X})(Y_i - \bar{Y})}{\sum (X_i - \bar{X})^2}$ siendo las medias $\bar{X} = \sum X_i / n$ y

$$\bar{Y} = \sum Y_i / n$$

Tomemos la primera parte: $\sum (X_i - \bar{X})(Y_i - \bar{Y})$ multiplicando:

$$\sum (X_i Y_i - \bar{X} Y_i - \bar{Y} X_i + \bar{X} \bar{Y})$$

$$\sum X_i Y_i - \bar{X} \sum Y_i - \bar{Y} \sum X_i + n \bar{X} \bar{Y}$$

$$\sum X_i Y_i - \frac{\sum X_i}{n} \sum Y_i - \frac{\sum Y_i}{n} \sum X_i + n \left(\frac{\sum X_i}{n} \right) \left(\frac{\sum Y_i}{n} \right)$$

$$\sum X_i Y_i - \frac{\sum X_i \sum Y_i}{n} - \frac{\sum Y_i \sum X_i}{n} + \frac{\sum X_i \sum Y_i}{n}$$

$$\sum X_i Y_i - \frac{\sum X_i \sum Y_i}{n} \quad \text{dividiendo y multiplicando por } n$$

$$\begin{aligned} & \sum X_i Y_i - \frac{\sum X_i \sum Y_i}{n} * \frac{n}{n} \\ & \sum X_i Y_i - \frac{n \sum X_i \sum Y_i}{n.n} \\ d_{XY} &= \sum X_i Y_i - n \bar{X} \bar{Y} \end{aligned}$$

Para hallar la covarianza dividir sobre n , queda:

$$S_{xy} = \frac{\sum X_i Y_i}{n} - \bar{X} \bar{Y} \text{ Covarianza}$$

Hallando la varianza poblacional de X:

Tomemos el denominador: $\sum (X_i - \bar{X})^2$ multiplicando:

$$\begin{aligned} & \sum (X_i - \bar{X})(X_i - \bar{X}) \\ & \sum (X_i^2 - 2\bar{X}X_i + \bar{X}^2) \\ & \sum X_i^2 - 2\bar{X} \sum X_i + n\bar{X}^2 \\ & \sum X_i^2 - 2 \frac{\sum X_i \sum X_i}{n} + n \frac{\sum X_i \sum X_i}{n^2} \\ & \sum X_i^2 - 2 \frac{\sum X_i \sum X_i}{n} + \frac{\sum X_i \sum X_i}{n} \\ & \sum X_i^2 - \frac{\sum X_i \sum X_i}{n} = \sum X_i^2 - \frac{(\sum X_i)^2}{n} \\ & \sum X_i^2 - \frac{\sum X_i \sum X_i}{n} * \frac{n}{n} \\ & \sum X_i^2 - \frac{n \sum X_i \sum X_i}{n.n} \\ S_{XX} &= \sum X_i^2 - n\bar{X}^2 \end{aligned}$$

La varianza se calcula dividiendo todo entre n :

$$S_X^2 = \frac{\sum X_i^2 - n\bar{X}^2}{n} = \frac{\sum X_i^2}{n} - \bar{X}^2$$

Varianza poblacional. Note que $S_X^2 = \sigma^2$

$$S_X^2 = \frac{\sum X_i^2 - n\bar{X}^2}{n - 1} = \frac{\sum X_i^2 - \frac{(\sum X_i)^2}{n}}{n - 1}$$

Varianza muestral.

Entonces:

$$b = \frac{S_{xy}}{S_X^2} = \frac{cov(X, Y)}{V(X)}$$

Hallando la varianza de Y

Tomando los desvíos: $\sum (Y_i - \bar{Y})^2$; podemos llegar a:

$$S_Y^2 = \frac{\sum Y_i^2 - n\bar{Y}^2}{n} = \frac{\sum Y_i^2}{n} - \bar{Y}^2$$

Nota: podemos observar que existen diferentes fórmulas para hallar b ,

Media aritmética de los estimadores

$E(Y_i) = \alpha + \beta X_i$; pero

$$\begin{aligned} E(\bar{Y}) &= E\left(\frac{\sum Y_i}{n}\right) = \frac{1}{n}E\left(\sum Y_i\right) = \frac{1}{n}\sum E(Y_i) = \\ &= \frac{1}{n}\sum (\alpha + \beta X_i) = \frac{1}{n}(n\alpha + \beta \sum X_i) = \alpha + \beta \frac{\sum X_i}{n} \\ E(\bar{Y}) &= \alpha + \beta \bar{X} \end{aligned}$$

i) Media aritmética de b

$$\begin{aligned} E(b) &= E\left[\frac{\sum (X_i - \bar{X})(Y_i - \bar{Y})}{\sum (X_i - \bar{X})^2}\right] = \frac{\sum (X_i - \bar{X})[E(Y_i) - E(\bar{Y})]}{\sum (X_i - \bar{X})^2} = \\ &= \frac{\sum (X_i - \bar{X})[\alpha + \beta X_i - \alpha - \beta \bar{X}]}{\sum (X_i - \bar{X})^2} = \frac{\beta \sum (X_i - \bar{X})^2}{\sum (X_i - \bar{X})^2} = \beta \end{aligned}$$

$$E(b) = \beta$$

ii) Media aritmética de a :

$$\begin{aligned} E(a) &= E(\bar{Y} - b\bar{X}) = E(\bar{Y}) - \bar{X}E(b) = \alpha + \beta \bar{X} - \bar{X}\beta = \alpha \\ E(a) &= \alpha \end{aligned}$$

f) Propiedades de la regresión lineal

Una vez que ya tenemos perfectamente definida \hat{Y} , (o bien \hat{X}) nos preguntamos las relaciones que hay entre la media y la varianza de esta y la de Y (o la de X). La respuesta nos la ofrece la siguiente proposición:

Proposición

En los ajustes lineales se conservan las medias, es decir $\hat{Y} = \bar{Y}$; $\hat{X} = \bar{X}$

En cuanto a la varianza, no necesariamente son las mismas para los verdaderos valores de las variables X e Y y sus aproximaciones \hat{X} y \hat{Y} , pues sólo se mantienen en un factor de r^2 , es decir,

$$\begin{aligned} S^2_{\hat{Y}} &= r^2 S^2_Y \\ S^2_{\hat{X}} &= r^2 S^2_X \end{aligned}$$

Demostración

Basta probar nuestra afirmación para la variable Y , ya que para X es totalmente análogo:

$$\bar{\hat{y}} = a + b\bar{x} = (\bar{y} - b\bar{x} + b\bar{x}) = \bar{y}$$

$$\begin{aligned} S^2_{\hat{Y}} &= b^2 S^2_X = \frac{S_{XY}^2}{S_X^2 \cdot S_X^2} \cdot S_X^4 \\ &= \frac{S_{XY}^2}{S_X^2 \cdot S_Y^2} \cdot S_Y^2 \\ &= \underbrace{\left[\frac{S_{XY}}{S_X \cdot S_Y} \right]^2}_r \cdot S_Y^2 \\ &= r^2 S_Y^2 \end{aligned}$$

Donde se ha utilizado la magnitud que denominamos *coeficiente de correlación*, r , y que ya definimos anteriormente como

$$r = \frac{S_{XY}}{S_X \cdot S_Y}$$

Observación

Como consecuencia de este resultado, podemos decir que *la proporción de varianza explicada por la regresión lineal es del $r^2 \cdot 100\%$.*

Nos gustaría tener que $r = 1$, pues en ese caso ambas variables tendrían la misma varianza, pero esto no es cierto en general. Todo lo que se puede afirmar, como sabemos, es que

$$-1 \leq r \leq 1$$

y por tanto $0 \leq S_{\hat{Y}}^2 \leq S_Y^2$

La cantidad que le falta a la **varianza de regresión**, $S_{\hat{Y}}^2$, para llegar hasta la varianza total de Y , S_Y^2 , es lo que se denomina **varianza residual**, que no es más que la varianza de $E = Y - \hat{Y}$, ya que

$$\begin{aligned} S_Y^2 &= \frac{1}{n} \sum_{i=1}^n (y_i - \bar{y})^2 \\ &= \frac{1}{n} \sum_{i=1}^n [\hat{y}_i - \bar{y} + e_i]^2 \\ &= \frac{1}{n} \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 + \frac{1}{n} \sum_{i=1}^n e_i^2 + \frac{1}{n} \sum_{i=1}^n (\hat{y}_i - \bar{y}) e_i \\ &= S_{\hat{Y}}^2 + S_E^2 + \underbrace{\frac{1}{n} \sum_{i=1}^n (\hat{y}_i - \bar{y}) e_i}_0 \end{aligned}$$

El tercer sumando se anula según las ecuaciones normales expresadas en la relación:

$$\begin{aligned} \sum_{i=1}^n (\hat{y}_i - \bar{y}) e_i &= \sum_{i=1}^n e_i [a + b x_i - (a + b \bar{x})] \\ &= b \sum_{i=1}^n e_i (x_i - \bar{x}) \\ &= b \sum_{i=1}^n e_i x_i - b \bar{x} \sum_{i=1}^n e_i \\ &= 0 \end{aligned}$$

Por ello

$$S_E^2 = S_Y^2 - S_{\hat{Y}}^2$$

Obsérvese que entonces la bondad del ajuste es

$$R_{Y|X}^2 = 1 - \frac{S_E^2}{S_Y^2} = 1 - (1 - r^2) = r^2$$

Para el ajuste contrario se define el error como $E = Y - \hat{Y}$, y su varianza residual es también proporcional a $1-r^2$:

$$S_E^2 = S_X^2 - S_{\hat{X}}^2 = S_X^2(1 - r^2)$$

y el coeficiente de determinación (que sirve para determinar la bondad del ajuste de X en función de Y) vale:

$$R_{X|Y}^2 = 1 - \frac{S_E^2}{S_X^2} = 1 - (1 - r^2) = r^2$$

lo que resumimos en la siguiente proposición:

Proposición

Para los ajustes de tipo lineal se tiene que los dos coeficientes de determinación son iguales a r^2 , y por tanto representan además la proporción de varianza explicada por la regresión lineal:

$$\boxed{R_{X|Y}^2 = r^2 = R_{Y|X}^2}$$

25. Bibliografía

- Amat, J. (2016). *Correlación lineal y Regresión lineal simple*.
https://raw.githubusercontent.com/JoaquinAmatRodrigo/Estadistica-con-R/master/PDF_format/24_Correlación_y_Regresión_lineal.pdf
- Baron, J. (1998). *Bioestadística: Métodos y Aplicaciones*.
<https://idoc.pub/documents/curso-de-estadistica2009-od4poyOkvwlp>
- Calcina, A. (2019). *Regresión Lineal simple MD*.
<https://es.scribd.com/document/600516829/Reresion-Lineal-Simple-MD-1>
- Camacho martinez, C. (2020). *ANALISIS DE DATOS EN PSICOLOGIA*.
<https://personal.us.es/vararey/adatos2/regcualitativas.pdf>
- Canavos, G. (1995). *Internet Archive*.
https://archive.org/stream/ProbabilidadYEstadistica.MetodosYEjercicios/EstadisticaYProbabilidad_djvu.txt
- Cybertesis. (2020). *Repositorio de Tesis Digitales*.
<http://cybertesis.unmsm.edu.pe/handle/cybertesis/859>
- Económicos. (2010). *Econometría*. <http://www.econometricos.com.ar/wp-content/uploads/2010/03/3-2012-IE-3.pdf>
- Gujarati, D. (2010). *Econometría* (5ta ed.). Centro universitario de ciencias Económico-administrativas.
- Hernandez, E. a. (2023). *Modelos de regresión con R*.
https://fhernanb.github.io/libro_regresion/
- Marin, J. (2020). *Estadística descriptiva*.
<https://halweb.uc3m.es/esp/Personal/personas/jmmarin/esp/EDescrip/tema6.pdf>
- Martinez, J. (1992). *La Seguridad social: Elementos, propiedades y relaciones* [Universidad de Complutense de madrid].
<https://eprints.ucm.es/id/eprint/3419/1/T18079.pdf>
- Morocho, D. (2015). Economía matemática. In *Economía matemática y modelos econométricos* (p. 389).
<https://danielmorochoruiz.files.wordpress.com/2015/09/economc3ada-matemc3aitica-y-modelos-econoc3b3micos.pdf>

Ruiz, F. (2020). *Econometría*. <https://idoc.tips/econometriavarelapdf-pdf-free.html>

Tareas, B. (2013). *Econometría*.

<https://www.buenastareas.com/ensayos/Econometria/7066172.html>

Trujillo, O. (2017). *Econometría*. <https://es.slideshare.net/orvilletrrivera/capitulo-de-introduccion-a-la-econometria>