

## COMPARACIÓN DE TRES MÉTODOS DE REGRESIÓN LINEAL USANDO PROCEDIMIENTOS DE SIMULACIÓN

JUAN C. TORRES C.\*

---

### Resumen

Cuando desea ajustarse un modelo lineal a un conjunto de datos, el método de regresión usualmente más empleado es el de mínimos cuadrados. Este método es óptimo si la distribución de los residuos es gaussiana. Existen casos en donde el supuesto de normalidad en los residuales no se cumple y se hace necesario el uso de métodos alternativos de regresión, como la regresión vía mínimas desviaciones absolutas (LAD) o la regresión no paramétrica basada en rangos, los cuales no requieren de supuestos distribucionales sobre los residuos y permiten obtener una mejor estimación de los parámetros del modelo de regresión.

*Palabras clave:* Modelo lineal; análisis de regresión; simulación.

### Abstract

When it's necessary to fit a lineal model to a data set, least squares regression method is usually used. This method is optimum if the residuals distribution is normal. When the assumption of residuals normality doesn't comply it's necessary to use alternative regression methods, as Least absolute deviations (LAD) or Non parametric regression based on ranks, which don't need the assumption about the residuals distribution and allow a better estimation of regression model parameters.

*Key words:* Linear model, regression analysis, simulation.

---

\*Estadístico egresado de la Universidad Nacional de Colombia y Consultor del Departamento Administrativo Nacional de Estadística DANE. Correo electrónico: jctorrres@dane.gov.co

## 1. Introducción

El objetivo de este artículo es determinar cuál de los tres métodos de regresión: Mínimo cuadrática, No paramétrica basada en rangos o Mínima desviación absoluta, presenta un mejor ajuste cuando la distribución de los errores aleatorios tiene forma diferente de la gaussiana en elongación y cuando hay presencia de observaciones atípicas u “outliers”. Para cumplir con el objetivo, se calcularon estimaciones de los coeficientes del modelo de regresión y sus medidas de ajuste para los tres métodos de regresión mediante procedimientos de simulación realizados con el paquete estadístico SAS.

## 2. Regresión Mínimo Cuadrática

Para el modelo dado por la expresión  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{e}$ , el método construye el estimador  $\hat{\boldsymbol{\beta}} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{Y}$ , que minimiza la suma de cuadrados de los errores. Para un conjunto de datos observados, cuando la expresión  $f(\mathbf{b}) = \sum (Y_i - \mathbf{x}_i\mathbf{b})^2$  se hace mínima, el vector de valores  $\mathbf{b}$  se conoce como la estimación mínimo cuadrática de  $\boldsymbol{\beta}$ . En la función anterior  $\mathbf{x}_i$  representa la  $i$ -ésima fila de la matriz  $\mathbf{X}$ . Se estima la varianza de la población de los errores con  $s^2 = \sum e_i^2 / (n - p - 1)$ , en donde los  $e_i$  son los errores obtenidos con los datos observados,  $n$  es el número de observaciones realizadas y  $p$  es el número de variables explicativas en el modelo.

## 3. Regresión no Paramétrica Basada en Rangos

El método se basa en el modelo lineal  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{e}$ . No existen fórmulas explícitas para el estimador  $\hat{\boldsymbol{\beta}}$ . Sin embargo, mediante un algoritmo iterativo es posible obtener el vector de estimación de  $\boldsymbol{\beta}$ . Para un conjunto de datos observados, el vector  $\boldsymbol{\beta} = [b_0 \ \mathbf{b}]$  se conoce como la estimación no paramétrica de  $\boldsymbol{\beta}$ , en donde el vector  $\mathbf{b}$  de tamaño  $1 \times p$ , minimiza la función

$$g(\mathbf{b}) = \sum \left[ \text{rango}(y_i - \mathbf{x}_i\mathbf{b}) - \frac{n+1}{2} \right] (y_i - \mathbf{x}_i\mathbf{b})$$

en donde  $\mathbf{b} = [b_1, b_2, \dots, b_p]'$  y  $\mathbf{x}_i = [x_{i1}, x_{i2}, \dots, x_{ip}]$  y la estimación no paramétrica de  $b_0$  se obtiene como la mediana de las diferencias de  $y_i - \mathbf{x}_i\mathbf{b}$ . Ya que no existen fórmulas explícitas para calcular los coeficientes estimados, Birkes y Dodge (1993) describen un algoritmo para estimarlos. De igual manera, los mismos autores describen un algoritmo para estimar la desviación estándar de

los errores, cuya notación usual es  $\tau$ . Además afirman que un estimador para la desviación estándar de los errores es  $\tilde{s} = \tau / 1.023$  si la distribución de los errores es normal.

## 4. Regresión Vía Mínima Desviación Absoluta

El método se basa en el modelo lineal  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{e}$ . No existen fórmulas explícitas para el estimador  $\hat{\boldsymbol{\beta}}$ . Sin embargo, mediante un algoritmo iterativo es posible obtener el vector de estimación de  $\boldsymbol{\beta}$  el cual es escogido de tal forma que la suma de los valores absolutos de los errores,  $\sum |e|$ , sea la más pequeña posible. Para un conjunto de datos observados, cuando la expresión  $f(\mathbf{b}) = \sum |y_i - \mathbf{x}_i \mathbf{b}|$  se hace mínima, el vector de estimación  $\mathbf{b}$  obtenido, es conocido como la estimación mínima desviación absoluta de  $\boldsymbol{\beta}$ . En la función anterior  $\mathbf{x}_i$  representa la  $i$ -ésima fila de la matriz  $\mathbf{X}$ . Como no existen fórmulas explícitas para calcular los coeficientes estimados, Birkes y Dodge (1.993) describen un algoritmo para estimar los coeficientes de regresión. La desviación estándar de los errores se estima con  $\tilde{s} = 1,483 * MAD$ , donde  $MAD$  es la mediana de las desviaciones absolutas con respecto a la mediana de los residuales, calculado con los residuos diferentes de cero.

## 5. Distribuciones $g$ y $h$ de Tukey

La familia de distribuciones  $g$  y  $h$  comprende una considerable variedad de distribuciones con características especiales en cuanto a asimetría y elongación, por lo cual resulta de gran utilidad cuando se desea simular datos que provengan de distribuciones con formas diferentes a la distribución normal.

### 5.1. Asimetría

Si  $Z$  es una variable normal estándar y  $g$  es una constante real, la variable aleatoria  $Y_g(Z)$ , definida como  $Y_g(Z) = (e^{gz} - 1)g^{-1}$  se dice que tiene la distribución  $g$  de Tukey para un valor dado de  $g$ . El parámetro  $g$  controla la magnitud y la dirección de la asimetría.

## 5.2. Elongación

Si  $Z$  es una variable aleatoria normal estándar y  $h$  es una constante real, la variable aleatoria  $Y_h(Z)$  dada por  $Y_h(Z) = Ze^{hZ^2/2}$  se dice que tiene distribución  $h$  de Tukey para un valor dado de  $h$ . El parámetro  $h$  controla la cantidad de la elongación de la distribución. Las distribuciones de la familia  $h$  de Tukey son simétricas y su valor esperado es cero.

## 6. Proceso de Simulación

Para el cálculo de las estimaciones del ajuste con errores aleatorios provenientes de la familia de distribuciones  $h$  de Tukey, se efectuaron 1000 simulaciones de modelos lineales con una y dos variables independientes, con tamaños de muestra 20 y 50, y usando cada uno de los tres métodos de regresión. Se utilizó el mismo procedimiento para errores provenientes de una distribución normal contaminada con un porcentaje de “outliers” dado. Todas las simulaciones fueron llevadas a cabo con el paquete estadístico SAS.

En el caso de la regresión simple, el modelo simulado fue  $\mathbf{Y} = \beta_0 + \beta_1 X_1 + \mathbf{e}$  donde  $\beta_0 = 4$  y  $\beta_1 = 5$  y la variable explicativa  $X_1$  se dejó fija en cada simulación y fue generada de una distribución uniforme  $U(2, 5)$  para los tamaños de muestra 20 y 50. Para la regresión múltiple, el modelo simulado fue  $\mathbf{Y} = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \mathbf{e}$  donde  $\beta_0 = 3$ ,  $\beta_1 = 4$ , y  $\beta_2 = 5$ . Las variables explicativas  $X_1$  y  $X_2$  se fijaron para cada simulación y fueron generadas de distribuciones  $U(1, 4)$  y  $U(2, 5)$  respectivamente, para los tamaños de muestra 20 y 50. Para la simulación de errores aleatorios con formas diferentes a la distribución normal se utilizaron las siguientes distribuciones:

Tabla 1: Distribuciones  $h$  de Tukey  
Simuladas con su varianza teórica

$h$	$V(e)$
0	1
0.2	2.1517
0.4	11.1803
0.6	Infinita
0.8	Infinita
1	Infinita

En la simulación de errores aleatorios con distribución normal y presencia de “outliers”, se utilizó el modelo contaminado  $f(x) = (1 - \epsilon)\phi(x) + \epsilon\phi(x/k)$  donde  $\phi(x)$  es la función de distribución normal estándar y  $k$  es igual a 3. Los valores de contaminación  $\epsilon$  que se utilizaron fueron 0.05, 0.10, 0.15 y 0.20, lo que equivale a contaminar la distribución  $N(0, 1)$  con los porcentajes de “outliers” 5%, 10%, 15% y 20%, generados de una distribución  $N(0, 3)$ , para cada uno de los modelos descritos anteriormente.

## 7. Criterios de Comparación.

Los criterios de comparación utilizados fueron los siguientes:

- Error cuadrático medio de los estimadores de los coeficientes del modelo (ECM)
- Error absoluto medio de los estimadores de los coeficientes del modelo (MAE)
- Promedios de los coeficientes de determinación estimados.

Cuando se simulan los errores de una distribución normal con un porcentaje de “outliers”, se compararon los promedios de las cantidades de “outliers” detectados que fueron generados para cada método de regresión, usando el método de los residuales estandarizados. Según Birkes y Dodge (1993), un residual estandarizado mayor a 2.5 en valor absoluto, puede catalogarse como “outlier”.

## 8. Resultados de las Simulaciones

Para el análisis de los resultados, se compararon el error cuadrático medio (ECM) y el valor absoluto medio (MAE) de los estimadores de los parámetros obtenidos con cada uno de los tres métodos de regresión: mínimos cuadrados, no paramétrica basada en rangos y LAD. Se exponen solamente los resultados para regresión simple ya que los correspondientes a regresión múltiple con dos variables son similares. Así mismo, se presentan los resultados para el error cuadrático medio (ECM) de los estimadores, ya que los resultados para el error absoluto medio (MAE) son similares (Ver Pulido y Torres, 1997). En cada una de las tablas presentadas los datos en letra negrilla son los menores valores de las medidas de comparación, los que poseen letra cursiva son los que siguen

en el orden de menor a mayor y los datos en letra normal corresponden a los mayores valores.

Tabla 2: Resultados de Regresión simple  $n = 20$

Distribución		$\beta_0$			$\beta_1$		
		MC	NP	LAD	MC	NP	LAD
normal		<b>1.1716</b>	<i>1.2739</i>	1.8169	<b>0.0448</b>	<i>0.0472</i>	0.0688
h	0.2	2.8079	<b>1.8567</b>	<i>2.2316</i>	0.1047	<b>0.0708</b>	<i>0.0869</i>
	0.4	12.2780	<b>2.4435</b>	<i>2.6841</i>	0.3529	<b>0.0972</b>	<i>0.1051</i>
	0.6	60.3540	<i>3.3256</i>	<b>3.2183</b>	2.6535	<i>0.1332</i>	<b>0.1269</b>
	0.8	828.1600	<i>4.2417</i>	<b>3.2424</b>	27.5080	<i>0.1738</i>	<b>0.1318</b>
	1	5432.562	<i>5.4006</i>	<b>3.6525</b>	174.5601	<i>0.2148</i>	<b>0.1440</b>
% outliers	5	<i>1.7075</i>	<b>1.5531</b>	2.1669	<i>0.0652</i>	<b>0.0588</b>	0.0827
	10	<i>2.1565</i>	<b>1.8126</b>	2.4468	<i>0.0835</i>	<b>0.0687</b>	0.0938
	15	2.7169	<b>1.9623</b>	<i>2.4529</i>	0.1058	<b>0.0755</b>	<i>0.0965</i>
	20	3.4506	<b>2.5211</b>	<i>3.0297</i>	0.1345	<b>0.0963</b>	<i>0.1164</i>

Tabla 3: Resultados de regresión simple  $n = 50$

Distribución		$\beta_0$			$\beta_1$		
		MC	NP	LAD	MC	NP	LAD
normal		<b>0.4407</b>	<i>0.483</i>	0.7151	<b>0.0154</b>	<i>0.0166</i>	0.0248
h	0.2	1.0208	<b>0.676</b>	<i>0.7780</i>	0.0356	<b>0.0235</b>	<i>0.0271</i>
	0.4	3.4151	<b>0.7573</b>	<i>0.7885</i>	0.1207	<b>0.0260</b>	<i>0.0271</i>
	0.6	17.2182	<i>0.9823</i>	<b>0.8817</b>	0.7438	<i>0.0341</i>	<b>0.0306</b>
	0.8	236.1500	<i>1.1016</i>	<b>0.8631</b>	7.7681	<i>0.0382</i>	<b>0.0295</b>
	1	1546.770	<i>1.3799</i>	<b>0.9745</b>	49.4998	<i>0.0485</i>	<b>0.0338</b>
% outliers	5	<i>0.5962</i>	<b>0.5575</b>	0.7769	<i>0.0207</i>	<b>0.0187</b>	0.0267
	10	<i>0.8799</i>	<b>0.6575</b>	0.8821	<i>0.0307</i>	<b>0.0226</b>	0.0309
	15	0.9968	<b>0.6949</b>	<i>0.9055</i>	0.0342	<b>0.0238</b>	<i>0.0316</i>
	20	1.1307	<b>0.7779</b>	<i>1.0061</i>	0.0391	<b>0.0271</b>	<i>0.0359</i>

## 8.1. Resultados de la Estimación de Parámetros

Los resultados se presentan en las Tablas 2 y 3. En ellas puede observarse que el ECM de los estimadores de los parámetros tiene un comportamiento similar para los tamaños de muestra 20 y 50.

### 8.1.1. Vector de errores con distribución normal

Observando las tablas 2 y 3, para los tamaños de muestra 20 y 50, el menor ECM de los estimadores de los parámetros se obtuvo con la regresión mínimo cuadrática y sigue en orden ascendente la regresión no paramétrica basada en rangos.

### 8.1.2. Vector de error con distribución $h$ de Tukey

Se observa en la Figura 1, realizada para el ECM de  $\beta_1$ , regresión simple y tamaño de muestra 50, que cuando el parámetro  $h$  de la distribución de los errores es igual 0.2, 0.4, 0.6, 0.8 y 1, el ECM de los estimadores calculados con la regresión de mínimos cuadrados aumenta aceleradamente, mientras que con la regresión LAD y la regresión no paramétrica basada en rangos es más estable.

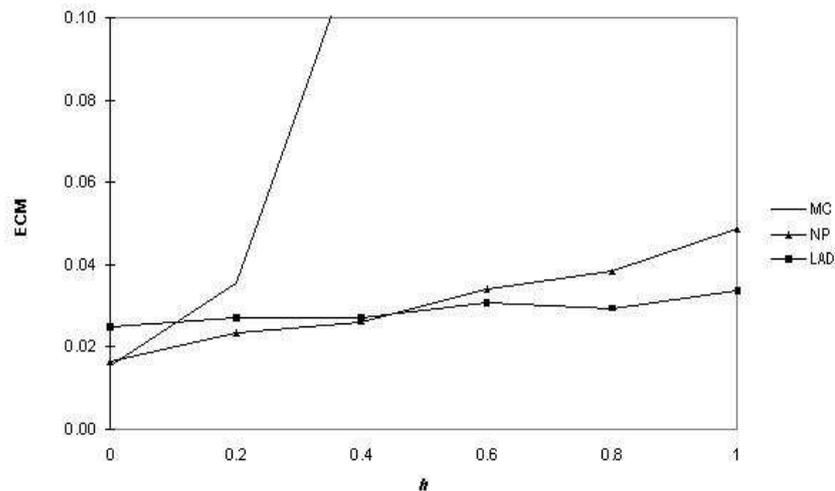


Figura 1. Error cuadrático medio. Estimador de  $\beta_1$ .  
Errores familia  $h$  de Tukey

En la figura también se observa que para distribuciones con parámetros  $h$  igual a 0.2 y 0.4 el menor ECM de los parámetros estimados fue obtenido con la regresión no paramétrica basada en rangos y le sigue la regresión LAD. Para distribuciones  $h$  igual a 0.6, 0.8 y 1 el menor ECM de los parámetros estimados fue obtenido con la regresión LAD y le sigue la regresión no paramétrica basada en rangos.

### 8.1.3. Distribución de errores normal con un porcentaje de “outliers”

Como se observa tanto en la Figura 2, realizada para  $\beta_1$ , con el modelo simple y tamaño de muestra 50, así como en la Tabla 2, el menor ECM de los estimadores de los parámetros para los porcentajes de “outliers” simulados, fue obtenido con la regresión no paramétrica basada en rangos.

En orden ascendente siguen las estimaciones obtenidas con la regresión mínimo cuadrática, con un 5 % y 10 % de “outliers” en la distribución de los errores, y cuando se tiene un 15 % y 20 % de “outliers” en la distribución de los errores, le sigue el ECM de los estimadores calculados vía regresión LAD.

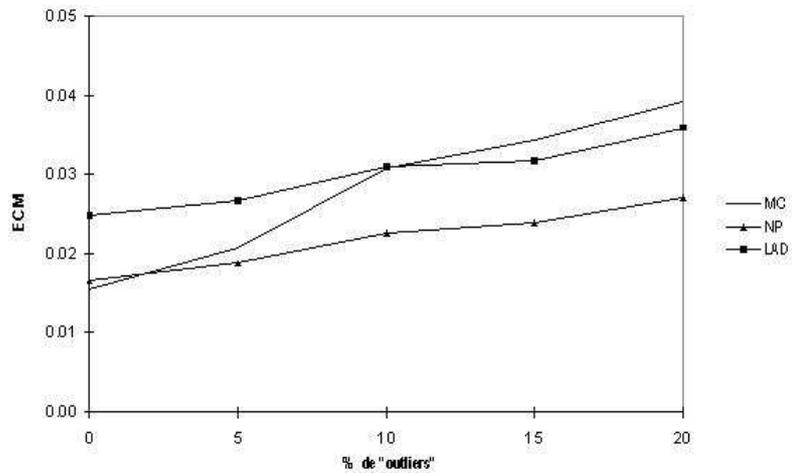


Figura 2. Error cuadrático Medio para  $\beta_1$ .  
Errores con porcentaje de “outliers”

## 8.2. Resultados de Generación y Detección de Outliers

Como se puede observar en la tabla 4, el método que más “outliers” detectó en promedio, usando el método de los residuales estandarizados fue la regresión LAD y le sigue la regresión no paramétrica basada en rangos, para los diferentes porcentajes de “outliers” generados.

Tabla 4: Resultados, Generación y Detección de Outliers  
Regresión Simple (Promedio)

		Generados	Detectados		
<b>n</b>	%		MC	NP	LAD
20	5	0.966	<b>0.223</b>	0.266	0.337
	10	1.981	<b>0.336</b>	0.462	0.583
	15	3.055	<b>0.440</b>	0.637	0.897
	20	3.994	<b>0.545</b>	0.755	1.112
50	5	2.473	<b>0.736</b>	0.838	0.910
	10	5.146	<b>1.212</b>	1.596	1.754
	15	7.603	<b>1.410</b>	2.152	2.444
	20	9.864	<b>1.589</b>	2.552	2.950

## 8.3. Resultados de la Medida de Ajuste $R^2$

En la Tabla 5 se puede observar que los valores de  $R^2$  son mayores con la regresión mínimo cuadrática y sigue la regresión no paramétrica basada en rangos, para las diferentes distribuciones y los tamaños de muestra usados en la simulación.

## 9. Conclusiones

- Cuando la distribución de los errores es simétrica, de colas pesadas, generada de las distribuciones  $h$  de Tukey, se obtienen mejores estimaciones de los parámetros del modelo con los métodos de regresión no paramétrica basada en rangos y LAD que con la regresión mínimo cuadrática. Si la distribución de los errores es simétrica de colas muy pesadas, las mejores estimaciones se obtienen con la regresión LAD.

Tabla 5: Valores de  $R^2$  para Regresión Simple (Promedio)

Distribución		20			50		
		MC	NP	LAD	MC	NP	LAD
<i>normal</i>	0	0.968	0.967	<b>0.966</b>	0.970	0.970	<b>0.970</b>
	0.2	0.935	0.932	<b>0.931</b>	0.939	0.938	<b>0.938</b>
	0.4	0.856	0.847	<b>0.845</b>	0.849	0.845	<b>0.845</b>
	<i>h</i> 0.6	0.728	0.707	<b>0.705</b>	0.677	0.667	<b>0.666</b>
	0.8	0.584	0.575	<b>0.543</b>	0.483	0.466	<b>0.4650</b>
	1	0.482	0.436	<b>0.432</b>	0.328	0.303	<b>0.302</b>
<i>%</i> <i>outliers</i>	5	0.957	0.955	<b>0.954</b>	0.959	0.958	<b>0.958</b>
	10	0.944	0.942	<b>0.941</b>	0.947	0.946	<b>0.946</b>
	15	0.931	0.928	<b>0.926</b>	0.936	0.935	<b>0.935</b>
	20	0.919	0.916	<b>0.914</b>	0.927	0.926	<b>0.925</b>

- Se puede observar la robustez de los métodos de regresión LAD y no paramétrica basada en rangos, frente a la regresión mínimo cuadrática, cuando la distribución de los errores presenta colas pesadas.
- Cuando la distribución de los errores es normal, con un porcentaje de “outliers” mayor o igual al 5%, se obtienen mejores estimaciones de los parámetros con la regresión no paramétrica y le sigue la regresión LAD.
- El empleo del coeficiente de determinación usual no es adecuado para comparar el ajuste de modelos con criterios de estimación diferentes, ya que para todas las distribuciones simuladas, éste fue mayor para la regresión mínimo cuadrática a pesar de que no siempre tuvo la mejor estimación de los parámetros del modelo de regresión.
- En cuanto a la detección de “outliers”, cuando la distribución de los errores es normal, el método que detectó más en promedio usando el criterio de los residuales estandarizados, fue la regresión LAD, seguido de la regresión no paramétrica basada en rangos.
- En general, se puede observar que ninguno de los tres métodos de regresión en estudio fue uniformemente mejor para la estimación de los parámetros. Si la distribución de los errores aleatorios simétrica de colas muy pesadas, es mejor estimar los parámetros con la regresión LAD,

mientras que en el caso de que los errores tengan distribución normal con altos porcentajes de “outliers” es preferible estimar los parámetros usando la regresión no paramétrica basada en rangos.

**Agradecimientos** El autor expresa sus agradecimientos al Doctor Jorge Martínez Collantes, por sus valiosos aportes en la elaboración de este documento. Este artículo se basó en el trabajo de grado del mismo nombre realizado por Elsa Leonor Pulido y Juan Carlos Torres en 1.997, para optar al título de Estadístico en la Universidad Nacional de Colombia.

## Referencias

- [1] BIRKES, D. y DODGE, Y. (1993). *Alternative methods of regression*. John Wiley & Sons, Inc.
- [2] PULIDO, E. y TORRES, J.(1.997). *Análisis comparativo de tres métodos de regresión: mínimos cuadrados, no paramétrica basada en rangos y mínima desviación absoluta, usando métodos de simulación*. Tesis de Grado. Universidad Nacional de Colombia.