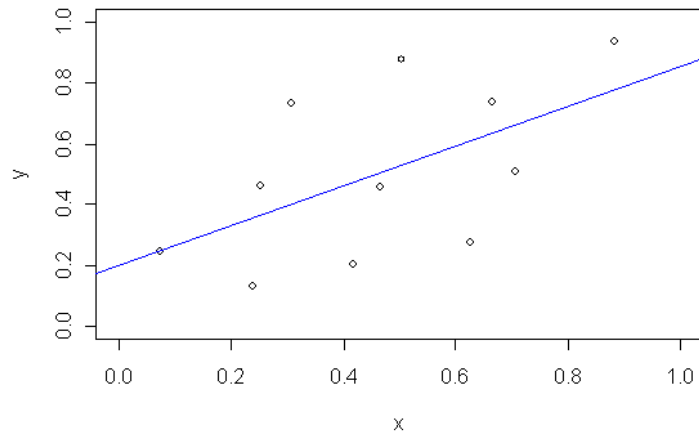


## REGRESIÓN LINEAL SIMPLE

Dado un conjunto de pares de datos  $(x_i, y_i)$ ,  $i = 1, \dots, n$ , se han desarrollado diversos métodos para ajustar una recta de la forma

$$y = a + bx,$$

al diagrama de dispersión de dichos datos.



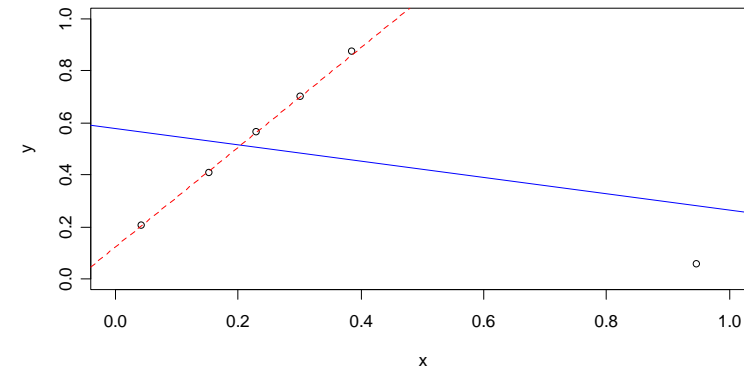
**Figura 1.** Diagrama de dispersión y recta ajustada

El más conocido y el más utilizado es el Método de Cuadrados Mínimos. La pendiente y la ordenada al origen se obtienen minimizando la suma de cuadrados de los residuos. Requiere cálculos algebraicamente sencillos y desarrollo matemático directo.

Desafortunadamente la recta ajustada por el método de cuadrados mínimos no es resistente. Un sólo punto puede tomar control del ajuste, como lo muestra la figura siguiente, y llevar a una interpretación absolutamente errónea de la relación entre  $x$  e  $y$ .

La figura muestra dos tipos de ajuste: por cuadrados mínimos y por el método resistente de **Suma de**

**cuadrados mínimos truncada.** La función `ltsreg`, por defecto, minimiza la suma de los cuadrados de aproximadamente la mitad de los residuos más pequeños.



**Figura 2.** Recta de Cuadrados Mínimos y Recta Resistente

Hemos obtenido el gráfico anterior mediante las siguientes instrucciones:

```
> plot(c(0,1),c(1,0),
type="n",xlab="x",ylab="y")
> xy<-locator(type="p") # marcamos los puntos en
el gráfico y sus coordenadas son guardadas en xy.
> xy
$x
[1] 0.0409885 0.1514274 0.2296549 0.3009800
0.3838091 0.9452066

$y
[1] 0.20550193 0.40878378 0.56776062
0.70067568 0.87789575 0.05955598

> abline(lm(y~x,xy),col=4)#recta de cuadr.
mín.
```

```
> library(MASS)
> abline(ltsreg(xy$x,xy$y),lty=2,col=2)
```

### Método de Cuadrados mínimos

Los coeficientes de la recta de cuadrados mínimos (CM) se eligen entre todos los posibles pares de valores aquellos que minimizan la suma de cuadrados de los residuos

$$S = \sum_{i=1}^n (y_i - (a + bx_i))^2 \quad (1)$$

Derivando (1) respecto de  $a$  y  $b$

$$\frac{\partial S}{\partial a} = -2 \sum_{i=1}^n (y_i - (a + bx_i))$$

$$\frac{\partial S}{\partial b} = -2 \sum_{i=1}^n (y_i - (a + bx_i)) x_i$$

e igualando a cero se obtienen los coeficientes de la recta estimada por CM:

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

$$a_{CM} = \bar{y} - b_{CM} \bar{x}$$

Ninguna otra recta tendrá, para el mismo conjunto de datos, una suma de cuadrados de los residuos tan baja como la obtenida por CM. En este sentido, CM garantiza la solución que mejor ajusta a ese conjunto de datos.

Se puede obtener fácilmente, que la recta de CM pasa por el punto  $(\bar{x}, \bar{y})$ .

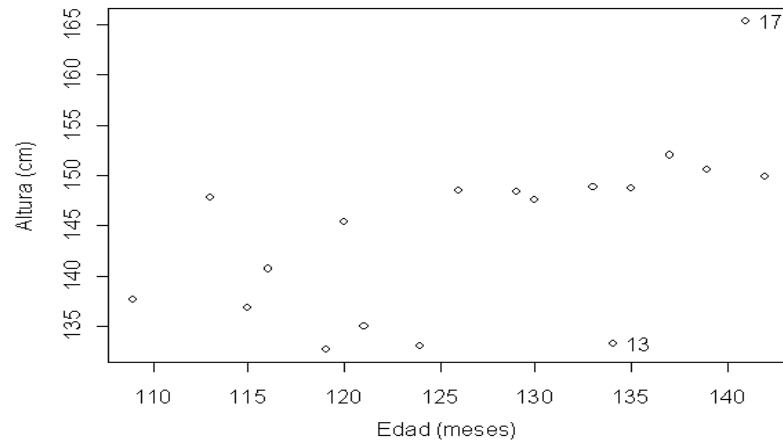
**EJEMPLO:** Edades y alturas (Greenberg 1953) correspondientes a 18 niños de una escuela (tabla 1) y graficamos la altura contra la edad en la figura 1.

Leemos los datos de un [archivo texto sin](#) la fila de nombres y luego asignamos nombres a las variables

```
> altedad <- read.table(file.choose())
> names(altedad) <- c("edad", "altura")
> plot(altedad$edad, altedad$altura,
      xlab="Edad
(meses)", ylab="Altura (cm) ")
> identify(altedad$edad, altedad$altura)
[1] 17 13
```

**Tabla 1: Edad y altura de niños en una escuela**

Niño	Edad (meses)	Altura (cm)
1	109	137.6
2	113	147.8
3	115	136.8
4	116	140.7
5	119	132.7
6	120	145.4
7	121	135.0
8	124	133.0
9	126	148.5
10	129	148.3
11	130	147.5
12	133	148.8
13	134	133.2
14	135	148.7
15	137	152.0
16	139	150.6
17	141	165.3
18	142	149.9



**Figura 3.** Diagrama de dispersión de altura vs. edad

Aunque los datos no siguen claramente una recta tampoco presentan un patrón notablemente curvo. De manera que una recta ajustada debería servir para resumir como aumenta la altura ( $y$ ) con la edad ( $x$ ) en ese grupo de niños.

Veamos que ocurre con el ajuste por cuadrados mínimos.

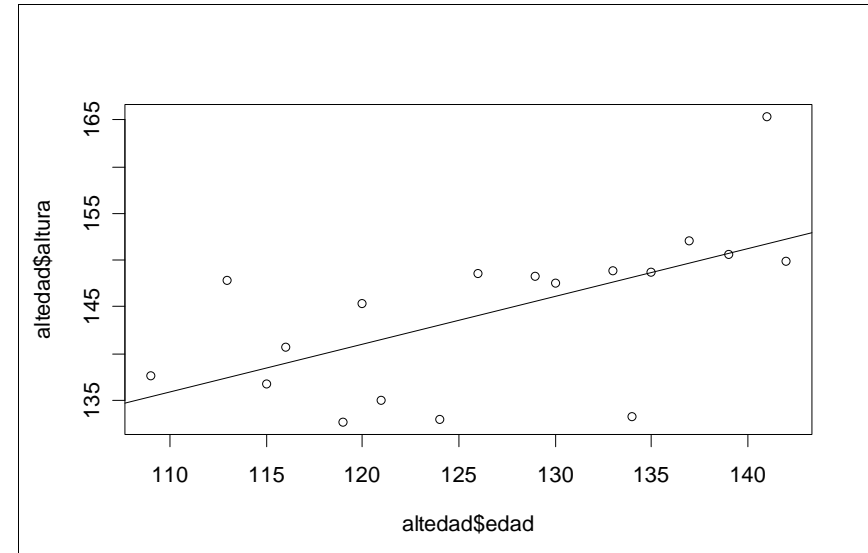
```
> plot(altedad$edad, altedad$altura)
> abline(lm(altedad$altura~altedad$edad))
> lm(altedad$altura~altedad$edad)
```

Call:

```
lm(formula = altedad$altura ~ altedad$edad)
```

Coefficients:

```
(Intercept) altedad$edad
79.6962      0.5113
```



**Figura 4.** Recta de Cuadrados Mínimos

Algunos datos pueden afectar al método de cuadrados mínimos. Sin embargo es útil ver que cuando el comportamiento de los puntos es razonable las rectas obtenidas por ajustes resistentes y por cuadrados mínimos son similares.

#### Comparación del ajuste por cuadrados mínimos y por cuadrados mínimos truncados

	Ordenada al origen	Pendiente
lm	79.6962	0.5113
ltsreg	127.6854	0.1615
ltsreg, quantile=12	89.7321	0.4429

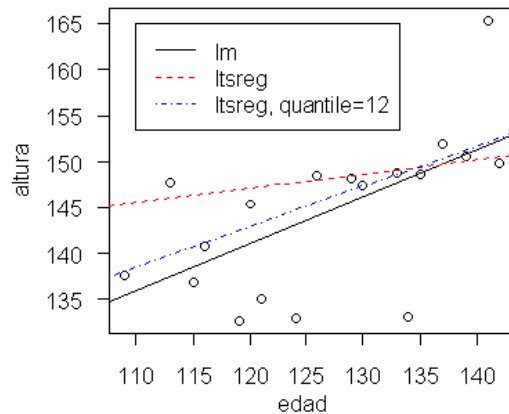


Figura 5. Comparación de tres ajustes

```
> library(MASS)
> plot(alteedad$edad,alteedad$altura, las=1,
      xlab="edad", ylab="altura")
> abline(lm(alteedad$altura~alteedad$edad))
>
> abline(ltsreg(alteedad$edad,alteedad$altura)
      ,col="red",lty= 2)
> abline(
  ltsreg(alteedad$edad,alteedad$altura,
    quantile=12),col="blue", lty=4 )
> legend(110,165,c("lm","ltsreg","ltsreg,
  quantile=12"),lty = c(1,2,4), col =
  c(1,2,4))
```

El parámetro "quantile" de la función ltsreg indica cuantos residuos minimiza toma por defecto

```
quantile = floor(n/2) + floor((p+1)/2)
```

En este caso, por defecto queda:

```
p = 2 (regresión lineal simple) y n = 18, luego
quantile = floor(18/2) + floor((2+1)/2) = 10
```

### Recta por medianas repetidas.

El procedimiento, propuesto por Siegel (1982), consiste en estimar la pendiente de la recta en dos etapas. En la primera etapa tomamos la mediana de las pendientes de las  $n-1$  rectas que pasan por un punto dado  $(x_i, y_i)$ ,  $i = 1, \dots, n$ ; y en la segunda etapa se toma la mediana de estas  $n$  pendientes. Es decir:

si definimos

$$b_{ij} = \frac{y_i - y_j}{x_i - x_j}$$

entonces la mediana de las pendientes de todas las rectas que pasan por  $(x_i, y_i)$  y algún otro punto será

$$\text{med}_{j \neq i} \{b_{ij}\}$$

luego a través de los puntos tendremos

$$b_{RM} = \text{med}_i \{ \text{med}_{j \neq i} \{b_{ij}\} \}$$

Para ajustar la ordenada al origen se calcula

$$a_i = y_i - b_{RM} x_i$$

y luego

$$a_{RM} = \text{med}_i \{a_i\}.$$

Siegel muestra que el método tiene punto de ruptura cercano a 1/2.

Veamos una derivación heurística para el caso  $n = 2k$ . En este caso el punto de ruptura exacto es  $(k-1)/n$ . Para ver esto supongamos que  $k-1$  de los datos son "salvajes" y que los restantes  $k+1$  son "buenos". Definimos

$$b_i = \text{med}_{j \neq i} \{b_{ij}\}$$

Si  $i_0$  indica un punto “bueno”,  $b_{i_0}$  está determinado por los restantes  $k$  puntos “buenos” y no por los  $k-1$  puntos “salvajes”. Por otro lado  $b_i$  para un punto “salvaje” debe ser “salvaje”. Exactamente  $k+1$  de los  $b_i$  son “buenos” y estas estimaciones de la pendiente determinan  $b_{RM}$ .

Cualquier número mayor de puntos salvajes” causaría la ruptura de  $b_{RM}$ . Como la ordenada al origen involucra únicamente una mediana simple, tolera  $k-1$  puntos salvajes entre  $2k$ . Tanto la pendiente como la ordenada al origen tienen el mismo punto de ruptura.

### EJEMPLO (cont) Cálculo de la pendiente y ordenada al origen por el método de medianas repetidas

La función “repmedians” está descrita en la práctica 8.

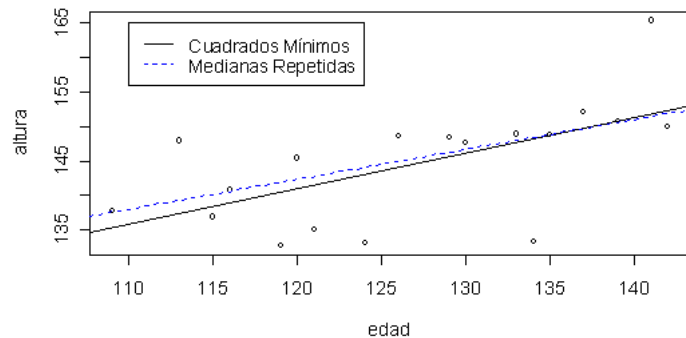


Figura 6

```
> attach(altedad)
> unlist(repmedians(edad, altura))
ord.origen  pendiente
90.4000000  0.4333333
> lm(altura ~ edad)[[1]] # Cuadrados Mínimos
(Intercept)      edad
 79.69623  0.5112868

> plot(edad, altura)
> abline(lm(altura ~ edad), lty = 1)
```

```
> abline(unlist(repmedians(edad, altura)),
lty=2,col=4)
> legend(110,165,c("Cuadrados
Mínimos","Medianas Repetidas"),lty =
c(1,2), col = c(1,4))
```

### Optativo. Recta resistente de tres grupos.

Para resumir el centro de un lote, el procedimiento resistente más simple es la mediana muestral. La técnica exploratoria que veremos para ajustar una recta (Tukey, 1970) deriva su resistencia de la mediana.

### Formación de tres grupos.

Ordenamos los valores de  $x$  de manera que  $x_1 \leq x_2 \leq \dots \leq x_n$ . Luego sobre la base de los datos ordenados dividimos los  $n$  datos  $(x_i, y_i)$  en tres grupos, el de la izquierda, el del centro y el de la derecha. Cuando no hay empates entre las  $x_i$ , la cantidad de puntos en los tres grupos depende del resto de dividir  $n$  por 3.

Ubicamos los puntos a los grupos de la siguiente manera:

Grupo	$n = 3k$	$n = 3k + 1$	$n = 3k + 2$
Izquierda	$k$	$k$	$k + 1$
Centro	$k$	$k + 1$	$k$
Derecha	$k$	$k$	$k + 1$

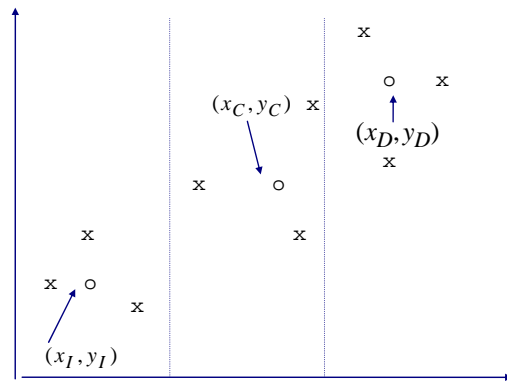
Los empates entre las  $x_i$  pueden impedir que se logre esta ubicación exactamente porque no separamos los empates. Todos los datos con el mismo valor de  $x$  van al mismo grupo.

### Puntos centrales de tres grupos (summary points).

Dentro de cada uno de los grupos formados, determinamos dos coordenadas de un punto central hallando primero la mediana de las  $x$ 's y luego la mediana de las  $y$ 's. Indicamos las coordenadas de los tres puntos centrales,  $I$  para izquierda,  $C$  para el centro y  $D$  para derecha:

$$(x_I, y_I), (x_C, y_C), (x_D, y_D).$$

La figura 7 muestra los puntos centrales en un ejemplo hipotético de nueve puntos. Como enfatiza la figura ninguno de los puntos centrales es necesariamente un dato, pues las medianas de las  $x$ 's y de las  $y$ 's se calculan en forma separada.



**Figura 7:** Los datos (x) y los puntos centrales (o) en un ejemplo hipotético.

Sin embargo todos los puntos centrales podrían ser datos, como ocurre frecuentemente cuando las  $x_i$ 's y las  $y_i$ 's tienen el mismo ordenamiento dentro de cada grupo.

Este método da una recta resistente. La mediana provee resistencia a valores salvajes en  $x$ ,  $y$  o ambos, siempre y

cuando el número de puntos en cada grupo no sea muy pequeño.

### Pendiente y ordenada al origen (slope and intercept)

Vamos a utilizar los puntos centrales para obtener la pendiente y la ordenada al origen de una recta

$$y = a + bx$$

Es usual utilizar la notación de  $y$  “sombrero” en la expresión de la recta ajustada para recordar que es la fuente de los valores ajustados tanto para valores  $x$ 's de los datos como en otros valores adecuados. Así tendremos:

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

La pendiente de la recta indica en cuantas unidades cambia  $y$  en respuesta a un cambio de una unidad en  $x$  y obtenemos esta información a partir de los puntos resumen izquierda y derecha:

$$b_0 = \frac{y_D - y_I}{x_D - x_I}.$$

De esta manera buscamos un balance entre

- (a) la ventaja de medir el cambio de  $y$  sobre un intervalo amplio de  $x$  y
- (b) la necesidad de tener suficientes datos en el grupo de la izquierda y el de la derecha para tener una resistencia adecuada.

Para cada punto central  $(x_I, y_I), (x_C, y_C), (x_D, y_D)$ , hay una recta que pasa por dicho punto y tiene como pendiente la pendiente ajustada  $b_0$ , el promedio de las ordenadas al origen de cada una de esas rectas es la ordenada al origen ajustada:

$$a_0 = \frac{1}{3}[(y_I - b_0 x_I) + (y_C - b_0 x_C) + (y_D - b_0 x_D)]$$

Nuevamente como los puntos centrales están basados en medianas  $a_0$  es resistente.

Para comparar, consideremos la pendiente y la ordenada al origen de la recta ajustada por cuadrados mínimos

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

$$a_{CM} = \bar{y} - b_{CM} \bar{x}$$

La imposibilidad de cualquier tipo de resistencia en los estimadores de cuadrados mínimos es evidente en la forma en que todos los datos entran en el cálculo de los coeficientes.

### Pendiente y valor central

Ajustar una recta en términos de la pendiente y la ordenada al origen es convencional pero generalmente artificial. La ordenada al origen da un valor de  $y$  cuando  $x = 0$ , que puede estar determinado muy imprecisamente y no tener importancia cuando los valores de  $x$  caen lejos del cero.

Generalmente es más útil realizar el ajuste en base a la pendiente y la ordenada correspondiente a un  $x$  que puede ser  $x = \bar{x}$  ó  $x = \text{mediana}(x_i)$  ó  $x = x_C$ , llamada **valor central**. Si por conveniencia trabajamos en  $x = x_C$ , entonces la recta inicial es:

$$\hat{y} = a + b_0(x - x_C) \quad (1)$$

Igual que antes consideramos las tres rectas, de la forma (1) con pendiente  $b_0$  que pasan por los tres puntos centrales respectivamente. Promediamos los valores  $a$  obtenidos para cada una de ellas:

para  $(x_I, y_I)$ ,  $a_I = y_I - b_0(x_I - x_C)$ ,

para  $(x_C, y_C)$ ,  $a_C = y_C$ ,

para  $(x_D, y_D)$ ,  $a_D = y_D - b_0(x_D - x_C)$ .

Por lo tanto la recta ajustada será

$$\hat{y} = a_0^* + b_0(x - x_C)$$

con  $b_0$  igual que antes y la ordenada en el valor central  $a_0^*$ , también llamado nivel, está dada por

$$a_0^* = \frac{1}{3}[a_I + a_C + a_D]$$

$$= \frac{1}{3}\{[y_I - b_0(x_I - x_C)] + y_C + [y_D - b_0(x_D - x_C)]\}.$$

### Residuos.

Una vez que se han obtenido la pendiente y el nivel para la recta ajustada, el paso inmediato siguiente es calcular el residuo para cada dato:

$$r_i = y_i - [a^* + b(x_i - x_C)]$$

Los residuos son la base de varios gráficos que permiten revelar una gran variedad de aspectos y patrones del ajuste.

Pero, en este caso, solamente necesitamos enfatizar una propiedad general de un conjunto de residuos, tanto en  $y$  versus  $x$  como en situaciones más complejas:

Sustituir los residuos en vez de los valores  $y$  originales (i.e. utilizar  $(x_i, r_i)$  en vez de  $(x_i, y_i)$ ,  $i = 1, \dots, n$ ) y repetir el ajuste lleva a un ajuste cero.

Para la recta esto significa que utilizar  $(x_i, r_i)$  como datos lleva a pendiente cero y nivel cero. En otras palabras, los residuos ya no contienen más información lineal para resumir, tal como ocurre al utilizar el método de cuadrados mínimos.

Los procedimientos resistentes frecuentemente requieren iteración.

La recta resistente de 3 grupos es un primer ejemplo. Si los residuos de la recta con pendiente  $b_0$  y nivel  $a_0^*$  no tienen pendiente cero y nivel cero, ajustamos una recta a ellos. La nueva pendiente y el nuevo nivel serán sustancialmente menores (en magnitud) que  $b_0$  y  $a_0^*$ . Por esta razón pensamos a  $b_0$  como un valor inicial para la pendiente y  $a_0^*$  como un valor inicial para el nivel (de allí el subíndice 0).

### Iteración

Generalmente esperamos que  $b_0$  y  $a_0^*$  necesiten alguna corrección. Ajustar una recta a los residuos de la recta inicial da las correcciones  $\delta_1$  y  $\gamma_1$  para la pendiente y el nivel respectivamente. Específicamente, utilizamos los residuos iniciales

$$r_i^{(0)} = y_i - [a_0^* + b_0(x_i - x_C)], \quad i = 1, \dots, n,$$

en lugar de  $y_i$  para repetir la mayoría de los pasos anteriores del proceso de ajuste. Las  $x_i$  no han cambiado, de manera que los grupos y las medianas de las  $x$ 's no varían durante el proceso iterativo.

La pendiente y el nivel corregidos son  $b_0 + \delta_1$  y  $a_0^* + \gamma_1$  y los nuevos residuos son

$$r_i^{(1)} = r_i^{(0)} - [\gamma_1 + \delta_1(x_i - x_C)], \quad i = 1, \dots, n.$$

Podríamos ahora intentar otra iteración. En general no sabemos si tenemos un conjunto adecuado de residuos hasta que verificamos que tienen un ajuste cero. En la práctica, continuamos las iteraciones hasta que la corrección a la pendiente es suficientemente pequeña (a lo sumo 1% ó 0.01%) del tamaño de  $b_0$ . Cada iteración agrega las correcciones de la pendiente y el nivel a los valores anteriormente modificados:

$$b_1 = b_0 + \delta_1, \quad b_2 = b_1 + \delta_2, \dots$$

y

$$a_1^* = a_0^* + \gamma_1, \quad a_2^* = a_1^* + \gamma_2, \dots$$

Las iteraciones son generalmente lo suficientemente pocas de manera de no llevar demasiado tiempo, la resistencia justifica el esfuerzo.

Para algunos conjuntos de datos las correcciones de la pendiente decrecen demasiado lentamente o después de algunos pasos pueden dejar de decrecer y en cambio oscilar entre dos valores con la misma magnitud y signo opuesto.

**Ejemplo:** continuamos con los datos de altura y edad



Como los datos no incluyen empates, los dividimos en tres grupos con 6 puntos cada uno. Los tres puntos centrales son:

$$(x_I, y_I) = (115.5, 139.15)$$

$$(x_C, y_C) = (127.5, 147.9)$$

$$(x_D, y_D) = (138, 150.25).$$

```
> sapply( altedad[1:6,] ,median)
  edad altura
115.50 139.15
> sapply( altedad[7:12,] ,median)
  edad altura
127.5  147.9
> sapply( altedad[13:18,] ,median)
  edad altura
138.00 150.25
```

Con `apply( altedad[13:18,] ,2, median)` se obtiene el resultado anterior.

```
> I <- apply( altedad[1:6,] ,2,median)
> C <- apply( altedad[7:12,] ,2,median)
> D <- apply( altedad[13:18,] ,2,
median)
```

La pendiente inicial es

```
> b0 <- (D[[2]]-I[[2]])/(D[[1]]-I[[1]])
>
> b0
[1] 0.4933333
```

$$b_0 = \frac{y_D - y_I}{x_D - x_I} = \frac{150.25 - 139.15}{138 - 115.5} = 0.4933.$$

y el nivel inicial es

$$a_0^* = \frac{1}{3}(145.07 + 147.90 + 145.07) = 146.01$$

La tabla 2 muestra los datos separados por grupo y los residuos de esta recta inicial.

**Tabla 2. Altura y edad de niños - residuos iniciales recta 3 grupos**

Niño	Edad (x)	Altura (y)	residuo	
			y - [146.01 + 0.4933(x- 127.5)]	
1	109	137.6		0.716
2	113	147.8		8.943
3	115	136.8		-3.044
4	116	140.7		0.363
5	119	132.7		-9.117
6	120	145.4		3.090
7	121	135.0		-7.804
8	124	133.0		-11.28
9	126	148.5		3.230
10	129	148.3		1.550
11	130	147.5		0.257
12	133	148.8		0.077
13	134	133.2		-16.02
14	135	148.7		-1.010
15	137	152.0		1.304
16	139	150.6		-1.083
17	141	165.3		12.63
18	142	149.9		-3.263

```
> altedad$altura -(146.01+
  0.4933*(altedad$edad-127.5))
[1] 0.71605 8.94285 -3.04375 0.36295
[5] -9.11695 3.08975 -7.80355 -11.28345
[9] 3.22995 1.55005 0.25675 0.07685
[13] -16.01645 -1.00975 1.30365 -1.08295
[17] 12.63045 -3.26285
```

Calculemos las correcciones a la pendiente y al nivel:

$$\delta_1 = \frac{-1.045 - 0.545}{138 - 115.5} = -0.0707$$

$$\gamma_1 = \frac{1}{3}(-0.30 + 0.07 - 0.30) = -0.14.$$

$\delta_1$  es sustancialmente menor que  $b_0$  pero aún no es despreciable. Dos iteraciones más nos llevan a una situación en que el proceso puede parar: La corrección

más reciente lleva e un cambio menor del 1% en la pendiente.

La recta resultante es

$$\hat{y} = 145.86 + 0.4285(x - 127.5)$$

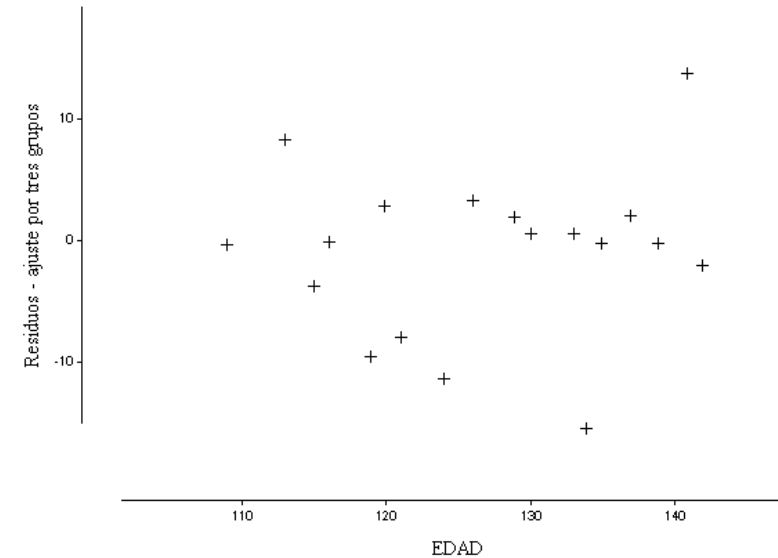
Los nuevos residuos son:

```
> altedad$altura
- 145.86-0.4285*(altedad$edad-127.5)
[1] -0.33275  8.15325 -3.70375 -0.23225
[5] -9.51775  2.75375 -8.07475 -11.36025
[9]  3.28275  1.79725  0.56875  0.58325
[13] -15.44525 -0.37375  2.06925 -0.18775
[17] 13.65525  -2.17325
```

En la figura 3 se destacan un punto bajo y otro alto que corresponden a los niños 13 y 17. Estos residuos son lo suficientemente extremos, si los analizáramos mediante un boxplot están fuera de los valores adyacentes.

También parecen bajos los residuos de tres niños con edades alrededor de 120 meses. Si tuviéramos más información podríamos intentar aprender porqué estos niños son bajos o altos de acuerdo a su edad. La distinción entre varones y mujeres podría ayudar.

Los puntos inusuales casi no tuvieron efecto en la recta que resume al conjunto de los datos. Una recta de cuadrados mínimos enfrenta más riesgo de distorsión debido a esos puntos.



**Figura 8.** Residuos de altura vs. edad después de un ajuste iterado por tres grupos