

# Estimación robusta: una introducción

Ricardo A. Maronna

Universidad de La Plata y Universidad de Buenos Aires

# 1 ¿Para qué necesitamos métodos robustos?

**Un ejemplo:** El siguiente dataset contiene 24 mediciones (ordenadas) del contenido de cobre en muestras de harina (en partes por millón).

2.20 2.20 2.40 2.40 2.50 2.70 2.80 2.90 3.03 3.03 3.10 3.37

3.40 3.40 3.40 3.50 3.60 3.70 3.70 3.70 3.70 3.77 5.28 28.95

El último valor 28.95 es llamativo. Podría deberse a un error de transcripción. Para ver su efecto, comparamos los valores de la media, el desvío standard y la mediana muestrales, con y sin él.

	con	sin
media	4.28	3.21
desvío	5.30	0.69
mediana	3.40	3.37

Vemos que el efecto en la media es notorio, en el desvío es espectacular, y en la mediana imperceptible.

## 1.1 Pensándolo desde los datos

Queremos ajustar un modelo a un conjunto de datos (“dataset”) y hacer inferencia usando métodos “clásicos” (medias y desvíos muestrales, regresión por mínimos cuadrados, etc.).

Puede haber algunas observaciones atípicas (“outliers”) que no parecen corresponder al conjunto y que pueden tener una desmedida influencia en los resultados.

El enfoque más antiguo (“diagnósticos”): detectar outliers, eliminarlos, y volver a estimar.

Usar un buen diagnóstico para outliers es claramente mucho mejor que no hacer nada, pero tiene sus inconvenientes:

- ¿Cuándo se justifica la eliminación?. La eliminación requiere una decisión subjetiva. ¿Cuando es una observación “lo bastante atípica” como para ser eliminada?

- El usuario o el autor de los datos puede pensar que “una observación es una observación” (o sea, las observaciones deberían hablar por sí mismas) y por lo tanto no estar de acuerdo con su eliminación.
- Dado que generalmente existe cierta incertidumbre sobre si una observación es realmente atípica, existe el riesgo de eliminar observaciones “buenas”, con lo que se subestima la variabilidad de los datos.
- Como los resultados dependen de las decisiones subjetivas del usuario, es difícil determinar el comportamiento estadístico del procedimiento completo.
- Last but not least: los valores atípicos no son necesariamente “malos”. ¡Pueden contener información valiosa!

Métodos robustos: automáticamente ajustan los datos "típicos"

Pueden ser descritos generalmente como métodos clásicos pesados, donde los pesos dependen del conjunto de datos, y dan menor peso a los outliers.

## **1.2 Pensándolo desde el modelo**

Enfoque clásico: creemos que el modelo vale exactamente.

Métodos "óptimos": generalmente Estimador de Máxima Verosimilitud (EMV) y Likelihood Ratio Tests (LRT)).

Pero en la vida real los modelos valen –en el mejor de los casos– sólo aproximadamente.

Si el modelo vale sólo aproximadamente, los métodos clásicos *no* son aproximadamente óptimos.

Métodos robustos: en vez de ser óptimos para un modelo, son aproximadamente óptimos en un "entorno" del modelo.

## 1.3 Un poco de historia

La media muestral, el desvío standard muestral, el estimador de Mínimos Cuadrados (EMC), los tests T, están justificados por ser óptimos para datos normales.

¿Por qué suponer normalidad?

Muchos datasets tienen distribución con forma de campana en el centro, pero con colas más pesadas que la normal.

En realidad Gauss (1820) justificó la suposición de normalidad porque ésta justificaba el uso del EMC, que era el único estimador de regresión calculable en esa época (y todavía mucho más tarde).



Puede parecer natural proceder como sigue:

- testear la hipótesis de que los datos son normales
- si no es rechazada, usar la media
- si no, usar la mediana

O, mejor aún, ajustar una distribución a los datos y usar el EMV correspondiente.

Pero esto tienen el inconveniente de que hacen falta muestras *muy grandes* para distinguir la distribución correcta, especialmente porque las *colas* –precisamente las regiones con menos datos– son las más influyentes.

## 2 El modelo de posición

Suponemos que el resultado  $x_i$  de cada observación depende del “valor verdadero”  $\mu$  del parámetro desconocido, y también de algún proceso de error. La suposición más simple es que el error actúa en forma aditiva, o sea

$$x_i = \mu + u_i \quad (i = 1, \dots, n)$$

donde los errores  $u_1, \dots, u_n$  son variables aleatorias.

Esto es el llamado *modelo de posición*.

Si las observaciones son réplicas independientes del mismo experimento bajo las mismas condiciones, se puede suponer que

- $u_1, \dots, u_n$  tienen la misma función de distribución  $F_0$
- $u_1, \dots, u_n$  son independientes.

Esto implica que  $x_1, \dots, x_n$  son independientes con la misma función de distribución:

$$F(x) = F_0(x - \mu),$$

y diremos que las  $x_i$  son variables aleatorias *i.i.d.* (independientes e idénticamente distribuidas).

La suposición de que no hay errores sistemáticos de medición se puede representar con:

- $u_i$  y  $-u_i$  tienen la misma distribución, y por lo tanto  $F_0(x) = 1 - F_0(-x)$  (“ $F_0$  es una distribución *simétrica* respecto de 0”).

Un *estimador*  $\hat{\mu}$  es una función de las observaciones:  $\hat{\mu} = \hat{\mu}(x_1, \dots, x_n) = \hat{\mu}(\mathbf{x})$

Una manera de medir la aproximación de  $\hat{\mu}$  a  $\mu$  es con el *error medio cuadrático* (*EMC*)

$$\text{EMC}(\hat{\mu}) = E(\hat{\mu} - \mu)^2$$

El EMC se puede descomponer como

$$\text{EMC}(\hat{\mu}) = \text{Var}(\hat{\mu}) + \text{Sesgo}(\hat{\mu})^2,$$

con

$$\text{Sesgo}(\hat{\mu}) = E\hat{\mu} - \mu.$$

Nótese que si  $\hat{\mu}$  es la media muestral y  $c$  es cualquier constante, entonces

$$\hat{\mu}(x_1 + c, \dots, x_n + c) = \hat{\mu}(x_1, \dots, x_n) + c$$

and

$$\hat{\mu}(cx_1, \dots, cx_n) = c\hat{\mu}(x_1, \dots, x_n)$$

Lo mismo ocurre con la mediana.

Estas dos propiedades se llaman respectivamente *equivariancia por posición* y *por escala*.

## 2.1 Alternativas a la normalidad

Una manera tradicional de representar datos “buenos” –o sea, sin outliers– es suponer que  $F_0$  es normal con media 0 y varianza desconocida  $\sigma^2$ , lo que implica

$$F = \mathcal{D}(x_i) = \mathbf{N}(\mu, \sigma^2),$$

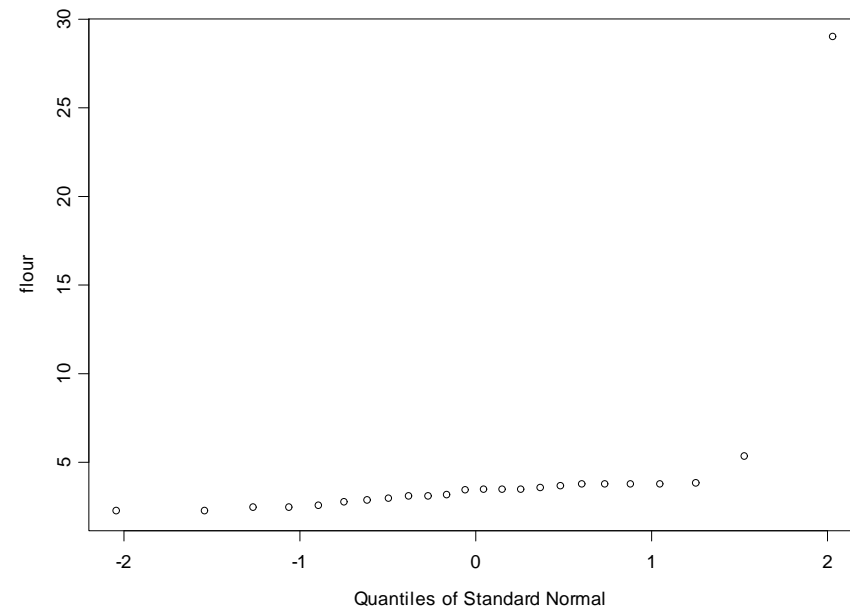
donde  $D(x)$  denota la distribución de la variable  $x$ , y  $\mathbf{N}(\mu, v)$  es la distribución normal con media  $\mu$  y varianza  $v$ .

Los métodos clásicos suponen que  $F$  pertenece a una familia de distribuciones *exactamente conocida*.

Si los datos fueran *exactamente* normales, la media muestral sería un estimador “óptimo”: es el de máxima verosimilitud, y minimiza el EMC entre los estimadores insesgados, y también entre los equivariantes.

Pero la realidad no suele comportarse siempre tan bien.

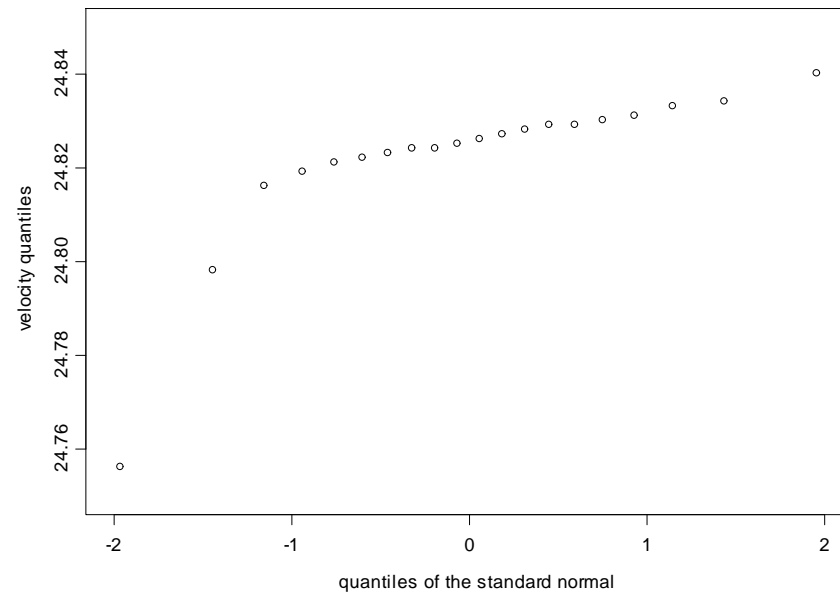
La siguiente figura muestra el Q-Q del dataset de contenido de cobre en harina mostrado al principio:



Cobre en harina: QQ-plot



Y el siguiente es de un dataset histórico de mediciones de la velocidad de la luz:



Velocidad de la luz: QQ-plot de los tiempos medidos

Vemos que la *mayoría* de los datos pueden ser descritos por una distribución normal, pero no la totalidad de ellos.

En este sentido, podemos considerar a  $F$  como *aproximadamente normal*, con la normalidad fallando en las colas.

Podemos entonces plantear nuestro objetivo inicial como: buscar estimadores

*que son “casi tan buenos” como la media muestral cuando  $F$  es exactamente normal*

*pero que son también “buenos” en algún sentido cuando  $F$  es sólo aproximadamente normal*

Para formalizar la idea de “normalidad aproximada” , podemos imaginar que una proporción  $1 - \epsilon$  de los datos es generada por el modelo normal, mientras que una proporción  $\epsilon$  es generada por un mecanismo desconocido, o sea

$$F = (1 - \epsilon)G + \epsilon H$$

donde  $G = N(\mu, \sigma^2)$  y  $H$  puede ser cualquier otra distribución; por ejemplo, otra normal con una varianza mayor y una media posiblemente distinta.

Esto se llama una *distribución normal contaminada*.

En general,  $F$  es llamada una *mezcla* de  $G$  y  $H$ , y se la llama una *mezcla normal* cuando tanto  $G$  como  $H$  son normales.

Si  $G$  y  $H$  tienen densidades  $g$  y  $h$ , respectivamente,  $F$  tiene densidad

$$f = (1 - \epsilon)g + \epsilon h.$$

Otro modelo para outliers son las llamadas *distribuciones con colas pesadas*, es decir, distribuciones cuya densidad tiende a cero más lentamente que las colas de la normal.

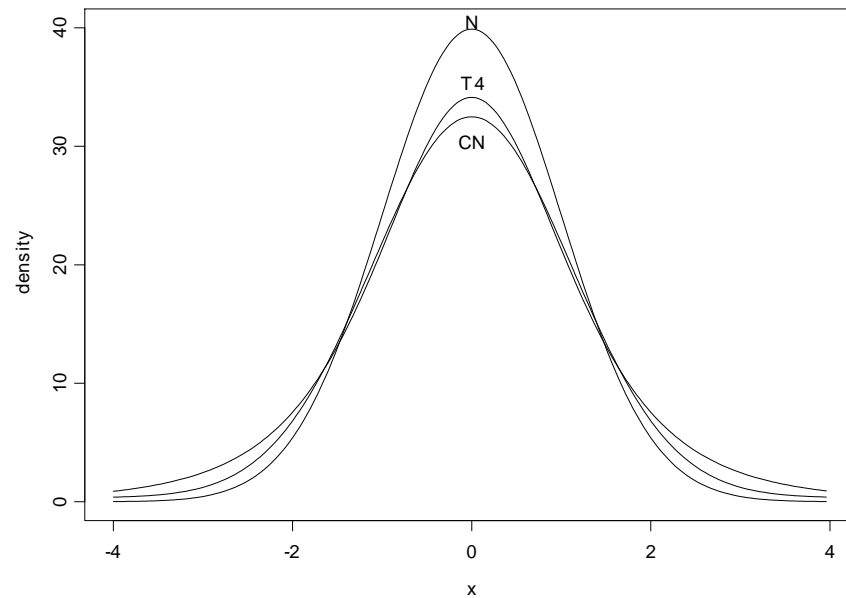
Un caso extremo es la llamada *distribución de Cauchy*, con densidad

$$f(x) = \frac{1}{\pi(1 + x^2)}.$$

La siguiente figura muestra las densidades de  $N(0, 1)$ , la de la distribución de Student con 4 grados de libertad, y la de la distribución normal contaminada

$$0.9N(0, 1) + 0.1N(0, 100)$$

indicadas por N, T4 y CN respectivamente.



Densidades de la Normal Standard (N), Student (T4), y normal contaminada (CN) escaladas para tener el mismo recorrido intercuartiles.

Para facilitar la comparación, las tres distribuciones están normalizadas para tener el mismo recorrido intercuartiles.

### 2.1.1 Eficiencia

Recordemos que si  $x \sim N(0, \sigma^2)$ , entonces la media muestral  $\bar{x} \sim N(\mu, \sigma^2/n)$  (donde “ $\sim$ ” significa “tiene distribución”).

Como veremos luego, la mediana muestral es aproximadamente  $N(\mu, 1.57\sigma^2/n)$ , de modo que para muestras normales, la mediana tiene un aumento de la varianza de 57% comparada con la media.

Decimos que la mediana tiene *baja eficiencia* respecto de la media, para datos normales.

La siguiente Tabla muestra para  $n$  “grande” el valor de  $n \times \text{varianza}$ , de la media y la mediana, para diferentes valores  $\tau$  en el modelo

$$F = (1 - \varepsilon)\mathbf{N}(\mu, 1) + \varepsilon\mathbf{N}(\mu, \tau^2)$$

$\varepsilon$	0.05		0.10	
$\tau$	$n\text{Var}(\bar{x})$	$n\text{Var}(\text{Med})$	$n\text{Var}(\bar{x})$	$n\text{Var}(\text{Med})$
3	1.40	1.68	1.80	1.80
4	1.75	1.70	2.50	1.84
5	2.20	1.70	3.40	1.86
6	2.75	1.71	4.50	1.87
10	5.95	1.72	10.90	1.90
20	20.9	1.73	40.90	1.92

Table 1: Varianzas ( $\times n$ ) ofde media y mediana para  $n$  grande

Se ve que los valores para la media aumentan rápidamente con  $\tau$ , mientras que los de la mediana se estabilizan.

En lo que sigue desarrollaremos estimadores que combinan la baja varianza de la media en la normal, con la robustez de la mediana bajo contaminación.

Po el momento nos limitaremos a distribuciones *simétricas*.



### 3 M-estimadores de posición

Ahora desarrollaremos una familia general de estimadores que contiene la media y la mediana como casos especiales.

#### 3.1 Generalizando Màxima Verosimilitud

Volvemos al modelo de posición:

$$x_i = \mu + u_i \text{ con } u_i \sim F_0.$$

Si conociéramos  $F_0$  exactamente, el EMV sería “óptimo”: tiene la menor varianza asintótica dentro de una clase “razonable” de estimadores.

Pero como la conocemos sólo aproximadamente, nuestro objetivo será encontrar estimadores que sean

(A) “casi óptimos” cuando  $F_0$  es exactamente normal,

*y también*

(B) “casi óptimos” cuando  $F_0$  es aproximadamente normal (por ejemplo, normal contaminada).

Si la densidad  $f_0 = F'_0$  es siempre positiva, dado que el logaritmo es una función creciente, el EMV puede escribirse como

$$\hat{\mu} = \arg \min_{\mu} \sum_{i=1}^n \rho(x_i - \mu)$$

donde

$$\rho = -\ln f_0.$$

Si  $F_0 = N(0, 1)$  es

$$f_0(x) = \frac{1}{\sqrt{2\pi}} \exp^{-x^2/2}$$

y (salvo una constante) es  $\rho(x) = x^2/2$ .

Por lo tanto el EMV es

$$\hat{\mu} = \arg \min_{\mu} \sum_{i=1}^n (x_i - \mu)^2,$$

lo que implica  $\hat{\mu} = \bar{x}$ .

Si  $F_0$  es la distribución “doble exponencial”

$$f_0(x) = \frac{1}{2}e^{-|x|}$$

resulta  $\rho(x) = |x|$ , y el EMV es

$$\hat{\mu} = \arg \min_{\mu} \sum_{i=1}^n |x_i - \mu|,$$

lo que implica  $\hat{\mu} = \text{mediana muestral}$ .

Si  $\rho$  es diferenciable, derivar respecto de  $\mu$  produce la “ecuación de estimación”

$$\sum_{i=1}^n \psi(x_i - \hat{\mu}) = 0$$

con

$$\psi = \rho'.$$

Si  $\psi$  es discontinua, la solución puede no existir, y en tal caso interpretamos la ecuación en el sentido de que el primer miembro cambia de signo en  $\mu$ .

Si  $\rho(x) = x^2/2$ , entonces  $\psi(x) = x$ ,  
y la ecuación de estimación resulta

$$\sum_{i=1}^n (x_i - \hat{\mu}) = 0.$$

cuya solución es  $\hat{\mu} = \bar{x}$ .

**Definición:** Dada una función  $\rho$ , un *M-estimador de posición* es una solución de

$$\hat{\mu} = \arg \min_{\mu} \sum_{i=1}^n \rho(x_i - \mu).$$

Si  $\rho$  es diferenciable, derivar respecto de  $\mu$  produce la “ecuación de estimación”

$$\sum_{i=1}^n \psi(x_i - \hat{\mu}) = 0$$

con

$$\psi = \rho'.$$

Esta familia contiene a los EMV, pero también estimadores que no son EMV para ninguna distribución.

Si  $\rho$  es convexa, entonces  $\psi$  es no decreciente.

Si  $\psi$  es estrictamente creciente, la solución es única, y si no, es todo un intervalo.

Si  $\rho$  no es convexa, hay soluciones múltiples.

Es fácil mostrar que son equivariantes por posición.

La media y la mediana son también equivariantes por escala, pero esto no se cumple en general para los M-estimadores. Más tarde se verá cómo se arregla esto.

Desde ahora se supondrá  $\rho$  par, y por lo tanto  $\psi$  impar.

## 3.2 La distribución de los M-estimadores

Para evaluar la performance de los M-estimadores, hay que calcular su distribución.

Salvo para la media y la mediana, no hay expresiones explícitas para la distribución de M-estimadores en muestras finitas, pero se pueden encontrar aproximaciones para muestras “grandes”.

Para una distribución  $F$ , sea  $\mu_0 = \mu_0(F)$

$$\mu_0 = \arg \min_{\mu} E_F \rho(x - \mu),$$

que cumple

$$E_F \psi(x - \mu_0) = 0.$$



Para la media muestral,  $\psi(x) = x$ , lo que implica  $\mu_0 = \mathbb{E}x$ , o sea, la media de  $F$ .

Para la mediana muestral es  $\mu_0 = \text{Med}(x)$ , o sea, la mediana de  $F$ .

En general, si  $F$  es simétrica,  $\mu_0$  coincide con el centro de simetría de  $F$ .

Se puede mostrar que cuando  $n \rightarrow \infty$ ,

$$\hat{\mu} \rightarrow_p \mu_0$$

donde “ $\rightarrow_p$ ” significa “tiende en probabilidad” (decimos que  $\hat{\mu}$  es “*consistente* para  $\mu_0$ ”), y que la distribución de  $\hat{\mu}$  es aproximadamente

$$N\left(\mu_0, \frac{v}{n}\right) \text{ con } v = \frac{\mathbb{E}_F \psi(x - \mu_0)^2}{(\mathbb{E}_F \psi'(x - \mu_0))^2}.$$

Si la distribución de un estimador  $\hat{\mu}$  es aproximadamente  $N(\mu_0, v/n)$  para  $n$  grande, decimos que  $\hat{\mu}$  es *asintóticamente normal*, con *valor asintótico*  $\mu_0$  y *varianza asintótica*  $v$ .

La *eficiencia asintótica* de  $\hat{\mu}$  bajo la distribución  $F$  es el cociente

$$\text{Efi}(\hat{\mu}) = \frac{v_0}{v},$$

donde  $v_0$  es la varianza asintótica del EMV y mide cuán cerca está  $\hat{\mu}$  del óptimo bajo  $F$ .

Para la media muestral es  $\psi' \equiv 1$  y por lo tanto  $v = \text{Var}(x)$ .

Para la mediana muestral, puede probarse que

$$v = \frac{1}{4f(\mu_0)^2},$$

donde  $f = F'$ ; y por lo tanto para  $F = N(0, 1)$  tenemos

$$v = \frac{2\pi}{4} = 1.571,$$

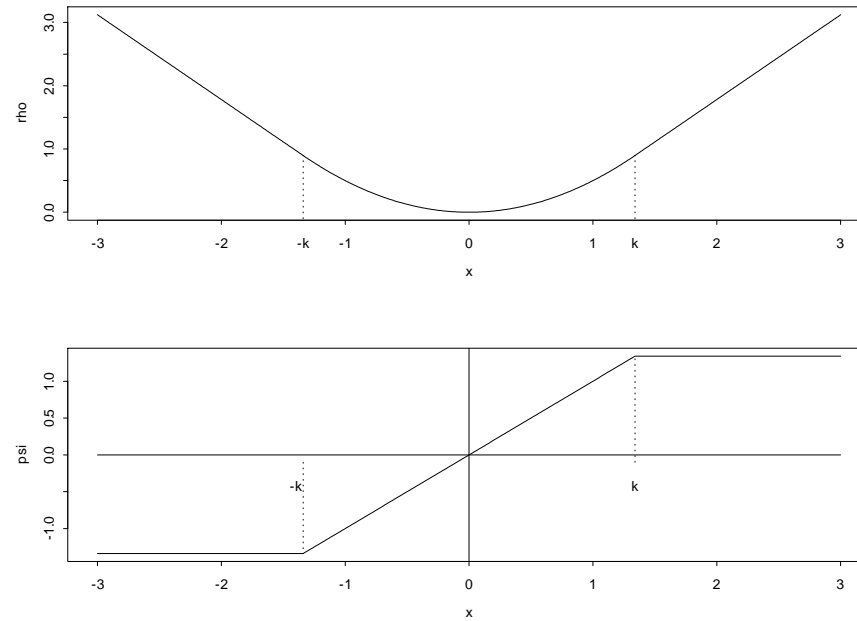
así que la eficiencia de la mediana es  $1/1.571=0.636$ .

Una familia de funciones  $\rho$  y  $\psi$  con propiedades importantes es la de *funciones de Huber* graficadas en la próxima figura:

$$\rho_k(x) = \begin{cases} x^2 & \text{si } |x| \leq k \\ 2k|x| - k^2 & \text{si } |x| > k \end{cases}$$

con derivada  $2\psi_k(x)$  donde

$$\psi_k(x) = \begin{cases} x & \text{si } |x| \leq k \\ \text{sgn}(x)k & \text{si } |x| > k \end{cases}$$



Funciones  $\rho$  y  $\psi$  de Huber

Se ve que  $\rho_k$  es cuadrática en un intervalo central, pero crece sólo linealmente fuera de él

Los M-estimadores correspondientes a los casos límites  $k \rightarrow \infty$  y  $k \rightarrow 0$  son respectivamente la media y la mediana (definimos  $\psi_0(x)$  como  $\text{sign}(x)$ ).

El valor de  $k$  se elige para obtener una varianza (y por lo tanto una eficiencia) deseada en la normal.

La siguiente tabla muestra la varianza asintótica del estimador bajo el modelo normal contaminado

$$F = (1 - \varepsilon)\text{N}(0, 1) + \varepsilon\text{N}(0, 10)$$

para diferentes valores de  $\varepsilon$  y  $k$ .

Aquí se ve el “regateo” entre robustez y eficiencia:

$k$	$\varepsilon = 0$	$\varepsilon = 0.05$	$\varepsilon = 0.10$
0	1.571	1.722	1.897
1.0	1.107	1.263	1.443
1.4	1.047	1.227	1.439
2.0	1.010	1.259	1.550
$\infty$	1.000	5.950	10.900

Table 2: Varanzas asintóticas del estimador de Huber

- Cuando  $k = 1.4$ , la varianza del M-estimador bajo la normal es sólo 4.7% mayor que la de  $\bar{x}$  (que corresponde a  $k = \infty$ ) y mucho menor que la de la mediana (que corresponde a  $k = 0$ ), mientras
- paa normales contaminadas el claramente menor que ambas.

### 3.3 Una visión intuitiva de los M-estimadores

Un M-estimador de posición se puede ver como una media pesada.

En la mayoría de los casos que interesan, es  $\psi(0) = 0$ , y  $\psi'(0)$  existe, de modo que  $\psi$  es aproximadamente lineal en el origen. Sea

$$W(x) = \begin{cases} \psi(x)/x & \text{si } x \neq 0 \\ \psi'(0) & \text{si } x = 0 \end{cases}$$

(la “función de peso”). Entonces la ecuación de estimación puede escribirse como

$$\sum_{i=1}^n W(x_i - \hat{\mu})(x_i - \hat{\mu}) = 0,$$

lo que implica

$$\hat{\mu} = \frac{\sum_{i=1}^n w_i x_i}{\sum_{i=1}^n w_i}, \text{ con } w_i = W(x_i - \hat{\mu}),$$

lo que expresa al estimador como una media pesada.

Notar que aunque esto parece una expresión explícita para  $\hat{\mu}$ , en realidad los pesos  $w_i$  dependen de  $\hat{\mu}$ .

Como en general  $W$  es decreciente, las observaciones recibirán pesos menores cuanto más alejadas estén del “centro” de los datos.

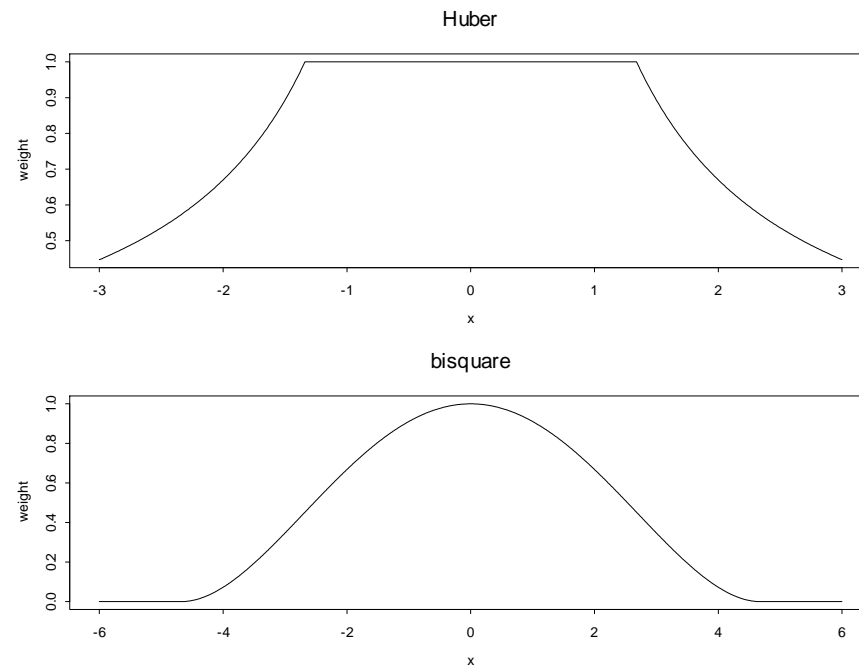
Además de su valor intuitivo, esta representación será útil para el cálculo numérico.

La función de peso correspondiente a la  $\psi$  de Huber es

$$W_k(x) = \min \left\{ 1, \frac{k}{|x|} \right\}$$



que se muestra en el panel superior de la próxima figura, junto con otra que trataremos luego.



Funciones de peso  $W$  Huber y bisquare

### 3.4 M-estimadores redescendientes

Es fácil mostrar que el EMV para la densidad de Student con  $\nu$  grados de libertad es

$$\psi(x) = \frac{x}{x^2 + \nu},$$

que tiende a 0 cuando  $x \rightarrow \infty$ .

Esto sugiere que para distribuciones simétricas con colas pesadas, es mejor usar una  $\psi$  "redescendiente" que tienda a 0 en infinito.

Esto implica que para  $x$  grande, la  $\rho$  correspondiente crece más lentamente que la de Huber, que es lineal para  $|x| > k$ .

Una elección popular es la familia *bisquare* (o *biweight*):

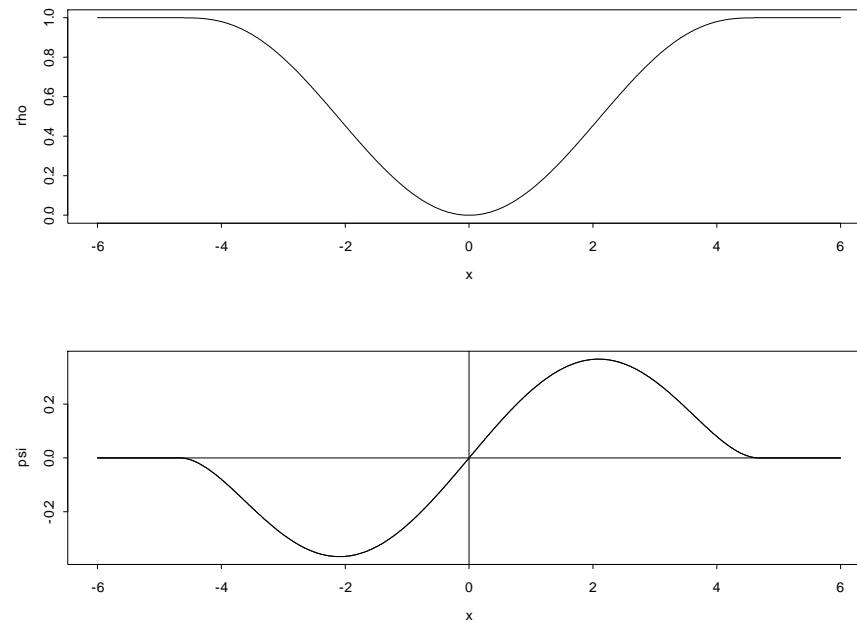
$$\rho(x) = \begin{cases} 1 - \left(1 - (x/k)^2\right)^3 & \text{if } |x| \leq k \\ 1 & \text{if } |x| > k \end{cases}$$

con derivada  $\rho'(x) = 6\psi(x)/k^2$  donde

$$\psi(x) = x \left(1 - \left(\frac{x}{k}\right)^2\right)^2 \mathbf{I}(|x| \leq k).$$

donde “I” es el indicador.

Estas funciones se muestran en la próxima figura.



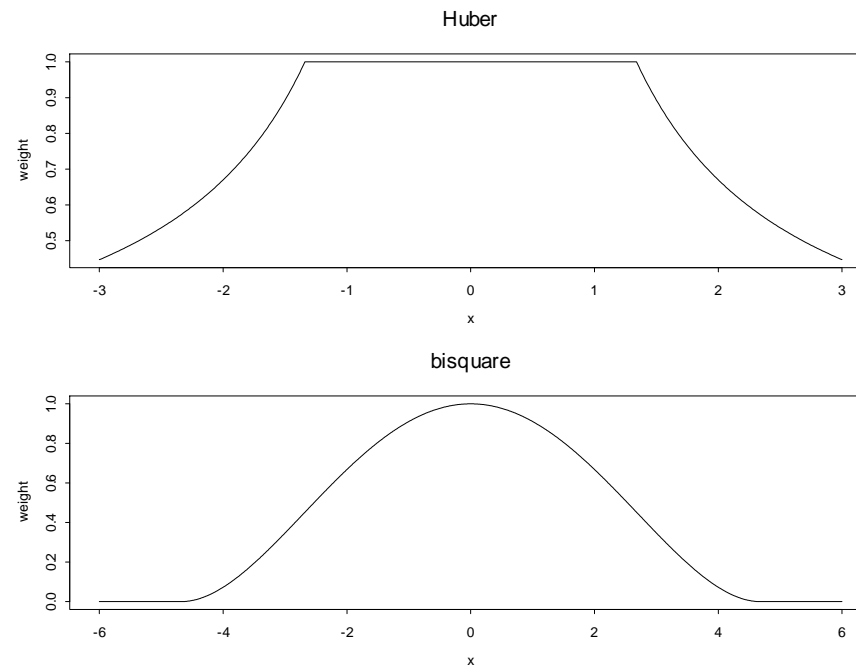
Funciones  $\rho$  y  $\psi$  para el estimador bisquare

Notar que  $\psi$  es diferenciable y se anula fuera de  $[-k, k]$ .

La correspondiente función de peso es

$$W(x) = \left(1 - \left(\frac{x}{k}\right)^2\right)^2 \mathbf{1}(|x| \leq k)$$

que volvemos a mostrar.



Funciones de peso  $W$  Huber y bisquare

Estos estimadores mejoran la resistencia a outliers grandes.

## 4 M-estimadores de posición con escala desconocida

Los estimadores definidos por

$$\hat{\mu} = \arg \min_{\mu} \sum_{i=1}^n \rho(x_i - \mu)$$

no son equivariantes por escala. Es natural exigir que si multiplico los datos por 100, el estimador se multiplique por 100.

Para arreglar este problema necesitamos un estimador robusto de dispersión.



Para fijar ideas, supongamos que queremos estimar  $\mu$  en el modelo de posición con

$$F = (1 - \varepsilon)\mathbf{N}(\mu, \sigma^2) + \varepsilon H.$$

Si  $\sigma$  fuera conocida, sería natural estimar  $\mu$  con

$$\hat{\mu} = \arg \min_{\mu} \sum_{i=1}^n \rho \left( \frac{x_i - \mu}{\sigma} \right).$$

Una idea intuitiva es usar

$$\hat{\mu} = \arg \min_{\mu} \sum_{i=1}^n \rho \left( \frac{x_i - \mu}{\hat{\sigma}} \right),$$

donde  $\hat{\sigma}$  es un estimador de dispersión calculado previamente.

Entonces  $\hat{\mu}$  resulta equivariante por escala.

Como  $\hat{\sigma}$  no depende de  $\mu$ , resulta que  $\hat{\mu}$  es solución de

$$\sum_{i=1}^n \psi \left( \frac{x_i - \hat{\mu}}{\hat{\sigma}} \right) = 0.$$

No podemos usar como  $\hat{\sigma}$  el desvío standard (DS), porque estaría afectado por los outliers. Se necesita un  $\hat{\sigma}$  *robusto*.

Al mismo tiempo, para datos normales  $\hat{\sigma}$  debería coincidir con el DS para mantener la eficiencia elegida.

Se define la MAD (median absolute deviation) como

$$\text{MAD}(x) = \text{Med}(|x - \text{Med}(x)|).$$

Se prueba que si  $x \sim N(\mu, \sigma^2)$  es  $\text{MAD}(x) = 0.675\sigma$ .

En consecuencia se define la MAD normalizada como

$$\text{MADN}(x) = \frac{\text{MAD}(x)}{0.675},$$

y entonces para  $x$  normal es  $\text{MADN}(x) = \sigma$ .

Este es un estimador muy robusto de dispersión. Para los datos de harina, la MAD con y sin el outlier es 0.35 y 0.34.

Elegimos entonces tomar  $\hat{\sigma}$  como la MADN de la muestra.

## 5 Cálculo numérico de los M-estimadores

Describiremos un procedimiento llamado *reponderación iterativa* basado en la presentación intuitiva que ya vimos.

La expresión de  $\hat{\mu}$  como media pesada

$$\hat{\mu} = \frac{\sum_{i=1}^n w_i x_i}{\sum_{i=1}^n w_i}, \text{ con } w_i = W(x_i - \hat{\mu})$$

sugiere un procedimiento iterativo.

Sean  $\hat{\sigma}_0$  un estimador robusto de dispersión (por ejemplo MADN) y  $\hat{\mu}_0$  un estimador inicial (por ejemplo la mediana). Ponemos  $m = 0$ .

Dado  $\hat{\mu}_m$  calcular

$$w_{m,i} = W \left( \frac{x_i - \hat{\mu}_m}{\hat{\sigma}} \right) \quad (i = 1, \dots, n)$$

y

$$\hat{\mu}_{m+1} = \frac{\sum_{i=1}^n w_{m,i} x_i}{\sum_{i=1}^n w_{m,i}}.$$

Si  $W(x)$  es acotada y es decreciente para  $x > 0$ , la sucesión  $\hat{\mu}_m$  converge a una solución

Si  $\psi$  es estrictamente creciente la solución es única, y el punto de partida  $\hat{\mu}_0$  sólo influye en la cantidad de iteraciones.

Si  $\psi$  es redescendiente,  $\hat{\mu}_0$  tiene que ser robusto para asegurar la convergencia a una solución “buena”.

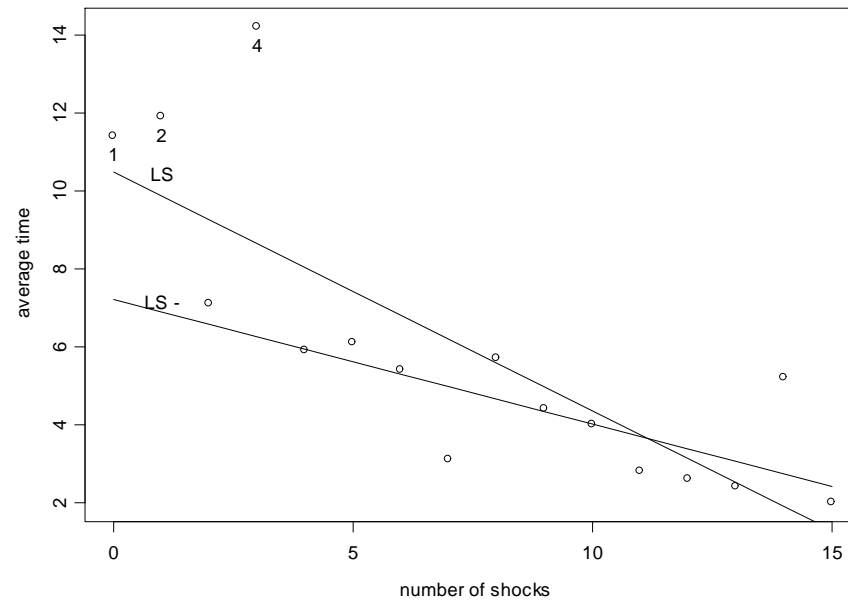
Para los datos de harina, los valores del M-estimador bisquare con eficiencia 0.95, con y sin el outlier, son 3.14 y 3.12, o sea que no es casi afectado por éste. La media sin el outlier era 3.40.

## 6 Regresión lineal (trailer)

El estimador clásico para regresión lineal es el de *mínimos cuadrados* (EMC).

Pese a sus virtudes, es muy sensible a outliers.

**Ejemplo:** Los datos corresponden a un experimento sobre la velocidad de aprendizaje de ratas. Para abreviar se omite la descripción detallada. El grafico muestra en la abscisa la cantidad de shocks, y en la ordenada el tiempo promedio de los intentos, para cada rata.



Ajuste del EMC con todos los datos, y omitiendo los puntos 1-2-4.

The Figure muestra los datos y las líneas ajustadas por MC para el m odelo

$$y_i = \beta_0 + \beta_1 x_i + u_i.$$



La relación entre las variables se ve como aproximadamente lineal, excepto para los 3 puntos en el ángulo superior izquierdo.

La línea del EMC no ajusta la mayoría de los datos, siendo un compromiso entre esos 3 puntos y el resto.

El ajuste sin los 3 puntos da una mejor representación de la mayoría de los datos, y muestra el carácter excepcional de 1-2-4.

Buscamos procedimientos que den un buen ajuste a la mayoría de los datos, sin ser perturbados por una pequeña proporción de outliers.

Consideremos en general un dataset de  $n$  observaciones  $(x_{i1}, \dots, x_{ip}, y_i)$ ,  $i = 1, \dots, n$ , donde

- $x_{i1}, \dots, x_{ip}$  son *predictores*, y
- $y_i$  es la *respuesta*.

Se supone que los datos cumplen el *modelo lineal*

$$y_i = \sum_{j=1}^p x_{ij}\beta_j + u_i, \quad i = 1, \dots, n$$

donde  $\beta_1, \dots, \beta_p$  son parámetros desconocidos que hay que estimar, y las  $u_i$  son variables aleatorias (“errores”).

Poniendo

$$\mathbf{x}_i = (x_{i1}, \dots, x_{ip})', \quad \beta = (\beta_1, \dots, \beta_p)',$$

el modelo puede ser escrito en forma más compacta como

$$y_i = \mathbf{x}_i' \beta + u_i$$

donde  $\mathbf{x}'$  is el traspuesto de  $\mathbf{x}$ .

Los *valores ajustados*  $\hat{y}_i$  y los *residuos*  $r_i$  correspondientes a un vector  $\beta$  se definen como

$$\hat{y}_i = \hat{y}_i(\beta) = \mathbf{x}_i' \beta \text{ y } r_i = r_i(\beta) = y_i - \hat{y}_i(\beta).$$

El EMC cumple las *ecuaciones normales*

$$\sum_{i=1}^n r_i(\hat{\beta}) \mathbf{x}_i = \mathbf{0}.$$

Trataremos los M-estimadores de regresión  $\hat{\beta}$  definidos como solución de

$$\sum_{i=1}^n \rho \left( \frac{r_i(\hat{\beta})}{\hat{\sigma}} \right) = \min .$$

Aquí  $\rho$  es una función como las que vimos en el modelo de posición, y  $\hat{\sigma}$  es un estimador auxiliar de escala, calculado previamente, necesario para hacer  $\hat{\beta}$  equivariante por escala.

Esta clase contiene a los EMV.

El EMC corresponde a  $\rho(t) = t^2$ .

Para  $\rho(t) = |t|$  tenemos el *estimador*  $L_1$  (o “LAD”: Least Absolute deviations”), que es el equivalente de la mediana para regresión.

Derivando respecto de  $\beta$  resulta el análogo de las ecuaciones normales:

$$\sum_{i=1}^n \psi \left( \frac{r_i(\hat{\beta})}{\sigma} \right) \mathbf{x}_i = \mathbf{0},$$

Recordemos que para M-estimadores de posición, estimábamos  $\sigma$  mediante la MAD.

Aquí el procedimiento equivalente es calcular el estimador  $L_1$ ,  $\hat{\beta}_{L_1}$  y los correspondientes residuos, y de allí obtener el análogo de la MAD normalizada tomando la mediana de los residuos no nulos:

$$\hat{\sigma} = \frac{1}{0.675} \text{Med}_i(|r_i| \left( \hat{\beta}_{L_1} \right) \mid r_i \neq 0).$$

El motivo para usar sólo los residuos no nulos es que el  $L_1$  produce al menos  $p$  residuos nulos, y por lo tanto usar todos los residuos puede producir una  $\sigma$  subestimada.

Nótese que  $L_1$  no necesita una escala auxiliar.

## 7 Interpretación intuitiva y cálculo numérico

Supongamos  $\psi$  diferenciable y definamos

$$W(x) = \frac{\psi(x)}{x}.$$

Entonces las ecuaciones de estimación pueden escribirse como

$$\sum_{i=1}^n w_i r_i \mathbf{x}_i = \sum_{i=1}^n w_i \mathbf{x}_i (y_i - \mathbf{x}_i' \hat{\beta}) = \mathbf{0}$$

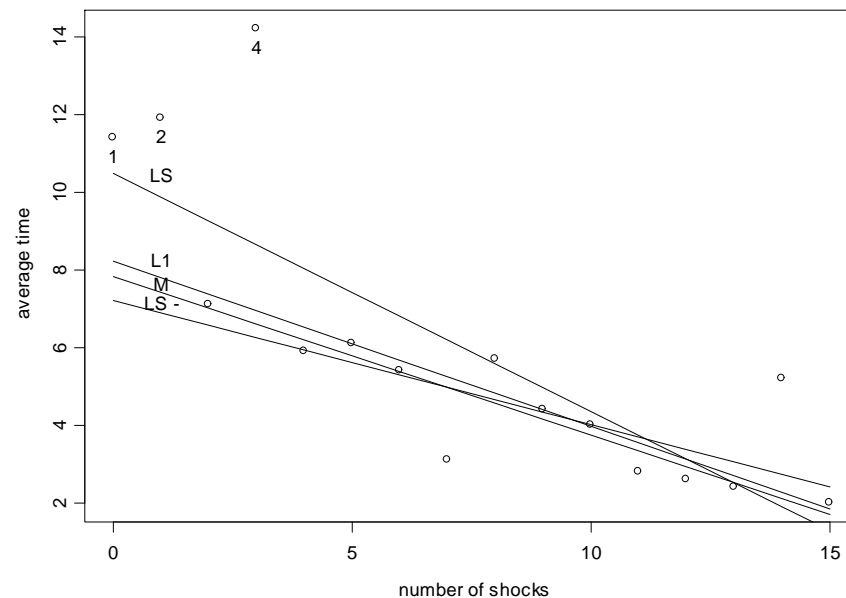
con  $w_i = W(r_i/\hat{\sigma})$ .

Esto muestra al estimador como un EMC ponderado.

Al mismo tiempo sugiere un algoritmo iterativo para su cómputo, análogo al que se mostró para posición.

Volvemos al ejemplo inicial. Se muestran los ajustes por EMC, EMC sin los outliers,  $L_1$ , y el M-estimador con la función bisquare que se definió para posición.

Se ve que el bisquare es el que más se parece al EMC sin outliers.





Ajustes por EMC (LS),  $L_1$ , M-estimador Bisquare (M) y EMC sin los outliers (LS-)

## Referencias

Maronna, R.A., Martin, R.D. Yohai, V.J. y Salibián-Barrera, M. (20019). *Robust Statistics: Theory and Methods (with R). Second Edition.* . John Wiley and Sons, New York.