

ESTIMACIÓN NO PARAMETRICA DE LA DENSIDAD

DANIELA RODRIGUEZ

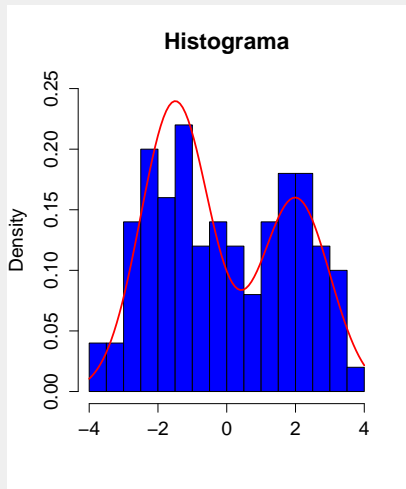
ESTIMADOR NO PARAMÉTRICO DE LA DENSIDAD

- X con densidad $f(x)$: queremos estimar $f(x)$
- X_1, \dots, X_n , i.i.d., $X_i \sim X$ donde $X_i \sim X$.
- Queremos estimar f sin asumir una determinada forma: sólo asumimos que f es suave.

ESTIMADOR NO PARAMÉTRICO DE LA DENSIDAD

- X con densidad $f(x)$: queremos estimar $f(x)$
- X_1, \dots, X_n , i.i.d., $X_i \sim X$ donde $X_i \sim X$.
- Queremos estimar f sin asumir una determinada forma: sólo asumimos que f es suave.
- La forma más sencilla: Histograma

EJEMPLO DATOS SIMULADOS



ESTIMACIÓN DE LA DENSIDAD: HISTOGRAMA

Sea X_1, \dots, X_n una muestra aleatoria con densidad desconocida f . El primer estimador de la densidad no paramétrico que se estudia es el histograma.

Construcción:

- Dividir el rango en intervalos o bins con origen x_0 y ancho h

$$B_j = [x_0 + (j-1)h, x_0 + jh)$$

- Contar la cantidad de observaciones n_j en B_j
- Normalizar

$$f_j = \frac{n_j}{nh}$$

- Graficar barras de altura f_j y ancho h sobre cada B_j

Formalmente el histograma puede escribirse como

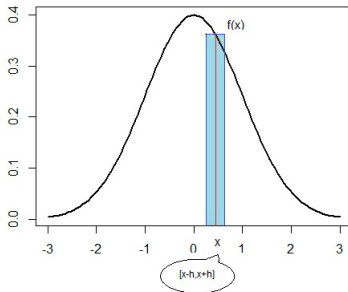
$$\hat{f}(x) = \frac{1}{nh} \sum_{i=1}^n \sum_j I(\mathbf{x}_i \in B_j) I(\mathbf{x} \in B_j)$$

DESVENTAJAS DEL HISTOGRAMA

- el estimador de la densidad depende del punto inicial de los bins: para un número de bins fijo, la forma puede cambiar moviendo la ubicación de los bins
- la densidad estimada no es suave, es *escalonada* y esto no es propio de la densidad sino de la herramienta de estimación
- por estas razones, el histograma es usado sólo para visualización
- aproxima a la densidad en el punto medio de cada bin

ESTIMADOR DE DENSIDAD POR NUCLEOS

- $\mathbb{P}(X \in (x-h, x+h)) = \int_{x-h}^{x+h} f(t) dt$
- Si h es pequeño y f continua en x ,



$$\int_{x-h}^{x+h} f(t) dt \approx 2hf(x)$$

JUNTANDO TODO...

- Si h es pequeño y f continua en x ,

$$\mathbb{P}(X \in (x - h, x + h)) \approx 2hf(x)$$

- Por la LGN X_1, \dots, X_n , i.i.d., $X_i \sim X$ donde $X_i \sim X$

$$\mathbb{P}(X \in (x - h, x + h)) \approx \frac{\#\{X_i \in (x - h, x + h)\}}{n}$$

- Entonces,

$$2hf(x) \approx \mathbb{P}(X \in (x - h, x + h)) \approx \frac{\#\{X_i \in (x - h, x + h)\}}{n}$$

$$f(x) \approx \frac{\#\{X_i \in (x - h, x + h)\}}{2hn}$$

$$\hat{f}(x) = \frac{\#\{X_i \in (x-h, x+h)\}}{2h n}$$

Notemos que

$$\hat{f}(x) = \frac{1}{2h n} \sum_{i=1}^n \mathcal{I}_{(x-h, x+h)}(X_i)$$

■ Estimador de Parzen

$$\hat{f}(x) = \frac{1}{nh} \sum_{i=1}^n \frac{1}{2} \mathcal{I}_{[-1,1]} \left(\frac{x - X_i}{h} \right)$$

PROPUESTA

$$\hat{f}(x) = \frac{\#\{X_i \in (x-h, x+h)\}}{2h n}$$

Notemos que

$$\hat{f}(x) = \frac{1}{2h n} \sum_{i=1}^n \mathcal{I}_{(x-h, x+h)}(X_i)$$

■ Estimador de Parzen

$$\hat{f}(x) = \frac{1}{nh} \sum_{i=1}^n \frac{1}{2} \mathcal{I}_{[-1,1]} \left(\frac{x - X_i}{h} \right)$$

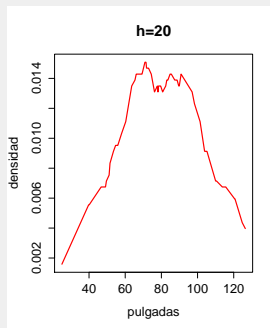
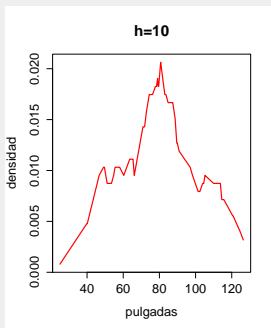
■ si $K(t) = \frac{1}{2} \mathcal{I}_{[-1,1]}(t) \Rightarrow$

$$\hat{f}(x) = \frac{1}{nh} \sum_{i=1}^n K \left(\frac{x - X_i}{h} \right)$$

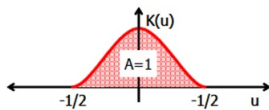
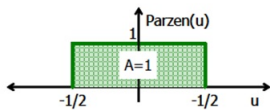
PROPUESTA

$$\blacksquare K(t) = \frac{1}{2}\mathcal{I}_{[-1,1]}(t) \Rightarrow \hat{f}(x) = \frac{1}{nh} \sum_{i=1}^n K\left(\frac{x-X_i}{h}\right)$$

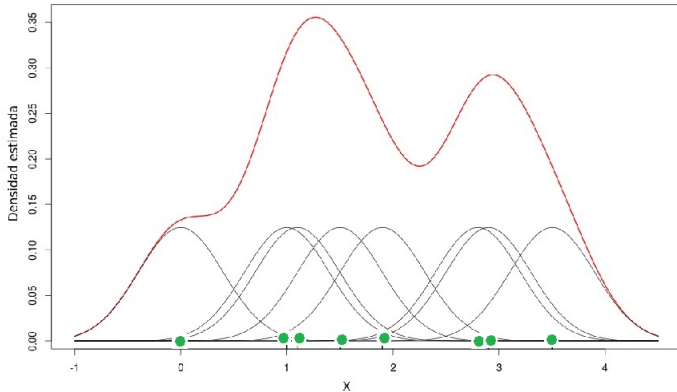
● K : núcleo ● h : ventana



NÚCLEOS

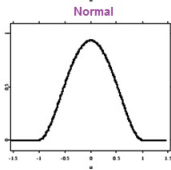
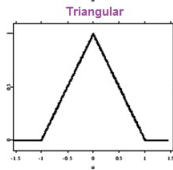
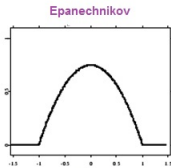
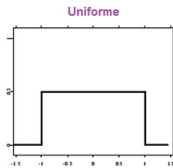


INTERPRETACIÓN DEL ESTIMADOR DE NÚCLEOS



TIPOS DE NÚCLEOS

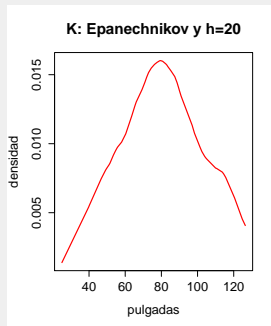
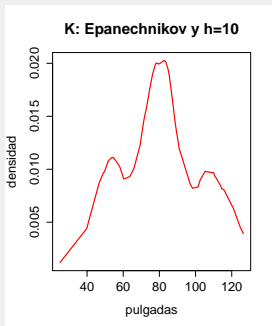
- Núcleo Rectangular: $K(t) = \frac{1}{2}\mathcal{I}_{[-1,1]}(t)$
- Núcleo Triangular: $K(t) = (1 - |t|)\mathcal{I}_{[-1,1]}(t)$
- Núcleo Gaussiano: $K(t) = \frac{1}{\sqrt{2\pi}}e^{-\frac{1}{2}t^2}$
- Núcleo Epanechnikov: $K(t) = \frac{3}{4}(1 - t^2)\mathcal{I}_{[-1,1]}(t)$



ESTIMADORES DE NÚCLEOS (ROSENBLATT-PARZEN)

$$\hat{f}(x) = \frac{1}{nh} \sum_{i=1}^n K\left(\frac{x-X_i}{h}\right)$$

- K núcleo: • $K \geq 0$ y • $\int K(x)dx = 1$.
- h : ventana o parámetro de suavizado
- Notemos que $\hat{f}(x)$ depende de n , del núcleo K y de h



Notemos que si $\int_{-\infty}^{+\infty} K(x)dx = 1$ y $K \geq 0$, entonces el estimador \hat{f} es también una función de densidad.

Pues,

$$\begin{aligned}\int_{-\infty}^{+\infty} \hat{f}(x)dx &= \int_{-\infty}^{+\infty} \sum_{i=1}^n \frac{1}{nh} K\left(\frac{x - \mathbf{x}_i}{h}\right) dx = \\ \frac{1}{n} \sum_{i=1}^n \int_{-\infty}^{+\infty} \frac{1}{h} K\left(\frac{x - \mathbf{x}_i}{h}\right) dx &= \frac{1}{n} \sum_{i=1}^n \int_{-\infty}^{+\infty} K(s) ds = 1\end{aligned}$$

PROPIEDADES: SESGO

i) f es 2-veces derivable tal que $\int f''(s)ds < \infty$.

ii) $\int K = 1$, $\int K(s)sds = 0$ y $\int K(s)s^2ds < \infty$.

Tenemos que

$$E \left[\widehat{f}(x) \right] = f(x) + \frac{h^2}{2} f''(x) \mu_2(K) + o(h^2)$$

si $h \rightarrow 0$ para cada x . Donde $\mu_2(K) = \int s^2 K(s)ds$

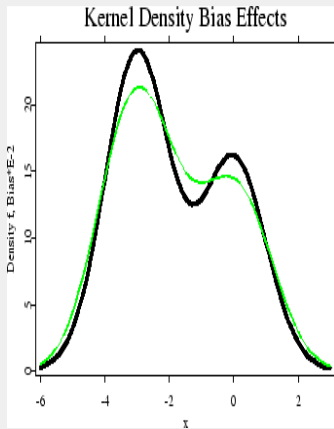
PROPIEDADES: SESGO

$$\begin{aligned}E[\widehat{f}(x)] &= E\left(\frac{1}{hn} \sum_{i=1}^n K\left(\frac{x - \mathbf{x}_i}{h}\right)\right) = \frac{1}{n} \sum_{i=1}^n E\left(\frac{1}{h} K\left(\frac{x - \mathbf{x}_i}{h}\right)\right) \\&= E\left(\frac{1}{h} K\left(\frac{x - \mathbf{x}_i}{h}\right)\right) = \frac{1}{h} \int_{-\infty}^{+\infty} K\left(\frac{x - u}{h}\right) f(u) du \\&= \int_{-\infty}^{+\infty} K(y) f(x - hy) dy \\&= \int_{-\infty}^{+\infty} K(y) \left[f(x) + f'(x)hy + \frac{f''(x)}{2}y^2h^2 + o(h^2) \right] dy \\&= f(x) \int_{-\infty}^{+\infty} K(y) dy + f'(x)h \int_{-\infty}^{+\infty} K(y)y dy + h^2 \frac{f''(x)}{2} \int_{-\infty}^{+\infty} K(y)y^2 dy \\&\quad + o(h^2)\end{aligned}$$

$$\text{Sesgo}(\hat{f}(x)) = h^2 \frac{f''(x)}{2} \int_{-\infty}^{+\infty} K(y)y^2 dy + o(h^2)$$

PROPIEDADES: SESGO

Estimador de densidad (en verde) y verdadera densidad (en negro).



Bajo las mismas hipótesis introducidas anteriormente, probaremos que la varianza del estimador es

$$\text{Var}(\hat{f}(x)) = \frac{1}{nh} \|K\|^2 f(x) + o\left(\frac{1}{nh}\right).$$

PROPIEDADES: VARIANZA

Sea $K_h(x) = \frac{1}{h}K(x/h)$, como las \mathbf{x}_i son i.i.d

$$\begin{aligned} \text{Var}(\hat{f}(x)) &= n^{-2} \text{Var}\left(\sum_{i=1}^n K_h(x - \mathbf{x}_i)\right) = n^{-2} \sum_{i=1}^n \text{Var}(K_h(x - \mathbf{x}_i)) \\ &= n^{-1} \text{Var}(K_h(x - \mathbf{x}_1)) \\ &= n^{-1} [E(K_h^2(x - \mathbf{x}_1)) - E^2(K_h(x - \mathbf{x}_1))] \\ &= n^{-1} [E(K_h^2(x - \mathbf{x}_1)) - (f(x) + o(h))^2] \end{aligned}$$

Usando los mismos argumentos que antes, es decir, cambio de variable y un desarrollo de Taylor tenemos que

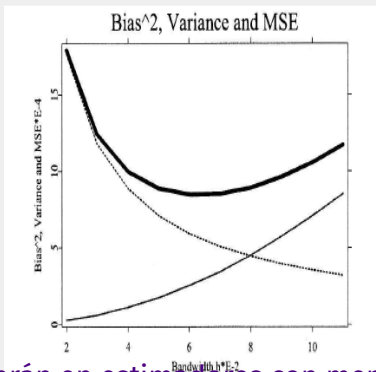
$$E(K_h^2(x - \mathbf{x}_1)) = h^{-1} \int K^2(s)f(x - hs)ds = h^{-1} \|K\|_2^2 f(x) + o(h).$$

De esta forma hemos calculado el error cuadrático medio del estimador (*ECM*) para cada x ,

$$ECM(\hat{f}(x)) = h^4 \frac{(f''(x))^2}{4} \mu_2^2(K) + o(h^4) + \frac{1}{nh} \|K\|^2 f(x) + o\left(\frac{1}{nh}\right)$$

PROPIEDADES: ECM

Compromiso entre sesgo y varianza. En la figura: Sesgo al cuadrado (línea sólida); varianza (línea punteada) y error cuadrático medio (línea sólida gruesa).



h pequeños derivarán en estimadores con menor sesgo mientras que al aumentar el ancho de banda lograremos disminuir la varianza.

SESGO Y VARIANZA DE $\hat{f}(x)$

- El sesgo es proporcional a $h^2 \Rightarrow$ elijamos h pequeña
- El sesgo depende de $f''(x)$ que mide la curvatura de f en x
- La varianza disminuye a medida que nh crece
- Para disminuir la varianza necesitamos h o n grandes.

La elección de la ventana es crucial.

- Una ventana h pequeña dará un estimador muy rugoso, con muchos picos y difícil de interpretar
- una ventana h grande sobre-suaviza al estimador de la densidad y enmascara estructuras de los datos.

PROPIEDADES: CONSISTENCIA

$$ECM(\hat{f}(x)) = h^4 \frac{(f''(x))^2}{4} \mu_2^2(K) + o(h^4) + \frac{1}{nh} \|K\|^2 f(x) + o\left(\frac{1}{nh}\right)$$

Si $h \rightarrow 0$ y $nh \rightarrow \infty$ tenemos que

$$\hat{f}(x) \xrightarrow{p} f(x)$$

para cada x .

Bajo ciertas condiciones de regularidad

- $h_n \rightarrow 0$
- $nh_n \rightarrow \infty$
- x tiene densidad f continua en x y dos veces diferenciable
- $K : \mathbb{R} \rightarrow \mathbb{R}$ es acotado, $\int K = 1$ y $\int u^2 K(u) > 0$ y con soporte compacto.

si $h_n = cn^{-1/5}$ entonces

$$\sqrt{nh}(\hat{f}(x) - f(x)) \xrightarrow{\mathcal{D}} N\left(\frac{c^{5/2}}{2}f''(x)\mu_2(K), f(x)\|K\|^2\right).$$

INTERVALOS Y BANDAS DE CONFIANZA

Luego resulta el siguiente intervalo de confianza de nivel aproximado $1 - \alpha$

$$\left[\hat{f}(x) - \frac{h^2}{2} f''(x) \mu_2(K) - z_{\alpha/2} \sqrt{\frac{f(x) \|K\|^2}{nh}}, \hat{f}(x) - \frac{h^2}{2} f''(x) \mu_2(K) + z_{\alpha/2} \sqrt{\frac{f(x) \|K\|^2}{nh}} \right]$$

si h es pequeña se puede despreciar el término que involucra a la derivada segunda y utilizar el siguiente intervalo

$$\left[\hat{f}(x) - z_{\alpha/2} \sqrt{\frac{\hat{f}(x) \|K\|^2}{nh}}, \hat{f}(x) + z_{\alpha/2} \sqrt{\frac{\hat{f}(x) \|K\|^2}{nh}} \right]$$

de lo contrario podemos estimar la derivada segunda, derivando un es de núcleos usando una ventana g .

INTERVALOS Y BANDAS DE CONFIANZA

Es importante notar que este intervalo es sólo para $f(x)$ y no para toda la densidad.

Para deducir bandas de confianza para toda la función es necesario emplear otras técnicas.

INTERVALOS Y BANDAS DE CONFIANZA

Bickel y Rosenblatt (1973) probaron el siguiente resultado: sea f una función de densidad definida sobre el $(0, 1)$, $h_n = n^{-\delta} \in (1/5, 1/2)$, entonces para todo $x \in (0, 1)$,

$$\lim_{n \rightarrow \infty} P \left(\hat{f}(x) - \sqrt{\frac{\frac{\hat{f}(x)\|K\|^2}{nh}}{\frac{z}{2\delta \log n} + d_n}} \leq f(x) \leq \hat{f}(x) + \sqrt{\frac{\frac{\hat{f}(x)\|K\|^2}{nh}}{\frac{z}{2\delta \log n} + d_n}} \right) = e^{-2e^{-z}}$$

donde $d_n = (2\delta \log n)^{1/2} + (2\delta \log n)^{-1/2} \log \left\{ \frac{\|K'\|_2}{2\pi\|K\|_2} \right\}$.

Luego si $z \approx 3.663$, se tiene que $e^{-2e^{-z}} = 1 - 0.05$.