

Profundidad de Datos Funcionales

Nicolas Maldonado Baracaldo

Andrés Felipe Patiño López

Noviembre 2021

Departamento de Matemáticas

Facultad de Ciencias

Universidad de los Andes

Motivación

Para una muestra aleatoria real X_1, X_2, \dots, X_n podemos definir:

Para una muestra aleatoria real X_1, X_2, \dots, X_n podemos definir:

- El k -ésimo estadístico de orden $X_{(k)}$.

Para una muestra aleatoria real X_1, X_2, \dots, X_n podemos definir:

- El k -ésimo estadístico de orden $X_{(k)}$.
- El vector de rangos $\mathcal{R} = (R_1, \dots, R_n)$.

Para una muestra aleatoria real X_1, X_2, \dots, X_n podemos definir:

- El k -ésimo estadístico de orden $X_{(k)}$.
- El vector de rangos $\mathcal{R} = (R_1, \dots, R_n)$.
- L-estadísticos. En particular, la media truncada.

A lo largo de este curso, ya hemos visto la utilidad del vector de rangos y los estadísticos de orden.

Profundidad: Una medida de profundidad en un conjunto de datos proporciona una noción de centralidad de los mismos. Es una función $D : \mathbb{R}^k \rightarrow [0, 1]$ que depende de la distribución de los datos (o de la empírica).

Profundidad: Una medida de profundidad en un conjunto de datos proporciona una noción de centralidad de los mismos. Es una función $D : \mathbb{R}^k \rightarrow [0, 1]$ que depende de la distribución de los datos (o de la empírica).

- Para finitos datos sin empates la profundidad induce estadísticos de orden.

Profundidad: Una medida de profundidad en un conjunto de datos proporciona una noción de centralidad de los mismos. Es una función $D : \mathbb{R}^k \rightarrow [0, 1]$ que depende de la distribución de los datos (o de la empírica).

- Para finitos datos sin empates la profundidad induce estadísticos de orden.
- La mediana será análoga al dato “más profundo” o “más central”. Es decir, el estadístico de orden $X_{(1)}$.

Profundidad: Una medida de profundidad en un conjunto de datos proporciona una noción de centralidad de los mismos. Es una función $D : \mathbb{R}^k \rightarrow [0, 1]$ que depende de la distribución de los datos (o de la empírica).

- Para finitos datos sin empates la profundidad induce estadísticos de orden.
- La mediana será análoga al dato “más profundo” o “más central”. Es decir, el estadístico de orden $X_{(1)}$.
- Los datos atípicos o los datos extremos serán los “menos profundos” o “menos centrales”. Por ejemplo, el estadístico de orden $X_{(n)}$.

Profundidad: Una medida de profundidad en un conjunto de datos proporciona una noción de centralidad de los mismos. Es una función $D : \mathbb{R}^k \rightarrow [0, 1]$ que depende de la distribución de los datos (o de la empírica).

- Para finitos datos sin empates la profundidad induce estadísticos de orden.
- La mediana será análoga al dato “más profundo” o “más central”. Es decir, el estadístico de orden $X_{(1)}$.
- Los datos atípicos o los datos extremos serán los “menos profundos” o “menos centrales”. Por ejemplo, el estadístico de orden $X_{(n)}$.
- Esto implica que las definiciones de profundidad son ligeramente diferentes en dimensiones > 1 .

Debido a la naturaleza y la libertad de movimiento en dimensiones más altas, existen una gran variedad de formas de extender las profundidades de datos. Algunas profundidades (Datos multivariados):

- **Tukey's Depth**
- **Simplicial Depth**
- **L1 Depth**

Objetivo: Extender estas definiciones a $X_1(t), \dots, X_n(t)$ datos funcionales i.i.d. F_t con trayectorias continuas.

Profundidad vía integración

Suponga que tenemos una profundidad D_t definida sobre \mathbb{R} para cada $t \in [0, 1]$ y $x(t) \in C[(0, 1)]$. Defina:

$$Z(t) = D_t(x(t)), \quad t \in [0, 1]$$

Suponga que tenemos una profundidad D_t definida sobre \mathbb{R} para cada $t \in [0, 1]$ y $x(t) \in C[(0, 1)]$. Defina:

$$Z(t) = D_t(x(t)), \quad t \in [0, 1]$$

La profundidad funcional (poblacional) de $x = x(t)$ es

$$I(x) = \int_0^1 Z(t) \, dt$$

Suponga que tenemos una profundidad D_t definida sobre \mathbb{R} para cada $t \in [0, 1]$ y $x(t) \in C[(0, 1)]$. Defina:

$$Z(t) = D_t(x(t)), \quad t \in [0, 1]$$

La profundidad funcional (poblacional) de $x = x(t)$ es

$$I(x) = \int_0^1 Z(t) \, dt$$

Obs:

- Decimos que x es profundo si $I(x)$ es grande.

Profundidad vía integración

Suponga que tenemos una profundidad D_t definida sobre \mathbb{R} para cada $t \in [0, 1]$ y $x(t) \in C[(0, 1)]$. Defina:

$$Z(t) = D_t(x(t)), \quad t \in [0, 1]$$

La profundidad funcional (poblacional) de $x = x(t)$ es

$$I(x) = \int_0^1 Z(t) \, dt$$

Obs:

- Decimos que x es profundo si $I(x)$ es grande.
- Si los valores de $D_t(x(t))$ son pequeños (grandes), los valores de $I(x)$ serán pequeños (grandes).

Dados $X_1(t), \dots, X_n(t)$ procesos estocásticos i.i.d. con trayectorias continuas definidas en $[0, 1]$, definimos

$$Z_i(t) = D_t(X_i(t)) \quad \text{e} \quad I_i = \int_0^1 Z_i(t) \, dt.$$

Ordenamos los valores I_k de forma decrecientes en profundidad. Si I_j está en la posición i , entonces $X_{(i)}(t) = X_j(t)$. Esto es $R_j = i$ y

$$X_{(R_j)}(t) = X_j(t).$$

Version muestral: En la definición de la profundidad, reemplazamos F_t por $F_{n,t}$ (la distribución empírica asociada) en la profundidad.

Dado $\alpha > 0$ (usualmente no superior 0.25). Podemos definir la media truncada a nivel α como el promedio de los $n - \lceil n\alpha \rceil$ datos más profundos.

Más precisamente, elegimos $\beta > 0$ tal que el intervalo $[\beta, \infty)$ contenga $n - \lceil n\alpha \rceil$ valores de $I_n(X_i)$ para $i = 1, \dots, n$. Entonces:

$$\hat{\mu}_n = \frac{\sum_{i=1}^n \mathbb{1}_{[\beta, \infty)}(I_n(X_i)) X_i}{\sum_{i=1}^n \mathbb{1}_{[\beta, \infty)}(I_n(X_i))}$$

Bajo dos condiciones H_1 y H_2 uno puede garantizar que

$$\hat{\mu}_n \rightarrow \mu \text{ c.s.}$$

Donde

$$\mu = \frac{E(\mathbf{1}_{[\beta, \infty)}(I(X_1))X_1)}{E(\mathbf{1}_{[\beta, \infty)}(I(X_1)))}$$

es la media poblacional recortada.

Profundidad vía bandas

El **grafo** de una curva $x \in C[a, b]$ está definido por:

$$G(x) = \{(t, x(t)) : t \in [a, b]\}$$

Y, la **banda** delimitada por j curvas x_1, \dots, x_j es:

$$B(x_1, \dots, x_j) = \left\{ (t, y) : t \in [a, b], \min_{k=1, \dots, j} x_k(t) \leq y \leq \max_{k=1, \dots, j} x_k(t) \right\}$$

Profundidad vía bandas

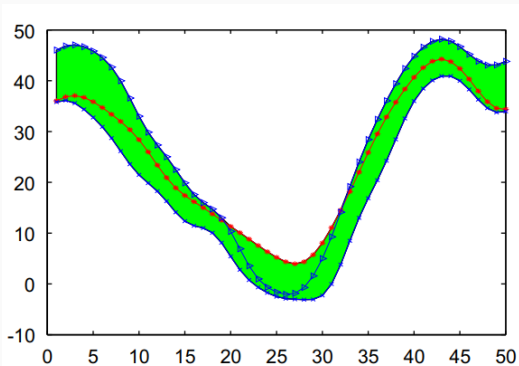


Fig. 1. Band defined by three curves.

La medida de profundidad **BD** para una curva x dado un conjunto de curvas x_1, \dots, x_n está dada por:

$$S_{n,J}(x) = \sum_{j=2}^J S_n^{(j)}(x)$$

donde

$$S_n^{(j)}(x) = \binom{n}{j}^{-1} \sum_{1 \leq i_1 \leq \dots \leq i_j \leq n} \mathbb{1}_{\{G(x) \subset B(x_{i_1}, \dots, x_{i_j})\}}(x)$$

EDOs: Dados $X_1(t), \dots, X_n(t)$ procesos estocásticos i.i.d.

EDOs: Dados $X_1(t), \dots, X_n(t)$ procesos estocásticos i.i.d.

- La mediana se define como

$$X_m = \operatorname{argmax} S_{n,J}(X_i)$$

EDOs: Dados $X_1(t), \dots, X_n(t)$ procesos estocásticos i.i.d.

- La mediana se define como

$$X_m = \operatorname{argmax} S_{n,J}(X_i)$$

- Puedo definir estadísticos de orden, rangos y medias truncadas con esta profundidad tal como antes.

Restricciones:

EDOs: Dados $X_1(t), \dots, X_n(t)$ procesos estocásticos i.i.d.

- La mediana se define como

$$X_m = \operatorname{argmax}_{S_{n,J}}(X_i)$$

- Puedo definir estadísticos de orden, rangos y medias truncadas con esta profundidad tal como antes.

Restricciones:

- Empates cuando $J = 2$.

EDOs: Dados $X_1(t), \dots, X_n(t)$ procesos estocásticos i.i.d.

- La mediana se define como

$$X_m = \operatorname{argmax}_{S_{n,J}}(X_i)$$

- Puedo definir estadísticos de orden, rangos y medias truncadas con esta profundidad tal como antes.

Restricciones:

- Empates cuando $J = 2$.
- Variabilidad y cálculo computacional intensivo cuando $J > 3$.

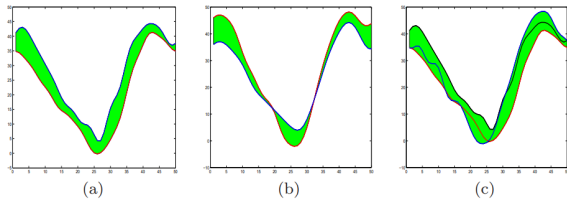


FIG 2. (a) Band determined by two curves, (b) Band determined by two curves, where the curves cross, and (c) Band determined by three curves.

Idea: Debilitar un poco la **BD** cambiando la función indicadora por la proporción de tiempo que una curva x vive “dentro” de un conjunto de curvas dadas.

Considere las curvas x, x_1, \dots, x_n y para $j \leq n$ defina

$$A(x; x_{i_1}, \dots, x_{i_j}) = \left\{ t \in [a, b] : \min_{k=1, \dots, j} x_{i_k}(t) \leq y \leq \max_{k=1, \dots, j} x_{i_k}(t) \right\}$$

Dados $x(t), x_1(t), \dots, x_n(t)$ definimos la profundidad generalizada de banda **BGD** por

$$GS_{n,J}(x) = \sum_{j=2}^J GS_n^{(j)}(x)$$

donde

$$GS_n^{(j)}(x) := \binom{n}{j}^{-1} \sum_{1 \leq i_1 \leq \dots \leq i_j \leq n} \frac{A(x; x_{i_1}, \dots, x_{i_j})}{b - a}$$

Profundidad vía “half-region”

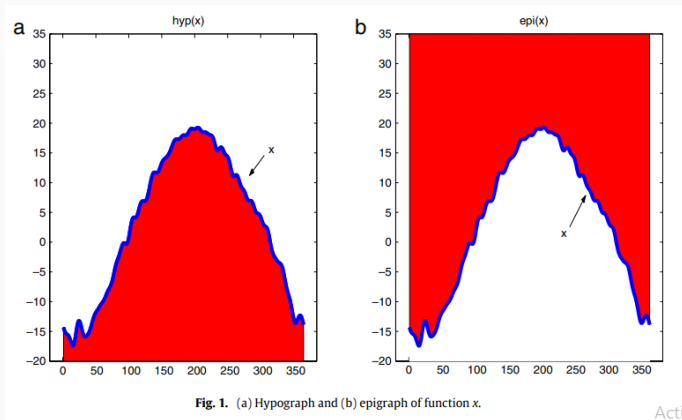
Grafo, hipografo, epigrafo

$$G(x) = \{(t, x(t)) : t \in [a, b]\}$$

$$\text{hyp}(x) = \{(t, y) \in [a, b] \times \mathbb{R} : y \leq x(t)\}$$

$$\text{epi}(x) = \{(t, y) \in [a, b] \times \mathbb{R} : y \geq x(t)\}$$

Profundidad vía “half-region”



Para una curva x defina:

$$DG_1(x) = P(G(x) \subset \text{hyp}(X)) = P(x(t) \leq X(t), t \in [a, b])$$

$$DG_2(x) = P(G(x) \subset \text{epi}(S)) = P(x(t) \geq X(t), t \in [a, b])$$

Entonces la profundidad half-region (poblacional) de x con respecto a P es definida por:

$$S_H(x) = \min\{DG_1(x), DG_2(x)\}$$

La profundidad half-region de x con respecto a x_1, \dots, x_n es:

$$S_{n,H}(x) = \min\{DG_{1,n}(x), DG_{2,n}(x)\}$$

donde

$$DG_{1,n}(x) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{\{G(x_i) \subset \text{hyp}(x)\}}(x)$$

$$DG_{2,n}(x) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{\{G(x_i) \subset \text{epi}(x)\}}(x).$$

Cómo en los ejemplos anteriores, podemos definir: mediana, estadísticos de orden, rangos y medias truncadas a nivel α .

- El dato que maximiza la profundidad es la mediana, es decir,

$$m_n = \operatorname{argmax} S_{n,H}(X)$$

- El dato que tiene la profundidad en la posición k -ésima es el estadístico de orden $X_{(k)}$

De nuevo tenemos un problema similar con los empates. Por tanto, podemos debilitar la definición y considerar la proporción de que la curva está en el epigrafo o en el hipografo. Es decir, para

$$SL(x) = \frac{1}{b-a} E[\lambda \{t \in [a, b] : x(t) \leq X(t)\}]$$

$$IL(x) = \frac{1}{b-a} E[\lambda \{t \in [a, b] : x(t) \geq X(t)\}]$$

se define la profundidad:

$$MS_H(x) = \min\{SL(x), IL(x)\}$$

Profundidad vía “half-region” generalizada

Con respecto a un conjunto de curvas x_1, \dots, x_n para

$$SL_n(x) = \frac{1}{n(b-a)} \sum_{i=1}^n \lambda\{t \in [a, b] : x(t) \leq x_i(t)\}$$

$$IL_n(x) = \frac{1}{n(b-a)} \sum_{i=1}^n \lambda\{t \in [a, b] : x(t) \geq x_i(t)\}$$

se define la profundidad

$$MS_{n,H}(x) = \min\{SL_n(x), IL_n(x)\}$$

Gracias
