

# Taller 1 Metodos No Parametricos

Mosquera Galvis, Liceth Cristina - lmosquerg@eafit.edu.co

Rios Naranjo, Johan Steward - jriosna1@eafit.edu.co

Estrada Pérez, Juan Diego - jestra15@eafit.edu.co

Cuscagua López, Juan Mauricio - jcuscagu@eafit.edu.co

Rozo Alzate Javier Arturo - jaroza@eafit.edu.co

Programa: Técnicas robustas y no paramétricas

Docente:

Henry Laniado Rodas

5 de septiembre de 2019

1. Grafique en un mismo plano las funciones de distribución empíricas de las temperaturas de cada año. Explique con base en las gráficas de las funciones, si es observable un efecto de cambio climático. Interprete los intervalos donde la distribución empírica de la temperaturas del año más reciente es mayor a la del año más antiguo.

El código en el anexo 11.1 se usa para graficar de manera iterativa la distribución empírica de las temperaturas. La figura 1 muestra las temperaturas para cada día del año a lo largo de 35 años. Por otro lado, la figura 2 muestra el comportamiento de las distribuciones empíricas para cada año.

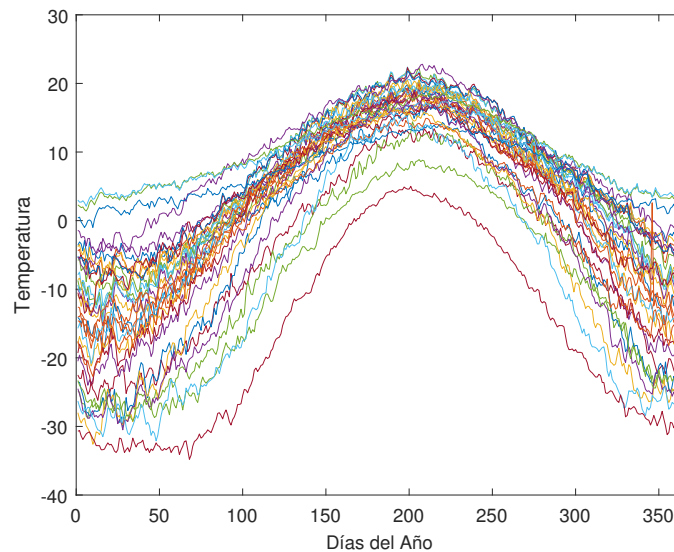


Figura 1: Temperaturas 35 años.

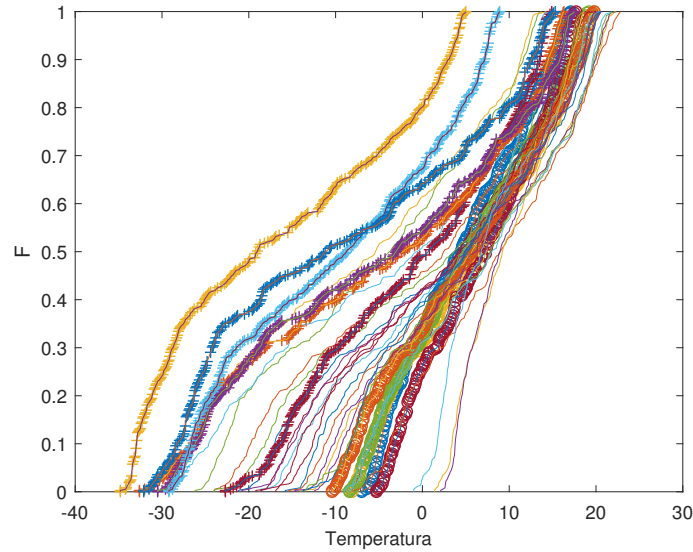


Figura 2: FDA Empíricas.

al comparar los 5 años mas antiguos (lineas punteadas en '+') con los 5 años mas recientes (lineas punteadas en 'o') se observa que las temperaturas en los primeros años eran menores a las temperaturas en los años mas recientes. En los primeros años se presentaban temperaturas entre -35 y 15. Este intervalo se ve más acotado en años posteriores donde hay temperaturas entre -10 y 20.

2. Calcule y grafique las bandas de confianza a un 95 % para la función de distribución empírica del primer año y último año. ¿Hay sectores solapados?

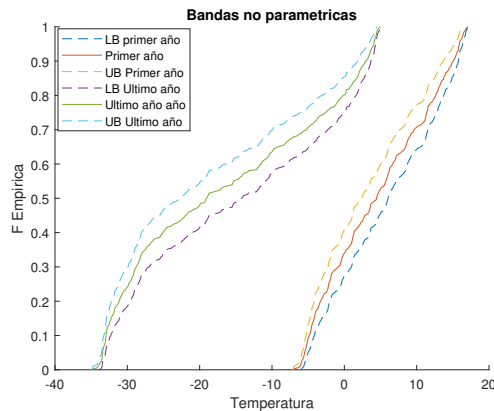
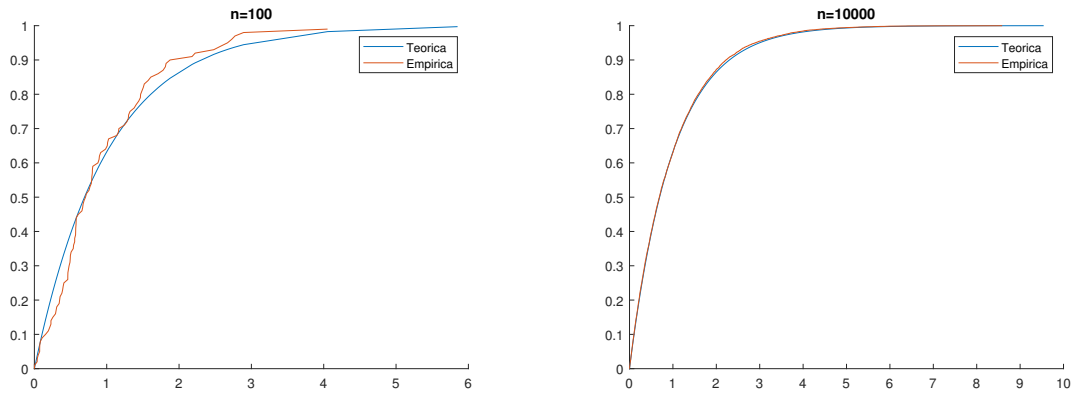


Figura 3: Bandas de confianza de funciones de distribución empíricas para el primer y último año a un 95 % de confianza.

Al comparar las funciones empíricas del primer y último año, se observa que no hay solapamiento entre las funciones o las bandas de confianza al 95 %.

3. Escriba y ejecute un código que permita visualizar el Teorema de Glivenko Cantelli para una distribución exponencial de parámetro 1.



(a) Distribución exponencial teórica vs estimada,  $n = 100$  (b) Distribución exponencial teórica vs estimada,  $n = 10000$

Figura 4: Comparación del ajuste ante el incremento de  $n$

<b>n</b>	$\max( F_n(x_i) - F(x_i) )$
100	0.11986
1000	0.018922
10000	0.0066632
100000	0.002496

Cuadro 1: Convergencia ante incremento de  $n$

El código adjunto en el anexo 11.3 permite calcular, para un tamaño  $n$  determinado, el valor  $\max(|F_n(x_i) - F(x_i)|)$  asociado al teorema de Glivenko Cantelli. Para diferentes valores de  $n$ , la tabla 1 muestra la convergencia a 0 de la expresión  $\max(|F_n(x_i) - F(x_i)|)$ . Las figuras 4a y 4b muestran cómo el ajuste a la distribución teórica mejora a medida que  $n$  incrementa.

4. Enuncie y demuestre la desigualdad de Jensen para funciones cóncavas.

Para cualquier función cóncava  $f$ , se cumple que

$$E[f(X)] \leq f(E[X])$$

**prueba:** Suponga que  $f$  es diferenciable. La función  $f$  es cóncava si, para cada  $x$  y  $y$  se cumple

$$f(x) \leq f(y) + (x - y)f'(y)$$

Sea  $x = X$  y  $y = E[X]$ , por lo tanto se puede escribir

$$f(X) \leq f(E[X]) + (X - E[X])f'(E[X])$$

Esta desigualdad se cumple para todo  $X$ , por lo que al tomar valor esperado a ambos lados, se tiene que

$$E[f(X)] \leq f(E[X]) + f'(E[X])E[(X - E[X])] = f(E[X])$$

**5. Suponga que  $X$  es una variable aleatoria exponencial de parámetro  $\beta$ . Calcule  $P(|X - \mu| > k\sigma)$  para  $k > 1$ . Compare el resultado con la cota obtenida de la desigualdad de Chebyshev.**

Dada la variable aleatoria  $X \sim \exp(\beta)$  y tomando  $\mu = \frac{1}{\beta}$ ,  $\sigma = \frac{1}{\beta}$  se tiene que:

$$\begin{aligned} P(|X - \mu| > k\sigma) &= 1 - P(|X - \mu| < k\sigma) \\ &= 1 - P(-k\sigma < X - \mu < k\sigma) \\ &= 1 - P(\mu - k\sigma < X < \mu + k\sigma) \\ &= 1 - P\left(\frac{1}{\beta} - k\frac{1}{\beta} < X < \frac{1}{\beta} + k\frac{1}{\beta}\right) \\ &= 1 - P\left(\frac{1}{\beta}(1 - k) < X < \frac{1}{\beta}(1 + k)\right) \\ &= 1 - (1 - e^{-\beta(\frac{1}{\beta}(k+1))}) \\ &= 1 - (1 - e^{-(k+1)}) \\ &= e^{-(k+1)} \end{aligned}$$

Al usar la desigualdad de Chebyshev se tiene que:

$$P(|X - \mu| > k\sigma) = \frac{\sigma^2}{k^2\sigma^2} = \frac{1}{k^2}$$

de esta manera, se tiene que  $e^{-(k+1)} \leq \frac{1}{k^2}$ . Por lo tanto, a medida que  $k$  aumenta, ambos valores tienden a converger, esto se ve al tomar límite en ambos lados de la desigualdad.

$$\lim_{k \rightarrow \infty} e^{-(k+1)} = 0 = \lim_{k \rightarrow \infty} \frac{1}{k^2}$$

**6. Demuestre que si  $X$  es Poisson de parámetro  $\lambda$  entonces  $P(x \geq 2\lambda) \leq \frac{1}{\lambda}$**

De la desigualdad de Chebyshev tenemos que

$$P(|x - \lambda| \geq \lambda) \leq \frac{\lambda}{\lambda^2} = \frac{1}{\lambda}$$

Por otro lado, se tiene que:

$$\begin{aligned} P(|X - \lambda| \geq \lambda) &= 1 - P(|X - \lambda| \leq \lambda) \\ &= 1 - P(-\lambda \leq x - \lambda \leq \lambda) \\ &= 1 - P(0 \leq x \leq 2\lambda) \\ &= 1 - P(x \leq 2\lambda) \\ &= P(x \geq 2\lambda) \end{aligned}$$

Por lo tanto,

$$P(x \geq 2\lambda) \leq \frac{1}{\lambda}$$

**7. Demuestre que convergencia en probabilidad está implicada por la convergencia en media cuadrática.**

Se debe probar que

$$\lim_{n \rightarrow \infty} P(|x_n - x| > \epsilon) = 0$$

Consideremos entonces,

$$\begin{aligned} P(|X_n - x| > \epsilon) &= P((X_n - x)^2 > \epsilon^2) \\ &\leq \frac{E[(x_n - x)^2]}{\epsilon^2} \end{aligned}$$

Esto se cumple debido a la desigualdad de Markov. Al tomar límites en ambos lados de la desigualdad

$$\begin{aligned} \lim_{n \rightarrow \infty} P(|X_n - x| > \epsilon) &\leq \lim_{n \rightarrow \infty} P((X_n - x)^2 > \epsilon^2) \\ &= 0 \end{aligned}$$

de esta manera

$$0 \leq \lim_{n \rightarrow \infty} P(|X_n - x| > \epsilon) \leq 0$$

por lo tanto

$$\lim_{n \rightarrow \infty} P(|X_n - x| > \epsilon) = 0$$

**8. Considere las temperaturas diarias del primer año. Calcule el intervalo de confianza para la temperatura mínima. Calcule el sesgo de  $T_{[1]}$  y la varianza**

En los códigos relacionados en el anexo 11.4, se relacionan los cálculos para el sesgo, el mínimo estimado y un intervalo de confianza usando la varianza de Jack-knife. Los resultados son los siguientes:

- Sesgo: 0.49863
- Mínimo Estimado: -7
- Intervalo de confianza: [-7.9773, -6.0227]

**9. Considere  $U_1, U_2, \dots, U_n$  una muestra de una distribución uniforme en el intervalo  $[0, 1]$ . Calcule la distribución teórica de  $U_{[1]}$ , su media y sesgo. Genere la muestra y utilice bootstrap para calcular la varianza de  $U_{[1]}$ . Calcule el sesgo por Jackknife y compárelo con el sesgo teórico**

Para encontrar el valor teórico del sesgo, hay que considerar  $P(\min\{U_1, \dots, U_n\} \leq t)$ :

$$\begin{aligned}
P(\min\{U_1, \dots, U_n\} \leq t) &= 1 - P(\min\{U_1, \dots, U_n\} \geq t) \\
&= 1 - P(u_1 > t \wedge \dots \wedge u_n > t) \\
&= 1 - [1 - t]^n \\
&= F(t)
\end{aligned}$$

por lo tanto, se tiene que  $f(t) = F'(t) = n(1 - t)^{n-1}$ . Por otro lado,

$$\begin{aligned}
E[\min\{u_1, \dots, u_n\}] &= \int_0^1 nt(1 - t)^{n-1} dt \\
&= - \int_1^0 n(1 - u)u^{n-1} du \\
&= n \int_0^1 u^{n-1} u^n du \\
&= n \left( \frac{u^n}{n} - \frac{u^{n+1}}{n+1} \right) \Big|_0^1 = 1 - \frac{n}{n+1}
\end{aligned}$$

Considerando  $\theta$  como el parámetro real, se tiene que el sesgo teórico es

$$E[\min\{u_1, \dots, u_n\}] - \theta = 1 - \frac{n}{n+1}$$

Usando el código del anexo 11.5, computacionalmente se obtienen los siguientes resultados:

- Varianza Bootstrap: 0.00019173
- Sesgo Jack-knife: 0.020724
- Sesgo Teórico: 0.0099014

## 10. Explique una forma no paramétrica y robusta de calcular las componentes principales. Aplique la técnica a las temperaturas del primer año y el último año. Compárelas con las componentes principales obtenidas de forma habitual.

Entre las principales técnicas de reducción de dimensionalidad, análisis de componentes principales es uno de los más comunes en implementación. Esta metodología usa bien sea la matriz de covarianzas o la matriz de correlaciones de tal forma que, al obtener los vectores propios de la matriz, se puede crear una combinación lineal de los datos en múltiples dimensiones a una proyección en una dimensión menor. En este orden de ideas, una versión robusta del análisis de componentes principales se daría al implementar una matriz de covarianza robusta la cual podría generar resultados más estables ante la presencia de datos atípicos.

La tabla 2 muestra los resultados del análisis de componentes principales paramétrico y no paramétrico para los datos originales. La tabla ?? muestra los mismos resultados para datos contaminados.

Paramétrico				No Paramétrico			
% Var	Valores Propios	Coeficientes		% Var	Valores Propios	Coeficientes	
96.8418	535.9271	0.6286	0.7777	98.5668	209.6023	0.4573	0.8893
3.1582	17.4777	0.7777	-0.6286	1.4332	3.0477	0.8893	-0.4573

Cuadro 2: Resultados de análisis de componentes principales sobre los datos originales.

Paramétrico				No Paramétrico			
% Var	Valores Propios	Coeficientes		% Var	Valores Propios	Coeficientes	
95.8948	322.2904	0.5721	0.8201	98.2432	198.9327	0.4674	0.8840
4.1052	13.7972	0.8201	-0.5721	1.7568	3.5573	0.8840	-0.4674

Cuadro 3: Resultados de análisis de componentes principales sobre los datos contaminados.

Se puede ver como incluir datos contaminados en la muestra altera la capacidad explicativa de los componentes principales. La implementación se encuentra en el código adjunto en el anexo 11.6

## 11. Anexos

### 11.1. Anexo punto 1

```

1 load('temperaturas.mat')
2 A= temperatura;
3 [rows,cols] = size(A);
4 for i = 1:cols
5     if i <= 5
6         C = A(:,i);
7         [f,x] = ecdf(C);
8         plot(x,f, 'o')
9         hold on
10    end
11    if i >= (cols-5)
12        C = A(:,i);
13        [f,x] = ecdf(C);
14        plot(x,f, '+')
15        hold on
16    end
17    C = A(:,i);
18    [f,x] = ecdf(C);
19    plot(x,f)
20    hold on
21 end
22
23 figure
24 for i = 1:cols
25     plot(A(:,i))
26     hold on
27 end

```

```
28 xlim([0 365])
```

## 11.2. Anexo punto 2

```
1 load('temperaturas.mat')
2
3 [f1,x1] = ecdf(temperatura(:,1));
4 [f35,x35] = ecdf(temperatura(:,end));
5
6 alfa = 0.05;
7
8 %% first year
9
10 error1 = sqrt(f1.*(1-f1)/length(f1)); %Error Estandar
11
12 %Calculo de las bandas semi-parametrico
13 L1 = f1-1.96.*error1;
14 U1 = f1+1.96.*error1;
15
16 %Calculo de las bandas No parametrico
17
18 eps = sqrt((1/(2*length(f1)))*log(2/alfa));
19 for i=1:length(f1)
20     L1np(i) = max([f1(i)-eps,0]);
21     U1np(i) = min([f1(i)+eps,1]);
22 end
23
24 %% last year
25
26 error35 = sqrt(f35.*(1-f35)/length(f35)); %Error Estandar
27
28 %Calculo de las bandas semi-parametrico
29 L35 = f35-1.96.*error35;
30 U35 = f35+1.96.*error35;
31
32 %Calculo de las bandas No parametrico
33 eps = sqrt((1/(2*length(f35)))*log(2/alfa));
34 for i=1:length(f35)
35     L35np(i) = max([f35(i)-eps,0]);
36     U35np(i) = min([f35(i)+eps,1]);
37 end
38
39 %% Plots
40
41 %Bandas Semiparametricas
42 figure
43 hold on
44 plot(x1,L1,'—')
45 plot(x1,f1)
46 plot(x1,U1,'—')
47
```



```

48 plot(x35,L35,'—')
49 plot(x35,f35)
50 plot(x35,U35,'—')
51 legend('LB primer a o ','Primer a o ','UB Primer a o ','LB Ultimo a o ','
    Ultimo a o a o ','UB Ultimo a o ')
52 ylabel('F Empirica')
53 xlabel('Temperatura')
54 ylim([0 1])
55 title('Bandas semi-parametricas')
56 hold off
57
58 %Bandas no parametricas
59 figure
60 hold on
61 plot(x1,L1,'—')
62 plot(x1,f1)
63 plot(x1,U1,'—')
64
65 plot(x35,L35,'—')
66 plot(x35,f35)
67 plot(x35,U35,'—')
68 legend('LB primer a o ','Primer a o ','UB Primer a o ','LB Ultimo a o ','
    Ultimo a o a o ','UB Ultimo a o ')
69 ylabel('F Empirica')
70 xlabel('Temperatura')
71 ylim([0 1])
72 title('Bandas no parametricas')
73 hold off

```

### 11.3. Anexo punto 3

```

1 n = 10000;
2 lambda = 1;
3 X = sort(-(1/lambda)*log(rand(n,1)/lambda));
4
5 F = (1-exp(-lambda*X));
6 [F_est, x_est] = ecdf(X);
7 F_est = F_est(1:end-1);
8 x_est = x_est(1:end-1);
9
10 figure
11 hold on
12 plot(X,F)
13 plot(x_est,F_est)
14 legend('Teorica','Empirica')
15 title(strcat('n=',num2str(n)))
16 hold off
17
18 GlivCant = max(abs(F_est-F));
19 disp(strcat('Glivenko-Cantelli: ',num2str(GlivCant)))

```

## 11.4. Anexo punto 8

```
1 load('temperaturas.mat')
2
3 n = 365;
4 X = temperatura(:,1);
5 minimo = min(X);
6
7 %Sesgo
8 Bias = Bias_min(X);
9 disp(strcat('Sesgo:', num2str(Bias)))
10
11 %Varianza
12 Ssq_tilde = 0;
13 for i=1:n
14     T_tilde_barr = 0;
15     for j=1:n
16         T_tilde_barr = T_tilde_barr + Ti_tilde(X,j);
17     end
18     T_tilde_barr = T_tilde_barr/n;
19     Ssq_tilde = Ssq_tilde + (Ti_tilde(X,i) - T_tilde_barr)^2;
20 end
21 Ssq_tilde = Ssq_tilde/(n-1);
22
23 Varianza = Ssq_tilde/n;
24 SE = sqrt(Varianza);
25
26 %Intervalo de Confianza:
27
28 Lx = minimo-1.96.*SE;
29 Ux = minimo+1.96.*SE;
30
31 disp(strcat('Minimo Estimado: ', num2str(minimo)))
32 disp(strcat('Intervalo de confianza (Usando varianza de Jackknife): ',
    num2str(Lx), ', ', num2str(Ux)))

1 function bias = Bias_min(X)
2 Tn_barr = 0;
3 [n y] = size(X);
4 Tn = min(X);
5
6 for i=1:n
7     Xi = X;
8     Xi(i) = [];
9     Ssq = min(Xi);
10    Tn_barr = Tn_barr + Ssq;
11 end
12
13 Tn_barr = Tn_barr/n;
14 bias = (n-1)*(Tn_barr-Tn);
15 end
```

## 11.5. Anexo punto 9

```
1  rng(19)
2  n=100;
3  U = rand(n,1);
4
5  %Varianza Bootstrap
6  B = 1000;
7  EB = bootstrp(B,@min,U);
8
9  Vboots = 0;
10 for b=1:B
11     Vboots = Vboots + (EB(b,1)-mean(EB))^2;
12 end
13
14 Vboots = Vboots/B;
15
16 %Sesgo
17 Bias = Bias_min(U);
18 Bias_teorico = 1-(n/(n+1));
19
20 disp(strcat('Varianza Bootstrap:',num2str(Vboots)))
21 disp(strcat('Sesgo Jackknife:',num2str(Bias)))
22 disp(strcat('Sesgo Teorico:',num2str(Bias_teorico)))
```

## 11.6. Anexo punto 10

```
1  %% Comedian y Datos
2  comedian = @(X,Y) median((X-median(X)).*(Y-median(Y)));
3  load('temperaturas.mat')
4  %% Componentes principales
5  X = [temperatura(:,1) temperatura(:,35)];
6
7  v = cov(X);
8  [Coef,latent,explained] = pcacov(v);
9
10 z1 = Coef(:,1);
11 z2 = Coef(:,2);
12
13 %Version robusta de la matriz de covarianzas
14
15 covR = zeros(2,2);
16 for i=1:2
17     for j=1:2
18         covR(i,j) = Comedian(X(:,i),X(:,j));
19     end
20 end
21
22 %PCA Robusto
23
24 [CoefR,latentR,explainedR] = pcacov(covR);
25
```

```

26 zr1 = CoefR(:,1);
27 zr2 = CoefR(:,2);
28
29 %% Contaminaci n de los datos
30 X = [X ; (X(28:35,:)+ 100) ];
31
32 v = cov(X);
33 [Coef,latent,explained] = pcacov(v);
34
35 z1 = Coef(:,1);
36 z2 = Coef(:,2);
37
38 %Version robusta de la matriz de covarianzas
39
40 covR = zeros(2,2);
41 for i=1:2
42     for j=1:2
43         covR(i,j) = Comedian(X(:,i),X(:,j));
44     end
45 end
46
47 PCA Robusto
48
49 [CoefR,latentR,explainedR] = pcacov(covR);

```