

ESTIMAREA NEDEPLASATĂ, CU DISPERSIE MINIMĂ, GENERALIZATĂ

Într-o problemă de estimare a unui parametru necunoscut, θ , dispunem de cele N eșantioane de date $\{x[0], x[1], \dots, x[N-1]\}$. Fiecare dintre acestea poartă informație despre parametrul necunoscut. Ne punem întrebarea dacă nu putem găsi un singur număr, T , care depinde de date și poartă toată informația despre θ

$$T(x[0], x[1], \dots, x[N-1]) = T(\mathbf{x})$$

Ca și datele, $T(\mathbf{x})$ este o variabilă aleatoare, numită "statistică". Vom considera modelul de semnal

$$x[n] = A + w[n]; \quad n = 0, 1, \dots, N-1; \quad \mathbf{w} \sim \mathcal{N}(0, \sigma^2 \mathbf{I}_u)$$

Am arătat că media eșantion este un estimator MVU eficient Pentru componenta continuă, A

$$\hat{A} = \bar{x} = \frac{1}{N} \sum_{n=0}^{N-1} x[n]$$

Pentru a estima A nu trebuie să cunoaștem toate valorile $x[n]$. Ne putem mărgini la cunoașterea unei singure valori, statistica $T(\mathbf{x})$

$$T(\mathbf{x}) = \sum_{n=0}^{N-1} x[n]$$

Ea se numește și "statistică suficientă"

1

Dacă pentru două seturi diferite de date

$$\mathbf{x}_1 \neq \mathbf{x}_2$$

dar pentru care valorile statisticii $T(\mathbf{x})$ sunt identice, adică

$$T(\mathbf{x}_1) = T(\mathbf{x}_2)$$

atunci valorile estimatului pentru θ determinate din cele două seturi de date sunt identice

Ne putem pune întrebarea firească "câte statistici suficiente există?"

Dacă ne referim tot la exemplul în care se estimează componenta continuă, A , cele N date măsurate sunt suficiente pentru măsurare (de altele nici nu dispunem). Prin urmare, mulțimea

$$S_1 = \{x[0], x[1], \dots, x[N-1]\}$$

constituie o statistică suficientă. În oricare situație, datele măsurate formează o statistică suficientă

Dar, în mod evident, o statistică suficientă o constituie și mulțimea

$$S_2 = \{x[0] + x[1], x[2] + x[3], \dots\}$$

La început am văzut că mulțimea cu un singur element

$$S_3 = \left\{ \sum_{n=0}^{N-1} x[n] \right\}$$

este și ea o statistică suficientă dar pe care o vom numi și "minimală" deoarece are numărul minim de elemente între statisticile suficiente posibile

Pentru exemplul luat în considerare, avem

$$p(\mathbf{x}; A) = \frac{1}{(\sqrt{2\pi}\sigma)^N} \exp\left\{-\frac{1}{2\sigma^2} \sum_{n=0}^{N-1} (x[n] - A)^2\right\}$$

Dacă datele \mathbf{x} măsurate determină o valoare fixă pentru statistica suficientă

$$T(\mathbf{x}) = \sum_{n=0}^{N-1} x[n] = T_0$$

atunci densitatea de repartiție condiționată de aceasta nu mai poate fi funcție de parametrul necunoscut, A .

$$p\left(\mathbf{x} \mid \sum_{n=0}^{N-1} x[n] = T_0; A\right)$$

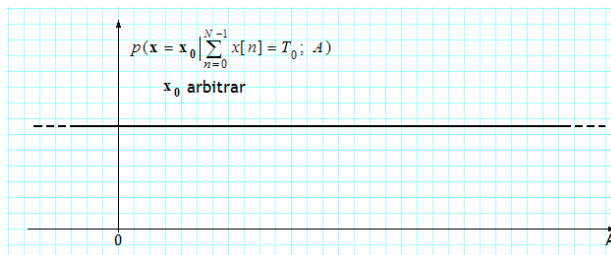
Dacă ar exista dependența de A , din date \mathbf{x} diferite, dar care ar da aceeași valoare a statisticii suficiente, am mai putea obține informații privind parametrul necunoscut, A . Dar atunci $T(\mathbf{x})$ nu ar fi statistică suficientă!

Pentru

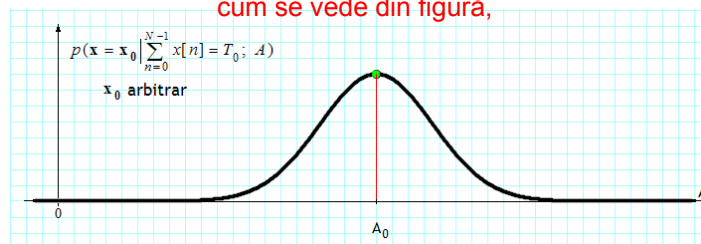
$$p\left(\mathbf{x} = \mathbf{x}_0 \mid \sum_{n=0}^{N-1} x[n] = T_0; A\right); \quad \mathbf{x}_0 \text{ fixat, dar arbitrar}$$

3

se arată în figură independența de valoarea parametrului necunoscut, A , dacă $T(\mathbf{x})$ este o statistică suficientă



Dacă există variație funcție de valoarea parametrului necunoscut, A , deși pentru diverse seturi de date se obține aceeași valoare a statisticii $T(\mathbf{x})$, așa cum se vede din figură,



atunci nu avem de-a face cu o statistică suficientă!

4

Se pune întrebarea “cum poate fi determinată statistica suficientă (eventual minimală)?” Răspunsul este dat de teorema de factorizare Neyman-Fisher (Noiman-Fișer), al cărui enunț îl dăm fără demonstrație

Teorema de factorizare Neyman-Fisher

Dacă densitatea de probabilitate a datelor, \mathbf{x} , $p(\mathbf{x}; \theta)$, dependentă de parametrul necunoscut θ , poate fi factorizată sub forma

$$p(\mathbf{x}; \theta) = g(T(\mathbf{x}), \theta)h(\mathbf{x})$$

în care $g(\cdot)$ este o funcție care depinde de datele \mathbf{x} numai prin intermediul statisticii $T(\mathbf{x})$ iar $h(\cdot)$ este o funcție numai de datele \mathbf{x} , nu și de statistica $T(\mathbf{x})$ sau de parametrul necunoscut θ , atunci $T(\mathbf{x})$ este o statistică suficientă pentru θ

Reciproc, dacă $T(\mathbf{x})$ este o statistică suficientă pentru θ , atunci se poate obține factorizarea de mai sus

5

Statistica suficientă pentru estimarea unui nivel continuu, în zgomot alb, gaussian

Dacă în expresia densității de repartiție corespunzătoare se dezvoltă pătratul de la exponent, putem obține forma cerută de teorema de factorizare

$$p(\mathbf{x}; A) = \underbrace{\frac{1}{(\sqrt{2\pi}\sigma)^N} \exp\left\{-\frac{1}{2\sigma^2} \left[NA^2 - 2A \sum_{n=0}^{N-1} x[n] \right] \right\}}_{g(T(\mathbf{x}), A)} \underbrace{\exp\left\{-\frac{1}{2\sigma^2} \sum_{n=0}^{N-1} x^2[n] \right\}}_{h(\mathbf{x})}$$

și deducem că statistica suficientă pentru estimarea componentei continue este

$$T(\mathbf{x}) = \sum_{n=0}^{N-1} x[n]$$

Statistica suficientă pentru estimarea puterii zgomotului alb, gaussian

Pentru zgomotul alb gaussian, cu puterea (dispersia) necunoscută

$$x[n] = w[n]; \quad n = 0, 1, \dots, N-1, \quad w[n] \sim \mathcal{N}(0, \sigma^2)$$

factorizarea este evidentă dacă scriem densitatea de repartiție sub forma

$$p(\mathbf{x}; \sigma^2) = \underbrace{\frac{1}{(2\pi\sigma^2)^{N/2}} \exp\left\{-\frac{1}{2\sigma^2} \sum_{n=0}^{N-1} x^2[n] \right\}}_{g(T(\mathbf{x}); \sigma^2)} \cdot \frac{1}{h(\mathbf{x})}$$

6

și rezultă că statistica suficientă pentru estimarea dispersiei zgomotului este

$$T(\mathbf{x}) = \sum_{n=0}^{N-1} x^2[n]$$

Problema estimării fazei unei sinusoide

Pentru modelul de semnal sinusoidal cu faza inițială, Φ , necunoscută, afectată de un zgomot alb, gaussian datele $x[n]$ au forma

$$x[n] = A \cos(2\pi f_0 n + \Phi) + w[n]; \quad n = 0, 1, \dots, N-1; \quad \mathbf{w} \sim \mathcal{N}(0, \sigma^2 \mathbf{I}_u)$$

iar densitatea de probabilitate a datelor \mathbf{x} are forma

$$p(\mathbf{x}; \Phi) = \frac{1}{(\sqrt{2\pi}\sigma)^N} \exp\left\{-\frac{1}{2\sigma^2} \sum_{n=0}^{N-1} [x[n] - A \cos(2\pi f_0 n + \Phi)]^2\right\}$$

Dezvoltăm pătratul de la exponent și obținem

$$\begin{aligned} & \sum_{n=0}^{N-1} [x[n] - A \cos(2\pi f_0 n + \Phi)]^2 \\ &= \sum_{n=0}^{N-1} x^2[n] - 2A \sum_{n=0}^{N-1} x[n] \cos(2\pi f_0 n + \Phi) + A^2 \sum_{n=0}^{N-1} \cos^2(2\pi f_0 n + \Phi) \\ &= \sum_{n=0}^{N-1} x^2[n] - 2A \cos \Phi \left[\sum_{n=0}^{N-1} x[n] \cos 2\pi f_0 n \right] \\ & \quad + 2A \sin \Phi \left[\sum_{n=0}^{N-1} x[n] \sin 2\pi f_0 n \right] + A^2 \sum_{n=0}^{N-1} \cos^2(2\pi f_0 n + \Phi) \end{aligned}$$

7

Cu notațiile

$$T_1(\mathbf{x}) = \sum_{n=0}^{N-1} x[n] \cos 2\pi f_0 n$$

$$T_2(\mathbf{x}) = \sum_{n=0}^{N-1} x[n] \sin 2\pi f_0 n$$

expresia anterioară devine

$$\begin{aligned} & \sum_{n=0}^{N-1} [x[n] - A \cos(2\pi f_0 n + \Phi)]^2 \\ &= \sum_{n=0}^{N-1} x^2[n] - 2AT_1(\mathbf{x}) \cos \Phi + 2AT_2(\mathbf{x}) \sin \Phi \\ & \quad + A^2 \sum_{n=0}^{N-1} \cos^2(2\pi f_0 n + \Phi) \end{aligned}$$

8

Densitatea de probabilitate a datelor \mathbf{x} se poate factoriza acum

$$p(\mathbf{x}; \Phi) = \frac{1}{(\sqrt{2\pi}\sigma)^N} \cdot \underbrace{\exp\left\{-\frac{1}{2\sigma^2}\left[A^2 \sum_{n=0}^{N-1} \cos^2(2\pi f_0 t + \Phi) - 2A \cos \Phi \cdot T_1(\mathbf{x}) + 2A \sin \Phi \cdot T_2(\mathbf{x})\right]\right\}}_{g(T_1(\mathbf{x}), T_2(\mathbf{x}), \Phi)} \cdot \underbrace{\exp\left\{-\frac{1}{2\sigma^2} \sum_{n=0}^{N-1} x^2[n]\right\}}_{h(\mathbf{x})}$$

Apar două statistici în loc de una. Teorema de factorizare Neyman-Fisher se poate extinde, după cum urmează

9

Extinderea teoremei de factorizare Neyman- Fisher pentru un grup de r statistici suficiente

Dacă densitatea de probabilitate a datelor, \mathbf{x} , dependentă de parametrul necunoscut θ , $p(\mathbf{x}; \theta)$ poate fi factorizată sub forma

$$p(\mathbf{x}; \theta) = g(T_1(\mathbf{x}), T_2(\mathbf{x}), \dots, T_r(\mathbf{x}); \theta) h(\mathbf{x})$$

atunci

$$\{T_1(\mathbf{x}), T_2(\mathbf{x}), \dots, T_r(\mathbf{x})\}$$

formează un grup de r statistici mutual suficiente pentru estimarea parametrului necunoscut, θ . Reciproca teoremei este adevărată.

În orice problemă de estimare, deoarece densitatea de repartiție a datelor se poate scrie sub forma

$$p(\mathbf{x}; \theta) = \underbrace{p(x[0], x[1], \dots, x[N-1]; \theta)}_{g(x[0], x[1], \dots, x[N-1]; \theta)} \cdot \underbrace{1}_{h(\mathbf{x})}$$

rezultă că, la limită, datele pot fi asimilate cu un grup de N statistici mutual suficiente pentru estimarea parametrului necunoscut

$$\{x[0], x[1], \dots, x[N-1]\}$$

10

Determinarea estimatorilor MVU plecând de la o statistică suficientă

Dacă am determinat o statistică suficientă, $T(\mathbf{x})$, pentru un parametru necunoscut, θ , se poate găsi un estimator MVU în două feluri, dintre care vom prezenta doar unul:

se caută o funcție $g(\cdot)$, astfel încât

$$\hat{\theta} = g(T(\mathbf{x}))$$

să fie un estimator nedeplasat pentru θ , adică

$$E\{\hat{\theta}\} = E\{g(T(\mathbf{x}))\} = \theta$$

Pentru exemplificare reluăm problema estimării componentei continue pentru care statistica suficientă este

$$T(\mathbf{x}) = \sum_{n=0}^{N-1} x[n]$$

Trebuie găsită funcția $g(x)$ pentru care să avem

$$E\left\{g\left(\sum_{n=0}^{N-1} x[n]\right)\right\} = A$$

11

Funcția $g(x)$ are forma evidentă

$$g(x) = \frac{x}{N}$$

astfel că estimatorul componentei continue A este, în mod evident

$$\hat{A} = \frac{1}{N} \sum_{n=0}^{N-1} x[n]$$

Dacă funcția $g(\cdot)$ este unică, statistica suficientă $T(\mathbf{x})$ se spune că este "completă"

Familia de repartiții exponențial-scalare, ce are forma

$$p(x; \theta) = \exp\{A(\theta)B(x) + C(x) + D(\theta)\}$$

are proprietatea de a genera statistici suficiente complete pentru parametrul necunoscut, θ . Repartiția gaussiană cu media μ necunoscută, repartiția Rayleigh cu dispersia necunoscută și repartiția exponențială cu parametrul λ necunoscut, fac toate parte din familia exponențial-scalară

12

Repartiția exponențială scalară poate fi factorizată conform teoremei Neyman-Fisher sub forma

$$p(x; \theta) = \underbrace{\exp\{A(\theta)B(x) + D(\theta)\}}_{g(T(x), \theta)} \underbrace{\exp\{C(x)\}}_{h(x)}$$

din care rezultă că statistica suficientă pentru parametrul necunoscut θ este

$$T(x) = B(x)$$

Considerăm că datele

$$\mathbf{x} = [x[0] \quad x[1] \quad \dots \quad x[N-1]]^T$$

sunt de tip IID, motiv pentru care densitatea de repartiție a vectorului \mathbf{x} este

$$\begin{aligned} p(\mathbf{x}; \theta) &= \prod_{n=0}^{N-1} p(x[n]; \theta) \\ &= \prod_{n=0}^{N-1} \left[\exp\{A(\theta)B(x[n]) + D(\theta)\} \exp\{C(x[n])\} \right] \\ &= \exp\left\{ A(\theta) \sum_{n=0}^{N-1} B(x[n]) + ND(\theta) \right\} \exp\left\{ \sum_{n=0}^{N-1} C(x[n]) \right\} \end{aligned}$$

Statistica suficientă și completă pentru parametrul necunoscut este, conform teoremei Neyman-Fisher

$$T(\mathbf{x}) = \sum_{n=0}^{N-1} B(x[n])$$

Prezentăm câteva exemple

1) Pentru o repartiție gaussiană de medie μ necunoscută, densitatea de repartiție este

$$\begin{aligned} p(x; \mu) &= \frac{1}{\sqrt{2\pi}\sigma} \exp\left\{ -\frac{1}{2\sigma^2} (x - \mu)^2 \right\} \\ &= \exp\left\{ \frac{\mu}{\sigma^2} x - \frac{\mu^2}{2\sigma^2} - \ln\sqrt{2\pi} - \ln\sigma \right\} \exp\left\{ -\frac{x^2}{2\sigma^2} \right\} \end{aligned}$$

Se observă că

$$B(x) = x$$

și deci statistica suficientă și completă, pentru cazul unui vector de date \mathbf{x} este, pentru medie

$$T_1(\mathbf{x}) = \sum_{n=0}^{N-1} x[n]$$

2) Pentru o repartiție Rayleigh de dispersie necunoscută, densitatea de repartiție este

$$p(x; \sigma^2) = \frac{x}{\sigma^2} \exp\left\{-\frac{x^2}{2\sigma^2}\right\} u(x)$$
$$= \exp\left\{-\frac{1}{2\sigma^2}x^2 - \ln \sigma^2\right\} \exp\{\ln[xu(x)]\}$$

în care $u(x)$ este treapta unitară

$$u(x) = \begin{cases} 1, & x > 0 \\ 0, & x < 0 \end{cases}$$

Prin identificare rezultă că

$$B(x) = x^2$$

și deci statistica suficientă și completă, pentru cazul unui vector de date \mathbf{x} este, pentru dispersie

$$T_2(\mathbf{x}) = \sum_{n=0}^{N-1} x^2[n]$$

15

3) Pentru o repartiție exponențială parametru λ necunoscut, densitatea de repartiție este

$$p(x; \lambda) = \lambda \exp\{-\lambda x\} u(x)$$
$$= \exp\{-\lambda x + \ln \lambda\} \exp\{\ln u(x)\}$$

Prin identificare rezultă că

$$B(x) = x$$

și deci statistica suficientă și completă, pentru cazul unui vector de date \mathbf{x} este, pentru parametrul λ

$$T_3(\mathbf{x}) = \sum_{n=0}^{N-1} x[n]$$

16

Pentru exemplul 1), am stabilit deja că funcția $g(x)$ este

$$g(x) = \frac{x}{N}$$

și deci estimatorul mediei devine

$$\hat{\mu} = \frac{1}{N} \sum_{n=0}^{N-1} x[n]$$

Pentru exemplul 2), vom căuta forma funcției $g(x)$. Media de ordinul doi a variabilei aleatoare Rayleigh este

$$\begin{aligned} E\{x^2\} &= \int_0^{\infty} x^2 \frac{x}{\sigma^2} \exp\left\{-\frac{x^2}{2\sigma^2}\right\} dx = 2\sigma^2 \int_0^{\infty} v e^{-v} dv \\ &= 2\sigma^2 v (-e^{-v}) \Big|_0^{\infty} + 2\sigma^2 \int_0^{\infty} e^{-v} dv = 2\sigma^2 \end{aligned}$$

Calculăm media statisticii suficiente pentru cazul 2). Avem

$$E\left\{\sum_{n=0}^{N-1} x^2[n]\right\} = \sum_{n=0}^{N-1} E\{x^2[n]\} = \sum_{n=0}^{N-1} \sigma^2 = 2N\sigma^2 \neq \sigma^2$$

Nu este greu de observat că

$$E\left\{\frac{1}{2N} \sum_{n=0}^{N-1} x^2[n]\right\} = \sigma^2$$

17

asa că funcția $g(x)$, pentru cazul 2) este, evident

$$g(x) = x/(2N)$$

Am obținut, pentru acest caz estimatorul

$$\hat{\sigma}^2 = \frac{1}{2N} \sum_{n=0}^{N-1} x^2[n]$$

Și pentru exemplul 3) trebuie căutată forma funcției $g(x)$. Media variabilei aleatoare exponențiale este

$$E\{x\} = \int_0^{\infty} x \lambda \exp\{-\lambda x\} dx = \lambda x \left(-\frac{1}{\lambda} e^{-\lambda x}\right) \Big|_0^{\infty} + \lambda \int_0^{\infty} \frac{1}{\lambda} e^{-\lambda x} dx = \frac{1}{\lambda}$$

și deci

$$E\left\{\sum_{n=0}^{N-1} x[n]\right\} = \sum_{n=0}^{N-1} E\{x[n]\} = \frac{N}{\lambda} \neq \lambda$$

Această expresie a mediei ne sugerează să luăm ca parametru necunoscut

$$\theta = 1/\lambda$$

Pentru noul parametru necunoscut, θ , avem densitatea de repartiție

$$\begin{aligned} p(x; \lambda) &= p(x; \theta) = \frac{1}{\theta} \exp\left\{-\frac{x}{\theta}\right\} u(x) \\ &= \exp\left\{-\frac{1}{\theta} \cdot x + \ln \frac{1}{\theta}\right\} \exp\{\ln u(x)\} \end{aligned}$$

18

Rezultă că

$$B(x) = x$$

Dacă ținem seama că media datelor $x[n]$ este, așa cum am stabilit mai înainte θ avem, pentru media statisticii suficiente din acest caz

$$E\{x[n]\} = \frac{1}{\lambda} = \theta$$

$$E\left\{\sum_{n=0}^{N-1} x[n]\right\} = \sum_{n=0}^{N-1} E\{x[n]\} = N\theta \neq \theta$$

Rezultă din

$$E\left\{\frac{1}{N} \sum_{n=0}^{N-1} x[n]\right\} = \theta$$

estimatorul MVU pentru θ ca fiind media eșantion a datelor

$$\hat{\theta} = \frac{1}{N} \sum_{n=0}^{N-1} x[n] = \frac{1}{\hat{\lambda}}$$

Ținând seama de relația dintre θ și λ rezultă estimatorul căutat

$$\hat{\lambda} = \frac{1}{\bar{x}} = \frac{1}{\frac{1}{N} \sum_{n=0}^{N-1} x[n]}$$

19

Cazul existenței unui grup de statistici mutual suficiente

Reluăm problema estimării fazei inițiale Φ a unui semnal sinusoidal, afectat de un zgomot alb, gaussian

$$x[n] = A \cos(2\pi f_0 n + \Phi) + w[n]; \quad n = 0, 1, \dots, N-1; \quad w[n] \sim \mathcal{N}(0, \sigma^2)$$

pentru care există două statistici mutual suficiente. Pentru a determina forma estimatorului MVU trebuie să găsim funcția $g(\cdot, \cdot)$, astfel încât media ei statistică să fie chiar Φ

$$E\{g(T_1(\mathbf{x}), T_2(\mathbf{x}))\} = \Phi$$

Prima statistică suficientă poate fi aproximată cu relația

$$\begin{aligned} T_1(\mathbf{x}) &= \sum_{n=0}^{N-1} x[n] \cos 2\pi f_0 n \\ &= \sum_{n=0}^{N-1} [A \cos(2\pi f_0 n + \Phi) + w[n]] \cos 2\pi f_0 n \\ &= \sum_{n=0}^{N-1} \frac{A}{2} [\cos \Phi + \cos(4\pi f_0 n + \Phi)] + \sum_{n=0}^{N-1} w[n] \cos 2\pi f_0 n \\ &\cong \frac{NA}{2} \cos \Phi + \sum_{n=0}^{N-1} w[n] \cos 2\pi f_0 n \end{aligned}$$

20

care, pentru un raport semnal/zgomot (SNR) mare devine

$$T_1(\mathbf{x}) \cong \frac{NA}{2} \cos \Phi; \quad \text{SNR} = \frac{A^2}{2\sigma^2} \gg 1$$

O formă asemănătoare se poate stabili și pentru a doua statistică suficientă

$$T_2(\mathbf{x}) \cong -\frac{NA}{2} \sin \Phi; \quad \text{SNR} = \frac{A^2}{2\sigma^2} \gg 1$$

Raportul celor două statistici suficiente

$$\frac{T_1(\mathbf{x})}{T_2(\mathbf{x})} \cong -\text{tg} \Phi; \quad \text{SNR} = \frac{A^2}{2\sigma^2} \gg 1$$

ne permite să estimăm faza inițială a sinusoidei dar numai în cazul unui raport semnal/zgomot mare

$$\hat{\Phi} \cong -\text{arctg} \frac{T_2(\mathbf{x})}{T_1(\mathbf{x})}; \quad \text{SNR} = \frac{A^2}{2\sigma^2} \gg 1$$

21

Mediile statistice ale celor două statistici suficiente se determină ușor

$$\begin{aligned} E\{T_1(\mathbf{x})\} &= \frac{NA}{2} \cos \Phi + \sum_{n=0}^{N-1} \underbrace{E\{w[n]\}}_{=0} \cos 2\pi f_0 n \\ &= \frac{NA}{2} \cos \Phi \end{aligned}$$

$$E\{T_2(\mathbf{x})\} = -\frac{NA}{2} \sin \Phi$$

Media statistică a estimatorului fazei, chiar și pentru SNR mare, nu este Φ !

$$E\left\{-\text{arctg} \frac{T_2(\mathbf{x})}{T_1(\mathbf{x})}\right\} \neq -\text{arctg} \frac{E\{T_2(\mathbf{x})\}}{E\{T_1(\mathbf{x})\}} = -\text{arctg} \frac{-\frac{NA}{2} \sin \Phi}{\frac{NA}{2} \cos \Phi} = \Phi$$

Prin urmare

$$E\{\hat{\Phi}\} \neq \Phi$$

și deci estimatorul introdus pentru fază nu este un estimator MVU. Acest fapt se datorează neliniarității funcției $\text{arctg}(x)$

22

O extindere pentru cazul mai multor parametri necunoscuți, organizați sub forma unui vector, θ de dimensiune p .

Cele r statistici se organizează sub forma unui vector

$$\mathbf{T}_{r \times 1}(\mathbf{x}) = [T_1(\mathbf{x}) \quad T_2(\mathbf{x}) \quad \dots \quad T_r(\mathbf{x})]^T$$

Teorema Neyman-Fisher se extinde și la acest caz. Astfel, dacă poate avea loc o factorizare de forma

$$p(\mathbf{x}; \theta) = g(\mathbf{T}(\mathbf{x}); \theta) h(\mathbf{x})$$

atunci vectorul $\mathbf{T}(\mathbf{x})$ este un vector statistică suficientă pentru parametrul vector θ

Reciproca teoremei este adevărată

23

Exemple de determinare a estimatorului MVU pentru cazul parametrului vector θ

Vom considera un exemplu de semnal sinusoidal, afectat de un zgomot alb, gaussian

$$x[n] = A \cos 2\pi f_0 n + w[n]; \quad n = 0, 1, \dots, N-1; \quad w[n] \sim \mathcal{N}(0, \sigma^2)$$

Parametrii necunoscuți sunt amplitudinea, frecvența (digitală) și puterea zgomotului. Cei trei parametri necunoscuți sunt componente ale vectorului

$$\theta = [A \quad f_0 \quad \sigma^2]^T$$

Vectorul de date are componentele gaussiene, IID, cu densitatea de repartiție

$$p(\mathbf{x}; \theta) = \frac{1}{(2\pi\sigma^2)^{\frac{N}{2}}} \exp \left\{ -\frac{1}{2\sigma^2} \sum_{n=0}^{N-1} (x[n] - A \cos 2\pi f_0 n)^2 \right\}$$

24

Paranteza de la exponent se dezvoltă sub forma

$$\sum_{n=0}^{N-1} (x[n] - A \cos 2\pi f_0 n)^2 = \sum_{n=0}^{N-1} x^2[n] - 2A \sum_{n=0}^{N-1} x[n] \cos 2\pi f_0 n + A^2 \sum_{n=0}^{N-1} \cos^2 2\pi f_0 n$$

Pentru determinarea estimatorului pentru vectorul θ vom proceda la o abordare graduală

1) Dacă se cunoaște frecvența digitală, vectorul parametrilor necunoscuți va avea doar două componente

$$\theta_1 = [A \quad \sigma^2]^T$$

Vom proceda la factorizarea Neyman-Fisher pentru acest caz de complexitate mai redusă. Avem

$$p(\mathbf{x}; \theta_1) = \frac{1}{(2\pi\sigma^2)^{\frac{N}{2}}} \exp \left\{ -\frac{1}{2\sigma^2} \left[\sum_{n=0}^{N-1} x^2[n] - 2A \sum_{n=0}^{N-1} x[n] \cos 2\pi f_0 n + A^2 \sum_{n=0}^{N-1} \cos^2 2\pi f_0 n \right] \right\} \cdot \underbrace{\frac{1}{h(\mathbf{x})}}_{g(\mathbf{T}(\mathbf{x}), \theta)}$$

25

Vectorul statisticilor suficiente are două componente (scalare)

$$\mathbf{T}(\mathbf{x}) = \begin{bmatrix} T_1(\mathbf{x}) = \sum_{n=0}^{N-1} x[n] \cos 2\pi f_0 n \\ T_2(\mathbf{x}) = \sum_{n=0}^{N-1} x^2[n] \end{bmatrix}$$

1.1) Vom arăta că prima componentă scalară a vectorului $\mathbf{T}(\mathbf{x})$ corespunde unei statistici suficiente pentru amplitudinea A . Într-adevăr, dacă numai parametrul A ar fi necunoscut, factorizarea Neyman-Fisher ar lua forma

$$p(\mathbf{x}; A) = \frac{1}{(2\pi\sigma^2)^{\frac{N}{2}}} \exp \left\{ -\frac{1}{2\sigma^2} \left[A^2 \sum_{n=0}^{N-1} \cos^2 2\pi f_0 n - 2A \sum_{n=0}^{N-1} x[n] \cos 2\pi f_0 n \right] \right\} \cdot \underbrace{\frac{1}{h(\mathbf{x})}}_{g(T'(\mathbf{x}); A)}$$

$$\exp \left\{ -\frac{1}{2\sigma^2} \sum_{n=0}^{N-1} x^2[n] \right\}$$

Din care rezultă că statistica suficientă pentru estimarea amplitudinii A este

$$T'(\mathbf{x}) = \sum_{n=0}^{N-1} x[n] \cos 2\pi f_0 n \equiv T_1(\mathbf{x})$$

statistică suficientă identică cu prima componentă a vectorului $\mathbf{T}(\mathbf{x})$ ²⁶

Va trebui să găsim acea funcție care face ca statistica suficientă pentru A să fie nedeplasată, adică să aibă media statistică egală cu A. Avem

$$E\{x[n]\} = A \cos 2\pi f_0 n$$

și apoi

$$\begin{aligned} E\{T'(\mathbf{x})\} &= \sum_{n=0}^{N-1} E\{x[n]\} \cos 2\pi f_0 n \\ &= A \sum_{n=0}^{N-1} \cos^2 2\pi f_0 n \end{aligned}$$

din care rezultă că

$$E\left\{\frac{T'(\mathbf{x})}{\sum_{n=0}^{N-1} \cos^2 2\pi f_0 n}\right\} = A$$

Prin urmare, estimatorul MVU pentru amplitudinea A a sinusoidei este

$$\hat{A} = \frac{\sum_{n=0}^{N-1} x[n] \cos 2\pi f_0 n}{\sum_{n=0}^{N-1} \cos^2 2\pi f_0 n}$$

27

1.2) Dacă nici A nici dispersia nu se cunosc, vom lua în considerare și a doua componentă a vectorului statistică suficientă, $\mathbf{T}(\mathbf{x})$. Media ei statistică este

$$\begin{aligned} E\{T_2(\mathbf{x})\} &= E\left\{\sum_{n=0}^{N-1} x^2[n]\right\} = \sum_{n=0}^{N-1} E\{x^2[n]\} \\ &= \sum_{n=0}^{N-1} E\left\{\left(A \cos 2\pi f_0 n + w[n]\right)^2\right\} \\ &= \sum_{n=0}^{N-1} E\left\{A^2 \cos^2 2\pi f_0 n + 2Aw[n] \cos 2\pi f_0 n + w^2[n]\right\} \\ &= \sum_{n=0}^{N-1} \left[A^2 \cos^2 2\pi f_0 n + 2AE\{w[n]\} \cos 2\pi f_0 n + E\{w^2[n]\}\right] \\ &= A^2 \sum_{n=0}^{N-1} \cos^2 2\pi f_0 n + N\sigma^2 \neq \sigma^2 \end{aligned}$$

28

Avem deja estimatorul pentru amplitudinea A. Momentul de ordinul doi al acestuia se poate calcula aplicând definiția

$$E\{\hat{A}^2\} = \frac{1}{\left(\sum_{n=0}^{N-1} \cos^2 2\pi f_0 n\right)^2} E\left\{\sum_{n=0}^{N-1} \sum_{m=0}^{N-1} x[n]x[m] \cos 2\pi f_0 n \cos 2\pi f_0 m\right\}$$

$$= \frac{1}{\left(\sum_{n=0}^{N-1} \cos^2 2\pi f_0 n\right)^2} \sum_{n=0}^{N-1} \sum_{m=0}^{N-1} E\{x[n]x[m]\} \cos 2\pi f_0 n \cos 2\pi f_0 m$$

În formulă intră corelația datelor, care se determină tot prin calculul direct, ținând seama de faptul că zgomotul $w[n]$ este alb și are deci eșantioanele necorelate.

Avem

$$E\{x[n]x[m]\} = E\{(A \cos 2\pi f_0 n + w[n])(A \cos 2\pi f_0 m + w[m])\}$$

$$= E\{A^2 \cos 2\pi f_0 n \cos 2\pi f_0 m\} + E\{Aw[m] \cos 2\pi f_0 n\}$$

$$+ E\{Aw[n] \cos 2\pi f_0 m\} + E\{w[n]w[m]\}$$

$$= A^2 \cos 2\pi f_0 n \cos 2\pi f_0 m + \sigma^2 \delta_{n,m}$$

29

care substituită în expresia momentului de ordinul doi al estimatorului amplitudinii conduce la relația

$$E\{\hat{A}^2\} = \frac{1}{\left(\sum_{n=0}^{N-1} \cos^2 2\pi f_0 n\right)^2}$$

$$\cdot \sum_{n=0}^{N-1} \sum_{m=0}^{N-1} (A^2 \cos 2\pi f_0 n \cos 2\pi f_0 m + \sigma^2 \delta_{n,m}) \cos 2\pi f_0 n \cos 2\pi f_0 m$$

$$= \frac{1}{\left(\sum_{n=0}^{N-1} \cos^2 2\pi f_0 n\right)^2} \left(A^2 \sum_{n=0}^{N-1} \cos^2 2\pi f_0 n \sum_{m=0}^{N-1} \cos^2 2\pi f_0 m + \sigma^2 \sum_{n=0}^{N-1} \cos^2 2\pi f_0 n \right)$$

$$= \frac{1}{\left(\sum_{n=0}^{N-1} \cos^2 2\pi f_0 n\right)^2} \left[A^2 \left(\sum_{n=0}^{N-1} \cos^2 2\pi f_0 n \right)^2 + \sigma^2 \sum_{n=0}^{N-1} \cos^2 2\pi f_0 n \right]$$

$$= A^2 + \frac{\sigma^2}{\sum_{n=0}^{N-1} \cos^2 2\pi f_0 n}$$

30

Din ultima relație rezultă că

$$A^2 = E\{\hat{A}^2\} - \frac{\sigma^2}{\sum_{n=0}^{N-1} \cos^2 2\pi f_0 n}$$

formă ce se substituie în $E\{T_2\}$. Obținem

$$\begin{aligned} E\{T_2(\mathbf{x})\} &= \left(E\{\hat{A}^2\} - \frac{\sigma^2}{\sum_{n=0}^{N-1} \cos^2 2\pi f_0 n} \right) \sum_{n=0}^{N-1} \cos^2 2\pi f_0 n + N\sigma^2 \\ &= E\left\{ \hat{A}^2 \sum_{n=0}^{N-1} \cos^2 2\pi f_0 n + (N-1)\sigma^2 \right\} \end{aligned}$$

Comparând cei doi membri ai egalității de mai sus, rezultă că

$$T_2(\mathbf{x}) = \hat{A}^2 \sum_{n=0}^{N-1} \cos^2 2\pi f_0 n + (N-1)\sigma^2$$

31

Se explicitează, din această ultimă relație dispersia

$$T_2'(\mathbf{x}) = \frac{1}{N-1} \left\{ T_2(\mathbf{x}) - \hat{A}^2 \sum_{n=0}^{N-1} \cos^2 2\pi f_0 n \right\} = \sigma^2$$

am definit astfel un estimator MVU pentru dispersie, deoarece

$$E\{T_2'(\mathbf{x})\} = \sigma^2$$

Grupăm cei doi estimatori MVU găsiți, pentru amplitudine și pentru dispersie, sub forma unui vector estimator cu două componente

$$\hat{\boldsymbol{\theta}}_1 = \begin{bmatrix} \hat{A} \\ \hat{\sigma}^2 \end{bmatrix} = \begin{bmatrix} \frac{\sum_{n=0}^{N-1} x[n] \cos 2\pi f_0 n}{\left(\sum_{n=0}^{N-1} \cos^2 2\pi f_0 n \right)^2} \\ \frac{1}{N-1} \left[\sum_{n=0}^{N-1} x^2[n] - \hat{A}^2 \sum_{n=0}^{N-1} \cos^2 2\pi f_0 n \right] \end{bmatrix}$$

32

Rezultatul obținut se poate aplica imediat, dacă se cunoaște frecvența digitală a sinusoidei. Pentru frecvență nulă se ajunge la modelul de semnal componentă continuă, în care nu se cunosc amplitudinea acesteia și puterea zgomotului alb

$$x[n] = A + w[n]; \quad n = 0, 1, \dots, N-1; \quad w[n] \sim \mathcal{N}(0, \sigma^2)$$

Vectorul parametrilor necunoscuți este

$$\boldsymbol{\theta}_1 = [A \quad \sigma^2]^T$$

Estimatorul vector MVU se obține din estimatorul stabilit anterior, punând valoarea zero pentru frecvența digitală. Obținem estimatorul vector

$$\hat{\boldsymbol{\theta}}_1 = \begin{bmatrix} \hat{A} \\ \hat{\sigma}^2 \end{bmatrix} = \begin{bmatrix} \frac{1}{N} \sum_{n=0}^{N-1} x[n] \\ \frac{1}{N-1} \left(\sum_{n=0}^{N-1} x^2[n] - N\hat{A}^2 \right) \end{bmatrix}$$

33

în care

$$\hat{A} = \frac{1}{N} \sum_{n=0}^{N-1} x[n] = \bar{x}$$

și

$$\hat{\sigma}^2 = \frac{1}{N-1} \left(\sum_{n=0}^{N-1} x^2[n] - N\bar{x}^2 \right)$$

Această ultimă relație se mai poate modifica

$$\begin{aligned} \frac{1}{N-1} \sum_{n=0}^{N-1} (x[n] - \bar{x})^2 &= \frac{1}{N-1} \left(\sum_{n=0}^{N-1} x^2[n] - 2\bar{x} \sum_{n=0}^{N-1} x[n] + \sum_{n=0}^{N-1} \bar{x}^2 \right) \\ &= \frac{1}{N-1} \left(\sum_{n=0}^{N-1} x^2[n] - 2\bar{x}N\bar{x} + N\bar{x}^2 \right) \\ &= \frac{1}{N-1} \left(\sum_{n=0}^{N-1} x^2[n] - N\bar{x}^2 \right) \end{aligned}$$

34

Estimatorul dispersiei poate fi utilizat și sub forma, în care se calculează puterea fluctuației în jurul mediei estimate, medierea fiind făcută prin împărțire cu N-1 și nu cu N

$$\hat{\sigma}^2 = \frac{1}{N-1} \sum_{n=0}^{N-1} (x[n] - \bar{x})^2$$

O formă posibilă a vectorului estimator, larg utilizată în calculele statistice este

$$\begin{bmatrix} \hat{A} \\ \hat{\sigma}^2 \end{bmatrix} = \begin{bmatrix} \bar{x} \\ \frac{1}{N-1} \sum_{n=0}^{N-1} (x[n] - \bar{x})^2 \end{bmatrix}$$

Semnal de tip componentă continuă necunoscută, afectată de un zgomot alb, gaussian cu puterea necunoscută. Reluare

$$x[n] = A + w[n]; \quad n = 0, 1, \dots, N-1; \quad w[n] \sim \mathcal{N}(0, \sigma^2)$$

Am stabilit în modelul semnalului sinusoidal cu frecvența cunoscută că vectorul statistică suficientă este

$$\mathbf{T}(\mathbf{x}) = \begin{bmatrix} T_1(\mathbf{x}) = \sum_{n=0}^{N-1} x[n] \cos 2\pi f_0 n \\ T_2(\mathbf{x}) = \sum_{n=0}^{N-1} x^2[n] \end{bmatrix}$$

35

Anulând frecvența digitală se obține vectorul statistică suficientă pentru cazul componentă continuă afectată de zgomot alb, gaussian

$$\mathbf{T}(\mathbf{x}) = \begin{bmatrix} \sum_{n=0}^{N-1} x[n] \\ \sum_{n=0}^{N-1} x^2[n] \end{bmatrix}$$

Mediind vectorul $\mathbf{T}(\mathbf{x})$ nu se obține vectorul θ adică componentele estimatorului nu sunt de tip MVU

$$E\{\mathbf{T}(\mathbf{x})\} = \begin{bmatrix} \sum_{n=0}^{N-1} E\{x[n]\} \\ \sum_{n=0}^{N-1} E\{x^2[n]\} \end{bmatrix} = \begin{bmatrix} \sum_{n=0}^{N-1} A \\ \sum_{n=0}^{N-1} (A^2 + \sigma^2) \end{bmatrix} = \begin{bmatrix} NA \\ N(A^2 + \sigma^2) \end{bmatrix} \neq \begin{bmatrix} A \\ \sigma^2 \end{bmatrix}$$

Prin mediere se stabilesc relațiile

$$E\left\{\frac{1}{N} T_1(\mathbf{x})\right\} = E\{\bar{x}\} = A$$

$$\frac{1}{N} E\{T_2(\mathbf{x})\} = A^2 + \sigma^2$$

36

Din a doua relație rezultă

$$E \left\{ \frac{1}{N} T_2(\mathbf{x}) - A^2 \right\} = \sigma^2$$

relație ce sugerează că un estimator “bun” pentru dispersie ar putea fi expresia

$$\frac{1}{N} T_2(\mathbf{x}) - \bar{x}^2$$

Media ei se determină cu

$$E \left\{ \frac{1}{N} T_2(\mathbf{x}) - \bar{x}^2 \right\} = E \left\{ \frac{1}{N} T_2(\mathbf{x}) \right\} - E \left\{ \bar{x}^2 \right\}$$

Media eșantion are repartiția normală

$$\bar{x} \sim N \left(A, \frac{\sigma^2}{N} \right)$$

Drept urmare momentul de ordinul doi al mediei eșantion este

$$E \left\{ \bar{x}^2 \right\} = A^2 + \frac{\sigma^2}{N}$$

37

Media estimatorului considerat a fi “bun” pentru dispersie devine deci

$$E \left\{ \frac{1}{N} T_2(\mathbf{x}) - \bar{x}^2 \right\} = \frac{1}{N} N \left(A^2 + \sigma^2 \right) - \left(A^2 + \frac{\sigma^2}{N} \right) = \frac{N-1}{N} \sigma^2$$

Din aceasta deducem că

$$E \left\{ \frac{1}{N-1} T_2(\mathbf{x}) - \frac{N}{N-1} \bar{x}^2 \right\} = \sigma^2$$

ceea ce înseamnă că estimatorul cu adevărat bun pentru dispersie este

$$\begin{aligned} \hat{\sigma}^2 &= \frac{1}{N-1} \sum_{n=0}^{N-1} x^2[n] - \frac{N}{N-1} \bar{x}^2 \\ &= \frac{1}{N-1} \sum_{n=0}^{N-1} (x[n] - \bar{x})^2 \end{aligned}$$

estimator găsit și mai înainte

38

Se specifică în literatura de specialitate că dacă datele $x[n]$ au o repartiție normală

$$x[n] \sim N(\mu, \sigma^2)$$

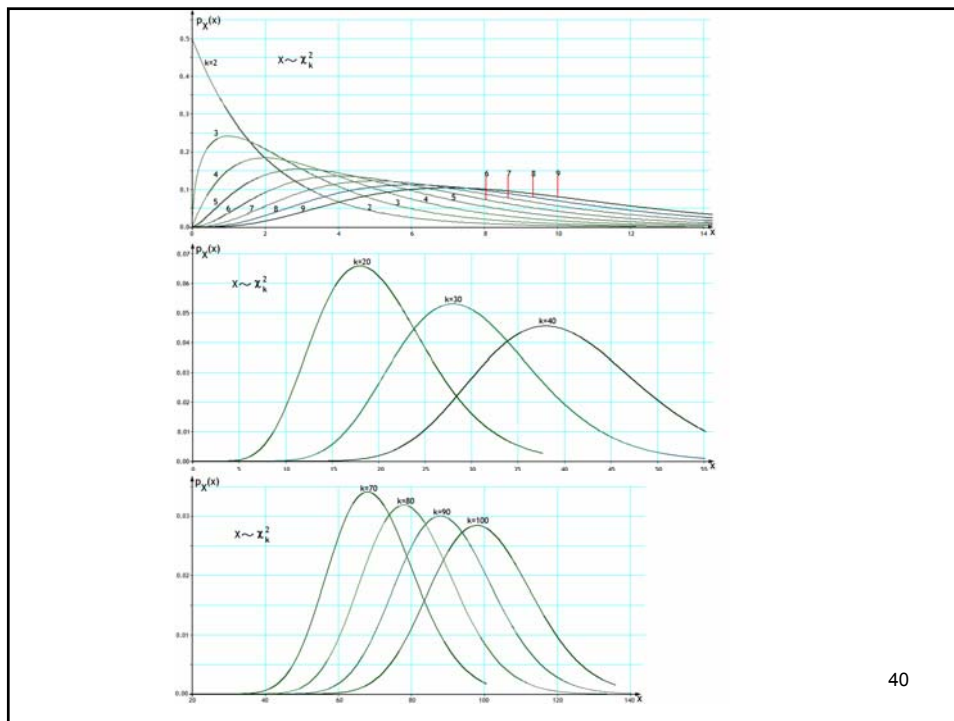
atunci variabila aleatoare

$$y = \sum_{n=0}^{N-1} \left(\frac{x[n] - \bar{x}}{\sigma} \right)^2 \sim \chi_{N-1}^2$$

are o repartiție "hi-pătrat" cu $k=N-1$ grade de libertate

În figura următoare se dau curbele densității de probabilitate pentru repartiții hi-pătrat cu $k=2,3..9$ grade de libertate și pentru repartiții hi-pătrat cu $k=70, 80, 90$ și 100 grade de libertate

39



40

Pe măsură ce crește numărul gradelor de libertate ale unei repartiții hi-pătrat, aceasta se apropie tot mai mult de o repartiție normală, așa cum se poate vedea și din figură



Media variabilei alatoare cu repartiție hi-pătrat, cu k grade de libertate este k iar dispersia ei este 2k

Din expresia estimatorului dispersiei

$$\hat{\sigma}^2 = \frac{1}{N-1} \sum_{n=0}^{N-1} (x[n] - \bar{x})^2 = \frac{\sigma^2}{N-1} \sum_{n=0}^{N-1} \left(\frac{x[n] - \bar{x}}{\sigma} \right)^2$$

se poate deduce că

$$\frac{N-1}{\sigma^2} \hat{\sigma}^2 = \sum_{n=0}^{N-1} \left(\frac{x[n] - \bar{x}}{\sigma} \right)^2 \sim \chi_{N-1}^2$$

41

Se poate scrie, ținând seama că dispersia repartiției hi-pătrat este 2(N-1) că

$$Disp \left\{ \frac{N-1}{\sigma^2} \hat{\sigma}^2 \right\} = \frac{(N-1)^2}{\sigma^4} Disp \left\{ \hat{\sigma}^2 \right\} = \underbrace{2(N-1)}_{Disp \{ \chi_{N-1}^2 \}}$$

din care rezultă că dispersia estimatorului pentru puterea zgomotului alb gaussian este

$$Disp \left\{ \hat{\sigma}^2 \right\} = \frac{2\sigma^4}{N-1}$$

Eșantioanele de zgomot fiind necorelate elementele din matricea de covarianță a estimatorului sunt nule, cu excepția celor de pe diagonala principală, care sunt dispersiile celor doi estimatori, dispersii determinate de noi. Avem deci

$$\mathbf{C}_{\hat{\theta}} = \begin{bmatrix} \frac{\sigma^2}{N} & 0 \\ 0 & \frac{2\sigma^4}{N-1} \end{bmatrix}$$

Anterior am stabilit că CRLB pentru puterea zgomotului alb este

$$CRLB_{\sigma^2} = \frac{2\sigma^4}{N}$$

42

Dispersia estimatorului pe care l-am stabilit este ușor mai mare decât CRLB. El este deci neeficient. Asimptotic însă poate fi considerat eficient

$$\frac{2\sigma^4}{N-1} > \frac{2\sigma^4}{N} = CRLB_{\sigma^2}$$

Trebuie să menționăm că se putea găsi un vector statistică suficientă și direct, pornind de la repartiția

$$p(\mathbf{x}; \theta) = \frac{1}{(2\pi\sigma^2)^{N/2}} \exp\left\{-\frac{1}{2\sigma^2} \sum_{n=0}^{N-1} (x[n] - A)^2\right\}$$

Suma de la exponent se dezvoltă după cum urmează

$$\begin{aligned} \sum_{n=0}^{N-1} (x[n] - A)^2 &= \sum_{n=0}^{N-1} [(x[n] - \bar{x}) + (\bar{x} - A)]^2 \\ &= \sum_{n=0}^{N-1} (x[n] - \bar{x})^2 - 2(\bar{x} - A) \sum_{n=0}^{N-1} (x[n] - \bar{x}) + \sum_{n=0}^{N-1} (\bar{x} - A)^2 \\ &= \sum_{n=0}^{N-1} (x[n] - \bar{x})^2 - 2(\bar{x} - A)(N\bar{x} - N\bar{x}) + \sum_{n=0}^{N-1} (\bar{x} - A)^2 \\ &= \sum_{n=0}^{N-1} (x[n] - \bar{x})^2 + N(\bar{x} - A)^2 \end{aligned}$$

43

Substituim această dezvoltare în expresia repartiției gaussiene și obținem factorizarea

$$p(\mathbf{x}; \theta) = \underbrace{\frac{1}{(2\pi\sigma^2)^{N/2}} \exp\left\{-\frac{1}{2\sigma^2} \left[\sum_{n=0}^{N-1} (x[n] - \bar{x})^2 + N(\bar{x} - A)^2 \right]\right\}}_{g(T(\mathbf{x}), \theta)} \cdot \frac{1}{h(\mathbf{x})}$$

din care se deduc imediat cele două statistici suficiente, pentru media A și pentru dispersie. Vectorul statisticilor suficiente este

$$T'(\mathbf{x}) = \begin{bmatrix} \bar{x} \\ \sum_{n=0}^{N-1} (x[n] - \bar{x})^2 \end{bmatrix}$$

44