

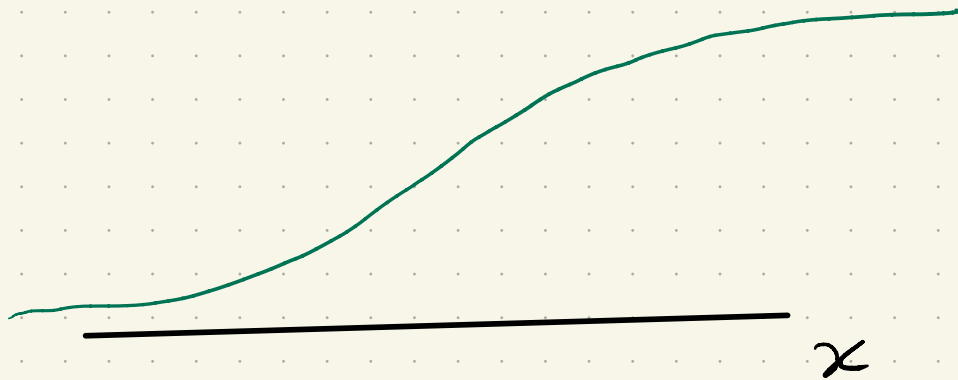
Introducción a la Teoría de la Información

Entropía diferencial.

Facultad de Ingeniería, UdelaR

Año ~~2020~~ 2024

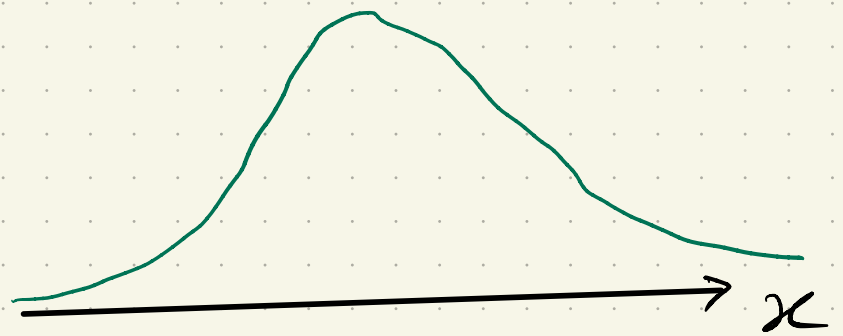
Distribución



$$F(x) = P(X \leq x)$$

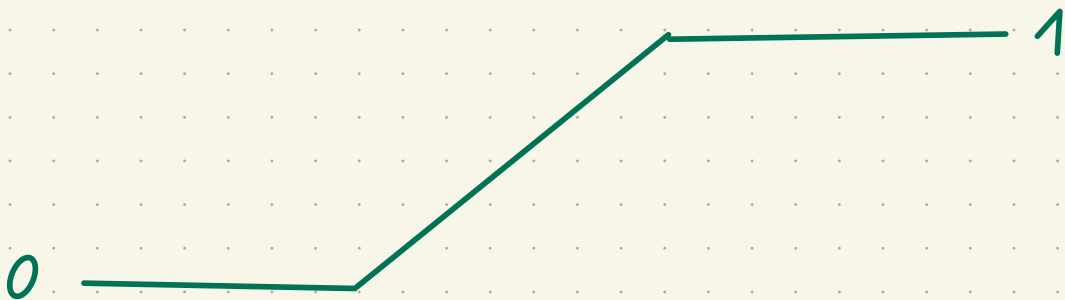
$F(x)$ creciente

Densidad



$$f(x) = F'(x)$$

$$\int_{-\infty}^{+\infty} f(x) dx = 1$$



Ej: uniforme



$$\sum p_i \log p_i$$

Definición (Entropía diferencial)

Para una variable aleatoria X con densidad de probabilidad $f(x)$ de soporte S , se define como:

$$h(X) = - \int_S f(x) \log f(x) dx \quad E_f[-\log f(x)]$$

- Para variables aleatorias continuas
- Similar a la entropía discreta
- Algunas diferencias importantes, hay que usarla con cuidado.
- Puede *no existir* para ciertas $f(x)$.



Ejemplo (Distribución uniforme en $(0, a)$)

$$h(X) = - \int_0^a \frac{1}{a} \log \frac{1}{a} dx = \log a$$

- Para $a < 1$, $h(X) = \log a$ es *negativa*

Ejemplos

una variable

Ejemplo (Distribución Normal $X \sim \phi(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-x^2/2\sigma^2}$)

$$h(X) = - \int \phi(x) \log \phi(x) = \frac{1}{2} \log(2\pi e \sigma^2) \text{ bits}$$

media $\mu = 0$
varianza σ

$$\mu = E_f(X)$$
$$\sigma^2 = E_f[(X - \mu)^2]$$

$$\int x f(x) dx$$

mejor usar \ln .

$$f(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{x^2}{2\sigma^2}}$$

$$\ln f(x) = \ln \frac{1}{\sqrt{2\pi}\sigma} - \frac{x^2}{2\sigma^2}$$

$$\log a^n = n \log a$$

$$E[x^2] = \sigma^2$$

$$-\int_{-\infty}^{+\infty} f(x) \ln f(x) dx = \frac{1}{2} \ln(2\pi\sigma^2) + \frac{1}{2} =$$

$$= \frac{1}{2} \ln(2\pi e \sigma^2) \text{ nats}$$

$$= \frac{1}{2} \log(2\pi e \sigma^2) \text{ bits.}$$

(cambio de base)

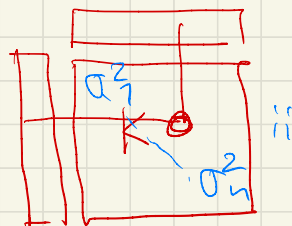
Multi-dimensional

$$f(x) = \frac{1}{\sqrt{2\pi} |K|} e^{-x^t K^{-1} x}$$

$$\text{Covarianz } K = E_f[x x^t]$$

$$K_{ij} = E[x_i x_j]$$

$$h(\bar{x}) = \frac{1}{2} \log((2\pi e)^n |K|)$$



$$\vec{X} = (x_1, x_2, \dots, x_n) \quad \text{Gaussverteilung}$$

$$E_f(\vec{x}) = \vec{\mu}$$

si son iid, $f(x_1, x_2, \dots, x_n) = f(x_1) f(x_2) \dots f(x_n)$

Teorema (AEP)

$X_1 \dots X_n$ variables *i.i.d.*, $X_i \sim f(x)$

$$-\frac{1}{n} \log f(X_1, X_2, \dots, X_n) \rightarrow E[-\log f(X)] = h(X)$$

Definición (Conjunto típico)

$$A_\epsilon^{(n)} = \left\{ (x_1, x_2, \dots, x_n) \in S^n : \left| -\frac{1}{n} \log f(x_1, x_2, \dots, x_n) - h(X) \right| \leq \epsilon \right\}$$

$$-\frac{1}{n} \log f(x_1, x_2, \dots, x_n) = -\frac{1}{n} \sum_i \log f(x_i) \rightarrow E_f[-\log f(x)] = h(x)$$

Propiedades

- $P\left(A_\epsilon^{(n)}\right) > 1 - \epsilon$ para n suf. grande
- $\text{Vol}\left(A_\epsilon^{(n)}\right) \leq 2^{n(h(X)+\epsilon)}$ para **todo** n
- $\text{Vol}\left(A_\epsilon^{(n)}\right) \geq (1 - \epsilon)2^{n(h(X)-\epsilon)}$ para n suf. grande

Las pruebas son análogas a las demostradas para el conjunto típico de una variable aleatoria discreta.

$$h(x) - \varepsilon \leq -\frac{1}{n} \log f(x_1 \dots x_n) \leq h(x) + \varepsilon$$

$$-h(x) + \varepsilon \geq \frac{1}{n} \log f(x_1 \dots x_n) \geq -h(x) - \varepsilon$$

$$2^{-n(h(x) + \varepsilon)} \geq f(x_1 \dots x_n) \geq 2^{-n(h(x) - \varepsilon)}$$

V.A. discrete - eodind

V.A. continua - Volumen

Def: $\text{Vol } A = \iiint_A dx_1 dx_2 \dots dx_n$

$$1 = \iint_S f(x_1, \dots, x_n) dx_1 \dots dx_n \geq \iint_{A_\varepsilon^{(n)}} f(x_1, \dots, x_n) dx_1 \dots dx_n \geq$$

$$\iint_{A_\varepsilon^{(n)}} 2^{-n(h(x) + \varepsilon)} dx_1 \dots dx_n = 2^{-n[h(x) + \varepsilon]} \text{Vol } A_\varepsilon^{(n)}$$

$$\text{Vol } A_\varepsilon^{(n)} \leq 2^{n[h(x) + \varepsilon]}$$

Algunas observaciones

?

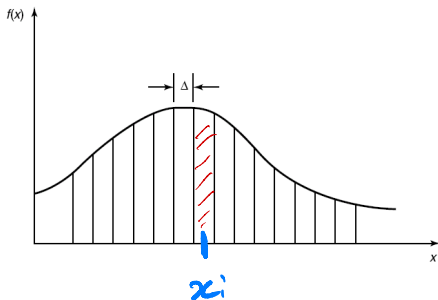
- El volumen del conjunto típico es aproximadamente 2^{nh} , o sea el volumen de un cubo n -dimensional de lado 2^h .
- Esto da una interpretación de $h(X)$ como el logaritmo del lado de un cubo, que tiene el mismo volumen que un subconjunto del soporte de X que concentra la mayoría de la probabilidad.

Entropía diferencial y Entropía discreta (1)

- Dividimos el soporte S de $X \sim f(x)$ en intervalos de ancho Δ
- Para f continua, por Teorema del Valor Medio podemos tomar en cada intervalo un punto x_i tal que

$$\Delta f(x_i) = \int_{i\Delta}^{(i+1)\Delta} f(x)dx$$

- Definimos $X^\Delta = x_i \Leftrightarrow i\Delta \leq X < (i+1)\Delta$



Entropía diferencial y Entropía discreta (2)

- $\Pr \{X^\Delta = x_i\} = \int_{i\Delta}^{(i+1)\Delta} f(x)dx = \Delta f(x_i)$
- X^Δ es una variable aleatoria discreta con entropía dada por:

$$\begin{aligned} H(X^\Delta) &= - \sum_{-\infty}^{\infty} p_i \log p_i \\ &= - \sum_{-\infty}^{\infty} \Delta f(x_i) \log (\Delta f(x_i)) \\ &= - \sum_{-\infty}^{\infty} \Delta f(x_i) \log (f(x_i)) - \sum_{-\infty}^{\infty} \Delta f(x_i) \log (\Delta) \\ &= - \sum_{-\infty}^{\infty} \Delta f(x_i) \log (f(x_i)) - \log \Delta \end{aligned}$$

Entropía diferencial y Entropía discreta (3)

Teorema ($h(X)$ vs $H(X^\Delta)$)

Si $f(x) \log f(x)$ es integrable Riemann,

$$H(X^\Delta) + \log \Delta \rightarrow h(f) = h(X) \text{ cuando } \Delta \rightarrow 0$$

- Si tomamos $\Delta = 2^{-n}$, vemos que $h(X) + n$ es aproximadamente la cantidad de bits promedio necesarios para codificar X con n bits de precisión.

Entropía diferencial y Entropía discreta (4)

Ejemplo ($X \sim U(0, a)$, $h(X) = \log a$)

- Si $X \sim U(0, 1)$, $h(X) = 0$ entonces $h(X) + n = n$ coincide con la cantidad necesaria de bits para describir X con n bits de precisión.
- $X \sim U(0, \frac{1}{8})$ Representando X en binario, los 3 bits a la derecha del punto quedan en 0, o sea que describir X con n bits de precisión requiere sólo $n - 3$ bits, que coincide con $h(X) + n = \log 1/8 + n$.

0,0001011...

Entropía diferencial conjunta y condicional

Definición (Entropía conjunta)

$$h(x, y) = - \int_S f(x, y) \log f(x, y) dx dy$$

soporte de los dos

Definición (Entropía condicional)

$$h(X|Y) = - \int f(x, y) \log f(x|y) dx dy$$

Regla de la cadena

De forma similar al caso discreto, obtenemos

$$h(X, Y) = h(X) + h(Y|X)$$

$h(Y) + h(X|Y)$

Entropía relativa e Información mutua

Definición (Distancia de Kullback-Leibler)

$$D(f||g) = \int f \log \frac{f}{g}$$

Notar que D es finita sólo si el soporte de f está contenido en el de g .

Definición (Información mutua)

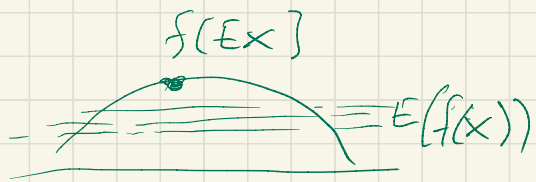
$$I(X; Y) = \int f(x, y) \log \frac{f(x, y)}{f(x)f(y)} dx dy$$

Por definición se obtiene, al igual que para el caso discreto, que

$$I(X; Y) = h(X) - h(X|Y) = h(Y) - h(Y|X)$$

$$I(X; Y) = D(f(x, y)||f(x)f(y))$$

$$D(f \parallel g) = \mathbb{E}_f \log \frac{f}{g} = \int f \log \frac{f}{g}$$



$$\rightarrow D(f \parallel g) = \int f \log \frac{g}{f} = \mathbb{E}_f \log \frac{g}{f} \leq \text{Jensen}$$

↑
concava

$$\leq \log \mathbb{E}_f \left(\frac{g}{f} \right) = \log \int \cancel{f} \frac{g}{\cancel{f}} = \log \int g = \log 1 = 0$$

La divergencia o distancia KL es positiva

Definición general de información mutua

Definición (Información mutua)

Sean \mathcal{X} , \mathcal{Y} los soportes de X, Y , respectivamente.

$$I(X; Y) = \sup_{\mathcal{P}, \mathcal{Q}} I([X]_{\mathcal{P}}; [Y]_{\mathcal{Q}})$$

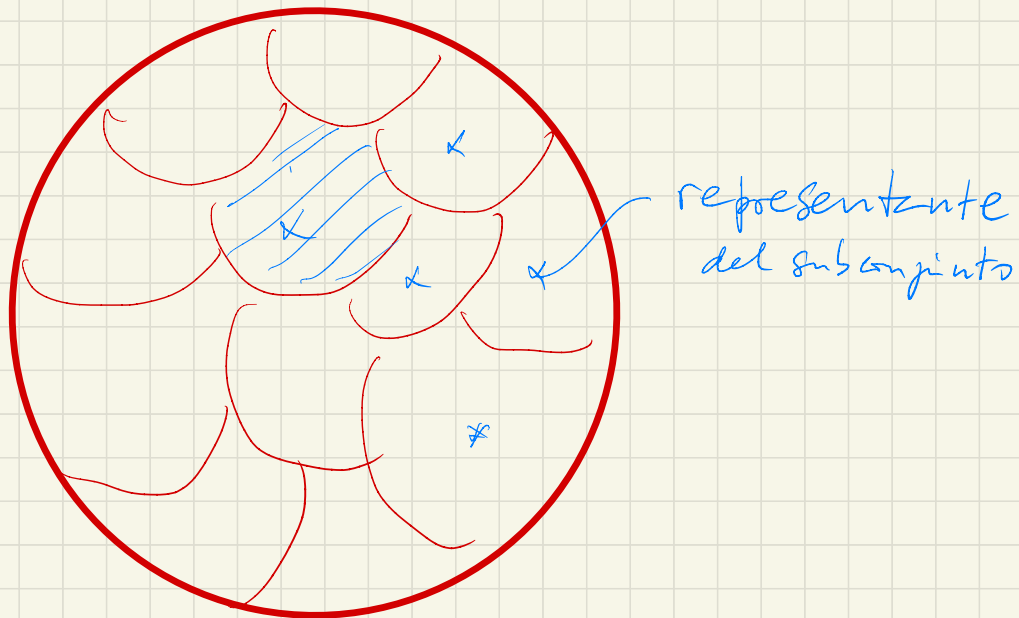
Donde $\mathcal{P} = \{P_i\}_{i=1 \dots N_p}$ y $\mathcal{Q} = \{Q_i\}_{i=1 \dots N_q}$ son particiones finitas de \mathcal{X} , \mathcal{Y} , respectivamente, y se cumple

$$\Pr \{[X]_{\mathcal{P}} = i\} = \Pr \{X \in P_i\}; \quad \Pr \{[Y]_{\mathcal{Q}} = i\} = \Pr \{Y \in Q_i\}.$$

- Esta definición aplica a variables discretas, continuas, y mezclas de uno y otro tipo.
- Si \mathcal{P}' refina \mathcal{P} , $[X]_{\mathcal{P}}$ es función de $[X]_{\mathcal{P}'}$, por lo tanto podemos escribir $[Y]_{\mathcal{Q}} \rightarrow [X]_{\mathcal{P}'} \rightarrow [X]_{\mathcal{P}}$ y por la desigualdad de procesamiento de datos

$$I([X]_{\mathcal{P}'}; [Y]_{\mathcal{Q}}) \geq I([X]_{\mathcal{P}}; [Y]_{\mathcal{Q}})$$

\mathcal{F}' refina \mathcal{F} si $\mathcal{F} \subset \mathcal{F}'$



Definición general de información mutua

- Cuando X, Y son discretas de alfabeto finito, para particiones \mathcal{P} y \mathcal{Q} suficientemente finas se cumple que

$$H([X]_{\mathcal{P}}) = H(X)$$

$$H([Y]_{\mathcal{Q}}) = H(Y)$$

$$H([X]_{\mathcal{P}}, [Y]_{\mathcal{Q}}) = H(X, Y)$$

$$\sup_{\mathcal{P}, \mathcal{Q}} I([X]_{\mathcal{P}}; [Y]_{\mathcal{Q}}) = H(X) + H(Y) - H(X, Y) = I(X; Y)$$

Definición general de información mutua

Variable y su discretización

- Cuando X continua, Y discreta de alfabeto finito, consideremos particiones \mathcal{P} y \mathcal{Q} a intervalos regulares de ancho Δ . Entonces,

$$\begin{aligned} I(X^\Delta; Y^\Delta) &= H(X^\Delta) + \log \Delta - (H(X^\Delta|Y^\Delta) + \log \Delta) \\ &= H(X^\Delta) + \log \Delta - \\ &\quad - \sum p(Y^\Delta = y_i) \left(H(X^\Delta|Y^\Delta = y_i) + \log \Delta \right) \end{aligned}$$

y_i

Para Δ suficientemente chico,

$$\begin{aligned} &= \underbrace{H(X^\Delta) + \log \Delta}_{\substack{\text{---} \\ \text{---}}} - \\ &\quad - \sum p(Y = y_i) \left(\underbrace{H(X^\Delta|Y = y_i) + \log \Delta}_{\substack{\text{---} \\ \text{---}}} \right) \\ \rightarrow & h(X) - \sum p(Y = y_i) \left(h(X|Y = y_i) \right) \\ &= \underbrace{h(X) - h(X|Y)} \end{aligned}$$

ruido de cuantificación

Definición general de información mutua

X e Y continuas

- Cuando X, Y continuas, también puede verse que $I(X^\Delta; Y^\Delta) \rightarrow h(X) + h(Y) - h(X, Y) = I(X; Y)$

Propiedades

1 $D(f||g) \geq 0$ con igualdad sii $f = g$ en casi todo punto

2 $I(X; Y) \geq 0$ con igualdad sii X, Y independientes $\rightarrow \cdot D$

3 $I(X; Y) = h(X) - h(X|Y) = h(Y) - h(Y|X)$

4 $h(X) \geq h(X|Y)$ con igualdad sii X, Y independientes

condicionar.....

5 $h(X_1, X_2, \dots, X_n) = \sum_{i=1}^n h(X_i | X_1, \dots, X_{i-1})$

cadena

6 $h(X_1, X_2, \dots, X_n) \leq \sum h(X_i)$

7 Si $X \rightarrow Y \rightarrow Z$, $I(X; Y) \geq I(X; Z)$

pro cesam. datos

8 $h(X + c) = h(X)$

9 $h(aX) = \log |a| + h(X)$

10 Para un vector \mathbf{X} , $h(A\mathbf{X}) = \log |\det A| + h(\mathbf{X})$

La distribución normal maximiza h para varianza dada

Teorema

Sea $\mathbf{X} \in \mathbb{R}^n$ con media cero y covarianza $\mathbf{K} = E[\mathbf{X}\mathbf{X}^t]$. Entonces $h(\mathbf{X}) \leq \frac{1}{2} \log [(2\pi e)^n |\mathbf{K}|]$ con igualdad sii \mathbf{X} tiene distribución normal.

Demostración.

Sea g densidad de probabilidad con media cero y covarianza \mathbf{K} .

Sea $\phi(\mathbf{x}) = \frac{1}{(\sqrt{2\pi})^n \sqrt{|\mathbf{K}|}} e^{-\frac{1}{2}\mathbf{x}^T \mathbf{K}^{-1} \mathbf{x}}$

$$\begin{aligned} D(g||\phi) &= \int g \log g/\phi = -h(g) - \int g \log \phi \\ &= -h(g) - E_g [\log \phi(\mathbf{X})] \\ &= -h(g) - E_\phi [\log \phi(\mathbf{X})] \text{ porque } E_\phi [X_i X_j] = E_g [X_i X_j] \\ &= -h(g) + h(\phi) \end{aligned}$$



La distribución normal maximiza h para varianza dada

$$E X^2 = P$$

Teorema

$$K = \sigma^2$$

Sea $\mathbf{X} \in \mathbb{R}^n$ con media cero y covarianza $\mathbf{K} = E[\mathbf{X}\mathbf{X}^t]$. Entonces $h(\mathbf{X}) \leq \frac{1}{2} \log[(2\pi e)^n |\mathbf{K}|]$ con igualdad sii \mathbf{X} tiene distribución normal.

Demostración.

Sea g densidad de probabilidad con media cero y covarianza \mathbf{K} . *en una var. $K \rightarrow \sigma^2$*

$$\text{Sea } \phi(\mathbf{x}) = \frac{1}{(\sqrt{2\pi})^n \sqrt{|\mathbf{K}|}} e^{-\frac{1}{2} \mathbf{x}^T \mathbf{K}^{-1} \mathbf{x}}$$

$$\int g \log g \, dx = E_g(\log g) = -h(g)$$

$$D(g||\phi) = \int g \log g / \phi = -h(g) - \int g \log \phi$$

$$= -h(g) - E_g[\log \phi(\mathbf{X})]$$

$$= -h(g) - E_\phi[\log \phi(\mathbf{X})] \text{ porque } E_\phi[X_i X_j] = E_g[X_i X_j] =$$

$$= -h(g) + h(\phi) \geq 0$$

log ϕ contiene $K^{-1}_{ij} X_i X_j$

$$\log e^A = \log e \ln e^A = \underline{\text{cte.}} \cdot A$$

Las páginas que siguen desarrollan y explican
la deducción de la entropía de la normal y
la demostración de que la normal maximiza
la entropía para una varianza dada.

Variance & entropia

$X \sim \mathcal{N}(0, \sigma)$ unidimensional

$$f(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-x^2/2\sigma^2}$$

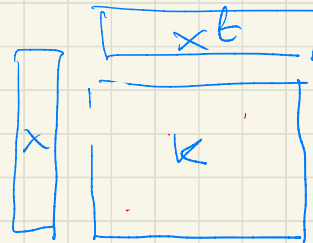
$$h(x) = \frac{1}{2} \log(2\pi e\sigma^2)$$

$\bar{X} \sim \mathcal{N}_n(\bar{0}, \kappa)$ multidimensional

covariance $\kappa = E[X X^t]$

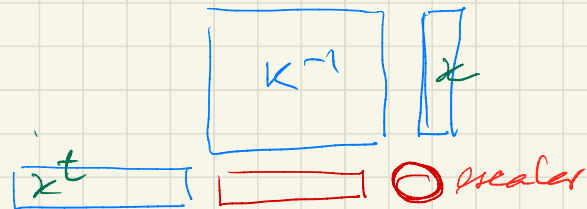
$$f(x) = \frac{1}{\sqrt{2\pi}^n \sqrt{|\kappa|}} e^{-\frac{x^t \kappa^{-1} x}{2}}$$

$\kappa = E$



$$\kappa_{ij} = E[x_i x_j]$$

$$h(\bar{X}) = \frac{1}{2} \log[(2\pi e)^n |\kappa|]$$



$$\kappa_{ij}^{-1} x_i x_j$$

Nota: κ es simétrica

$$\xrightarrow{1 \text{ dim}} x \quad f(x) \quad E_f(x) \quad E_f(x^2) \quad \sigma^2 = E x^2 - (E x)^2$$

$$n \text{ dim } \vec{x} \quad f(x_1 \dots x_n) \quad E(x) = \vec{\mu} \quad K = E_f[x x^t]$$

$$\rightarrow \phi(\vec{x}) = \frac{1}{\sqrt{(2\pi)^n |K|}} e^{-\frac{1}{2}(x^t K^{-1} x)} \quad \text{cov. } K$$

$$E_\phi[x_i x_j] = K_{ij} = E_f[x_i x_j] \quad \text{Matrix Covarianza}$$

Mayorante medida

Para aligerar, suponemos spq que restamos k medida

$$\bar{\mu} = E[\bar{x}]$$

$$\phi(\bar{x}) = \frac{1}{\sqrt{(2\pi)^n |k|}} e^{-\frac{1}{2} x^t k^{-1} x}$$

$$h(\phi) = \int_{-\infty}^{+\infty} \phi(x) \left[\underbrace{\frac{1}{2} x^t k^{-1} x}_{\textcircled{A}} - \underbrace{\frac{1}{2} \ln((2\pi)^n |k|)}_{= \frac{1}{2} \ln(2\pi)^n |k|} \right] dx$$

$$\textcircled{A} = \frac{1}{2} E \left[\sum_{i,j} x_i (k^{-1})_{ij} x_j \right] =$$

$$= \frac{1}{2} E \left[\sum_{i,j} x_j x_i (k^{-1})_{ij} \right] = \frac{1}{2} \sum_{i,j} E(x_j x_i) (k^{-1})_{ij} =$$

$$= \frac{1}{2} \sum_{i,j} k_{ji} (k^{-1})_{ij} = \frac{1}{2} \sum_j (k k^{-1})_{jj} = \frac{1}{2} n = \frac{1}{2} n \ln e$$

↑
traza de J

$$h(\phi) = \frac{1}{2} \ln \left[(2\pi e)^n |k| \right] \text{ nats}$$

Producto de matrices

En general, $C = A \times B$

$$C_{ij} = \sum_k A_{ik} B_{kj}$$

este caso

$$\begin{aligned} \sum_{ij} k_{ji} [k^{-1}]_{ij} &= \sum_j \sum_i k_{ji} (k^{-1})_{ij} = \\ &= \sum_j (k k^{-1})_{jj} = \sum_j I_{jj} = n \quad (\text{traza}) \end{aligned}$$

Las que siguen son cálculos auxiliares o
aplicados.

$$x \sim f(x)$$

$$\sigma^2 = E_f[x^2]$$

$$H(f) \leq H(\phi)$$

$$\phi(x) = \frac{1}{\sqrt{2\pi} \sigma} e^{-x^2/2\sigma^2}$$

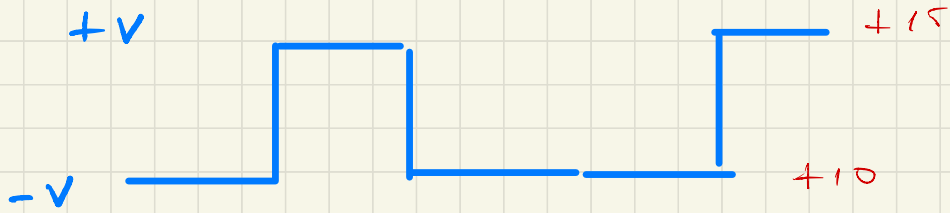
$$E_\phi[x^2] = \sigma^2$$

$$D(f \parallel \phi) \geq 0 \quad D(f \parallel \phi) = \int f \log\left(\frac{f}{\phi}\right) =$$

$$= \underbrace{\int f \log(f)}_{-H(f)} - \underbrace{\int f \log \phi}_{E_f(\log \phi)} = -H(f) + H(\phi) \geq 0$$

$$E_\phi(\log \phi)$$

$$-H(\phi)$$



X, Y v.a. $X \sim U[0, a]$ $Y \sim U[0, b]$

? soporte (X, Y)

? se desea cuantificar X con n bits precisión n
 Y m

1) Dibujar soporte y su cuantificación

2) Expresión para las v.d. X^Δ Y^Δ

3) Relacionar $H(X^\Delta, Y^\Delta)$ con $H(X, Y)$

4) Cuántos dígitos para (X, Y) con $n+m$ bits

Uniforme dimensional 2

$$Z \text{ v.a.c.} \sim \underbrace{\frac{e^{-z/\lambda}}{\lambda}}_{z \in \mathbb{R}^+}$$

$$0 \text{ meriti } A e^{-z/\lambda}$$

$$1 = \int_0^{+\infty} A e^{-z/\lambda} dz = A \left(-\lambda e^{-z/\lambda} \right) \Big|_0^{+\infty} \\ = \Delta \lambda \cdot 1 \Rightarrow A = \frac{1}{\lambda}$$

$$1) h(z)$$

$$2) W \text{ v.a.} \geq 0 \quad W \sim f \text{ di media } \lambda$$

$$\text{dem. } h(W) \leq h(z)$$

$$\text{media } E(z) = \int_0^{+\infty} z \frac{e^{-z/\lambda}}{\lambda} dz = -z e^{-z/\lambda} \Big|_0^{+\infty} - \int_0^{+\infty} 1 e^{-z/\lambda} dz \\ = 0 - \lambda e^{-z/\lambda} \Big|_0^{+\infty} = \lambda$$

$$h(z) \text{ nats} = \int \frac{e^{-z/\lambda}}{\lambda} \ln \left(\frac{1}{\lambda} e^{-z/\lambda} \right) dz = \int \frac{e^{-z/\lambda}}{\lambda} \left(-\frac{z}{\lambda} + \ln \lambda \right)$$

$$= \ln \lambda + \frac{1}{\lambda} \lambda = \ln e \lambda \text{ nats} \quad h(z) \text{ bits} = \log(\lambda e)$$

Exponential

Exponencial

comparer un w de densité $f(w)$

$$D(f \parallel g) \geq 0 \quad D(f \parallel g) = E_f \left[\log \left(\frac{f}{g} \right) \right] = -h(f) - E_f(\log g)$$

Construis le exponentiel $g(z) = \frac{1}{\lambda} e^{-z/\lambda}$ avec $\lambda = E_w(w)$

Es decir, ¿qué medida?

$$E_f(\log g) = E_f(\text{des} + \text{terminos lineales}) =$$

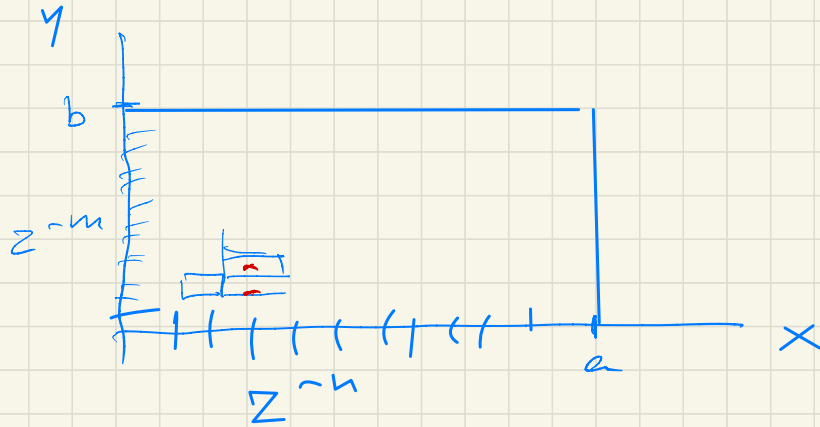
$$= E_f[A + Bz] = E_g[\log g] = h(g)$$

← cambio esencial

Entonces $h(g) \geq h(f)$

Mejora a v.a. de ¿qué medida

Uniformes



$$h(x, y) = h(x) + h(y)$$

$$H(x^\Delta, y^\Delta) = H(x^\Delta) + H(y^\Delta)$$

$$H(x^\Delta, y^\Delta) = \underbrace{h(x, y)} + n + m$$

$$h(x) = \log a$$

$$h(x) = -\int_0^a \frac{1}{a} \log \frac{1}{a} dx$$

$$H(x^\Delta) = h(x) - \log 2^{-n} = h(x) + n$$

$$H(y^\Delta) = h(y) + m$$

Son independientes.

