

Elementos de Teoría de la Información y Codificación

Última actualización: 30/08/2002

INT. TEORIA DE LA INFORMACION

$$I(E) \approx f\left(\frac{1}{p}\right)$$

$$\begin{aligned} & \text{si } p_A > p_B \\ & \Rightarrow I_A < I_B \end{aligned}$$

Si los sucesos A y B son simultáneos e independientes:

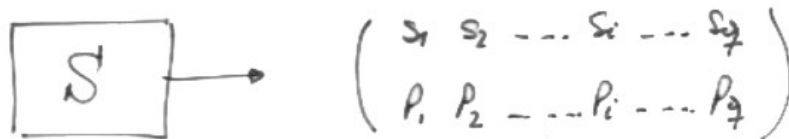
$$I = I_A + I_B \quad I = f\left(\frac{1}{p_A}\right) + f\left(\frac{1}{p_B}\right)$$

$$p_{A \cdot B} = p_A \cdot p_B$$

$$\therefore \boxed{I = \log \frac{1}{p}}$$

$$I = \log_2 \frac{1}{p} \text{ [bit]}$$

FUENTE de Información de MEMORIA NULA



MEM. NULA \Leftrightarrow s_i ESTÁN IND. INDEPENDIENTES.

ENTROPÍA

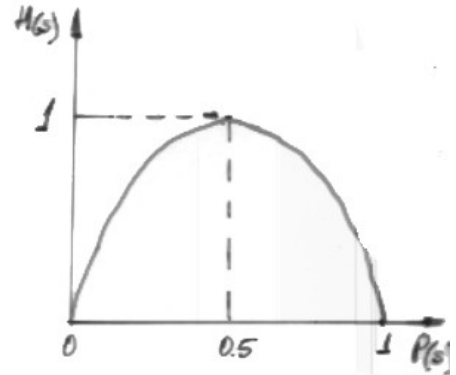
$$H(s) \triangleq \sum_s P(s) I(s) = - \sum P(s) \log P(s)$$

Cantidad media de información por símbolo

Ej

$$\boxed{S} \rightarrow \begin{pmatrix} s_1 & s_2 & s_3 \\ \frac{1}{2} & \frac{1}{4} & \frac{1}{4} \end{pmatrix} \quad H(s) = \frac{1}{2} \log 2 + 2 \left(\frac{1}{4} \log 4 \right) = 1.6 \text{ bits/smb.}$$

Fuente binaria



$$H(s) = - \sum_{i=1}^q p_i \log_2 \frac{1}{p_i}$$

$$H(s) \leq \log_2 q$$

$$H(s) = \log_2 q \iff p_i = \frac{1}{q} \forall i$$

$H(s)$ MAX si todos los símbolos son equiprobables

EXTENSION DE UNA FUENTE DE MEMORIA NULA

Toda una fuente binaria podría generar de 2 a los bits, tres, etc. \Rightarrow un nuevo alfabeto.



la nueva secuencia:

$$v_i = \{s_{i1}, s_{i2}, \dots, s_{in}\} \quad P(v_i) = \prod_{j=1}^n p_{ij}$$

la extensión de orden "n", s^n

$$H(s^n) = n H(s)$$

Ex

$$S \rightarrow S = \begin{pmatrix} s_1 & s_2 & s_3 \\ 1/2 & 1/4 & 1/4 \end{pmatrix}$$

$$H(S) = 1.5 \frac{\text{bits}}{\text{simb.}}$$

Extensión de 2º orden

$$\Rightarrow \boxed{S^2} \rightarrow S^2 = \begin{pmatrix} s_1 s_1 & s_1 s_2 & s_1 s_3 & \dots & s_3 s_3 \\ p_1 p_1 & p_1 p_2 & \dots & p_3 p_3 \end{pmatrix}$$

$$H(S^2) = H(S^2) = 2 H(S) = 3 \frac{\text{bits}}{\text{simb.}}$$

$$H(S^2) = \sum_{\langle S^2 \rangle} P(\langle S^2 \rangle) \log \frac{1}{P(\langle S^2 \rangle)}$$

Fuente de INFORMACION DE MARKOV

Si depende de "m" símbolos precedentes.
"Fuente de MARKOV de orden m".

$$P(s_i / s_{j1}, s_{j2} \dots s_{jm}) \quad i = 1, 2, \dots, q$$

\exists "q" símbolos \neq

$\Rightarrow \exists$ "q^m" estados.

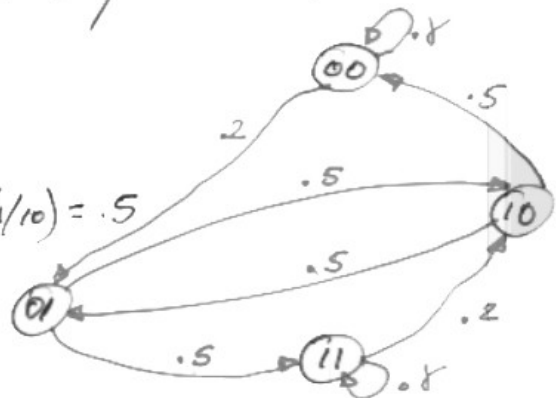
Ex: Una fuente de MARKOV 2º ORDEN con
 $S = \{0, 1\}$ siendo las prob. condicionales

$$P(0/00) = P(1/11) = .8$$

$$P(1/00) = P(0/11) = .2$$

$$P(0/01) = P(0/10) = P(1/01) = P(1/10) = .5$$

- transiciones.



"ESTADÍSTICA DEL LENGUAJE"

Tomando INGLÉS 26 LETRAS + 1 ESPACIO

- Aproximación cero:

Fuente de memoria nula y todos los simb.
equi probables.

$$H(s) = \log 27 = 4.75 \frac{\text{bits}}{\text{simb.}}$$

- Primera Aproximación:

a) $P(A) = .0642$ letra 1961

$P(Z) = 0.0005$

$P(\text{espacio}) = 0.1859$

$$H(s) = - \sum_{27} P_i \log P_i = 4.03 \frac{\text{bits}}{\text{simb.}}$$

b) Usando fuente de MARKOV de 1º orden con simb.
de probs. condic. bien elegidos. (Prob # 42)

$$H(s) = - \sum_{S^2} P(i,j) \log P(i,j) = 3.32 \frac{\text{bits}}{\text{simb.}}$$

- Segunda Aproximación

a) Shannon (51), sobre 1000000 libro de TEXTO.

$$\text{Obtengo } H(s) = 3.1 \frac{\text{bits}}{\text{simb.}}$$

Por otros métodos Shannon $0.6 < H(s) < 1.3$

Mejor aproximación hay que tomar palabras como
símbolos.

COMPOS



Blogue

Se le origina secuencia fite.

No singular

Toda sus palabras \neq

s_1	0
s_2	00
s_3	11
s_4	01

Pero si llega 0011 es singular.

Univoco

Se dice univocamente decodificable \Leftrightarrow su extensión de orden n es "no singular" para cualquier valor finito de n .

Infinito

Puedo decodif. los palabras de la secuencia sin precisar los símbolos que la producen.

Ej.

	códigos		
	A	B	Φ
s_1	00	0	0
s_2	01	10	01
s_3	10	110	011
s_4	11	1110	0111

código coma

$S \rightarrow 10110101101$

Con Φ no distingue
si es 0 o 01

A y B son instantáneos.

Prefixos: la s_4 0111, los pref. son:
0, 01, 011, 0111

"Condición necesaria y suficiente para que un código sea instantáneo es que ninguna palabra del código comience con el prefijo de otra"

Inecuación de Kraft

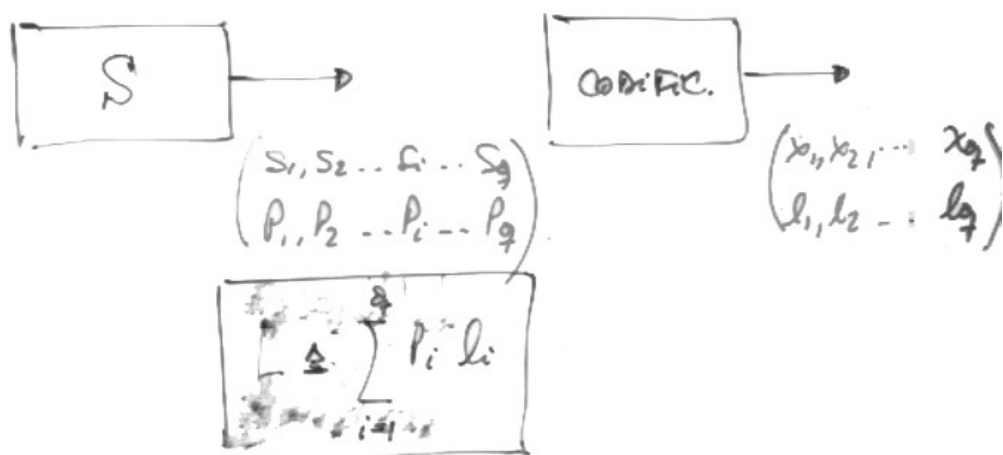
En un alfabeto binario

$$\sum_{i=1}^n 2^{-l_i} \leq 1$$

código inst.

Longitud Media de un Código

Sea un código bloque



Def: Todo un código unívoco, este código será compacto (respecto a S') si su L es \leq a la L de todos los códigos unívocos que pueden aplicarse a la misma fuente y el mismo alfabeto.

En base a la desigualdad de Kraft.

$$\boxed{H(s) \leq L}$$

Muestra la relación que \exists entre la definición de información y una cantidad (L) que no depende de la definición.

NOTA:

L el menor su valor mínimo \Leftrightarrow puede elegirse la long de los palabras l_i , $l_i = \log_2 \frac{1}{P_i}$

$$P_i = \left(\frac{1}{2}\right)^{k_i} \quad k_i: \text{entero.}$$

Se toman $l_i = k_i$ para diseñar un código de L min.

Ex

$$\boxed{S} \rightarrow \begin{pmatrix} S_1 & S_2 & S_3 & S_4 \\ 1/4 & 1/4 & 1/4 & 1/4 \end{pmatrix}$$

$$H(S) = - \sum_{i=1}^4 P_i \log P_i = 2 \frac{\text{bits}}{\text{simb.}}$$

Es imposible codificar esta sin haber medido un código binario univoco de $L < 2 \frac{\text{bits}}{\text{simb.}}$

un código compacto $P_i = \frac{1}{4} \quad P_i = \left(\frac{1}{2}\right)^2$

$$l_i = k_i = 2$$

S_1 00
 S_2 01
 S_3 10
 S_4 11

$$L = 2 \frac{\text{bits}}{\text{simb.}}$$

*

Frases Original		Frases Reducidas.	
Símb.	Pi		
S ₁	.4	<u>1</u>	→ .4 <u>1</u> → .4 <u>1</u> → .6 0
S ₂	.2	<u>01</u>	→ .2 <u>01</u> → .4 <u>001</u> → .4 1
S ₃	.2	<u>000</u>	→ .2 <u>000</u> → .2 <u>01</u>
S ₄	.1	<u>0010</u>	→ .2 <u>001</u>
S ₅	.1	<u>0011</u>	

* Construcción de Códigos Compactos - Código HUFFMAN.

TECNICAS DE COMPRESION

Se pueden hacer {
no función del conjunto de símbolos
frecuencia relativa
Contexto donde aparecen.

COMPRESION DEPENDIENTE DE LA FRECUENCIA

- HUFFMAN (52)

Pero no alcanza en muchos casos el límite teórico (H_0) Prob: algo de 2.5 bits con 3 bits

Huffman lo que hace c/simb. codifica en forma indep. y después lo transmite.

- COMPRESION ARITMETICA

Von Prob. 2.5 bits con 3 bits

No codifica a los símbolos en forma independiente.

COMPRESION DEPENDIENTE CONTEXTO

RLL

Compresión de datos { Esquemas basados en "diccionarios"
Métodos estadísticos.

En sistemas pequeños basados en diccionarios es el más popular.

Pero la combinación de codific. aritmética y técnicas de "modelado" dan una buena performance.

Sistema basado en diccionarios pero reemplaza un grupo de símbolos del `TEXTinput` con código de Long Fitt. (El más conocido LZW: Lempel - Zip - Welch)

- Modelado basado Contexto Finito

Las probabilidades se calculan en función del "contexto" en que los símbolos aparecen.

Modelo orden 0 se calcula los probs. en forma "independiente" de cada símbolo previo.

En orden -0 una sola tabla de frecuencias.
orden -1 256 tablas \neq .
orden -2 65.536

Modelo Adaptativo

U solamente 5% pero si el carácter de contexto previo ^{orden} puede llegar a 95%.

Como el modelo crea líneas nuevas, pero la memoria necesaria es exponencial. !!
oo

Es por ello que la solución modelo "ADAPTATIVO"



Motor SEMI ADAPTATIVO:

Use un modelo \neq para cada texto.
Se envía el decodificador al motor.