

RI - CHAPTER 2 RANKING WITH INDEXES

Elias - γ (gamma) Code

To encode K :

$$K_d = \lfloor \log_2 K \rfloor, \quad K_d \text{ es el mayor de que cumple } 2^{K_d} \leq K$$

$$K_r = K - 2^{K_d}$$

Codificación:

- K_d en unario
- K_r en binario usando K_d dígitos

Bits que requiere:

$$\begin{array}{l} \lfloor \log_2 K \rfloor + 1 \text{ bits para } K_d \text{ en unario} \\ \lfloor \log_2 K \rfloor \text{ bits para } K_r \text{ en binario} \\ \hline \text{Total } 2 \lfloor \log_2 K \rfloor + 1 \text{ bits} \end{array}$$

Elias δ (Delta) Code

To encode K :

$$K_d = \lfloor \log_2 K \rfloor$$

$$K_r = K - 2^{K_d}$$

K_{d+1} se codifica en Elias gamma
 K_r se codifica en binario usando
 K_d dígitos

$$K_{dd} = \lfloor \log_2 (K_{d+1}) \rfloor \rightarrow \text{en binario}$$

$$K_{dr} = (K_{d+1}) - 2^{K_{dd}} \rightarrow \text{en binario usando } K_{dd} \text{ dígitos}$$

↓
Empezar en la slide del libro!

Ej: $K = 1023$

$$K_d = \lfloor \log_2 K \rfloor = \underline{9}$$

$$(2^9 = 512, 2^{10} = 1024)$$

$$K_r = K - 2^{K_d} = 1023 - 512 = \underline{511}$$

$K_{d+1} = 10$ se codifica en gamma

$K_r = 511$ se codifica en binario usando $K_d = 9$ dígitos

$$K_{dd} = \lfloor \log_2 (K_{d+1}) \rfloor = \underline{3} \quad 2^3 = 8 \leq 9$$

$$K_{dr} = (K_{d+1}) - 2^{K_{dd}} = 10 - 2^3 = \underline{2}$$

1110 010 111111111

$K_{dd} = 3$ en binario

$K_{dr} = 2$ en binario usando $K_{dd} = 3$ dígitos

$K_r = 511$ en binario usando $K_d = 9$ dígitos

Shannon - Teoría de la Información

Cantidad de información de un símbolo s_i que se emite con probabilidad P_i $I(s_i) = -\log_2 P_i$ $-I(s_i)$

• Si $P_i = 1 \Rightarrow I(s_i) = 0$

• Si $P_i \rightarrow 0 \Rightarrow I(s_i) \rightarrow \infty$

$\| P_i = \frac{1}{2}$

Entropía o cantidad de información de una distribución de probabilidad

$$H(P) = - \sum_{i=1}^n P_i \log_2 P_i = \sum_{i=1}^n P_i I(s_i)$$

↓
promedio de la cantidad de información de los eventos
por ponderando con sus probabilidades

Ej: $P = \{ \text{cara}, \text{cruz} \} = \{ \frac{1}{2}, \frac{1}{2} \}$

$$H(P) = -\frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{2} \log_2 \frac{1}{2} = -\frac{1}{2} - \frac{1}{2} = \boxed{1 \text{ bit}}$$

↳ con 2 eventos equiprobables, la cantidad de información de saber el resultado por adelantado es 1 bit.

Ej: moneda trucada 99% cara

$$P = \{ \frac{99}{100}, \frac{1}{100} \}$$

$$H(P) = 0.08 \text{ bits}$$

La cantidad de información de saber la salida por adelantado es mucho menor.

la longitud ^{ideal} del código de un símbolo x y su probabilidad $Pr(x)$ idealmente deben seguir la relación de Shannon

$$\boxed{l_x = -\log_2 Pr(x)} \quad \left| \quad \boxed{Pr(x) = 2^{-l_x}} \right|$$

$Pr(x)=1 \Rightarrow$ Sólo un símbolo posible $l_x=0$

$Pr(x) \rightarrow 0 \Rightarrow l_x \rightarrow \infty$

• N símbolos con probabilidad uniforme, i.e.,

$$Pr(x) = \frac{1}{N} \quad l_x = -\log_2 \frac{1}{N}, \text{ es decir,}$$

longitud etc para todos los símbolos. \rightarrow Codificación binaria con un tamaño fijo y el mismo para todos los valores posibles a codificar

• Codificación unaria

$l_x = x$ (la longitud del código es igual al valor que se codifica).

$$Pr(x) = 2^{-x}$$

Codificación ideal para una distribución de probabilidad

$$\hookrightarrow \frac{1}{2}, \frac{1}{4}, \frac{1}{8}, \frac{1}{16} \dots \hookrightarrow$$

El código 0 ocurre el 50% de las veces, el código 10, 25% veces, 110, 12.5% de las veces, etc.

Codificació gamma γ

$$Pr(x) = 2^{-lx} = 2^{-(1+2\log_2 x)} = \frac{1}{2x^2}$$

$$\frac{1}{2^{(1+2\log_2 x)}} = \frac{1}{2 \cdot 2^{2\log_2 x}} = \frac{1}{2 \cdot 2^{\log_2 x^2}} = \frac{1}{2x^2}$$

Distribució probabilitat ideal

$$\left\{ \frac{1}{2}, \frac{1}{8}, \frac{1}{18}, \dots \right\}$$

El valor $x=1$ (codificat 10) $\frac{1}{2}$ veges, valor $x=2$ (codificat 100)

$\frac{1}{8}$ veges, valor $x=3$ (codificat 101) $\frac{1}{18}$ veges

Reserve longituds de codificació larges per a probabilitats més baixes que en un cas.

Codificació delta δ

$$Pr(x) = 2^{-lx} = 2^{-(1+2\log_2 \log_2 x + \log_2 x)} = \frac{1}{2x(\log x)^2}$$

Per a valors petits els delta codes són més llargs que els gamma codes, però per a valors grans de contrari.

$x = 1.000.000 \rightarrow \gamma \text{ code } 39 \text{ bits}$
 $\rightarrow \delta \text{ code } 28 \text{ bits}$

Compression Example

(1, 2, [1, 7]) (2, 3, [6, 17, 197]) (3, 1, [1])

positions
doc-id

D-gaps Para doc-id 7 positions:

(1, 2, [1, 6]) (1, 3, [6, 11, 180]) (1, 1, [1])

compress with 4-byte

81 82 81 86 81 83 86 8B 01 B4 81 81 81

0000 0001

1000 1011

0000 0001 1011 0100

1011 0100

128+

32

16

4

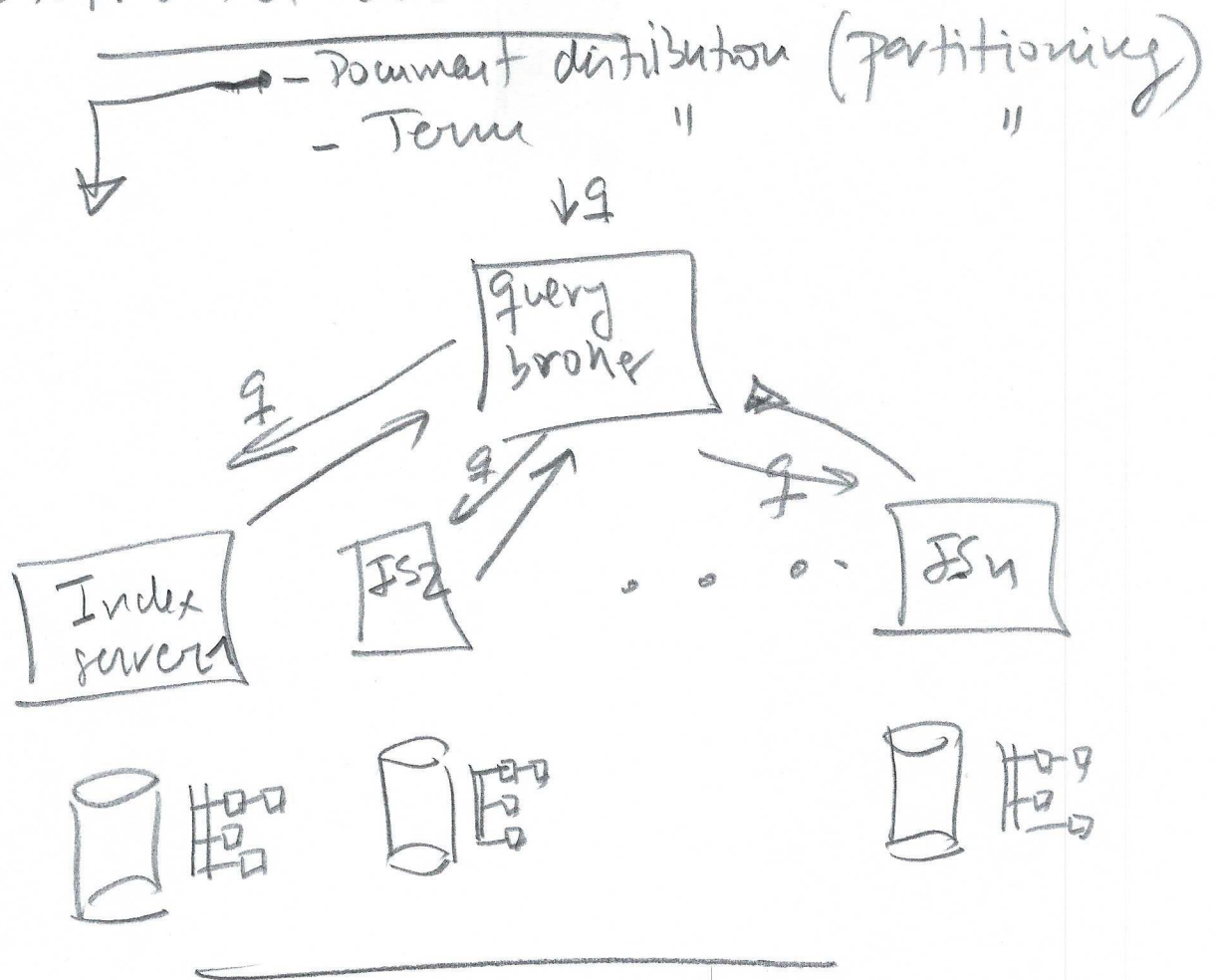
180

Skip Pointers

D-gaps 9

5, 6, 6, 4, 5, 8, ...

Distributed Evaluation



Term distribution: For a query of k terms and n index servers, the number of disjuncts necessary will be $O(k)$, but in document distribution will be $O(k \cdot n)$ because each query term has to be processed in the n servers (index servers).

Term Distribution, although is more complex and requires more messages between index servers. For a single query, it is more efficient in EIS.

- Federated (Pure Distributed) Search.

Los docs y por tanto índices creados por esos docs están en sitios distintos por razones de privacidad, leyes, etc.

Parsons en Document Distribution: El broker no tiene las estadísticas globales y no tiene acceso a los colecciones. Puede hacer

- Pedir a los sitios las estadísticas de sus colecciones o hacer query sampling.
- Resource (file) selection
- Enviar lo query solo a los sitios seleccionados
- Fusionar los resultados.

- Metasearch

Unión de scores o de rankings (posiciones) de diversos search engines a los que puedes acceder por un API

Técnicas de combinación.
Combortz, Berdelouet, etc.

Ver Aslam, Montague. Models for Metasearch, ACM SIGIR 2001

Caching de queries

Ley de Zipf (Potencia negativa)

$$f(u) = C u^{-\alpha} = C \cdot \frac{1}{u^\alpha}$$

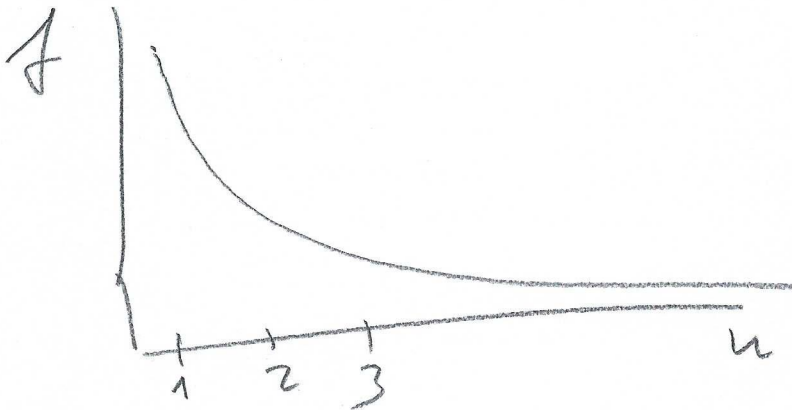
↑ frecuencia o número de queries

n, rango por frecuencia

$n=1$ query más frecuente

$n=2$ 2^a query más frecuente

etc.



Ley exponencial que favorece caching

