

Costos de Búsqueda en Árboles Binarios de Búsqueda

Introducción

Ya hemos definido recursivamente un árbol binario de búsqueda y hemos planteado una manera de deducir el esfuerzo medio de localización exitosa a priori, usando para ello las funciones i , I y E . Ahora analizaremos otras formas de hacerlo.

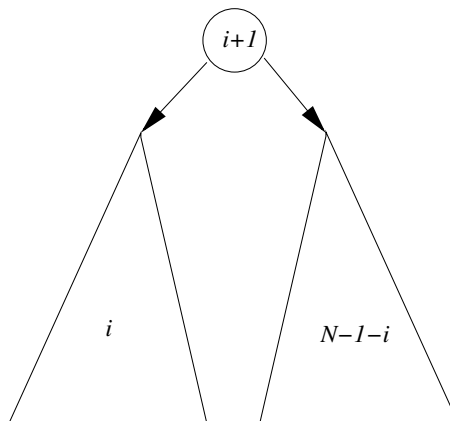
Tenemos una secuencia de N elementos con los que armaremos un árbol binario de búsqueda y queremos calcular el esfuerzo medio de buscar en él exitosamente uno de esos N elementos. Además, por simplicidad, vamos a suponer que una sucesión de N valores del conjunto se puede ver como isomorfa a una sucesión de los valores de $1 \dots N$.

Vamos a obtener el esfuerzo medio a priori, y por ello estamos pensando en obtener el esfuerzo medio sobre todos los elementos a buscar y también el medio respecto del árbol a construir. Por lo tanto, el universo a promediar tiene $N! \times N$ datos ($N!$ secuencias y N elementos a ser buscados).

Vamos a medir los costos en *cantidad de celdas consultadas* y asumiremos que cualquiera de las secuencias de entrada es igualmente probable de ocurrir y que podemos buscar con éxito, con igual probabilidad, cualquiera de los N elementos.

Si tenemos N elementos, consultamos la raíz (1 consulta) y la probabilidad de tener éxito en esa primera consulta es $\frac{1}{N}$ y no necesito hacer nada más (0 consultas adicionales). Pero con probabilidad $\frac{(N-1)}{N}$ ¹ fracasaremos en esa primera pregunta y tendremos que seguir buscando. A su vez, la probabilidad de tener que buscar en un determinado subárbol está relacionada con la cantidad de elementos que éste contiene.

Entonces, en un árbol cuya raíz es el $i + 1$ -ésimo elemento (en la sucesión ordenada de los N elementos) tendría el siguiente árbol:



si el elemento $i + 1$ se encuentra como raíz, habrán i elementos en el subárbol izquierdo y $N - 1 - i$ elementos en el derecho, y la fórmula en ese caso sería:

¹Recordar que la suma de probabilidades sobre todo el espacio debe ser igual a 1.



$$C(N)_i = 1 + \left(\frac{1}{N} 0 + \frac{i}{N} C(i) + \frac{(N-1-i)}{N} C(N-1-i) \right)$$

Pero, no conocemos realmente qué elemento es el que está como raíz y por lo tanto debemos promediar sobre todos los i posibles. Como son isoprobables la probabilidad de cada i es $\frac{1}{N}$. Así llegamos a la siguiente fórmula, que es la recurrencia inicial para obtener el esfuerzo medio a priori de localización exitosa.

Fórmula recurrente inicial:

$$C(N) = 1 + \frac{1}{N} \sum_{i=0}^{N-1} \left(\frac{i}{N} C(i) + \frac{(N-1-i)}{N} C(N-1-i) \right) \quad (1)$$

$$= 1 + \frac{2}{N^2} \sum_{i=0}^{N-1} i C(i) \quad (2)$$

$$C(1) = 1 \quad (3)$$

Se puede observar en (1) que al desarrollar la sumatoria cada $\frac{i}{N} C(i)$ aparece dos veces; y por lo tanto se llega a la ecuación (2).

La última igualdad no es imprescindible, porque la recurrencia la reconstruye para $N = 1$; es decir que si en (1) o en (2) reemplazo por el valor de $N = 1$, obtengo 1 como resultado.

El desarrollo de esta fórmula se puede realizar por uno de dos caminos posibles: en forma *algebraica* o usando *funciones generatrices*.

Camino Algebraico

El primer cambio a realizar sobre la ecuación (2) se debe a que como aparece a la derecha $i C(i)$, entonces queremos lograr que a la izquierda aparezca $N C(N)$. Entonces quedaría:

$$N C(N) = N + \frac{2}{N} \sum_{i=0}^{N-1} i C(i) \quad (4)$$

y ahora por simplicidad de escritura rebautizamos $N C(N)$ por $D(N)$, y entonces obtenemos la siguiente ecuación:

$$D(N) = N + \frac{2}{N} \sum_{i=0}^{N-1} D(i) \quad (5)$$

ahora conviene hacer desaparecer $N - 1$ y tener N como límite superior de la sumatoria. Entonces, usamos que:

$$D(N+1) = (N+1) + \frac{2}{N+1} \sum_{i=0}^N D(i) \quad (6)$$



en este punto, en cada una de las dos ecuaciones anteriores ((5) y (6)) eliminamos los denominadores haciendo:

$$N D(N) = N^2 + 2 \sum_{i=0}^{N-1} D(i)$$

$$(N + 1) D(N + 1) = (N + 1)^2 + 2 \sum_{i=0}^N D(i)$$

podemos restar miembro a miembro las dos últimas ecuaciones para eliminar las sumatorias y quedaría:

$$(N + 1) D(N + 1) - N D(N) = (N + 1)^2 - N^2 + 2 D(N)$$

ahora resolviendo, despejamos $D(N + 1)$ y nos queda:

$$\begin{aligned} (N + 1) D(N + 1) &= (N + 1)^2 - N^2 + 2 D(N) + N D(N) \\ &= (N + 1)^2 - N^2 + (N + 2) D(N) \\ &= N^2 + 2 N + 1 - N^2 + (N + 2) D(N) \\ &= 2 N + 1 + (N + 2) D(N) \end{aligned}$$

$$D(N + 1) = \frac{N + 2}{N + 1} D(N) + \frac{2N + 1}{N + 1} \quad (7)$$

Sabemos que $C(1) = 1$ y que $D(1) = 1$ por ser búsqueda exitosa. Pero queremos que la última fórmula hable de N y no de $N + 1$, entonces bajo el $N + 1$ obteniendo:

$$D(N) = \frac{N + 1}{N} D(N - 1) + \frac{2N - 1}{N} \quad (8)$$

Si reemplazamos en (8) (la última fórmula) $D(N - 1)$, dejando a $D(N)$ ahora en función de $D(N - 2)$ y luego reemplazamos $D(N - 2)$ en función de $D(N - 3)$, tendríamos:

$$\begin{aligned} D(N) &= \frac{N + 1}{N} \left(\frac{N}{N - 1} D(N - 2) + \frac{2N - 3}{N - 1} \right) + \frac{2N - 1}{N} \\ &= \frac{N + 1}{N} \frac{N}{N - 1} \left(\frac{N - 1}{N - 2} D(N - 3) + \frac{2N - 5}{N - 2} \right) + \frac{N + 1}{N} \frac{2N - 3}{N - 1} + \frac{2N - 1}{N} \\ &= \frac{N + 1}{N} \frac{N}{N - 1} \frac{N - 1}{N - 2} D(N - 3) + \frac{N + 1}{N} \frac{N}{N - 1} \frac{2N - 5}{N - 2} + \frac{N + 1}{N} \frac{2N - 3}{N - 1} + \frac{2N - 1}{N} \quad (9) \end{aligned}$$

Por otra parte, si observamos el primer término de la ecuación (9) podemos deducir la forma que tendrá ese término cuando se alcance el $D(1)$:

$$\frac{N + 1}{N} \frac{N}{N - 1} \frac{N - 1}{N - 2} \frac{N - 2}{N - 3} \dots \frac{3}{2} D(1)$$

si simplificamos quedarían el mayor numerador y el menor denominador, es decir:

$$\frac{N + 1}{2} D(1) \quad (10)$$



Podemos simplificar las fracciones que aparecen en los términos de la última fórmula, quedando:

$$D(N) = \frac{N+1}{N-2} D(N-3) + \frac{N+1}{N-1} \frac{2N-5}{N-2} + \frac{N+1}{N} \frac{2N-3}{N-1} + \frac{2N-1}{N}$$

ahora podemos multiplicar y dividir el último término por $\frac{N+1}{N+1}$ para que aparezcan términos similares:

$$D(N) = \frac{N+1}{N-2} D(N-3) + \frac{N+1}{N-1} \frac{2N-5}{N-2} + \frac{N+1}{N} \frac{2N-3}{N-1} + \frac{2N-1}{N} \frac{N+1}{N+1} \quad (11)$$

Recapitulando, de (10) y (11), tenemos la siguiente ecuación:

$$D(N) = \frac{N+1}{N+1} \frac{2N-1}{N} + \frac{N+1}{N} \frac{2N-3}{N-1} + \frac{N+1}{N-1} \frac{2N-5}{N-2} + \dots + \frac{N+1}{4} \frac{5}{3} + \frac{N+1}{3} \frac{3}{2} + \frac{N+1}{2}$$

Para ver mejor la fórmula podemos reescribirla como:

$$D(N) = \frac{N+1}{3} \frac{3}{2} + \frac{N+1}{4} \frac{5}{3} + \dots + \frac{N+1}{N-1} \frac{2N-5}{N-2} + \frac{N+1}{N} \frac{2N-3}{N-1} + \frac{N+1}{N+1} \frac{2N-1}{N} + \frac{N+1}{2} \quad (12)$$

así, nos podemos dar cuenta que en todos los términos, salvo el último, el patrón es $\frac{(2i-1)}{i(i+1)}$ y que podemos sacar factor común de ellos a $(N+1)$, además si los abreviamos usando una sumatoria vemos que iría desde 2 a N . Entonces llegamos así a la fórmula:

$$D(N) = (N+1) \left(\sum_{i=2}^N \frac{2i-1}{i(i+1)} \right) + \frac{N+1}{2} \quad (13)$$

Otra manera de encontrar la forma de los términos y resolver la sumatoria sería:

$$D(N) = (N+1) \underbrace{\left(\frac{(2N-1)}{(N+1)N} + \frac{(2N-3)}{N(N-1)} + \frac{(2N-5)}{(N-1)(N-2)} + \dots + \frac{3}{3 \cdot 2} \right)}_{S(N)} + \frac{(N+1)}{2} \quad (14)$$

sacamos factor común $(N+1)$ y llamamos $S(N)$ a la parte de la sumatoria que debemos resolver.

Ahora tenemos que resolver $S(N)$ en (14):

$$\begin{aligned} S(N) &= \left(\frac{(2N-1)}{(N+1)N} + \frac{(2N-3)}{N(N-1)} + \frac{(2N-5)}{(N-1)(N-2)} + \dots + \frac{3}{3 \cdot 2} \right) \\ &= \sum_{j=1}^{N-1} \left(\frac{2N - (2j-1)}{(N+2-j)(N+1-j)} \right) \end{aligned} \quad (15)$$

si ahora llamamos en (15) $k = (N-j)$ vemos que k tomaría los valores desde $(N-1)$ a 1, o lo que es lo mismo desde 1 a $(N-1)$, así tendríamos:

$$\begin{aligned} S(N) &= \sum_{k=N-1}^1 \frac{(2k+1)}{(k+2)(k+1)} \\ &= \sum_{k=1}^{N-1} \frac{(2k+1)}{(k+2)(k+1)} \end{aligned} \quad (16)$$



ahora podemos volver a renombrar la variable haciendo que $i = (k - 1)$ en (16), y redefiniendo los límites para i vemos que irían desde 2 a N , así tenemos que:

$$S(N) = \sum_{i=2}^N \frac{(2i-1)}{i(i+1)} \quad (17)$$

ahora volviendo a $D(N)$ y reemplazando en (14) la fórmula de $S(N)$, obtenida en (17), llegaríamos a la misma expresión que teníamos en (13):

$$D(N) = (N+1) \left(\sum_{i=2}^N \frac{2i-1}{i(i+1)} \right) + \frac{N+1}{2} \quad (18)$$

Podemos verificar la fórmula obteniendo el resultado para los valores conocidos. Viendo los árboles, tenemos que $D(1) = 1$ y $D(2) = 3$; y reemplazando en la fórmula anterior, para $N = 1$ y $N = 2$, obtenemos también que $D(1) = 1$ y $D(2) = 3$.

Para continuar con el método algebraico podemos usar que una división de polinomios, en la que el numerador es de menor orden que el denominador, se puede descomponer como suma de factores:

$$\begin{aligned} \frac{(2i-1)}{i(i+1)} &= \frac{a}{i} + \frac{b}{(i+1)} \\ &= \frac{a(i+1) + bi}{i(i+1)} = \frac{(a+b)i + a}{i(i+1)} \end{aligned}$$

y resolviendo esta última ecuación tenemos que:

$$a + b = 2 \quad (19)$$

$$a = -1$$

$$b = 3 \quad (20)$$

Así, reemplazando en (18) los coeficientes obtenidos en (19) y (20) y tendríamos:

$$\begin{aligned} D(N) &= (N+1) \left(\sum_{i=2}^N \frac{-1}{i} + \frac{3}{(i+1)} \right) + \frac{N+1}{2} \\ &= (N+1) \left(\sum_{i=2}^N \frac{-1}{i} + \sum_{i=2}^N \frac{3}{(i+1)} \right) + \frac{N+1}{2} \\ &= (N+1) \left(\frac{1}{2} - \sum_{i=2}^N \frac{1}{i} + 3 \sum_{i=2}^N \frac{1}{(i+1)} \right) \quad (21) \end{aligned}$$

Ahora podemos observar que la primera sumatoria “casi” corresponde a la serie armónica H_N y la segunda “casi” correspondería a la serie armónica H_{N+1} . A dichas series no se les conoce el resultado, a pesar de no ser siempre divergentes, pero se conoce su aproximación.

$$H_N = \sum_{i=1}^N \frac{1}{i} \simeq \ln N + \gamma + O\left(\frac{1}{N}\right)$$



la constante γ es la llamada *constante de Euler* y es $\gamma = 0,5772156649\dots$ y además a medida que N crece el último término es despreciable. Así, normalmente usamos $\ln N$ como una aproximación a H_N .

Ahora reescribimos nuestra fórmula (21) usando lo anterior:

$$D(N) = (N + 1) \left(\frac{1}{2} - (H_N - 1) + 3(H_{N+1} - \frac{1}{2} - 1) \right) \quad (22)$$

Ahora podemos volver a reemplazar $D(N)$, por lo que era $N C(N)$, y luego resolver lo más posible:

$$\begin{aligned} N C(N) &= (N + 1) \left(\frac{1}{2} - (H_N - 1) + 3(H_{N+1} - \frac{1}{2} - 1) \right) \\ C(N) &= \frac{N + 1}{N} \left(\frac{1}{2} - (H_N - 1) + 3(H_{N+1} - \frac{1}{2} - 1) \right) \\ &= \frac{N + 1}{N} \left(\frac{1}{2} - H_N + 1 + 3(H_N + \frac{1}{N+1} - \frac{1}{2} - 1) \right) \\ &= \frac{N + 1}{N} \left(\frac{1}{2} - H_N + 1 + 3H_N + \frac{3}{N+1} - \frac{3}{2} - 3 \right) \\ &= \frac{N + 1}{N} \left(\frac{1}{2} + 1 + 2H_N + \frac{3}{N+1} - \frac{3}{2} - 3 \right) \\ &= \frac{N + 1}{N} \left(2H_N + \frac{1}{2} + 1 - \frac{3}{2} - 3 \right) + \frac{3}{N+1} \frac{N + 1}{N} \\ &= \frac{N + 1}{N} \left(2H_N + \frac{1}{2} + 1 - \frac{3}{2} - 3 \right) + \frac{3}{N} \\ &= \frac{N + 1}{N} (2H_N - 3) + \frac{3}{N} \\ &= 2 \frac{N + 1}{N} H_N - 3 \frac{(N + 1)}{N} + \frac{3}{N} \\ &= 2 \frac{N + 1}{N} H_N - \frac{3N - 3}{N} + \frac{3}{N} \\ &= 2 \frac{N + 1}{N} H_N - \frac{3N}{N} \\ &= 2 \frac{N + 1}{N} H_N - 3 \\ &\simeq 2 \frac{(N + 1)}{N} \ln N - 3 \end{aligned} \quad (23)$$

ahora, para que el logaritmo de N quede en base 2, hacemos transformación de la base quedando:

$$C(N) \simeq 2 \ln 2 \frac{(N + 1)}{N} \log N - 3 \quad (24)$$

$$C(N) \simeq 1,38 \frac{(N + 1)}{N} \log N - 3 \quad (25)$$



Como se puede observar, la ecuación (23) obtenida aquí es igual a la obtenida en el apunte de *árboles Binarios Ordenados* (ver ecuación (43) del mencionado apunte).

Usando Funciones Generatrices

Primero a la ecuación planteada en (2) se le puede hacer un cambio de variable $F(N) = N C(N)$.

$$\frac{F(N)}{N} = 1 + \frac{2}{N^2} \sum_{i=0}^{N-1} F(i) \quad (26)$$

$$F(1) = 1 \quad (27)$$

$$F(0) = 0 \quad (28)$$

Para poder eliminar la sumatoria debemos dejarla con coeficiente constante:

$$NF(N) = N^2 + 2 \sum_{i=0}^{N-1} F(i) \quad (29)$$

la reescribimos para $N - 1$:

$$(N - 1)F(N - 1) = (N - 1)^2 + 2 \sum_{i=0}^{N-2} F(i) \quad (30)$$

la restamos miembro a miembro:

$$NF(N) - (N - 1)F(N - 1) = N^2 - (N - 1)^2 + 2F(N - 1) \quad (31)$$

$$NF(N) = (N + 1)F(N - 1) + 2N - 1 \quad (32)$$

Para acercarnos a la función generatriz lo multiplicamos por z^N y lo sumamos desde 1 (para no tener un argumento negativo en la $F()$) hasta ∞ . Finalmente pasamos la N a i más habitual como subíndice y escribimos el argumento de $F()$ como subíndice.

$$\sum_{i=1}^{\infty} i F_i z^i = \sum_{i=1}^{\infty} (i + 1) F_{i-1} z^i + \sum_{i=1}^{\infty} (2i - 1) z^i \quad (33)$$

La primera sumatoria se puede extender a 0 tanto por el factor i como por el valor de F_0 . En las sumatorias del segundo miembro saco una z y cambio la variable para empezar en 0.

$$\sum_{i=0}^{\infty} i F_i z^i = z \sum_{i=0}^{\infty} (i + 2) F_i z^i + z \sum_{i=0}^{\infty} (2i + 1) z^i \quad (34)$$

En lo que sigue D será el operador de derivada. Se consigue introducir los factores i con el operador

$z D$. Por otra parte $\sum_{i=0}^{\infty} z^i = \frac{1}{(1 - z)}$

$$zDF(z) = z[zDF(z) + 2F(z)] + z \left[2zD \frac{1}{1 - z} + \frac{1}{1 - z} \right] \quad (35)$$



Podemos suprimir un factor z común a todos los términos

$$DF(z) = [zDF(z) + 2F(z)] + \left[2zD\frac{1}{1-z} + \frac{1}{1-z} \right] \quad (36)$$

$$(1-z)DF(z) = 2F(z) + \frac{2z}{(1-z)^2} + \frac{1}{1-z} \quad (37)$$

$$DF(z) = \frac{2}{(1-z)} F(z) + \frac{2z}{(1-z)^3} + \frac{1}{(1-z)^2} \quad (38)$$

$$DF(z) - \frac{2}{(1-z)} F(z) = \frac{2z}{(1-z)^3} + \frac{1}{(1-z)^2} \quad (39)$$

Y esto que hemos obtenido es una ecuación diferencial lineal. La manera de resolverlo es hallar primero la solución de la ecuación homogénea, o sea, con término independiente nulo.

$$DF(z) - \frac{2}{(1-z)} F(z) = 0 \quad (40)$$

$$DF(z) = \frac{2}{(1-z)} F(z) \quad (41)$$

$$\frac{dF(z)}{dz} = \frac{2}{(1-z)} F(z) \quad (42)$$

$$d \ln F(z) = d \ln(1-z)^{-2} \quad (43)$$

$$\ln F(z) = \ln \frac{1}{(1-z)^2} + K \quad (44)$$

$$F(z) = \frac{K}{(1-z)^2} \quad (45)$$

La solución particular se obtiene suponiendo que la constante multiplicativa es una función

$$F(z) = \frac{K(z)}{(1-z)^2} \quad (46)$$

$$DF(z) = \frac{DK(z)}{(1-z)^2} + \frac{2K(z)}{(1-z)^3} \quad (47)$$

$$= \frac{2}{(1-z)} \frac{K(z)}{(1-z)^2} + \frac{2z}{(1-z)^3} + \frac{1}{(1-z)^2} \quad (48)$$

Igualando los segundos miembros y suprimiendo el término común se tiene

$$\frac{DK(z)}{(1-z)^2} + \frac{2K(z)}{(1-z)^3} = \frac{2}{(1-z)} \frac{K(z)}{(1-z)^2} + \frac{2z}{(1-z)^3} + \frac{1}{(1-z)^2} \quad (49)$$

$$\frac{DK(z)}{(1-z)^2} = \frac{2z}{(1-z)^3} + \frac{1}{(1-z)^2} \quad (50)$$

$$DK(z) = \frac{2z}{1-z} + 1 \quad (51)$$

$$= \frac{2}{1-z} - 1 \quad (52)$$

$$K(z) = -2 \ln(1-z) - z \quad (53)$$

$$F(z) = \frac{-2 \ln(1-z)}{(1-z)^2} - \frac{z}{(1-z)^2} \quad (54)$$

¿Cómo sacar de aquí los coeficientes F_n ?

Multiplicar por $\frac{1}{(1-z)}$ es lo mismo que hace sumatorias parciales, o sea:

$$\frac{1}{1-z} \sum_{i=0}^{\infty} a_i z^i = \sum_{j=0}^{\infty} \left(\sum_{k=0}^j a_k \right) z^j \quad (55)$$

Desarrollando según Mac Laurin se puede deducir:

$$-\ln(1-z) = \frac{z}{1} + \frac{z^2}{2} + \frac{z^3}{3} + \dots \quad (56)$$

Las sumatorias parciales de estos coeficientes son H_n .

$$\frac{-\ln(1-z)}{1-z} = \sum_{i=0}^{\infty} H_i z^i \quad (57)$$

Pero tenemos $1-z$ con exponente 2, o sea, tenemos que volver a sumar.

$$\frac{-\ln(1-z)}{(1-z)^2} = \frac{1}{1-z} \sum_{i=0}^{\infty} H_i z^i \quad (58)$$

$$= \sum_{j=0}^{\infty} \left(\sum_{k=0}^j H_k \right) z^j \quad (59)$$

$$\sum_{k=0}^j H_k = \sum_{k=0}^j \sum_{i=1}^k \frac{1}{i} \quad (60)$$

$$= \sum_{i=1}^k \frac{j+1-i}{i} \quad (61)$$

$$= (j+1)H_j - j \quad (62)$$

$\frac{z}{(1-z)}$ engendra la sucesión 0, 1, 1, 1, ... Su segunda suma engendra la sucesión 0, 1, 2, 3, ..., con lo cual tenemos:

$$F_n = 2(n+1)H_n - 2n - n \quad (63)$$

$$C(N) = 2 \frac{(N+1)}{N} H_N - 3 \quad (64)$$

$$C(N) \simeq 2 \ln N \frac{(N+1)}{N} - 3 \quad (65)$$

$$C(N) \simeq 2 \ln 2 \frac{(N+1)}{N} \log N - 3 \quad (66)$$

$$C(N) \simeq 1,38 \frac{(N+1)}{N} \log N - 3 \quad (67)$$

El camino que hemos seguido es mixto, empezamos con un trabajo algebraico y en cierto punto pasamos a funciones generatrices. Podríamos haberlo hecho desde el comienzo convirtiendo (1) a funciones generatrices.



Como reconocemos que los dos términos de la sumatoria dan lo mismo arrancamos de (2). Para no tener integrales, sacamos las N que dividen multiplicando todo por N^2 .

$$N^2 C(N) = N^2 + 2 \sum_{i=0}^{N-1} i C(i) \quad (68)$$

Un factor N se introduce con $z D$, el número 1 con $\frac{1}{(1-z)}$. Como la sumatoria llega hasta $N - 1$ es una función de $N - 1$ y no de N por lo tanto sobra una z y se construye una sumatoria multiplicando por $\frac{1}{(1-z)}$.

Aplicando todo esto tenemos:

$$z D z D C(z) = z D z D \frac{1}{1-z} + 2 \frac{z}{1-z} z D C(z) \quad (69)$$

Conclusiones

Finalmente, cabe destacar que son resultados equivalentes:

- La fórmula (65) obtenida usando funciones generatrices a partir de la fórmula (2),
- La fórmula (23) obtenida usando el camino algebraico y
- La fórmula (43), del apunte de *árboles Binarios Ordenados*, obtenida usando las funciones i , I y E y la relación existente entre esfuerzos de búsqueda exitosa y búsqueda que fracasa.

