

1. Exercise 5.7 (a), (b) and (c)

(a) $\rho(A \dots B, k) \rightarrow$ 1st interval in $A \dots B$ ending at or after k

$\rho(A \dots B, k) \equiv$

if $k == \infty$ then

return $[\infty, \infty]$

if $k == -\infty$ then

return $[-\infty, -\infty]$

$[u, v] \leftarrow \rho(A, k) \dots$ 1st interval in A ending at/after k

if $[u, v] == [\infty, \infty]$ then

return $[\infty, \infty]$

$[u', v'] \leftarrow \rho(B, v+1) \dots$ 1st interval in B ending at/after $v+1$

if $[u, v] == [\infty, \infty]$ then

return $[\infty, \infty]$

$[u'', v''] \leftarrow \rho'(A, u'-1)$

return $[u'', v']$

(b) $\tau(A \Delta B, K)$

(First interval in $A \Delta B$ starting at/after K)

$\Rightarrow \tau(A \Delta B, K) \equiv$

if $K == \infty$ then
return $[\infty, \infty]$

if $K == -\infty$ then
return $[-\infty, \infty]$

$[u, v] \leftarrow \tau(A, K) \dots$ (get 1st interval in A
if $[u, v] == [\infty, \infty]$ then starting at/after K)
return $[\infty, \infty]$

~~if $A' \neq \emptyset$~~
 $[u', v'] \leftarrow \tau(B, K) \dots$ get 1st interval in B
starting at/after K)

// Now we find then interval
having both the intervals

if $u' < u$ && $v' < v$ then // overlap case-1
return $[u', v']$

if $u > u'$ && $v > v'$ then // overlap case-2
return $[u, v']$

if $u > u'$ && $v' > v$ then // $[u', v'] \in [u, v]$
return $[u, v]$

if $v' > v$ && $u' > v$ then // disjoint intervals.
return $\tau(A \Delta B, v)$

(c) $\tau(A \nabla B, K)$

(First interval in $A \nabla B$ starting at/after K)

$\tau(A \nabla B, K) \equiv$

if $K = \infty$ then

return $[\infty, \infty]$

if $K = -\infty$ then

return $[-\infty, -\infty]$

$[u, v] \leftarrow \tau(A, K)$

if $[u, v] = [\infty, \infty]$ then

return $[\infty, \infty]$

$[u', v'] \leftarrow \tau(B, K)$

// Now we find the first interval
having either A or B's interval

if $u > u' \&\& v' > v$ then // $[u', v'] \in [u, v]$

return $[u, v]$

if $u < u' \&\& v < v'$ then // overlap case 1

return $[u, v]$

if $u' < u \&\& v' < v$ then // overlap case 2

return $[u', v']$

if $u' < u$ and $v' < u'$ then

return $[u', v']$.

2. $Pr["a"] = 0.55$ & $Pr["b"] = 0.45$

When we group symbols in block of size 2, we get

$$Pr["aa"] = Pr["a"] \times Pr["a"] = 0.55 \times 0.55 = 0.3025$$

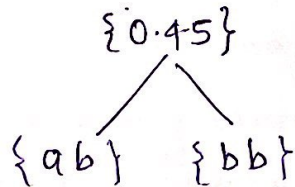
$$Pr["ab"] = Pr["a"] \times Pr["b"] = 0.55 \times 0.45 = 0.2475$$

$$Pr["ba"] = Pr["b"] \times Pr["a"] = 0.45 \times 0.55 = 0.2475$$

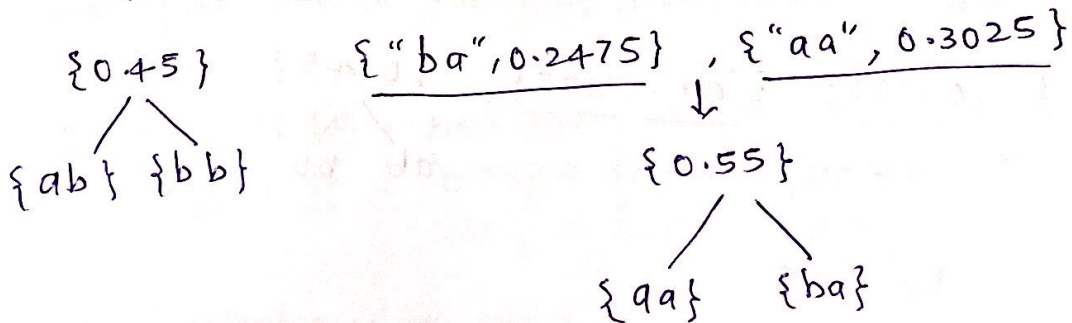
$$Pr["bb"] = Pr["b"] \times Pr["b"] = 0.45 \times 0.45 = 0.2025$$

to build Huffman tree using above distribution we get.

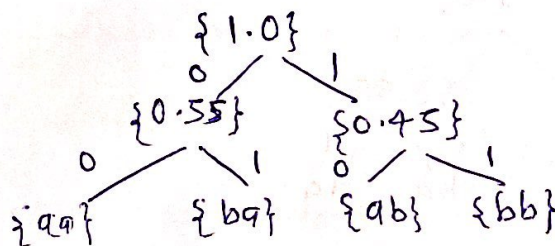
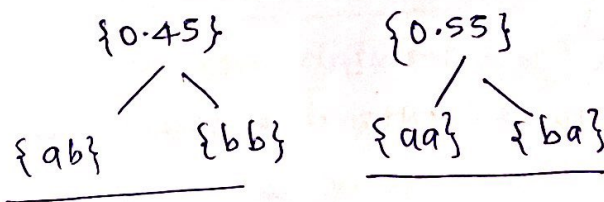
Step 1: $\{ \text{"bb"}, 0.2025 \}$, $\{ \text{"ab"}, 0.2475 \}$, $\{ \text{"ba"}, 0.2475 \}$, $\{ \text{"aa"}, 0.3025 \}$



Step 2:



Step 3:



Huffman code for 'aa' = 00, 'ba' = 01, 'ab' = 10, 'bb' = 11

3. Exercise 6.2 - show that the γ code is prefix-free

Solution - γ code is a non parametric gap compression code used to compress the posting list under the assumption that Gaps between entries are relatively short.

A γ -code has two parts a selector and body.

If we want to encode k in γ code we need $2\lceil \log_2(k) \rceil + 1$ bits i.e., $\log_2(k)$ bits for selector and $\log_2(k)$ bits for the body.

To prove that γ code is prefix free we consider following 2 cases.

(a) When length of selectors is same

When we have 2 messages (Gap numbers) which need to be converted into γ code and having same length for selector field, we can assume that they are prefix free because the body field of those 2 numbers will be different.

Eg: $k_1 = 6 \rightarrow \underline{001} \underline{10}$
 $k_2 = 7 \rightarrow \underline{001} \underline{11}$ } Here even when the selector is of same size γ code is prefix free because body is different for every number.

(b) When length of selectors is not same

When selectors have different lengths, bodies of those codes will also differ in their respective lengths. selectors for vice codes are themselves prefix free making the γ -code prefix free as a whole

Eg: $k_1 = 17 \rightarrow \underline{00001} \underline{0001}$
 $k_2 = 11 \rightarrow \underline{0001} \underline{011}$ } Here, selector field for both k_1 and k_2 is different and prefix free.

Eg: $k_1 = 4 \rightarrow \underline{001} \underline{00}$
 $k_2 = 12 \rightarrow \underline{0001} \underline{100}$ } As the selector increases in size, it follows prefix free property making γ -code prefix free.

4. What is the expected number of bits per codeword when using a rice code with parameter $M=2^6$ to compress a geometrically distributed posting list for term T with $N_T/N=0.05$?

Solution:

$$|E| = \sum_{k=1}^{\infty} \text{Pr}[\Delta=k] \cdot L(\Delta=k) \quad \dots \dots \dots (1)$$

where, $|E|$ - expected no. of bits/codeword

$\text{Pr}[\Delta=k]$ - geometric distribution with $\Delta=k$

$L(\Delta=k)$ - Length of codeword with gap k ($\Delta=k$)

we know that,

$$\text{Pr}[\Delta=k] = (1-p)^{k-1} \cdot p$$

since we know the distribution is $N_T/N = p = 0.05$

$$\text{Pr}[\Delta=k] = (1-N_T/N)^{k-1} \cdot N_T/N$$

$$= (1-0.05)^{k-1} \cdot 0.05$$

$$\text{Pr}[\Delta=k] = (0.95)^{k-1} \cdot 0.05 \quad \dots \dots \dots (2)$$

and,

$$L(\Delta=k) = \left\lfloor \frac{k-1}{M} \right\rfloor + 1 + \lceil \log_2(M) \rceil$$

$$= \left\lfloor \frac{k-1}{2^6} \right\rfloor + 1 + \lceil \log_2(2^6) \rceil$$

$$L(\Delta=k) = 7 + \left\lfloor \frac{k-1}{2^6} \right\rfloor \quad \dots \dots \dots (3)$$

Substituting (2) and (3) in original equation (1) we get,

$$|E| = \sum_{k=1}^{\infty} 0.05 (0.95)^{k-1} \cdot \left(7 + \left\lfloor \frac{k-1}{2^6} \right\rfloor \right)$$

$$= 0.35 \sum_{k=1}^{\infty} (0.95)^{k-1} + 0.00078125 \cdot \sum_{k=1}^{\infty} (k-1)(0.95)^{k-1}$$

$$|E| = 0.35 \sum_{k=1}^{\infty} (0.95)^{k-1} + 7.8125 \times 10^{-4} \sum_{k=1}^{\infty} (0.95)^{k-1} (k-1)$$

$$= 0.35 \cdot \left(\frac{1}{1-x} \right) + 7.8125 \times 10^{-4} \cdot \frac{d}{dx} \left(\frac{1}{1-x} \right)$$

where $x = 0.95$ &

$$\sum_{k=1}^{\infty} x^k = \frac{1}{1-x}$$

$$= \frac{0.35}{1-0.95} + \frac{7.8125 \times 10^{-4}}{(1-0.95)^2}$$

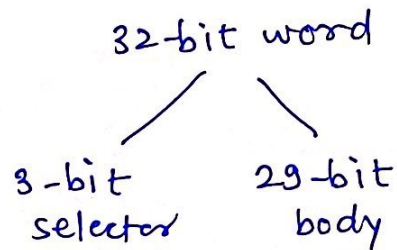
$$|E| = 7.3125$$

Thus, expected no. of bits per codeword is 7.3125 bits.

(Reference - Homework 4 (2012))

5. Comp up with a variation on simple-g that uses only a 3-bit selector rather than 4 bit selector.

Solution - If we use simple-g with a 3 bit selector on Δ -values with 32-bit integer format, we get 29 bits for the body and 3 bits for selector



Following table shows all possible ways to divide the 29bit body into chunks of equal size

selector	0	1	2	3	4	5	6	7
Number of Δ 's (1	2	3	4	7	9	14	29
Bits per Δ	29	14	9	7	4	3	2	1
Unused bits/word	0	1	2	1	1	2	1	0

Case when above distribution is better -

As, shown in the table simple-g with 3 bit selector can encode 29 1-bit numbers in a single word (32-bit) as compared to simple-g with 4 bit selector which can only accommodate 28 1-bit numbers.

Case when simple-g is better:

If a Gap distribution has varying integers, (i.e., more number of selectors are useful), then simple-g with 4-bit selector would be better than simple-g with 3-bit selector