

Content Based compression :-

Prediction with Partial Match (PPM) :-

- ∴ it is the best known content Based algorithm.
- ∴ it has been not as popular as the Ziv-Lempel based algorithms mainly because of the faster execution speed of latter algo.
- ∴ the idea of PPM algorithm is very simple.
- ∴ we would like to use large contents to determine the probability of the (content) symbol being encoded.
- ∴ instead of estimating these probabilities ahead of time, we can reduce the burden by estimating the probabilities as the coding proceeds.
- ∴ this way we need to store only those contents that have occurred in the sequence being coded.

Algorithm:-

- ∴ if the symbol to be encoded has not previously been encountered in this content an escape symbol is encoded and the algorithm attempts to use the next smaller content.
- ∴ if the symbol has not occurred in this content either, the size is further reduced.

∴ This process continues until either (31)
 we obtain a content that has previously
 been encountered with this symbol or
 we arrive at the conclusion that the
 symbol has not been encountered previously
 in any content.

The Escape Symbol:-

∴ Count of one for the escape symbol, thus
 inflating the total count in each content
 by 1. This is called as PPM.

e.g. Suppose in a given sequence 'a' occurs
 10 times in the content of probability 'l'.
 occurs 9 times and '0' occurs ~~3~~ 3 times
 in the same content.

∴ In method A, we assign a count of 'one'
 to the escape symbol, resulting in a
 total count of 23, which is one more than
 the no. of times prob has occurred

Counts using Method A		
content	symbol	count
Prob	a	10
	l	9
	0	3
Total Count	(ESC)	<u>23</u>

Counts using Method B		
content	symbol	count
Prob	a	9
	l	8
	0	2
<esc>		3
Total Count		<u>22</u>

Count using method c

content	symbol	count
Prob	a	10
	l	9
	o	3
	<ESC>	3
Total count		25

- = In a second method known as method 'B', we reduce the count of each of the symbols a, l, o by one and give the escape symbol a count of three resulting in the counts.
- = A variant of method 'B' appropriately named method 'c' was proposed by Moffat.
- = In method 'c', the count assigned to the escape symbol is the number of symbols that have occurred in that content.
- = In this content, method 'c' is similar to method B. The difference comes in the fact that instead of taking this from the counts of individual symbols, the total count is inflated by this amount
- ~~method~~ = while there is some variation in the performance depending on the characteristics of the data being encoded of the three methods for assigning counts to the escape symbol on an avg, method 'c' seems to provide best performance.

#) Length of Content :-

- A longer maximum length will usually result in a higher probability if the symbol to be encoded has a non-zero count with respect to that content.
- However, a long maximum length also means a higher probability of long sequence of escapes which in turn can increase the no. of bits used to encode the sequence.
- If we plot the compression performance vs. max. content length, we see an initial sharp increase in performance until some value of max length, followed by a steady drop as the max. length is further reduced.
- The value at which we see a downturn in performance changes depending on the characteristic of the source sequence.

#) The Exclusion Principle :-

- The basic idea behind arithmetic coding is the division of the unit interval into sub-intervals, each of which represents a particular letter.

- ∴ the smaller the subinterval, the more bits are required to distinguish it from other subintervals.
 - ∴ if we can reduce the no. of symbols to be represented, the no. of subintervals goes down as well.
 - ∴ this in turn means that the size of the subintervals increase, leading to a reduction in the no. of bits required for encoding.
 - ∴ the exclusion principle used in the PPM provides this kind of reduction in rate.
-

The Burrows - Wheeler Transform (BWT)

- The BWT algorithm also uses the content of the symbol being encoded in a different way for lossless compression.
- BWT algorithm requires that the entire (compression) sequence to be coded be available to the encoder before the coding takes place.

BWT Algorithm :-

- Given a sequence of length N , we create $N-1$ other sequences where each of these $N-1$ sequences is a cyclic shift of the original sequence.
- These N sequences are arranged in lexicographic order.
- The encoder then transmits the sequence of length N created by taking the last letter of each sorted cyclically shifted sequence.
- This sequence of last letters and the position of the original sequence in the sorted list are coded and sent to the decoder.

- This information is sufficient to recover the original sequence at decoder side.
- We start with an sequence of length N' and end with a representation that contains $N+1$ elements. However, this sequence has a structure that makes it highly amenable to compression.
- In particular we will use a method of coding called move-to-front (mtf) which is particularly effective on the type of structure exhibited by the sequence L .

Move to front Coding :-

- A coding scheme that takes advantage of long runs of identical symbols is the move-to front (mtf) coding.
- In this coding scheme, we start with some initial listing of the source alphabet.
- The symbol at the top of the list is assigned the number 0, the next one is assigned 1 and so on.
- The first time a particular symbol occurs, the number corresponding to its place in the list

is transmitted, then it is moved to the top of the list.

- If we have a run of these symbols, we transmit a sequence of 0's. This way long runs of different symbols get transformed to a large number of 0's.

example (BWT) :-

eg Lets encode the sequence
this is the

- Start with all the cyclic permutations of the given sequence. As here $N = 11$ so there are 11 permutations.

Permutations of this is the											
0	t	h	i	s	p	i	s	p	t	h	e
1	h	i	s	p	i	s	p	t	h	e	t
2	i	s	p	i	s	p	t	h	e	t	h
3	s	p	i	s	p	t	h	e	t	h	i
4	p	i	s	p	t	h	e	t	h	i	s
5	i	s	p	t	h	e	t	h	i	s	p
6	s	p	t	h	e	t	h	i	s	p	i
7	p	t	h	e	t	h	i	s	p	i	s
8	t	h	e	t	h	i	s	p	i	s	p
9	h	e	t	h	i	s	p	i	s	p	t
10	e	t	h	i	s	p	i	s	p	t	h

- Now sort these sequences in lexicographic (dictionary) order —

sequences sorted into lexicographic Order

(38)

	s	s	i	t	n	e	t	n	i	s
0	b	i	s	b	t	n	e	t	n	i
1	b	+	h	e	+	h	i	s	b	s
2	e	+	h	i	s	b	i	s	b	h
3	h	e	+	h	i	s	s	b	s	t
4	h	i	s	b	i	s	s	b	t	h
5	i	s	b	i	s	b	e	t	s	b
6	i	s	b	t	n	b	t	h	s	i
7	s	b	i	s	b	t	h	e	b	i
8	s	b	t	h	e	t	h	i	s	b
9	t	h	e	t	h	i	s	b	t	s
10	t	h	i	s	b	i	s	b	t	e

- L

→ the sequence of last letters L in this case is
L = sshtthb*iibe*

- the original sequence appears as sequence no. 10 in the sorted list.
- so encoding of the sequence consists of the sequence L and the index value 10.

⇒ Decoding Procedure :- [continue with previous example.]

- Now lets decode the original sequence by using the sequence 'L' and index to the original sequence in the sorted list.
- All the elements of the initial sequence are contained in L, just need to figure out the permutation that will recover the original seq.

① Obtain the permutation is to generate the sequence F consisting of first element of each row. Therefore the sequence F is simply the sequence L in lexicographic order. so

$$F = \underline{\text{K} \text{K} \text{e} \text{h} \text{h} \text{i} \text{i} \text{s} \text{t} \text{t}}$$

Now use F and L to decode the sequence.

If we know that the original sequence is in the k^{th} row then we can begin unravelling the original sequence starting with k^{th} element of F

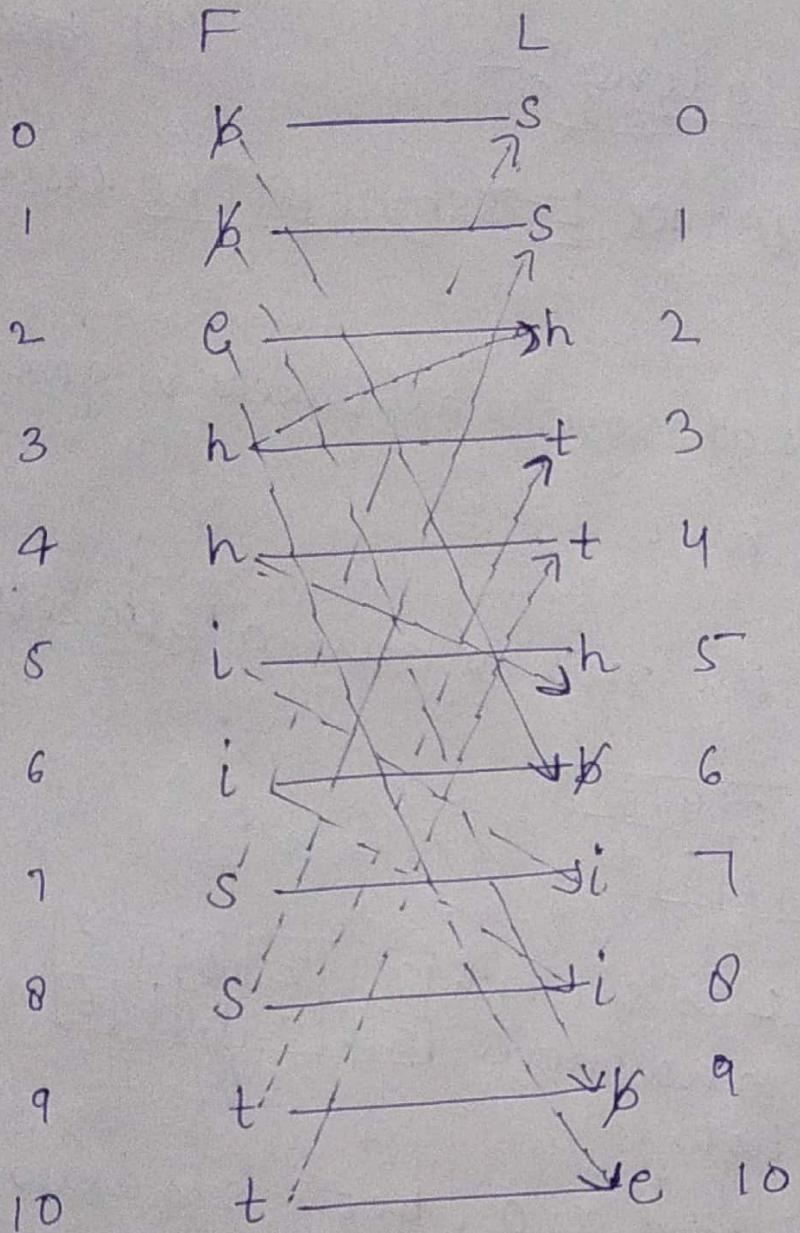
$$F = \begin{bmatrix} \text{K} \\ \text{K} \\ \text{e} \\ \text{h} \\ \text{i} \\ \text{i} \\ \text{s} \\ \text{s} \\ \text{t} \\ \text{t} \end{bmatrix}$$

$$L = \begin{bmatrix} \text{s} \\ \text{s} \\ \text{h} \\ \text{t} \\ \text{t} \\ \text{h} \\ \text{K} \\ \text{i} \\ \text{i} \\ \text{K} \\ \text{e} \end{bmatrix}$$

The original seq is seq no. 10, so the first letter in of the original seq is $\boxed{F[10] = t}$.

To find the letter following t we look for t in the array L.

There are 2 t's in L so the t in F that was we are working is lower of 2 t's so we consider lower of 2 t's in L.



(40)

\therefore hence lowest
+ is at $L[4]$.

\therefore therefore, the
next letter in
our recommended
sequence is
 $F[4] = h$.

\therefore the reconstructed
sequence is th .

\therefore To find the next
letter, we look
in L .

\therefore Repeat this
process.

\equiv

decoding sequence

L index $2[4] 1[5] \quad L[7] [0] [6] [8] [1] [9] [3] [2] [10]$
 $t \quad h \quad .i \quad s \quad p \quad i \quad s \quad p \quad + \quad h \quad e$
 F index $F[10] [4] \quad F[5] [7] [0] [6] [8] [1] [9] [3] [2]$

Hence the reconstructed sequence is

this p is $p + h e$

\equiv

Example - MTF Coding :-

(41)

~~(42)~~

e.g. encode the sequence $L = \text{sshtttkike}$ using MTF coding.

B Let's assume that source alphabets is given by $A = \{\text{k, e, h, i, s, t}\}$

① first we assign the nos to our alphabets

0	1	2	3	4	5
s	e	h	i	s	t

② first letter of the L is 's' which is encoded as 0. and move s to the top of list.

0	1	2	3	4	5
s	k	e	h	i	t

③ the next s is encoded as 0, bcoz s is already at the top of the list, hence ~~not~~ encode. seq is 40.

④ The next letter is h, which is encoded as 3, & move h at the top of list, seq becomes 403.

0	1	2	3	4	5
h	s	k	e	t	

⑤ the next letter is t - 5, move t at top of list hence seq - 4035

0	1	2	3	4	5
t	h	s	k	e	

⑥ continue the procedure for further 9/P seq, hence the encoded seq is 40350135015

*) '0' means ^{previous} letter is repeat. //