

- Arithmetic Coding :- Generating variable length codes through arithmetic coding method.
- Coding Rate :- Average nos of bits used to represent a symbol from a source.
- Entropy :- lowest rate at which the source can be coded.
- P_{max} :- Probability of most frequently occurring symbol

eg consider a source that puts out iid letters from the alphabet $A = \{a_1, a_2, a_3\}$ with the probability model $P(a_1) = 0.95$, $P(a_2) = 0.02$ and $P(a_3) = 0.03$. Entropy for this source is 0.335 bits /symbol. Average length is = 1.05 bits /symbol.

<u>symbol</u>	<u>codeword</u>
a_1	0
a_2	11
a_3	10

- the difference b/w the average code length and entropy is 0.715 bits /symbol is redundancy

(2)

- which is 213% of the entropy.
- Another approach is extended alphabet approach.
 - We can see that it is more efficient to generate codeword for group or sequences of symbols rather than generating a separate codeword for each symbol in a sequence.
 - However, this approach becomes impractical when we try to obtain Huffman codes for long sequence.

~~III~~ Mathematics

- #### # Arithmetic Coding
- In arithmetic coding, a unique identifier or tag is generated for the sequence to be encoded.
- The tag corresponds to binary functions which becomes the binary code for the sequence.
 - A unique arithmetic code can be generated for a sequence of length 'm' without the need of generating codewords for all sequence of length 'm'.

II Coding a Sequence :- In order to distinguish a sequence of symbol from another sequence of symbols, we need a unique tag with unique identifier.

- ∴ One possible set of tags for representing sequences of symbols are the no's in the unit interval $[0, 1]$. For this we need a function that will map sequences of symbols into unit interval.
- ∴ A function that maps random variables and sequences of random variables, into the unit interval is the cumulative distribution function (cdf) of the random variable associated with the source.
- ∴ cdf will be used in developing in arithmetic codes.

IIa Generating a Tag :-

- ∴ Tag works by reducing the size of the interval in which tag resides as more and more elements of the sequence are received.

- We start by first dividing the unit interval into subintervals of the form
- $[F_n(i-1), F_n(i)]$, $i=1, 2 \dots m$, as the minimum value of cdf is '0' and maximum value is '1' this exactly partitions the unit interval.
- We associate the subinterval $[F_n(i-1), F_n(i)]$ with symbol a_i . The 'j'th interval corresponding to the symbol a_j is given by

$$\left[\frac{F_n(k-1) + F_n(j-1)}{(F_n(k) - F_n(k-1))}, \frac{F_n(k-1) + F_n(j)}{(F_n(k) - F_n(k-1))} \right]$$

eg :- consider a 3-letter alphabet $A = \{a_1, a_2, a_3\}$ with $P(a_1) = 0.7$, $P(a_2) = 0.1$, $P(a_3) = 0.2$, also $F_n(1) = 0.7$, $F_n(2) = 0.8$ and $F_n(3) = 1$

Sol :- the partition in which the tag resides depends on the first symbol of the sequence being encoded. for eg - if the first symbol is a_1 , the tag lies in the interval $[0.0, 0.7]$; if the symbol is a_2 , than tag lies in the interval $[0.8, 1.0]$.

Once the interval containing, the tag has been determined, the rest of the unit interval

is discarded, and this restricted interval is again divided in some proportions as the original interval. ⑤

• Suppose the first symbol was a_1 , the tag would be contained in the subinterval $[0.0, 0.7]$.

• This subinterval is then subdivided in exactly the same proportions as the original intervals, yielding the subintervals $[0.0, 0.49], [0.49, 0.56]$ and $[0.56, 0.7]$.

• The first partition as before corresponds to a_1 , the second partition to a_2 . The tag value is then restricted to lie in the interval $[0.49, 0.56] \dots$ and so on.

⇒ Generating a Tag :- Define a random variable $X(a_i) = i$, suppose we wish to encode the sequence 1321. From the probability model, we know that

$$F_n(k) = 0, k \leq 0, F_n(1) = 0.8, F_n(2) = 0.82, F_n(3) = 1.$$

$$F_n(k) = 1, k > 3$$

initializing $l^{(0)}$ to 1, and $U^{(0)}$ to 0, the first element of sequence 1 results in the

(6)

following update

$$l^{(1)} = 0 + (1-0)0 = 0$$

$$u^{(1)} = 0 + (1-0)(0 \cdot 8) = 0 \cdot 8$$

i.e., the tag is contained in the interval $[0, 0 \cdot 8]$.

\Rightarrow the second element of the sequence is 3,
using the update equation we get

$$l^{(2)} = 0 + (0 \cdot 8 - 0) F_x^{(2)} = 0 \cdot 8 \times 0 \cdot 8^2 = 0 \cdot 656$$

$$u^{(2)} = 0 + (0 \cdot 8 - 0) F_x^{(3)} = 0 \cdot 8 \times 1 \cdot 0 = 0 \cdot 8$$

Therefore, the interval containing the tag
for the sequence #3 is $[0 \cdot 656, 0 \cdot 8]$.

\Rightarrow the third element 2, then

$$l^{(3)} = 0 \cdot 656 + (0 \cdot 8 - 0 \cdot 656) F_x^{(1)} *$$

$$= 0 \cdot 656 + 0 \cdot 144 \times 0 \cdot 8 = 0 \cdot 7712$$

$$u^{(3)} = 0 \cdot 656 + (0 \cdot 8 - 0 \cdot 656) F_x^{(2)}$$

$$= 0 \cdot 656 + 0 \cdot 144 \times 0 \cdot 82 = 0 \cdot 77408.$$

and the interval for the tag is $[0 \cdot 7712, 0 \cdot 77408]$.

\Rightarrow continuing with the last element, the
upper and lower interval for sequence 4th Benet
is $l^{(4)} = 0 \cdot 7712$ & $u^{(4)} = 0 \cdot 773504$.

(7)

and the Tag for the sequence 1321
can be generated as -

$$\bar{T}_n(1321) = \frac{0.7712 + 0.773504}{2}$$

$$\boxed{\bar{T}_n(1321) = \underline{0.772352}}$$

Reciphering the Tag :-

[Reverse procedure of Generating a tag]

⇒ The tag is useless unless we can also decipher it with minimal computational cost.

eg:- The tag value obtained in previous eg.
is 0.772352. The interval containing this tag value is a subset of every interval obtained in the encoding process.

For decoding procedure, we decode the element in the sequence in such a way that the upper and lower limits $u^{(k)}$ & $l^{(k)}$ will always contain the tag value for each k.

We start with $l^{(0)}=0$ & $u^{(0)}=1$. After decoding the first element of the sequence n_1 the upper & lower limits become -

$$l^{(1)} = 0 + (1-0) f_n(x_1 - 1) = f_n(x_1 - 1)$$

$$u^{(1)} = 0 + (1-0) f_n(x_1) = f_n(x_1)$$

ie the interval containing the tag is ~~is~~)

$$[f_n(x_1 - 1), f_n(x_1)]$$

* Now choose or find the value of n_1 for which 0.772352 lies in the interval $[f_n(x_1 - 1), f_n(x_1)]$, if $x_1 = 1$ the interval is $[0, 0.8]$, if $x_1 = 2$ the interval is $[0.8, 1.0]$, if $x_1 = 3$ the value of interval is $[0.82, 1.0]$,
 \Rightarrow As 0.772352 lies in the interval $[0.8, 1.0]$
we take $x_1 = 1$.

\Rightarrow Now repeat this procedure for 2nd, 3rd & 4th elements.

$$l^{(2)} = 0 + (0.8 - 0) f_n(x_2 - 1) = 0.8 f_n(x_2 - 1)$$

$$u^{(2)} = 0 + (0.8 - 0) f_n(x_2) = 0.8 f_n(x_2)$$

2

$$l^{(3)} = 0.656 + (0.8 - 0.656) f_n(x_3 - 1)$$

$$u^{(3)} = 0.656 + (0.8 - 0.656) f_n(x_3)$$

2

$$l^{(4)} = 0.7712 + (0.77408 - 0.7712) f_n(x_4 - 1)$$

$$u^{(4)} = 0.7712 + (0.77408 - 0.7712) f_n(x_4)$$

Hence the Algorithm for decipher the tag. (9)

1. Initialize $l^{(0)} = 0$ & $u^{(0)} = 1$
2. For each k find $t^* = (\text{tag} - l^{(k-1)}) / (u^{(k-1)} - l^{(k-1)})$
3. Find the value of x_k for which $F_n(x_k-1) \leq t^* < F_n(x_k)$
4. update $u^{(k)}$ & $l^{(k)}$
5. Continue until the entire sequence has been decoded.

⇒ There are two ways to know when the entire sequence has been decoded.

① The decoder may know the length of the sequence

② A particular symbol is denoted as an end-of-transmission symbol, the decoding of this symbol would bring the decoding process to a close.

example :- Encode "BILL GATES" using Arithmetic coding Technique

(1)

chr.	freq.
A	0.1
B	0.1
E	0.1
G	0.1
I	0.1
L	0.2
S	0.1
T	0.1

contd.

(10)

<u>Character</u>	<u>Probability</u>	<u>Range</u>
Space	0.1	[0.00, 0.10)
A	0.1	[0.10, 0.20)
B	0.1	[0.20, 0.30)
E	0.1	[0.30, 0.40)
G	0.1	[0.40, 0.50)
I	0.1	[0.50, 0.60)
L	0.2	[0.60 - 0.80)
S	0.1	[0.80 - 0.90)
T	0.1	[0.90 - 1.00)

<u>char</u>	<u>low</u>	<u>high</u>
B	0.0	1.0
B	0.2	0.3
I	0.25	0.26
L	0.256	0.258
L Space	0.2572	0.2576
Space	0.25720	0.25724
G	0.257216	0.257220
A	0.2572164	0.2572168
T	0.25721676	0.25721688
E	0.257216772	0.257216776
S	0.2572167752	0.2572167756

hence the final low value - 0.2572167752 will uniquely encode the name "BILL GATES." //

Arithmetic Coding Vs Huffman Coding -

- ⇒ Arithmetic coding requires less m/o as symbol representation is calculated on the fly.
- ⇒ It is more suitable for high performance models where there are confident predictions.
- ⇒ In large collections of text and images, Huffman coding is likely to be used for text, and its redundancy is quite small.
- ⇒ ~~Arithmetic~~ Arithmetic coding is used for the images and much better than Huffman coding unless a blocking technique is used.
- ⇒ Huffman decoding is generally faster than arithmetic decoding.
- ⇒ In arithmetic coding it is not easy to start decoding in the middle of the stream, while in Huffman coding we can use.
- ⇒ A major advantage of arithmetic coding over Huffman coding is the ability to separate the modeling and coding aspects of the compression approach.

Generating a Binary Code :-

- = We want a binary code that will represent (x) in a unique and different way in a efficient manner.
- = 'Tags' forms a unique representation for the sequence. This means that the binary representation of the tag forms a unique binary code for all sequence.

Uniqueness and Efficiency of the Arithmetic Code

- = $\bar{T}_x(x)$ is all numbers in the interval $[0, 1]$. A binary code for $\bar{T}_x(x)$ can be obtained by taking the binary representation of this no. and truncating it to $l(x) = \lceil \log \frac{1}{P(x)} \rceil + 1$ bits.

e.g. consider a source A that generates letters from an alphabet of size four

$A = \{a_1, a_2, a_3, a_4\}$ with Probability

$$P(a_1) = \frac{1}{2}, P(a_2) = \frac{1}{4}, P(a_3) = \frac{1}{8}, P(a_4) = \frac{1}{8}$$

\Rightarrow A Binary code for four letter alphabets

(13)

Symbol	f_n	\bar{T}_n	In Binary	$\lceil \log \frac{1}{P_n} \rceil + 1$	Code
1.	0.5	0.25	0.010	2	01
2.	0.75	0.625	0.101	3	101
3.	0.875	0.8125	0.1101	4	1101
4.	1.0	0.9375	0.1111	4	1111

Here we can see that, a code obtained in this fashion is a unique decodable code.

=

Dictionary Techniques :-

- Here we will look at techniques that incorporate the structure in the data in order to increase the amount of compression.
 - These techniques - both static and dynamic or Adaptive build a list of commonly occurring patterns and encode these patterns by transmitting their index in the list.
- ⇒ Introduction :- some application, the O/P of the source consists of recurring patterns. A classic eg is a text source in which certain patterns or words occur constantly. Some patterns do not occur or if they do, occur with rarity. So techniques for handling these 2 data is different.

i) static Approach

ii) Adaptive Approach.

I) static Approach :- (Dictionary)

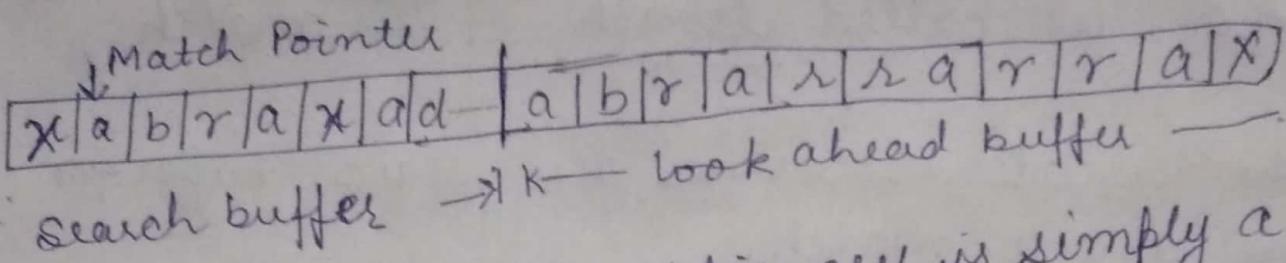
⇒ This approach is best if we have prior knowledge about the source.

- (15)
- if the task were to compress student records at university, a static dictionary approach may be the ~~best~~ best.
 - the reason is - we know ahead of time that certain words such as "Name" and "Student-id" are going to appear in almost all the records.
 - These schemes would work well only for the application and data they were designed for.

II.) Adaptive Dictionary :-

- Most of the adaptive dictionary-based techniques have their roots into landmark papers by Jacob Ziv and Abraham Lempel, in 1977 and 1978.
- These papers provide two different approaches to adaptively building dictionaries and each approach has given rise to number of variations.
- The approach based on 1977 papers are said to belong to LZ77 family, while approaches based on the 1978 papers are said to belong to LZ78 or LZ2 family.

IIa) The LZ77 Approach :-



- In LZW approach, dictionary is simply a portion of the previously encoded sequence.
- The sliding window has the input sequence, encoder examines the sequence.
- Window has 2 parts :-
 - 1. Search Buffer :- contains portion of recently encoded sequence
 - 2. Lookahead Buffer :- contains next sequence of the buffer to be encoded
- To encode the sequence in the lookahead buffer, the encoder moves a search pointer back through the search buffer until it encounters the first match to the first symbol in the lookahead buffer.
- offset - It is the distance of the pointer from the lookahead buffer.
- length of the match - the number of consecutive symbols in the search buffer that match consecutive symbols in the lookahead buffer, starting with the first

symbol, is called length of the match.

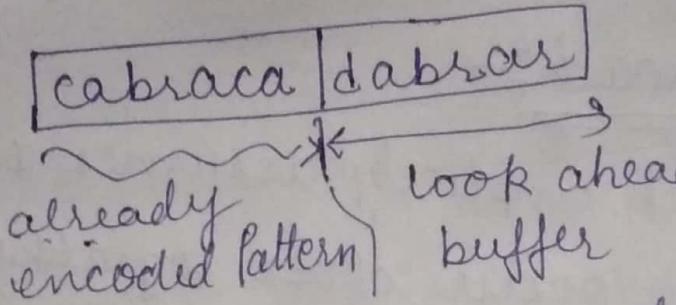
- The encoder searches the search buffer for the longest match.
- Once the longest match has been found, the encoder encodes it with a triple $\langle o, l, c \rangle$, where ' o ' is the offset, ' l ' is the length of the match and ' c ' is the codeword corresponding to the symbol in the look-ahead buffer that follows the match.
- If the size of the search buffer is ' S ', the size of the window (search and look ahead buffer) is ' w ' and the size of the source alphabet is ' A ', then the number of bits needed to code the triple using fixed length code is $\lceil \log_2 S \rceil + \lceil \log_2 w \rceil + \lceil \log_2 A \rceil$

$$\boxed{\lceil \log_2 S \rceil + \lceil \log_2 w \rceil + \lceil \log_2 A \rceil}$$

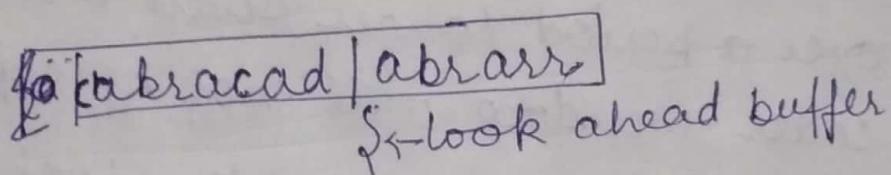
Eg:- using the LZ77 approach-

Suppose the sequence to be encoded is
cabracadabrarraarrad.

and suppose the length of window is 13, the size of the look ahead buffer is 6.

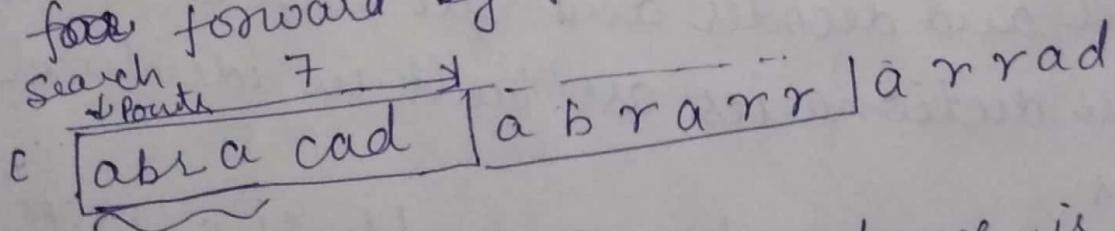


sol.: we look back in the already encoded portion of the window to find a match for d. There is no match, so we transmit the triple $\langle 0, 0, c(d) \rangle$. $c(d)$ is the code for character d.



now we find a match for \underline{a} at an offset of 2. The length of this match is 1.

so now we encode the string 'abra' with the triple $\langle 7, 4, c(r) \rangle$ and move the window forward by five characters.



$$l=4$$

As we see that, LZW scheme is very simple adaptive scheme that requires no prior knowledge of the source and seems to require no assumptions about the characteristics of the source.

IIb The LZ 78 approach:-

- The LZ-77 approach implicitly assumes that like patterns will occur close together.
- It makes use of this structure by using the recent part of the sequence as the dictionary for encoding.
- However this means that any pattern that occurs over a period longer than that covered by the code's window will not be captured.
- The LZ-78 algorithm solves this problem by dropping the reliance on the search buffer and keeping an explicit dictionary.
- This dictionary has to be built at both the encoder and decoder and care must be taken that the dictionaries are built in identical manner.
- The input are coded as a double $\langle i, c \rangle$ with ' i ' being an index corresponding to the dictionary entry that was longest match to the input, and ''c'' being the code for the character in the input following the match portion of the input. '0' is used in case of no match.

eq :- (LZ-78 approach)

20

The sequence is -

The sequence is
wabbawabbawabbawabbawabbawoobwoobwoobwoob²

Initial dictionary $\{ \}$ while LZ-78 algo has
Index Entry $\{ \}$ due ability to capture
 1 w patterns and hold them
 2 a indefinitely. It also has a
 3 b serious drawback, dictionary
 keeps growing without bound.

Development of Dictionary

Encoder	Index	Entry	Encoder	Index	Entry
$\langle 0, c(w) \rangle$	1	w	$\langle 1, c(0) \rangle$	15	w0
$\langle 0, c(a) \rangle$	2	a	$\langle 14, c(w) \rangle$	16	opw
$\langle 0, c(b) \rangle$	3	b	$\langle 13, c(0) \rangle$	17	oo
$\langle 3, c(a) \rangle$	4	ba			
$\langle 0, c(\beta) \rangle$	5	β			
$\langle 1, c(a) \rangle$	6	wa			
$\langle 3, c(b) \rangle$	7	bb			
$\langle 2, c(\beta) \rangle$	8	$\alpha\beta$			
$\langle 6, c(b) \rangle$	9	wab wab			
$\langle 4, c(b) \rangle$	10	ba β			
$\langle 9, c(b) \rangle$	11	wabb			
$\langle 8, c(w) \rangle$	12	$\alpha\beta w$			
$\langle 0, c(0) \rangle$	13	o			
$\langle 13, c(b) \rangle$	14	$\delta\beta$			

IIc) ZZW Encoding :-

- Let us take the previous sequence -

• Let us take the previous -
"wabba wabba wabba wabba kooookwookwook"

\therefore The alphabet for the source is = { p, a, b, o, w }

- the LZW dictionary initially looks like

<u>Index</u>	<u>Entry</u>
1	b
2	a
3	b
4	c
5	w.

⇒ The LZW dictionary for encoding for the given sequence - "wabba - - - ... oo"

<u>order</u>	<u>entry</u>	<u>order</u>	<u>entry</u>
1	b	13	bba
2	a	14	abw
3	b	15	wabb
4	o	16	baþ
5	w	17	bwa
6	wa	18	abb
7	ab	19	bablo
8	bb	20	w0
9	ba	21	oo
10	ab	22	oþ
11	þw	23	þwo
12	wab	24	ooþ
		25	þwooo

Hence our encoder O/P sequence is -

"5 2 3 3 2 1 6 0 1 0 1 2 9 1 1 7 1 6 5 4 4 1 1 2 1 2 3 4".

II(d) LZW-Decoding :-

- the O/P sequence of LZW-encoding procedure is -

- first ~~also~~ take the encoded sequence -

"5 2 3 3 2 1 6 0 1 0 1 2 9 1 1 7 1 6 5 4 4 1 1 2 1 2 3 4"

- Now constructing ^{11th index of} LZW dictionary while decoding - .

<u>Index</u>	<u>Entry</u>
1	∅
2	a
3	b
4	o
5	w
6	wa
7	ab
8	bb
9	ba
10	ab∅
11	∅...

and so on -

- ∴ the index value '5' corresponds to the letter w;
- so we decode w as the first element of our sequence.
- ∴ the next decode input is 2, which is corresponds to letter a.
- ∴ we decode an a and concatenate it with our current pattern and become the pattern 'wa'. As this does not exist with dictionary, so we add it at the 6th element of the dictionary and start new pattern begining with letter 'a'
- ∴ the next four inputs 3 3 2 1 correspond to the letters b b a ϕ and generate the dictionary entries ab, bb, ba, $\alpha\phi$
- ∴ the next input is '6', which is the index of the pattern 'wa'. Therefore we decode 'aw' and an 'a'; we first concatenate 'w' to the existing pattern, which is ϕ and form the pattern ϕw .
- ∴ similarly do for further sequence.

III Applications of Dictionary Techniques

- ⇒ Levy Welch's article publication has been steadily increasing number of applications that use some variant of the LZ78 algorithm.
- ⇒ Among the LZ78 variant by far, the most popular is the LZW algorithm.
- ⇒ There are two best known applications of LZW is -
 - GIF and V.42 bits.

File Compression - Unix Compress :-

- ⇒ the unix compress command is one of the earlier applications of LZW.
- ⇒ the size of the dictionary is adaptive, we start with a dictionary of size 512 i.e. codeword are 9 bits long.
- ⇒ Once the dictionary is doubled to 1024 entries the codeword transmitted is now 10 bits.
- ⇒ The size of the dictionary is progressively doubled as it fills up.
- ⇒ The size of the dictionary
- ⇒ The maximum size of the codeword bmax, can be set up by the user to b/w 9 & 16 with

16 bits being the default. After ~~some~~⁵⁰⁰⁰ entries, ~~compress~~⁽²⁾ becomes a static dictionary coding technique.

- At this point the algorithm monitors the compression ratio. If the compression ratio falls below a threshold, the dictionary is flushed and the dictionary building process is restarted.

Image Compression - The Graphics Interchange Format (GIF) :-

- The GIF was developed by CompuServe Information Service to encode graphical images.
- It is another implementation of LZW algorithm and is very similar to the compress command.
- The compressed image is stored with the first byte being the min. number of bits ~~16~~⁸ b/pixel in the original image.
- The binary number 2^b is defined to be the clear code. This code is used to reset all compression & decompression parameter to a start-up state.

- the initial size of the dictionary is 2^{b+1} , when this fills up, the dictionary size is doubled as was done in compress algorithm.
- the codeword from the LZW algorithm are stored in blocks of character.
- the characters are 8 bits long and maximum block size is 255.
- Each block is preceded by a header that contains the block size.
- GIF has become quite popular for encoding all kinds of images both computer generated and natural generated images.
- even if we account for the extra overhead in the GIF files, for these images GIF barely holds its own even with simple arithmetic coding of the original pixels.
- therefore, for images like these, the straight forward dictionary coding is done.

(27)

Image compression — Portable N/w Graphics (PNG) —

- The PNG standard is one of the first standards to be collaboratively developed over the internet.
- The community decided that a patent free replacement for GIF should be developed, and within 3 months, PNG took place.
- Unlike GIF, the compression algorithm used in PNG is based on LZ-77.
- In particular it is based on the deflate implementation of LZ-77.
- This implementation allows for match length of b/w 3 & 258. At each step the encoder examines 3 bytes. If it cannot find a match of 3 bytes, it puts the first byte and examine next 3 bytes.
- An initial estimate of the pixel is first made by adding the pixel to the left and the pixel that is above and subtracting the pixel to the upper left. Then the pixel that is closer to the initial estimate is taken as the estimate.

→ Comparison of PNG with GIF and arithmetic coding :-

Image name (sample)	PNG	GIF	Arithmetic code of Pixel value	Arithmetic code of Pixel difference
Sena	31,577	51,085	53,431	31,847
Sensin	34,488	60,649	58,306	37,126
Earth	26,995	34,276	38,248	32,137
Omaha	50,185	61,580	58,061	51,393

Comparison over Modems - V.42 bits :-

- The ITU-T Recommendation V.42 bits is a compression standard devised for use over a telephone network along with error correcting procedures described in CCITT recommendation V.42.
- This algorithm is used in modems connecting computers to remote users.
- The reason for the existence of two modes is that at times the data being transmitted do not have repetitive structure and therefore cannot be compressed using the

LZW algorithm.

- The V.42 bits recommendation suggest periodic testing of the O/P of the compression algorithm to see if data expansion is taking place.
- In the compressed mode, the system uses LZW compression with a variable size dictionary.
- The initial dictionary size is negotiated at the time a link is established b/w the transmitter and receiver.
- When the numbers of entries in the dictionary exceed a prearranged threshold C_3 , the encoder size is incremented by 1 bit.
At the same time, the threshold C_3 is also doubled.

Note :

- Dictionary-based algo. are being used to compress all kinds of data; ~~however~~
- This approach is most useful when structural constraints restrict the frequently occurring patterns to a small subset of all possible patterns.