# CSE 326: Data Structures

## Dictionaries for Data Compression

Neva Cherniavsky
Summer 2006

---

# Why compress files?

---

# What is a file?

---

# Data Compression

| original X | → | Encoder | compressed Y | → | Decoder | decompressed X' → |

- **Lossless** compression  X = X'
- **Lossy** compression  X != X'
- **Compression Ratio**  |X|/|Y|
  - Where |X| is the # of bits in X.

---

# Dictionary Coding

- Does not use statistical knowledge of data.
- Encoder: As the input is processed develop a dictionary and transmit the index of strings found in the dictionary.
- Decoder: As the code is processed reconstruct the dictionary to invert the process of encoding.
- Examples: LZW, LZ77, Sequitur
- Applications: Unix Compress, gzip, GIF

---

# LZW Encoding Algorithm

```
Repeat
  find the longest match w in the dictionary
  output the index of w
  put wa in the dictionary where a was the
      unmatched symbol
```

## LZW Encoding Example (1)

Dictionary

0  a
1  b

a b a b a b a b a

## LZW Encoding Example (2)

Dictionary

0  a
1  b
2  ab

<u>a</u> b a b a b a b a
0

## LZW Encoding Example (3)

Dictionary

0  a
1  b
2  ab
3  ba

<u>a</u> <u>b</u> a b a b a b a
0 1

## LZW Encoding Example (4)

Dictionary

0  a
1  b
2  ab
3  ba
4  aba

<u>a</u> <u>b</u> <u>a</u> b a b a b a
0 1 2

## LZW Encoding Example (5)

Dictionary

0  a
1  b
2  ab
3  ba
4  aba
5  abab

<u>a</u> <u>b</u> <u>a</u> <u>b</u> <u>a</u> <u>b</u> <u>a</u> b a
0 1 2    4

## LZW Encoding Example (6)

Dictionary

0  a
1  b
2  ab
3  ba
4  aba
5  abab

<u>a</u> <u>b</u> <u>a</u> <u>b</u> <u>a</u> <u>b</u> <u>a</u> <u>b</u> <u>a</u>
0 1 2    4    3

## LZW Decoding Algorithm

- Emulate the encoder in building the dictionary.
  Decoder is slightly behind the encoder.

```
initialize dictionary;
decode first index to w;
put w? in dictionary;
repeat
    decode the first symbol s of the index;
    complete the previous dictionary entry with s;
    finish decoding the remainder of the index;
    put w? in the dictionary where w was just decoded;
```

13

## LZW Decoding Example (1)

Dictionary

0  a
1  b
2  a?

$\underline{0}$ 1 2 4 3 6
a

14

## LZW Decoding Example (2a)

Dictionary

0  a
1  b
2  ab

0 $\underline{1}$ 2 4 3 6
a b

15

## LZW Decoding Example (2b)

Dictionary

0  a
1  b
2  ab
3  b?

$\underline{0}$ 1 2 4 3 6
a b

16

## LZW Decoding Example (3a)

Dictionary

0  a
1  b
2  ab
3  ba

0 1 $\underline{2}$ 4 3 6
a b a

17

## LZW Decoding Example (3b)

Dictionary

0  a
1  b
2  ab
3  ba
4  ab?

$\underline{0}$ 1 $\underline{2}$ 4 3 6
a b ab

18

## LZW Decoding Example (4a)

Dictionary

0 a
1 b
2 ab
3 ba
4 aba

0 1 2 4 3 6
a  b  ab  a

19

## LZW Decoding Example (4b)

Dictionary

0 a
1 b
2 ab
3 ba
4 aba
5 aba?

0 1 2 4 3 6
a  b  ab  aba

20

## LZW Decoding Example (5a)

Dictionary

0 a
1 b
2 ab
3 ba
4 aba
5 abab

0 1 2 4 3 6
a  b  ab  aba  b

21

## LZW Decoding Example (5b)

Dictionary

0 a
1 b
2 ab
3 ba
4 aba
5 abab
6 ba?

0 1 2 4 3 6
a  b  ab  aba  ba

22

## LZW Decoding Example (6a)

Dictionary

0 a
1 b
2 ab
3 ba
4 aba
5 abab
6 bab

0 1 2 4 3 6
a  b  ab  aba  ba  b

23

## LZW Decoding Example (6b)

Dictionary

0 a
1 b
2 ab
3 ba
4 aba
5 abab
6 bab
7 bab?

0 1 2 4 3 6
a  b  ab  aba  ba  bab

24

## Decoding Exercise

Base Dictionary

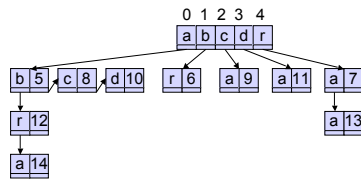0 1 4 0 2 0 3 5 7

```
0 a
1 b
2 c
3 d
4 r
```

25

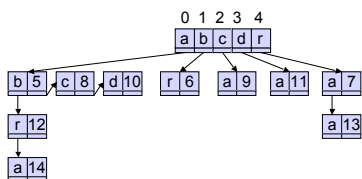## Trie Data Structure for Encoder's Dictionary

- Fredkin (1960)

```
0 a   9 ca
1 b   10 ad
2 c   11 da
3 d   12 abr
4 r   13 raa
5 ab  14 abra
6 br
7 ra
8 ac
```
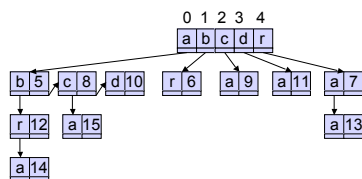


26

## Encoder Uses a Trie (1)



a b r a c a d a b r a a b r a c a d a b r a
0 1 4 0 2 0 3 5 7  12

27

## Encoder Uses a Trie (2)



a b r a c a d a b r a a b r a c a d a b r a
0 1 4 0 2 0 3 5 7  12  8

28

## Decoder's Data Structure

- Simply an array of strings

```
0 a   9 ca
1 b   10 ad
2 c   11 da
3 d   12 abr
4 r   13 raa
5 ab  14 abr?
6 br
7 ra
8 ac
```

0 1 4 0 2 0 3 5 7 12  8 ...
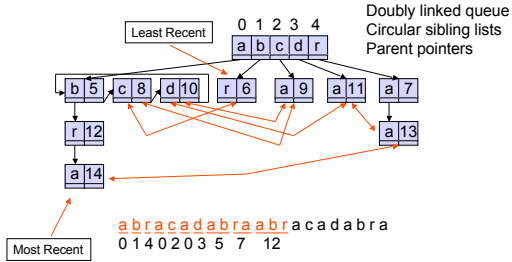a b r a c a d ab ra abr

29

## Bounded Size Dictionary

- Bounded Size Dictionary
  - n bits of index allows a dictionary of size $2^n$
  - Doubtful that long entries in the dictionary will be useful.
- Strategies when the dictionary reaches its limit.
  1. Don't add more, just use what is there.
  2. Throw it away and start a new dictionary.
  3. Double the dictionary, adding one more bit to indices.
  4. Throw out the least recently visited entry to make room for the new entry.
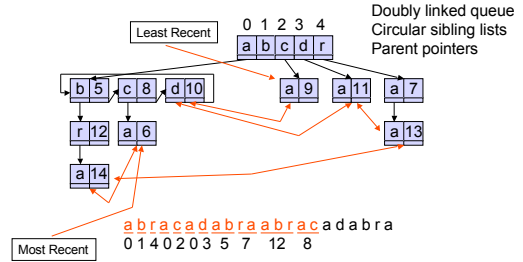
30

## Implementing the LRV Strategy



Least Recent
0 1 2 3 4
a b c d r

Doubly linked queue
Circular sibling lists
Parent pointers

b 5  c 8  d 10  r 6  a 9  a 11  a 7
r 12
a 14
a 13

Most Recent

a b r a c a d a b r a a b r a c a d a b r a
0 1 4 0 2 0 3 5 7   12

31

## Implementing the LRV Strategy



Least Recent
0 1 2 3 4
a b c d r

Doubly linked queue
Circular sibling lists
Parent pointers

b 5  c 8  d 10  a 9  a 11  a 7
r 12  a 6
a 14
a 13

Most Recent

a b r a c a d a b r a a b r a c a d a b r a
0 1 4 0 2 0 3 5 7   12  8

32

## Notes on LZW

- Extremely effective when there are repeated patterns in the data that are widely spread.
- Negative: Creates entries in the dictionary that may never be used.
- Applications:
  - Unix compress, GIF, V.42 bis modem standard

33

## LZ77

- Ziv and Lempel, 1977
- Dictionary is implicit
- Use the string coded so far as a dictionary.
- Given that $x_1x_2...x_n$ has been coded we want to code $x_{n+1}x_{n+2}...x_{n+k}$ for the largest k possible.

34

## Solution A

- If $x_{n+1}x_{n+2}...x_{n+k}$ is a substring of $x_1x_2...x_n$ then $x_{n+1}x_{n+2}...x_{n+k}$ can be coded by <j,k> where j is the beginning of the match.
- Example

  abababab babababababababab....
      coded

  abababab babababa babababab....
          <2,8>

35

## Solution A Problem

- What if there is no match at all in the dictionary?

  abababab cabababababababab....
      coded

- Solution B. Send tuples <j,k,x> where
  - If k = 0 then x is the unmatched symbol
  - If k > 0 then the match starts at j and is k long and the unmatched symbol is x.

36

## Solution B

- If $x_{n+1}x_{n+2}...x_{n+k}$ is a substring of $x_1x_2...x_n$ and $x_{n+1}x_{n+2}...x_{n+k}x_{n+k+1}$ is not then $x_{n+1}x_{n+2}...x_{n+k}x_{n+k+1}$ can be coded by
$$<j,k, x_{n+k+1} >$$
where j is the beginning of the match.
- Examples

<u>abababab</u> cababababababababab....

<u>abababab</u> <u>c</u> <u>ababababab</u> ababab....
<0,0,c>  <1,9,b>

37

## Solution B Example

<u>a</u> babababababababababababab.....
<0,0,a>

<u>a b</u> abababababababababababab.....
<0,0,b>

<u>a b</u> <u>aba</u> babababababababababab.....
<1,2,a>

<u>a b</u> <u>aba</u> <u>babab</u> abababababababab.....
<2,4,b>

<u>a b</u> <u>aba</u> <u>babab</u> <u>ababababab</u> bab.....
<1,10,a>

38

## Surprise Code!

<u>a</u> babababababababababababab$
<0,0,a>

<u>a b</u> abababababababababababab$
<0,0,b>

<u>a b</u> <u>ababababababababababababab$</u>
<1,22,$>

39

## Surprise Decoding

<0,0,a><0,0,b><1,22,$>

| | |
|---|---|
| <0,0,a> | a |
| <0,0,b> | b |
| <1,22,$> | a |
| <2,21,$> | b |
| <3,20,$> | a |
| <4,19,$> | b |
| ... | |
| <22,1,$> | b |
| <23,0,$> | $ |

40

## Surprise Decoding

<0,0,a><0,0,b><1,22,$>

| | |
|---|---|
| <0,0,a> | a |
| <0,0,b> | b |
| <1,22,$> | a |
| <2,21,$> | b |
| <3,20,$> | a |
| <4,19,$> | b |
| ... | |
| <22,1,$> | b |
| <23,0,$> | $ |

41

## Solution C

- The matching string can include part of itself!
- If $x_{n+1}x_{n+2}...x_{n+k}$ is a substring of
$$x_1x_2...x_n x_{n+1}x_{n+2}...x_{n+k}$$
that begins at $j \leq n$ and $x_{n+1}x_{n+2}...x_{n+k}x_{n+k+1}$ is not then $x_{n+1}x_{n+2}...x_{n+k}x_{n+k+1}$ can be coded by
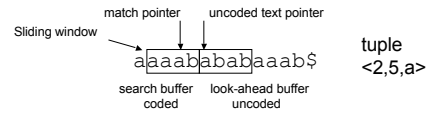$$<j,k, x_{n+k+1} >$$

42

## In Class Exercise

- Use Solution C to code the string
  - aaaabaaabaabab$

## Bounded Buffer – Sliding Window

- We want the triples <j,k,x> to be of bounded size. To achieve this we use bounded buffers.
  - Search buffer of size s is the symbols $x_{n-s+1}...x_n$ j is then the offset into the buffer.
  - Look-ahead buffer of size t is the symbols $x_{n+1}...x_{n+t}$
- Match pointer can start in search buffer and go into the look-ahead buffer but no farther.

Sliding window

match pointer    uncoded text pointer

a aaab abab aaab$

search buffer    look-ahead buffer
coded            uncoded

tuple
<2,5,a>

## Search in the Sliding Window

|  | offset | length |
|---|---|---|
| a aaab abab aaab$ | 1 | 0 |
| a aaab abab aaab$ | 2 | 1 |
| a aaab abab aaab$ | 2 | 2 |
| a aaab abab aaab$ | 2 | 3 |
| a aaab abab aaab$ | 2 | 4 |
| a aaab abab aaab$ | 2 | 5 |

tuple
<2,5,a>

## Coding Example

s = 4, t = 4, a = 3

tuple

| aaaa babababaaab$ | <0,0,a> |
| a aaab ababaaab$ | <1,3,b> |
| aaaab abab aaab$ | <2,5,a> |
| aaaabab abaa ab$ | <4,2,$> |

## Coding the Tuples

- Simple fixed length code

$$\lceil \log_2(s+1) \rceil + \lceil \log_2(s+t+1) \rceil + \lceil \log_2 a \rceil$$

s = 4, t = 4, a = 3

tuple    fixed code
<2,5,a>  010 0101 00

- Variable length code using adaptive Huffman or arithmetic code on Tuples
  - Two passes, first to create the tuples, second to code the tuples
  - One pass, by pipelining tuples into a variable length coder

## Zip and Gzip

- Search Window
  - Search buffer 32KB
  - Look-ahead buffer 258 Bytes
- How to store such a large dictionary
  - Hash table that stores the starting positions for all three byte sequences.
  - Hash table uses chaining with newest entries at the beginning of the chain. Stale entries can be ignored.
- Second pass for Huffman coding of tuples.
- Coding done in blocks to avoid disk accesses.

## Example

12

```
aaaabababaaabaaaababababaaabba$
```



aba →

Offset = 12 − 8 = 4
Length = 5
Tuple = <4,5,a>

49

---

## Example

18

```
aaaabababaaabaaaababababaaabba$
```



bab →

No match
Tuple = <0,0,b>

50

---

## Notes on LZ77

- Very popular especially in unix world
- Many variants and implementations
  - Zip, Gzip, PNG, PKZip,Lharc, ARJ
- Tends to work better than LZW
  - LZW has dictionary entries that are never used
  - LZW has past strings that are not in the dictionary
  - LZ77 has an implicit dictionary.  Common tuples are coded with few bits.

51

---

## Huffman Coding

- Huffman (1951)
- Uses frequencies of symbols in a string to build a variable rate prefix code.
  - Each symbol is mapped to a binary string.
  - More frequent symbols have shorter codes.
  - No code is a prefix of another.
- Example:

a  0
b  100
c  101
d  11



52

---

## Variable Rate Code Example

- Example:  a  0, b  100, c  101, d  11
- Coding:
  - aabddcaa = 16 bits
  - 0 0 100 11 11 101 0 0= 14 bits
- Prefix code ensures unique decodability.
  - 00100111110100
  - a a b d d c a a

53

---
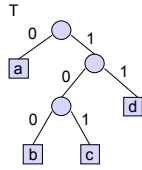
## Cost of a Huffman Tree

- Let $p_1, p_2, \ldots, p_m$ be the probabilities for the symbols $a_1, a_2, \ldots, a_m$, respectively.
- Define the cost of the Huffman tree T to be

$$C(T) = \sum_{i=1}^{m} p_i r_i$$

where $r_i$ is the length of the path from the root to $a_i$.
- C(T) is the expected length of the code of a symbol coded by the tree T.   C(T) is the bit rate of the code.

54

---

## Example of Cost

- Example:  a  1/2, b  1/8, c  1/8, d  1/4



$$C(T) = 1 \times 1/2 + 3 \times 1/8 + 3 \times 1/8 + 2 \times 1/4 = 1.75$$
$$\quad\quad\quad a \quad\quad\quad b \quad\quad\quad c \quad\quad\quad d$$

55

## Huffman Tree

- Input: Probabilities $p_1, p_2, \dots, p_m$ for symbols $a_1, a_2, \dots, a_m$, respectively.
- Output: A tree that minimizes the average number of bits (bit rate) to code a symbol. That is, minimizes
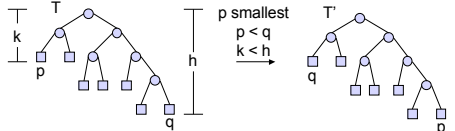
$$HC(T) = \sum_{i=1}^{m} p_i r_i \quad \text{bit rate}$$

where $r_i$ is the length of the path from the root to $a_i$. This is the Huffman tree or Huffman code

56

## Optimality Principle 1

- In a Huffman tree a lowest probability symbol has maximum distance from the root.
  - If not exchanging a lowest probability symbol with one at maximum distance will lower the cost.
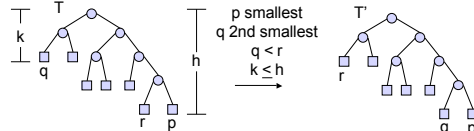


p smallest
$p < q$
$k < h$

$$C(T') = C(T) + hp - hq + kq - kp = C(T) - (h-k)(q-p) < C(T)$$

57

## Optimality Principle 2

- The second lowest probability is a sibling of the the smallest in some Huffman tree.
  - If not, we can move it there not raising the cost.



p smallest
q 2nd smallest
$q < r$
$k \leq h$

$$C(T') = C(T) + hq - hr + kr - kq = C(T) - (h-k)(r-q) \leq C(T)$$

58

## Optimality Principle 3

- Assuming we have a Huffman tree T whose two lowest probability symbols are siblings at maximum depth, they can be replaced by a new symbol whose probability is the sum of their probabilities.
  - The resulting tree is optimal for the new symbol set.



p smallest
q 2nd smallest

$$C(T') = C(T) + (h-1)(p+q) - hp - hq = C(T) - (p+q)$$

59

## Optimality Principle 3 (cont')

- If T' were not optimal then we could find a lower cost tree T''. This will lead to a lower cost tree T''' for the original alphabet.



$$C(T''') = C(T'') + p + q < C(T') + p + q = C(T) \quad \text{which is a contradiction}$$

60

10

## Recursive Huffman Tree Algorithm

1. If there is just one symbol, a tree with one node is optimal. Otherwise
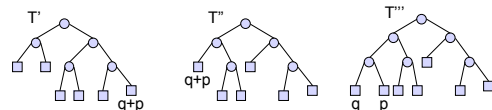2. Find the two lowest probability symbols with probabilities p and q respectively.
3. Replace these with a new symbol with probability p + q.
4. Solve the problem recursively for new symbols.
5. Replace the leaf with the new symbol with an internal node with two children with the old symbols.

61

## Iterative Huffman Tree Algorithm

```
form a node for each symbol aᵢ with weight pᵢ;
insert the nodes in a min priority queue ordered by probability;
while the priority queue has more than one element do
   min1 := delete-min;
   min2 := delete-min;
   create a new node n;
   n.weight := min1.weight + min2.weight;
   n.left := min1;
   n.right := min2;
   insert(n)
return the last node in the priority queue.
```

62

## Example of Huffman Tree Algorithm (1)

- P(a) =.4, P(b)=.1, P(c)=.3, P(d)=.1, P(e)=.1



63

## Example of Huffman Tree Algorithm (2)



64

## Example of Huffman Tree Algorithm (3)



65

## Example of Huffman Tree Algorithm (4)



66

11

## Huffman Code



average number of bits per symbol is
.4 x 1 + .1 x 4 + .3 x 2 + .1 x 3 + .1 x 4 = 2.1

a  0
b  1110
c  10
d  110
e  1111

67

---

## Optimal Huffman Code vs. Entropy

- P(a) =.4, P(b)=.1, P(c)=.3, P(d)=.1, P(e)=.1

Entropy

H = -(.4 x $\log_2$(.4) + .1 x $\log_2$(.1) + .3 x $\log_2$(.3)
  + .1 x $\log_2$(.1) + .1 x $\log_2$(.1))
  = 2.05 bits per symbol

Huffman Code

HC = .4 x 1 + .1 x 4 + .3 x 2 + .1 x 3 + .1 x 4
  = 2.1 bits per symbol
    pretty good!

68

---

## In Class Exercise

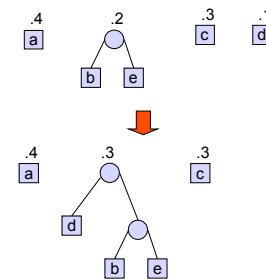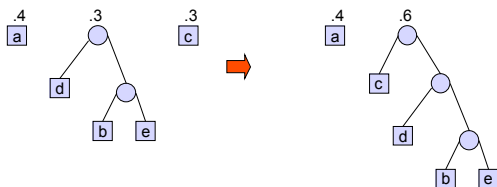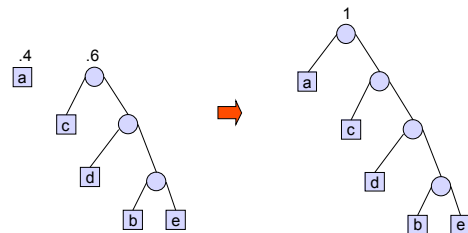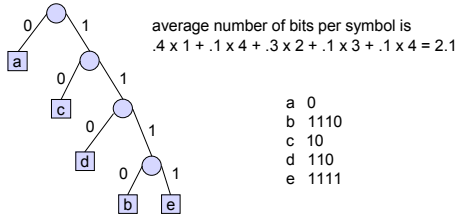- P(a) = 1/2, P(b) = 1/4, P(c) = 1/8, P(d) = 1/16, P(e) = 1/16
- Compute the Huffman tree and its bit rate.
- Compute the Entropy
- Compare
- Hint: For the tree change probabilities to be integers: a:8, b:4, c:2, d:1, e:1.  Normalize at the end.

69

---

## Quality of the Huffman Code

- The Huffman code is within one bit of the entropy lower bound.

$$H \leq HC \leq H + 1$$

- Huffman code does not work well with a two symbol alphabet.
  - Example: P(0) = 1/100, P(1) = 99/100
  - HC = 1 bits/symbol



  - H = -((1/100)*$\log_2$(1/100) + (99/100)$\log_2$(99/100))
    = .08 bits/symbol

70

---

## Powers of Two

- If all the probabilities are powers of two then

$$HC = H$$

- Proof by induction on the number of symbols.

  Let $p_1 \leq p_2 \leq ... \leq p_n$ be the probabilities that add up to 1

  If n = 1 then HC = H (both are zero).

  If n > 1 then $p_1 = p_2 = 2^{-k}$ for some k, otherwise the sum cannot add up to 1.

  Combine the first two symbols into a new symbol of probability $2^{-k} + 2^{-k} = 2^{-k+1}$.

71

---

## Powers of Two (Cont.)

By the induction hypothesis

$$HC(p_1 + p_2, p_3, ..., p_n) = H(p_1 + p_2, p_3, ..., p_n)$$
$$= -(p_1 + p_2)\log_2(p_1 + p_2) - \sum_{i=3}^{n} p_i \log_2(p_i)$$
$$= -2^{-k+1}\log_2(2^{-k+1}) - \sum_{i=3}^{n} p_i \log_2(p_i)$$
$$= -2^{-k+1}(\log_2(2^{-k}) + 1) - \sum_{i=3}^{n} p_i \log_2(p_i)$$
$$= -2^{-k}\log_2(2^{-k}) - 2^{-k}\log_2(2^{-k}) - \sum_{i=3}^{n} p_i \log_2(p_i) - 2^{-k} - 2^{-k}$$
$$= -\sum_{i=1}^{n} p_i \log_2(p_i) - (p_1 + p_2)$$
$$= H(p_1, p_2, ..., p_n) - (p_1 + p_2)$$

72

## Powers of Two (Cont.)

By the previous page,

$$HC(p_1 + p_2, p_3, \ldots, p_n) = H(p_1, p_2, \ldots, p_n) - (p_1 + p_2)$$

By the properties of Huffman trees (principle 3),

$$HC(p_1, p_2, \ldots, p_n) = HC(p_1 + p_2, p_3, \ldots, p_n) + (p_1 + p_2)$$

Hence,

$$HC(p_1, p_2, \ldots, p_n) = H(p_1, p_2, \ldots, p_n)$$

73

---

## Extending the Alphabet

- Assuming independence P(ab) = P(a)P(b), so we can lump symbols together.
- Example: P(0) = 1/100, P(1) = 99/100
  - P(00) = 1/10000, P(01) = P(10) = 99/10000, P(11) = 9801/10000.



HC = 1.03 bits/symbol (2 bit symbol)
= .515 bits/bit

Still not that close to H = .08 bits/bit

74

---

## Quality of Extended Alphabet

- Suppose we extend the alphabet to symbols of length k then

$$H \le HC \le H + 1/k$$

- Pros and Cons of Extending the alphabet
  - + Better compression
  - - $2^k$ symbols
  - - padding needed to make the length of the input divisible by k

75

---

## Huffman Codes with Context

- Suppose we add a one symbol context. That is in compressing a string $x_1 x_2 \ldots x_n$ we want to take into account $x_{k-1}$ when encoding $x_k$.
  - New model, so entropy based on just independent probabilities of the symbols doesn't hold. The new entropy model (2nd order entropy) has for each symbol a probability for each other symbol following it.
  - Example: {a,b,c}

|      |   | next |    |    |
|------|---|------|----|----|
|      |   | a    | b  | c  |
|      | a | .4   | .2 | .4 |
| prev | b | .1   | .9 | 0  |
|      | c | .1   | .1 | .8 |

76

---

## Multiple Codes

|      |   | next |    |    |
|------|---|------|----|----|
|      |   | a    | b  | c  |
|      | a | .4   | .2 | .4 |
| prev | b | .1   | .9 | 0  |
|      | c | .1   | .1 | .8 |

Code for first symbol
a 00
b 01
c 10



a b b a c c

00 00 0 1 01 0

77

---

## Complexity of Huffman Code Design

- Time to design Huffman Code is O(n log n) where n is the number of symbols.
  - Each step consists of a constant number of priority queue operations (2 deletemin's and 1 insert)

78

# Approaches to Huffman Codes

1. Frequencies computed for each input
   – Must transmit the Huffman code or frequencies as well as the compressed input
   – Requires two passes
2. Fixed Huffman tree designed from training data
   – Do not have to transmit the Huffman tree because it is known to the decoder.
   – H.263 video coder
3. Adaptive Huffman code
   – One pass
   – Huffman tree changes as frequencies change

79