

CSE 589 Applied Algorithms Autumn 2001

Golomb Coding
Arithmetic Coding
LZW
Sequitur

Binary Golomb Codes

- Binary source with 0's much more frequent than 1's.
- Variable-to-variable length code
- Prefix code.
- Golomb code of order 4

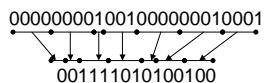
input	output
0000	0
0001	100
001	101
01	110
1	111

CSE 589 - Lecture 6 - Autumn 2001

2

Example

input	output
0000	0
0001	100
001	101
01	110
1	111



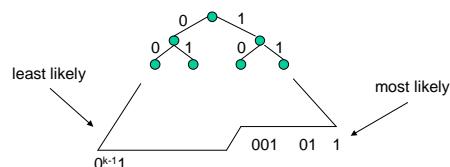
CSE 589 - Lecture 6 - Autumn 2001

3

Constructing a Binary Golomb Code

- Let p be the probability of 0
- Choose the order k such that

$$p^{k+1} \leq 1 - p^k < p^{k-1}$$



CSE 589 - Lecture 6 - Autumn 2001

4

Example

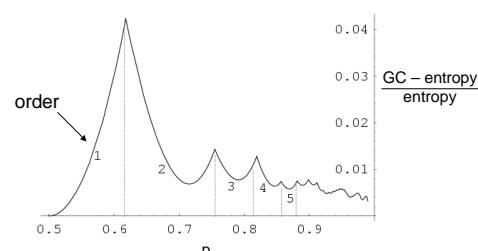
$$\begin{aligned} p &= .89 \\ p^7 &= .442 \\ 1 - p^6 &= .503 \\ p^5 &= .558 \\ \text{Order } 6 \end{aligned}$$

input	output
000000	0
000001	1000
00001	1001
0001	1010
001	1011
01	110
1	111

CSE 589 - Lecture 6 - Autumn 2001

5

Comparison of GC with Entropy



CSE 589 - Lecture 6 - Autumn 2001

6

Notes on Binary Golomb Code

- Very effective as a binary run-length code
- As p goes to 1 effectiveness near entropy
- There are adaptive Golomb codes, so that the order can change according to observed frequency of 0's.

CSE 589 - Lecture 6 - Autumn
2001

7

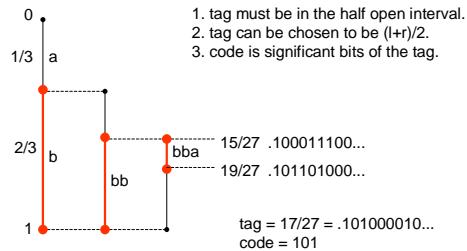
Arithmetic Coding

- Huffman coding works well for larger alphabets and gets to within one bit of the entropy lower bound. Can we do better. Yes
- Basic idea in arithmetic coding:
 - represent each string x of length n by an interval $[l, r]$ in $[0, 1]$. We assume n is known.
 - The width $r-l$ of the interval $[l, r]$ represents the probability of x occurring.
 - The interval $[l, r]$ can itself be represented by any number, called a tag, within the half open interval.
 - The k significant bits of the tag $.t_1 t_2 t_3 \dots$ is the code of x . That is, $\dots t_1 t_2 t_3 \dots t_k 000 \dots$ is in the interval $[l, r]$.

CSE 589 - Lecture 6 - Autumn
2001

8

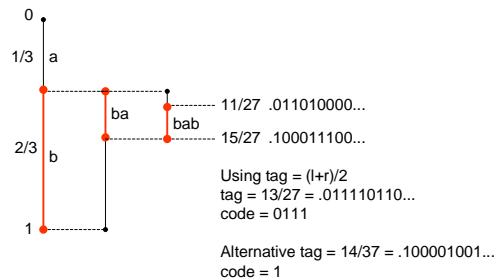
Example of Arithmetic Coding (1)



CSE 589 - Lecture 6 - Autumn
2001

9

Some Tags are Better than Others



CSE 589 - Lecture 6 - Autumn
2001

10

Example of Codes

		tag = $(l+r)/2$	code
0	aa	0/27 .000000000...	000001001...
a	aab	1/27 .000010010...	000100110...
a	aba	3/27 .000101100...	001001100...
a	abb	5/27 .000101110...	010000101...
b	baa	9/27 .010101010...	010111110...
b	bab	11/27 .011010000...	010111111...
b	bba	15/27 .100011100...	011101111...
b	bbb	19/27 .101101000...	101000010...
1		27/27 .111111111...	.110110100...
			.95 bits/symbol
			.92 entropy lower bound

CSE 589 - Lecture 6 - Autumn
2001

11

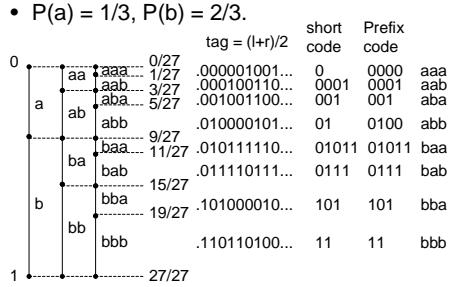
Code Generation from Tag

- If binary tag is $.t_1 t_2 t_3 \dots = (r+l)/2$ in $[l, r]$ then we want to choose k to form the code $t_1 t_2 \dots t_k$.
- Short code:
 - choose k to be as small as possible so that $l \leq .t_1 t_2 \dots t_k 000 \dots < r$.
 - Good if n is known.
- Guaranteed code:
 - choose $k = \lceil -\log_2(r-l) \rceil + 1$
 - $l \leq .t_1 t_2 \dots t_k b_1 b_2 b_3 \dots < r$ for any bits $b_1 b_2 b_3 \dots$
 - for fixed length strings provides a good prefix code.
 - example: [.000000000..., .000010010...], tag = .000001001...
Short code: 0
Guaranteed code: 000001

CSE 589 - Lecture 6 - Autumn
2001

12

Guaranteed Code Example



CSE 589 - Lecture 6 - Autumn 2001

13

Arithmetic Coding Algorithm

- $P(a_1), P(a_2), \dots, P(a_m)$
- $C(a_i) = P(a_1) + P(a_2) + \dots + P(a_{i-1})$
- Encode $x_1 x_2 \dots x_n$

```
Initialize l := 0 and r := 1;
for i = 1 to n do
    w := r - l;
    l := l + w * C(x);
    r := l + w * P(x);
    t := (l+r)/2;
    choose code for the tag
```

CSE 589 - Lecture 6 - Autumn 2001

14

Arithmetic Coding Example

- $P(a) = 1/4, P(b) = 1/2, P(c) = 1/4$
- $C(a) = 0, C(b) = 1/4, C(c) = 3/4$
- abca

symbol	w	l	r
a	1	0	1/4
b	1/4	1/16	3/16
c	1/8	5/32	6/32
a	1/32	5/32	21/128

w := r - l;
 $l := l + w C(x);$
 $r := l + w P(x)$

tag = $(5/32 + 21/128)/2 = 41/256 = .00101001\dots$
 $l = .00101000\dots$
 $r = .00101010\dots$
code = 00101
prefix code = 00101001

CSE 589 - Lecture 6 - Autumn 2001

15

Arithmetic Coding Exercise

- $P(a) = 1/4, P(b) = 1/2, P(c) = 1/4$
- $C(a) = 0, C(b) = 1/4, C(c) = 3/4$
- bbbb

symbol	w	l	r
b	1	0	1
b	1	0	1
b	1	0	1
b	1	0	1

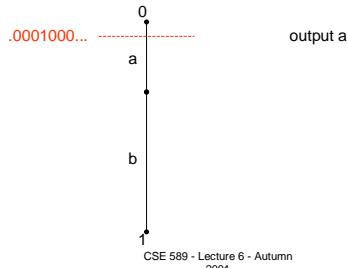
tag =
 $l =$
 $r =$
code =
prefix code =

CSE 589 - Lecture 6 - Autumn 2001

16

Decoding (1)

- Assume the length is known to be 3.
- 0001 which converts to the tag .0001000...

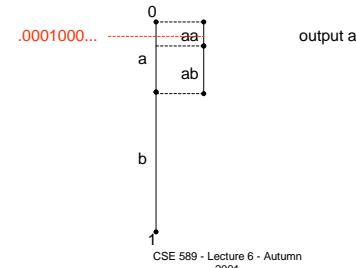


CSE 589 - Lecture 6 - Autumn 2001

17

Decoding (2)

- Assume the length is known to be 3.
- 0001 which converts to the tag .0001000...

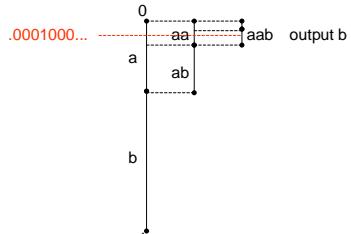


CSE 589 - Lecture 6 - Autumn 2001

18

Decoding (3)

- Assume the length is known to be 3.
- 0001 which converts to the tag .0001000...



CSE 589 - Lecture 6 - Autumn
2001

19

Arithmetic Decoding Algorithm

- $P(a_1), P(a_2), \dots, P(a_m)$
- $C(a_i) = P(a_1) + P(a_2) + \dots + P(a_{i-1})$
- Decode $b_1 b_2 \dots b_m$, number of symbols is n .

```
Initialize l := 0 and r := 1;
t := .b1b2...bm000...
for i = 1 to n do
    w := r - l;
    find j such that l + w * C(aj) ≤ t < l + w * (C(aj) + P(aj))
    output aj;
    l := l + w * C(aj);
    r := l + w * P(aj);
```

CSE 589 - Lecture 6 - Autumn
2001

20

Decoding Example

- $P(a) = 1/4, P(b) = 1/2, P(c) = 1/4$
- $C(a) = 0, C(b) = 1/4, C(c) = 3/4$
- 00101

tag = .00101000... = 5/32			
w	l	r	output
1/32	0	1/4	a
1/8	1/16	3/16	b
5/32	5/32	6/32	c
21/128	5/32	21/128	a

CSE 589 - Lecture 6 - Autumn
2001

21

Practical Arithmetic Coding

- Scaling:
 - By scaling we can keep l and r in a reasonable range of values so that $w = r - l$ does not underflow.
 - The code can be produced progressively, not at the end.
 - Complicates decoding some.
- Integer arithmetic coding avoids floating point altogether.

CSE 589 - Lecture 6 - Autumn
2001

22

Notes on Arithmetic Coding

- Arithmetic codes come close to the entropy lower bound.
- Grouping symbols is effective for arithmetic coding.
- Arithmetic codes can be used effectively on small symbol sets. Advantage over Huffman.
- Context can be added so that more than one probability distribution can be used.
 - The best coders in the world use this method.
- There are very effective adaptive arithmetic coding methods.

CSE 589 - Lecture 6 - Autumn
2001

23

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, Unix Compress, gzip, GIF

CSE 589 - Lecture 6 - Autumn
2001

24

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
```

CSE 589 - Lecture 6 - Autumn
2001

25

LZW Encoding Example (1)

Dictionary

0 a
1 b

a b a b a b a

CSE 589 - Lecture 6 - Autumn
2001

26

LZW Encoding Example (2)

Dictionary
0 a
1 b
2 ab

a b a b a b a
0

CSE 589 - Lecture 6 - Autumn
2001

27

LZW Encoding Example (3)

Dictionary
0 a
1 b
2 ab
3 ba

a b a b a b a
0 1

CSE 589 - Lecture 6 - Autumn
2001

28

LZW Encoding Example (4)

Dictionary
0 a
1 b
2 ab
3 ba
4 aba

a b a b a b a
0 1 2

CSE 589 - Lecture 6 - Autumn
2001

29

LZW Encoding Example (5)

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

a b a b a b a b a
0 1 2 4

CSE 589 - Lecture 6 - Autumn
2001

30

LZW Encoding Example (6)

Dictionary

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

a b a b a b a

0 1 2 4 3

CSE 589 - Lecture 6 - Autumn
2001

31

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;
```

CSE 589 - Lecture 6 - Autumn
2001

32

LZW Decoding Example (1)

Dictionary

0 a
1 b
2 a?

0 1 2 4 3 6

a

CSE 589 - Lecture 6 - Autumn
2001

33

LZW Decoding Example (2a)

Dictionary

0 a
1 b
2 ab

0 1 2 4 3 6

a b

CSE 589 - Lecture 6 - Autumn
2001

34

LZW Decoding Example (2b)

Dictionary

0 a
1 b
2 ab
3 b?

0 1 2 4 3 6

a b

CSE 589 - Lecture 6 - Autumn
2001

35

LZW Decoding Example (3a)

Dictionary

0 a
1 b
2 ab
3 ba

0 1 2 4 3 6

a b a

CSE 589 - Lecture 6 - Autumn
2001

36

LZW Decoding Example (3b)

Dictionary

0 a
1 b
2 ab
3 ba
4 ab?

0 1 2 4 3 6
a b ab

CSE 589 - Lecture 6 - Autumn
2001

37

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

CSE 589 - Lecture 6 - Autumn
2001

38

LZW Decoding Example (4b)

Dictionary

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

0 1 2 4 3 6
a b ab aba

CSE 589 - Lecture 6 - Autumn
2001

39

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

CSE 589 - Lecture 6 - Autumn
2001

40

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

CSE 589 - Lecture 6 - Autumn
2001

41

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

CSE 589 - Lecture 6 - Autumn
2001

42

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

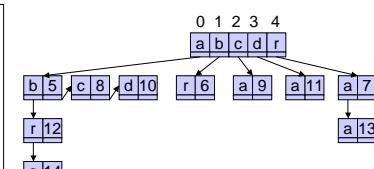
CSE 589 - Lecture 6 - Autumn
2001

43

Trie Data Structure for Encoder's Dictionary

- Fredkin (1960)

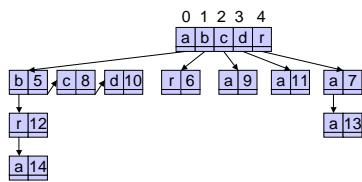
0	a	9	ca
1	b	10	ad
2	c	11	da
3	d	12	abr
4	r	13	raa
5	ab	14	abra
6	br		
7	ra		
8	ac		



CSE 589 - Lecture 6 - Autumn
2001

44

Encoder Uses a Trie (1)

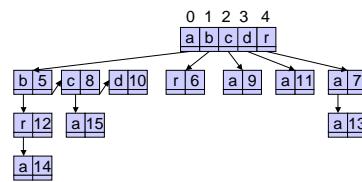


abracadabraabracadabra
0 1 4 0 2 0 3 5 7 12

CSE 589 - Lecture 6 - Autumn
2001

45

Encoder Uses a Trie (2)



abracadabraabracadabra
0 1 4 0 2 0 3 5 7 12 8

CSE 589 - Lecture 6 - Autumn
2001

46

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	ab?
6	br		
7	ra		
8	ac		

0 1 4 0 2 0 3 5 7 12 ...
abracadabra abr

CSE 589 - Lecture 6 - Autumn
2001

47

Notes on Dictionary Coding

- Extremely effective when there are repeated patterns in the data that are widely spread. Where local context is not as significant.
 - text
 - some graphics
 - program sources or binaries
- Variants of LZW are pervasive.

CSE 589 - Lecture 6 - Autumn
2001

48

Sequitur

- Nevill-Manning and Witten, 1996.
- Uses a context-free grammar (without recursion) to represent a string.
- The grammar is inferred from the string.
- If there is structure and repetition in the string then the grammar may be very small compared to the original string.
- Clever encoding of the grammar yields impressive compression ratios.
- Compression plus structure!

CSE 589 - Lecture 6 - Autumn
2001

49

Context-Free Grammars

- Invented by Chomsky in 1959 to explain the grammar of natural languages.
- Also invented by Backus in 1959 to generate and parse Fortran.
- Example:
 - terminals: b, e
 - nonterminals: S, A
 - Production Rules: $S \rightarrow SA$, $S \rightarrow A$, $A \rightarrow bSe$, $A \rightarrow be$
 - S is the start symbol

CSE 589 - Lecture 6 - Autumn
2001

50

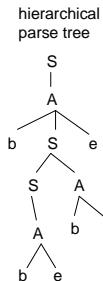
Context-Free Grammar Example

- $S \rightarrow SA$
- $S \rightarrow A$
- $A \rightarrow bSe$
- $A \rightarrow be$

Example: b and e matched as parentheses

derivation of bbebee

S
 A
 bSe
 $bAAe$
 $bbeAe$
 $bbebee$



CSE 589 - Lecture 6 - Autumn
2001

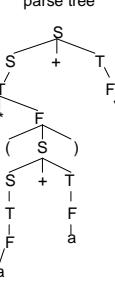
51

Arithmetic Expressions

- $S \rightarrow S + T$
- $S \rightarrow T$
- $T \rightarrow T * F$
- $T \rightarrow F$
- $F \rightarrow a$
- $F \rightarrow (S)$

derivation of $a^*(a+a)+a$

S
 $S+T$
 $T+T$
 T^*F+T
 $F+F+T$
 a^*F+T
 $a^*(S)+F$
 $a^*(S+F)+T$
 $a^*(T+F)+T$
 $a^*(F+F)+T$
 $a^*(a+F)+T$
 $a^*(a+a)+T$
 $a^*(a+a)+F$
 $a^*(a+a)+a$



CSE 589 - Lecture 6 - Autumn
2001

52

Sequitur Principles

- Digram Uniqueness:
 - no pair of adjacent symbols (digram) appears more than once in the grammar.
- Rule Utility:
 - Every production rule is used more than once.
- These two principles are maintained as an invariant while inferring a grammar for the input string.

CSE 589 - Lecture 6 - Autumn
2001

53

Sequitur Example (1)

bbebeebebebbeebee

$S \rightarrow b$

CSE 589 - Lecture 6 - Autumn
2001

54

Sequitur Example (2)

bbeebeebbebbbebee

$S \rightarrow bb$

CSE 589 - Lecture 6 - Autumn
2001

55

Sequitur Example (3)

bbeeebeebbebbbebee

$S \rightarrow bbe$

CSE 589 - Lecture 6 - Autumn
2001

56

Sequitur Example (4)

bbebbeebbebbbebee

$S \rightarrow bbeb$

CSE 589 - Lecture 6 - Autumn
2001

57

Sequitur Example (5)

bbebebebebbbebee

$S \rightarrow bbebe$

Enforce digram uniqueness.
be occurs twice.
Create new rule A $\rightarrow be$.

CSE 589 - Lecture 6 - Autumn
2001

58

Sequitur Example (6)

bbebebebebbbebee

$S \rightarrow bAA$
 $A \rightarrow be$

CSE 589 - Lecture 6 - Autumn
2001

59

Sequitur Example (7)

bbebebebebbbebee

$S \rightarrow bAAe$
 $A \rightarrow be$

CSE 589 - Lecture 6 - Autumn
2001

60

Sequitur Example (8)

bbebeebebbbebee

$S \rightarrow bAAeb$
 $A \rightarrow be$

CSE 589 - Lecture 6 - Autumn
2001

61

Sequitur Example (9)

bbebeebebebbee

$S \rightarrow bAAe$
 $A \rightarrow be$

Enforce digram uniqueness.
be occurs twice.
Use existing rule $A \rightarrow be$.

CSE 589 - Lecture 6 - Autumn
2001

62

Sequitur Example (10)

bbebeebebbbeee

$S \rightarrow bAAeA$
 $A \rightarrow be$

CSE 589 - Lecture 6 - Autumn
2001

63

Sequitur Example (11)

bbebeebebebbee

$S \rightarrow bAAeAb$
 $A \rightarrow be$

CSE 589 - Lecture 6 - Autumn
2001

64

Sequitur Example (12)

bbebeebebebbee

$S \rightarrow bAAeAbe$
 $A \rightarrow be$

Enforce digram uniqueness.
be occurs twice.
Use existing rule $A \rightarrow be$.

CSE 589 - Lecture 6 - Autumn
2001

65

Sequitur Example (13)

bbebeebebebbee

$S \rightarrow bAAeAA$
 $A \rightarrow be$

Enforce digram uniqueness.
AA occurs twice.
Create new rule $B \rightarrow AA$.

CSE 589 - Lecture 6 - Autumn
2001

66

Sequitur Example (14)

bbebeebebebbebee

S -> bBeB
A -> be
B -> AA

CSE 589 - Lecture 6 - Autumn
2001

67

Sequitur Example (15)

bbebeebebebebbee

S -> bBeBb
A -> be
B -> AA

CSE 589 - Lecture 6 - Autumn
2001

68

Sequitur Example (16)

bbebeebebebbbebee

S -> bBeBbb
A -> be
B -> AA

CSE 589 - Lecture 6 - Autumn
2001

69

Sequitur Example (17)

bbebeebebebbbee

S -> bBeBb_{be}
A -> _{be}
B -> AA

Enforce digram uniqueness.
be occurs twice.
Use existing rule A ->be.

CSE 589 - Lecture 6 - Autumn
2001

70

Sequitur Example (18)

bbebeebebebbbee

S -> bBeBbA
A -> be
B -> AA

CSE 589 - Lecture 6 - Autumn
2001

71

Sequitur Example (19)

bbebeebebebbbee

S -> bBeBbAb
A -> be
B -> AA

CSE 589 - Lecture 6 - Autumn
2001

72

Sequitur Example (20)

bbebeebebebbeee

S -> bBeBbA**be**
A -> **be**
B -> AA

Enforce digram uniqueness.
be occurs twice.
Use existing rule A -> be.

CSE 589 - Lecture 6 - Autumn
2001

73

Sequitur Example (21)

bbebeebebebbeee

S -> bBeBb**AA**
A -> be
B -> **AA**

Enforce digram uniqueness.
AA occurs twice.
Use existing rule B -> AA.

CSE 589 - Lecture 6 - Autumn
2001

74

Sequitur Example (22)

bbebeebebebbeee

S -> bBeBbB
A -> be
B -> AA

Enforce digram uniqueness.
bB occurs twice.
Create new rule C -> bB.

CSE 589 - Lecture 6 - Autumn
2001

75

Sequitur Example (23)

bbebeebebebbeee

S -> CeBC
A -> be
B -> AA
C -> bB

CSE 589 - Lecture 6 - Autumn
2001

76

Sequitur Example (24)

bbebeebebebbeee

S -> **CeBCe**
A -> be
B -> AA
C -> bB

Enforce digram uniqueness.
Ce occurs twice.
Create new rule D -> Ce.

CSE 589 - Lecture 6 - Autumn
2001

77

Sequitur Example (25)

bbebeebebebbeee

S -> DBD
A -> be
B -> AA
C -> bB
D -> Ce

Enforce rule utility.
C occurs only once.
Remove C -> bB.

CSE 589 - Lecture 6 - Autumn
2001

78

Sequitur Example (26)

bbebeebebebbeebee

S -> DBD
A -> be
B -> AA
D -> bBe

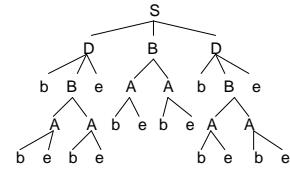
CSE 589 - Lecture 6 - Autumn
2001

79

The Hierarchy

bbebeebebebbeebee

S -> DBD
A -> be
B -> AA
D -> bBe



Is there compression? In this small example, probably not.

CSE 589 - Lecture 6 - Autumn
2001

80

Sequitur Algorithm

```

Input the first symbol s to create the production S -> s;
repeat
    match an existing rule:
        A -> ....XY...
        B -> XY          →   A -> ....B....
        B -> XY          →   B -> XY
    create a new rule:
        A -> ....XY...   →   A -> ....C....
        B -> ....XY...   →   B -> ....C....
        C -> XY
    remove a rule:
        A -> ....B....
        B -> X1X2...Xk →   A -> .... X1X2...Xk ...
    input a new symbol:
        S -> X1X2...Xk →   S -> X1X2...Xks
until no symbols left
  
```

CSE 589 - Lecture 6 - Autumn
2001

81

Complexity

- The number of non-input sequitur operations applied $\leq 2n$ where n is the input length.
- Amortized Complexity Argument**
 - Let $s =$ the sum of the right hand sides of all the production rules. Let $r =$ the number of rules.
 - We evaluate $2s - r$.
 - Initially $2s - r = 1$ because $s = 1$ and $r = 1$.
 - $2s - r \geq 0$ at all times because each rule has at least 1 symbol on the right hand side.
 - $2s - r$ increases by 2 for every input operation.
 - $2s - r$ decreases by at least 1 for each non-input sequitur rule applied.

CSE 589 - Lecture 6 - Autumn
2001

82

Sequitur Rule Complexity

- Diagram Uniqueness - match an existing rule.

$$\begin{array}{rcl} A \rightarrow \dots XY\dots & \longrightarrow & A \rightarrow \dots B\dots \\ B \rightarrow XY & \longrightarrow & B \rightarrow XY \end{array} \quad \begin{array}{ccccc} s & r & & & 2s - r \\ -1 & 0 & & & -2 \end{array}$$

- Diagram Uniqueness - create a new rule.

$$\begin{array}{rcl} A \rightarrow \dots XY\dots & \longrightarrow & A \rightarrow \dots C\dots \\ B \rightarrow \dots XY\dots & \longrightarrow & B \rightarrow \dots C\dots \\ C \rightarrow XY & \longrightarrow & C \rightarrow XY \end{array} \quad \begin{array}{ccccc} s & r & & & 2s - r \\ 0 & 1 & & & -1 \end{array}$$

- Rule Utility - Remove a rule.

$$\begin{array}{rcl} A \rightarrow \dots B\dots \\ B \rightarrow X_1X_2\dots X_k & \longrightarrow & A \rightarrow \dots X_1X_2\dots X_k\dots \end{array} \quad \begin{array}{ccccc} s & r & & & 2s - r \\ -1 & -1 & & & -1 \end{array}$$

CSE 589 - Lecture 6 - Autumn
2001

83

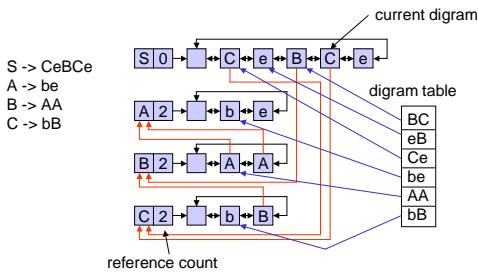
Linear Time Algorithm

- There is a data structure to implement all the sequitur operations in constant time.
 - Production rules in an array of doubly linked lists.
 - Each production rule has reference count of the number of times used.
 - Each nonterminal points to its production rule.
 - diagrams stored in a hash table for quick lookup.

CSE 589 - Lecture 6 - Autumn
2001

84

Data Structure Example



CSE 589 - Lecture 6 - Autumn
2001

85

Basic Encoding a Grammar

Grammar	S -> DBD A -> be B -> AA C -> bB	Symbol Code	b 000 e 001 S 010 A 011 B 100 D 101 # 110
Grammar Code			
D	B	D	#
101	100	101	110
b	e	#	A
000	001	011	011
AA	BB	010	001
39 bits			

$$|\text{Grammar Code}| = (s + r - 1) \lceil \log_2(r + a + 1) \rceil$$

r = number of rules

s = sum of right hand sides

a = number in original symbol alphabet

CSE 589 - Lecture 6 - Autumn
2001

86

Better Encoding of the Grammar

- Nevill-Manning and Witten suggest a more efficient encoding of the grammar that resembles LZ77.
 - The first time a nonterminal is sent, its right hand side is transmitted instead.
 - The second time a nonterminal is sent the new production rule is established with a pointer to the previous occurrence sent along with the length of the rule.
 - Subsequently, the nonterminal is represented by the index of the production rule.

CSE 589 - Lecture 6 - Autumn
2001

87

Compression Quality

- Neville-Manning and Witten 1997

	size	compress	gzip	sequitur	PPMC
bib	111261	3.35	2.51	2.48	2.12
book	768771	3.46	3.35	2.82	2.52
geo	102400	6.08	5.34	4.74	5.01
obj2	246814	4.17	2.63	2.68	2.77
pic	513216	0.97	0.82	0.90	0.98
progc	38611	3.87	2.68	2.83	2.49

Files from the Calgary Corpus

Units in bits per character (8 bits)

compress - based on LZW

gzip - based on LZ77

PPMC - adaptive arithmetic coding with context

CSE 589 - Lecture 6 - Autumn
2001

88

Notes on Sequitur

- Very new and different from the standards.
- Yields compression and hierarchical structure simultaneously.
- With clever encoding is competitive with the best of the standards.
- Practical linear time encoding and decoding.
- Alternatives
 - Off-line algorithms – (i) find the most frequent digram, (ii) find the longest repeated substring

CSE 589 - Lecture 6 - Autumn
2001

89