

CSE 490 GZ Introduction to Data Compression Winter 2002

Sequitur

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 490gz - Lecture 9 - Winter 2002

2

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
 - non-terminals: S, A
 - Production Rules:
 $S \rightarrow SA$, $S \rightarrow A$, $A \rightarrow bSe$, $A \rightarrow be$
 - S is the start symbol

CSE 490gz - Lecture 9 - Winter 2002

3

Context-Free Grammar Example

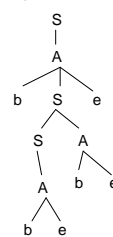
- $S \rightarrow SA$
 $S \rightarrow A$
 $A \rightarrow bSe$
 $A \rightarrow be$

derivation of bbebee

S
A
bSe
bSAe
bAAe
bbeAe
bbebee

Example: b and e matched as parentheses

hierarchical parse tree



CSE 490gz - Lecture 9 - Winter 2002

4

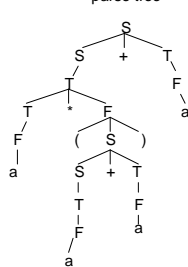
Arithmetic Expressions

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

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

parse tree



CSE 490gz - Lecture 9 - Winter 2002

5

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 490gz - Lecture 9 - Winter 2002

6

Sequitur Example (1)

bbebebebebebebe

$S \rightarrow b$

Sequitur Example (2)

bbeebebebebebe

$S \rightarrow bb$

Sequitur Example (3)

bbeeebebebebe

$S \rightarrow bbe$

Sequitur Example (4)

bbebeebebebebe

$S \rightarrow bbeb$

Sequitur Example (5)

bbebeeebebebebe

$S \rightarrow bbebe$

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

Sequitur Example (6)

bbebeeebebebebe

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

Sequitur Example (7)

bbebeebebbebee

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

CSE 490gz - Lecture 9 - Winter 2002

13

Sequitur Example (8)

bbebeebebbebee

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

CSE 490gz - Lecture 9 - Winter 2002

14

Sequitur Example (9)

bbebeebebbebee

$S \rightarrow bAAebe$
 $A \rightarrow be$

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

CSE 490gz - Lecture 9 - Winter 2002

15

Sequitur Example (10)

bbebeebebbebee

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

CSE 490gz - Lecture 9 - Winter 2002

16

Sequitur Example (11)

bbebeebebbebee

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

CSE 490gz - Lecture 9 - Winter 2002

17

Sequitur Example (12)

bbebeebebbebee

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

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

CSE 490gz - Lecture 9 - Winter 2002

18

Sequitur Example (13)

bbebeebbebeeb

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

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

Sequitur Example (14)

bbebeebbebeeb

$S \rightarrow bBeB$
 $A \rightarrow be$
 $B \rightarrow AA$

Sequitur Example (15)

bbebeebbebeeb

$S \rightarrow bBeBb$
 $A \rightarrow be$
 $B \rightarrow AA$

Sequitur Example (16)

bbebeebbebeeb

$S \rightarrow bBeBbb$
 $A \rightarrow be$
 $B \rightarrow AA$

Sequitur Example (17)

bbebeebbebeeb

$S \rightarrow bBeBbbe$
 $A \rightarrow be$
 $B \rightarrow AA$

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

Sequitur Example (18)

bbebeebbebeeb

$S \rightarrow bBeBbA$
 $A \rightarrow be$
 $B \rightarrow AA$

Sequitur Example (19)

bbebeebbebebeee

S → bBeBbAb
A → be
B → AA

CSE 490gz - Lecture 9 - Winter 2002

25

Sequitur Example (20)

bbebeebbebebee

S → bBeBbA**be**
A → **be**
B → AA

Enforce digram uniqueness.
be occurs twice.
Use existing rule A → be.

CSE 490gz - Lecture 9 - Winter 2002

26

Sequitur Example (21)

bbebeebbebebee

S → bBeBb**AA**
A → be
B → **AA**

Enforce digram uniqueness.
AA occurs twice.
Use existing rule B → AA.

CSE 490gz - Lecture 9 - Winter 2002

27

Sequitur Example (22)

bbebeebbebebee

S → **bBeBbB**
A → be
B → AA

Enforce digram uniqueness.
bB occurs twice.
Create new rule C → bB.

CSE 490gz - Lecture 9 - Winter 2002

28

Sequitur Example (23)

bbebeebbebebee

S → CeBC
A → be
B → AA
C → bB

CSE 490gz - Lecture 9 - Winter 2002

29

Sequitur Example (24)

bbebeebbebebee

S → **CeBCe**
A → be
B → AA
C → bB

Enforce digram uniqueness.
Ce occurs twice.
Create new rule D → Ce.

CSE 490gz - Lecture 9 - Winter 2002

30

bbebeebbebee

Enforce rule utility.
C occurs only once.
Remove $C \rightarrow bB$.

CSF 490az - Lecture 9 - Winter 2002

31

bbebebebebebebe

$$\begin{aligned} S &\rightarrow DBD \\ A &\rightarrow be \\ B &\rightarrow AA \\ D &\rightarrow bBe \end{aligned}$$

CSF 49007 - Lecture 9 - Winter 2002

32

bbeeebeebbee

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

CSE 490az - Lecture 9 - Winter 2002

33

Input the first symbol s to create the production $S \rightarrow s$;
repeat

 match an existing rule:

$A \rightarrow \dots XY \dots$ $A \rightarrow \dots B \dots$
 $B \rightarrow XY$ $B \rightarrow XY$

 create a new rule:

$A \rightarrow \dots XY \dots$ $A \rightarrow \dots C \dots$
 $B \rightarrow \dots XY \dots$ $B \rightarrow \dots C \dots$
 $C \rightarrow XY$

 remove a rule:

$A \rightarrow \dots B \dots$
 $B \rightarrow X_1 X_2 \dots X_k$ $A \rightarrow \dots X_1 X_2 \dots X_k \dots$

 input a new symbol:

$S \rightarrow X_1 X_2 \dots X_k$ $S \rightarrow X_1 X_2 \dots X_k s$

until no symbols left

CSE 490az - Lecture 9 - Winter 2002

34

- 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 490az - Lecture 9 - Winter 2002

35

- Digram Uniqueness - match an existing rule.

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

- Digram Uniqueness - create a new rule.

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

- Rule Utility - Remove a rule.

$$\begin{array}{l} A \rightarrow \dots B \dots \\ B \rightarrow X_1 X_2 \dots X_k \end{array} \longrightarrow \begin{array}{l} A \rightarrow \dots X_1 X_2 \dots X_k \dots \end{array} \quad \begin{array}{ccc} s & r & 2s-r \\ -1 & -1 & -1 \end{array}$$

CSE 490az - Lecture 9 - Winter 2002

36

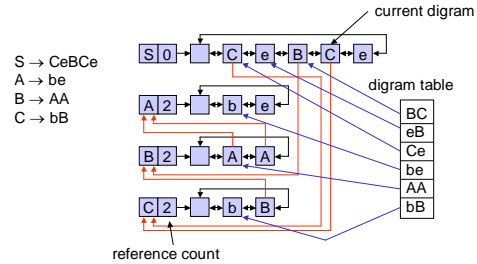
Linear Time Algorithm

- 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 non-terminal points to its production rule.
 - digrams stored in a hash table for quick lookup.

CSE 490gz - Lecture 9 - Winter 2002

37

Data Structure Example



CSE 490gz - Lecture 9 - Winter 2002

38

Basic Encoding a Grammar

Grammar	$S \rightarrow DBD$	Symbol Code	b 000	No code for S needed
	$A \rightarrow be$		e 001	
	$B \rightarrow AA$		A 010	
	$D \rightarrow bBe$		B 011	
			D 100	
			# 101	

Grammar Code

D B D # b e # A A # b B e
 101 100 101 110 000 001 110 011 011 110 000 100 001 39 bits

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

r = number of rules

s = sum of right hand sides

a = number in original symbol alphabet

CSE 490gz - Lecture 9 - Winter 2002

39

Better Encoding of the Grammar

- Nevill-Manning and Witten suggest a more efficient encoding of the grammar that uses LZ77 ideas.
 - Send the right hand side of the S production.
 - The first time a non-terminal is sent, its right hand side is transmitted instead.
 - The second time a non-terminal is sent as a tuple $\langle i, j, k \rangle$ which says the right hand side starts occurs in production i , at position j and is k long. A new production rule is then added to a dictionary.
 - Subsequently, the non-terminal is represented by the index of the production rule.

CSE 490gz - Lecture 9 - Winter 2002

40

Compression Quality

- Neville-Manning and Witten 1997

	size	comp	gzip	sequitur	PPMC	bzip2	First Second Third
bib	111261	3.35	2.51	2.48	2.12	1.98	
book	768771	3.46	3.35	2.82	2.52	2.42	
geo	102400	6.08	5.34	4.74	5.01	4.45	
obj2	246814	4.17	2.63	2.68	2.77	2.48	
pic	513216	0.97	0.82	0.90	0.98	0.78	
progc	38611	3.87	2.68	2.83	2.49	2.53	

Files from the Calgary Corpus
 Units in bits per character (8 bits)
 compress - based on LZ78 and LZW
 gzip - based on LZ77
 PPMC - adaptive arithmetic coding with context
 bzip2 - Burrows Wheeler transform

CSE 490gz - Lecture 9 - Winter 2002

41

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 490gz - Lecture 9 - Winter 2002

42

Other Grammar Based Methods

- YK Algorithm
 - Kieffer, Yang 2000
 - Like Sequitur, but does not allow different non-terminals to generate the same string
 - Slower, but has some better theoretical properties
- Longest Match
- Most frequent digram
- Match producing the best compression