

CSEP 590 Data Compression Autumn 2007

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!

CSEP 590 - Lecture 5 - Autumn 2007

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

CSEP 590 - Lecture 5 - Autumn 2007

3

Context-Free Grammar Example

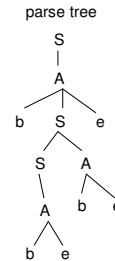
- $S \rightarrow SA$
- $S \rightarrow A$
- $A \rightarrow bSe$
- $A \rightarrow be$

derivation of bbebee

S
A
bSe
bSAe
bAAe
bbeAe
bbebee

Example: b and e matched as parentheses

hierarchical
parse tree



CSEP 590 - Lecture 5 - Autumn 2007

4

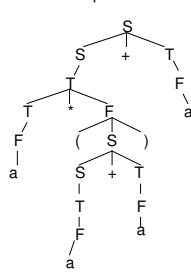
Arithmetic Expressions

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

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

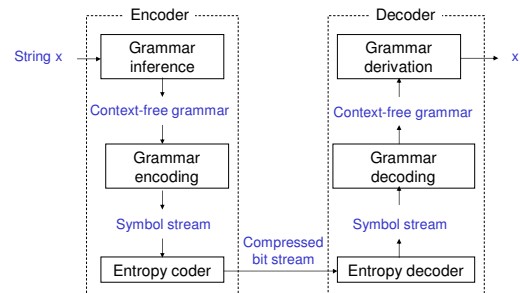
parse tree



CSEP 590 - Lecture 5 - Autumn 2007

5

Overview of Grammar Compression



CSEP 590 - Lecture 5 - Autumn 2007

6

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.

CSEP 590 - Lecture 5 - Autumn 2007

7

Sequitur Example (1)

bbeebbebbee

$S \rightarrow b$

CSEP 590 - Lecture 5 - Autumn 2007

8

Sequitur Example (2)

bbeebbebbee

$S \rightarrow bb$

CSEP 590 - Lecture 5 - Autumn 2007

9

Sequitur Example (3)

bbeebbebbee

$S \rightarrow bbe$

CSEP 590 - Lecture 5 - Autumn 2007

10

Sequitur Example (4)

bbeebbebbee

$S \rightarrow bbeb$

CSEP 590 - Lecture 5 - Autumn 2007

11

Sequitur Example (5)

bbeebbebbee

$S \rightarrow bbebe$

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

CSEP 590 - Lecture 5 - Autumn 2007

12

Sequitur Example (6)

bbeeebebebebe

S → bAA
A → be

Sequitur Example (7)

bbeeebebebebe

S → bAAe
A → be

Sequitur Example (8)

bbeeebebebebe

S → bAAeb
A → be

Sequitur Example (9)

bbeeebebebebe

S → bAAebe
A → be

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

Sequitur Example (10)

bbeeebebebebe

S → bAAeA
A → be

Sequitur Example (11)

bbeeebebebebe

S → bAAeAb
A → be

Sequitur Example (12)

bbebebebebebe

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

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

Sequitur Example (13)

bbebebebebebe

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

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

Sequitur Example (14)

bbebebebebebe

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

Sequitur Example (15)

bbebebebebbebe

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

Sequitur Example (16)

bbebebebebbebe

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

Sequitur Example (17)

bbebebebebbebe

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

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

Sequitur Example (18)

bbeeeebbee

S → bBeBbA
A → be
B → AA

Sequitur Example (19)

bbeeeebbee

S → bBeBbAb
A → be
B → AA

Sequitur Example (20)

bbeeeebbee

S → bBeBbA**be** Enforce digram uniqueness.
A → **be** be occurs twice.
B → AA Use existing rule A → be.

Sequitur Example (21)

bbeeeebbee

S → bBeBb**AA** Enforce digram uniqueness.
A → be AA occurs twice.
B → **AA** Use existing rule B → AA.

Sequitur Example (22)

bbeeeebbee

S → **bBeBbB** Enforce digram uniqueness.
A → be bB occurs twice.
B → AA Create new rule C → bB.

Sequitur Example (23)

bbeeeebbee

S → CeBC
A → be
B → AA
C → bB

Sequitur Example (24)

bbeebbeebbeeb

$S \rightarrow CeBCe$ Enforce digram uniqueness.
 $A \rightarrow be$ Ce occurs twice.
 $B \rightarrow AA$ Create new rule $D \rightarrow Ce$.
 $C \rightarrow bB$

CSEP 590 - Lecture 5 - Autumn 2007

31

Sequitur Example (25)

bbeebbeebbeeb

$S \rightarrow DBD$ Enforce rule utility.
 $A \rightarrow be$ C occurs only once.
 $B \rightarrow AA$ Remove $C \rightarrow bB$.
 $C \rightarrow bB$
 $D \rightarrow Ce$

CSEP 590 - Lecture 5 - Autumn 2007

32

Sequitur Example (26)

bbeebbeebbeeb

$S \rightarrow DBD$
 $A \rightarrow be$
 $B \rightarrow AA$
 $D \rightarrow bBe$

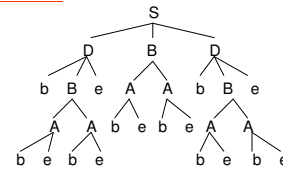
CSEP 590 - Lecture 5 - Autumn 2007

33

The Hierarchy

bbeebbeebbeeb

$S \rightarrow DBD$
 $A \rightarrow be$
 $B \rightarrow AA$
 $D \rightarrow bBe$



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

CSEP 590 - Lecture 5 - Autumn 2007

34

Sequitur Algorithm

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

CSEP 590 - Lecture 5 - Autumn 2007

35

Exercise

Use Sequitur to construct a grammar for $aaaaaaaaa = a^{10}$

CSEP 590 - Lecture 5 - Autumn 2007

36

Complexity

- The number of non-input sequitur operations applied $< 2n$ where n is the input length.
- Since each operation takes constant time, sequitur is a linear time algorithm

CSEP 590 - Lecture 5 - Autumn 2007

37

Amortized Complexity Argument

- Let m = # of non-input sequitur operations.
Let n = input length. Show $m \leq 2n$.
- 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 $r = 1$.
- $2s - r > 0$ at all times because each rule has at least 1 symbol on the right hand side.

CSEP 590 - Lecture 5 - Autumn 2007

38

Sequitur Rule Complexity

- Digram Uniqueness - match an existing rule.

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

- Digram Uniqueness - create a new rule.

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

- Rule Utility - Remove a rule.

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

CSEP 590 - Lecture 5 - Autumn 2007

39

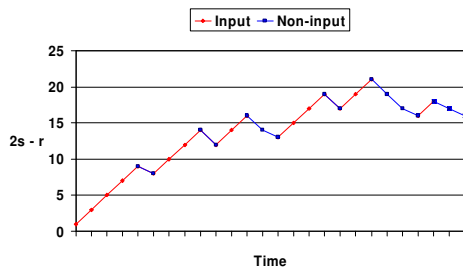
Amortized Complexity Argument

- $2s - r \geq 0$ 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.
- n = number of input symbols
 m = number of non-input operations
- $2n - m \geq 0$. $m \leq 2n$.

CSEP 590 - Lecture 5 - Autumn 2007

40

Amortized Complexity Argument



CSEP 590 - Lecture 5 - Autumn 2007

41

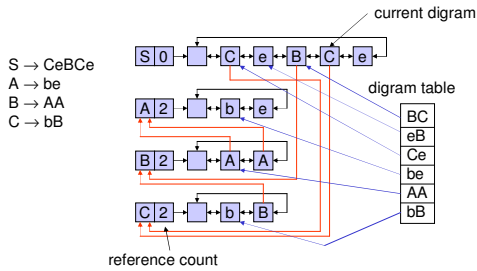
Linear Time Algorithm

- 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.
 - Digrams stored in a hash table for quick lookup.

CSEP 590 - Lecture 5 - Autumn 2007

42

Data Structure Example



CSEP 590 - Lecture 5 - Autumn 2007

43

Basic Encoding a Grammar

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

Grammar Code

$D B D \# b e \# A A \# b B e$
 100 011 100 101 000 001 101 010 010 101 000 011 001 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

CSEP 590 - Lecture 5 - Autumn 2007

44

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 nonterminal is sent, its right hand side is transmitted instead.
 - The second time a nonterminal is sent as a triple $[i, j, d]$, which says the right hand side starts at position j in production rule i and is d long. A new production rule is then added to a dictionary.
 - Subsequently, the nonterminal is represented by the index of the production rule.

CSEP 590 - Lecture 5 - Autumn 2007

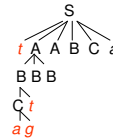
45

Transmission Example

$S \rightarrow tAABCa$
 $A \rightarrow BBB$
 $B \rightarrow Ct$
 $C \rightarrow ag$

$T = \text{Transmitted}$
 $T \text{ tagt}$
 $X_0 \text{ tagt}$

$l_0 = 4$



CSEP 590 - Lecture 5 - Autumn 2007

46

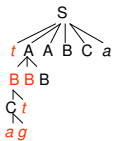
Transmission Example

$S \rightarrow tAABCa$
 $A \rightarrow BBB$
 $B \rightarrow Ct$
 $C \rightarrow ag$

$T = \text{Transmitted}$
 $T \text{ tagt}[0, 1, 3]$

$X_0 \text{ t} X_1 X_1$
 $X_1 \text{ agt}$

$l_0 = 3$
 $l_1 = 3$



CSEP 590 - Lecture 5 - Autumn 2007

47

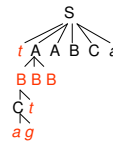
Transmission Example

$S \rightarrow tAABCa$
 $A \rightarrow BBB$
 $B \rightarrow Ct$
 $C \rightarrow ag$

$T = \text{Transmitted}$
 $T \text{ tagt}[0, 1, 3] 1$

$X_0 \text{ t} X_1 X_1 X_1$
 $X_1 \text{ agt}$

$l_0 = 4$
 $l_1 = 3$



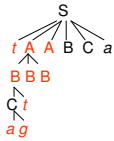
CSEP 590 - Lecture 5 - Autumn 2007

48

Transmission Example

S → tAABCa
A → BBB
B → Ct
C → ag

T = Transmitted
T tagt[0, 1, 3] 1 [0, 1, 3]
X₀ tX₂X₂ l₀ = 3
X₁ agt l₁ = 3
X₂ X₁X₁ l₂ = 3



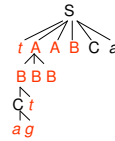
CSEP 590 - Lecture 5 - Autumn 2007

49

Transmission Example

S → tAABCa
A → BBB
B → Ct
C → ag

T = Transmitted
T tagt[0, 1, 3] 1 [0, 1, 3] 1
X₀ tX₂X₂X₁ l₀ = 4
X₁ agt l₁ = 3
X₂ X₁X₁ l₂ = 3



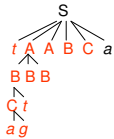
CSEP 590 - Lecture 5 - Autumn 2007

50

Transmission Example

S → tAABCa
A → BBB
B → Ct
C → ag

T = Transmitted
T tagt[0, 1, 3] 1 [0, 1, 3] 1 [1, 0, 2]
X₀ tX₂X₂X₁X₃ l₀ = 5
X₁ X₃f l₁ = 2
X₂ X₁X₁ l₂ = 3
X₃ ag l₃ = 2



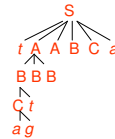
CSEP 590 - Lecture 5 - Autumn 2007

51

Transmission Example

S → tAABCa
A → BBB
B → Ct
C → ag

T = Transmitted
T tagt[0, 1, 3] 1 [0, 1, 3] 1 [1, 0, 2] a
X₀ tX₂X₂X₁X₃a l₀ = 6
X₁ X₃f l₁ = 2
X₂ X₁X₁ l₂ = 3
X₃ ag l₃ = 2



CSEP 590 - Lecture 5 - Autumn 2007

52

Kieffer-Yang Improvement

- Kieffer and Yang developed a theoretical framework for studying these types of grammars in 2000.

– KY is universal; it achieves entropy in the limit

- Add to sequitur Reduction Rule 5:

S → AB		S → AA	
A → CD		A → CD	
B → aE	⇒	B → aE	Adding this constraint
C → ab		C → ab	makes sequitur
D → cd		D → cd	universal.
E → bD		E → bD	

$\langle A \rangle = \langle B \rangle = abcd$

CSEP 590 - Lecture 5 - Autumn 2007

53

Compression Quality

- Neville-Manning and Witten 1997

	size	comp	gzip	sequitur	PPMC	bzip2
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
prog	38611	3.87	2.68	2.83	2.49	2.53

■ = First; ■ = Second; ■ = Third.

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
bzip2 - Burrows-Wheeler block sorting

CSEP 590 - Lecture 5 - Autumn 2007

54

Notes on Sequitur

- Yields compression and hierarchical structure simultaneously.
- With clever encoding is competitive with the best of the standards.
- The grammar size is not close to approximation algorithms
 - Upper = $O((n/\log n)^{3/4})$; Lower = $\Omega(n^{1/3})$. (Lehman, 2002)
- *But!* Practical linear time encoding and decoding.

CSEP 590 - Lecture 5 - Autumn 2007

55

Other Grammar Based Methods

- Longest Match
- Most frequent digram
- Match producing the best compression

CSEP 590 - Lecture 5 - Autumn 2007

56