

Diff

Leonardo Gresta Paulino Murta

leomurta@ic.uff.br

Exercise

- Conceive an algorithm for identifying the differences among two files without a common ancestor (2-way diff)

A
B
C
D
E
F

A
C
D
E
G
F

Possible Solution

- Identify the longest common subsequence among both files

A
B
C
D
E
F

A
C
D
E
G
F

Possible Solution

- Subtract the longest common sequence from both sides to identify what was added/removed

	A	
- B		
	C	
	D	
	E	
		+ G
	F	

However, how can we find the longest common sequence?

- Possible solution:
 - Generate all subsequences for one of the files
 - Check, for each generated subsequence, if it is also a subsequence of the other file
- Problem:
 - Complexity = $O(2^n n)$

LCS

- Problem characteristics
 - Can be divided into subproblems
 - The subproblems can repeat during recursion (leading to redundant computation)
- LCS algorithm
 - Longest Common Subsequence
 - Used both in bioinformatics and diff program
 - Adopts Dynamic Programming technique
 - Complexity = $O(n^2)$

LCS

- Considering the following sequences
 - $X_i = (x_1, x_2, \dots, x_i)$
 - $Y_j = (y_1, y_2, \dots, y_j)$
- Algorithm

$$LCS(X_i, Y_j) = \begin{cases} \emptyset & \text{if } i = 0 \vee j = 0 \\ (LCS(X_{i-1}, Y_{j-1}), x_i) & \text{if } x_i = y_j \\ \text{longest}(LCS(X_i, Y_{j-1}), LCS(X_{i-1}, Y_j)) & \text{if } x_i \neq y_j \end{cases}$$

LCS

- It can be computed in a bottom-up fashion
 - Using a matrix with all elements of one sequence in the line and all elements of the other sequence in the column
 - Computing line 1 and column 1, then line 2 and column 2, and so on
 - Storing in each cell the length of the sequence and the path to the cells that belong to the LCS

Longest Common Subsequence

line & column = 0

	\emptyset	A	C	D	E	G	F
\emptyset	0	$\leftarrow 0$	$\leftarrow 0$	$\leftarrow 0$	$\leftarrow 0$	$\leftarrow 0$	$\leftarrow 0$
A	$\uparrow 0$						
B	$\uparrow 0$						
C	$\uparrow 0$						
D	$\uparrow 0$						
E	$\uparrow 0$						
F	$\uparrow 0$						

Longest Common Subsequence

line & column = 1

	\emptyset	A	C	D	E	G	F
\emptyset	0	$\leftarrow 0$	$\leftarrow 0$	$\leftarrow 0$	$\leftarrow 0$	$\leftarrow 0$	$\leftarrow 0$
A	$\uparrow 0$	$\nwarrow 1$	$\leftarrow 1$	$\leftarrow 1$	$\leftarrow 1$	$\leftarrow 1$	$\leftarrow 1$
B	$\uparrow 0$	$\uparrow 1$					
C	$\uparrow 0$	$\uparrow 1$					
D	$\uparrow 0$	$\uparrow 1$					
E	$\uparrow 0$	$\uparrow 1$					
F	$\uparrow 0$	$\uparrow 1$					

Longest Common Subsequence

line & column = 2

	\emptyset	A	C	D	E	G	F
\emptyset	0	$\leftarrow 0$	$\leftarrow 0$	$\leftarrow 0$	$\leftarrow 0$	$\leftarrow 0$	$\leftarrow 0$
A	$\uparrow 0$	$\nwarrow 1$	$\leftarrow 1$	$\leftarrow 1$	$\leftarrow 1$	$\leftarrow 1$	$\leftarrow 1$
B	$\uparrow 0$	$\uparrow 1$	$\leftarrow 1 \uparrow$	$\leftarrow 1 \uparrow$	$\leftarrow 1 \uparrow$	$\leftarrow 1 \uparrow$	$\leftarrow 1 \uparrow$
C	$\uparrow 0$	$\uparrow 1$	$\nwarrow 2$				
D	$\uparrow 0$	$\uparrow 1$	$\uparrow 2$				
E	$\uparrow 0$	$\uparrow 1$	$\uparrow 2$				
F	$\uparrow 0$	$\uparrow 1$	$\uparrow 2$				

Longest Common Subsequence

line & column = 3

	\emptyset	A	C	D	E	G	F
\emptyset	0	$\leftarrow 0$	$\leftarrow 0$	$\leftarrow 0$	$\leftarrow 0$	$\leftarrow 0$	$\leftarrow 0$
A	$\uparrow 0$	$\nwarrow 1$	$\leftarrow 1$	$\leftarrow 1$	$\leftarrow 1$	$\leftarrow 1$	$\leftarrow 1$
B	$\uparrow 0$	$\uparrow 1$	$\leftarrow 1 \uparrow$	$\leftarrow 1 \uparrow$	$\leftarrow 1 \uparrow$	$\leftarrow 1 \uparrow$	$\leftarrow 1 \uparrow$
C	$\uparrow 0$	$\uparrow 1$	$\nwarrow 2$	$\leftarrow 2$	$\leftarrow 2$	$\leftarrow 2$	$\leftarrow 2$
D	$\uparrow 0$	$\uparrow 1$	$\uparrow 2$	$\nwarrow 3$			
E	$\uparrow 0$	$\uparrow 1$	$\uparrow 2$	$\uparrow 3$			
F	$\uparrow 0$	$\uparrow 1$	$\uparrow 2$	$\uparrow 3$			

Longest Common Subsequence

line & column = 4

	\emptyset	A	C	D	E	G	F
\emptyset	0	$\leftarrow 0$	$\leftarrow 0$	$\leftarrow 0$	$\leftarrow 0$	$\leftarrow 0$	$\leftarrow 0$
A	$\uparrow 0$	$\nwarrow 1$	$\leftarrow 1$	$\leftarrow 1$	$\leftarrow 1$	$\leftarrow 1$	$\leftarrow 1$
B	$\uparrow 0$	$\uparrow 1$	$\leftarrow 1 \uparrow$	$\leftarrow 1 \uparrow$	$\leftarrow 1 \uparrow$	$\leftarrow 1 \uparrow$	$\leftarrow 1 \uparrow$
C	$\uparrow 0$	$\uparrow 1$	$\nwarrow 2$	$\leftarrow 2$	$\leftarrow 2$	$\leftarrow 2$	$\leftarrow 2$
D	$\uparrow 0$	$\uparrow 1$	$\uparrow 2$	$\nwarrow 3$	$\leftarrow 3$	$\leftarrow 3$	$\leftarrow 3$
E	$\uparrow 0$	$\uparrow 1$	$\uparrow 2$	$\uparrow 3$	$\nwarrow 4$		
F	$\uparrow 0$	$\uparrow 1$	$\uparrow 2$	$\uparrow 3$	$\uparrow 4$		

Longest Common Subsequence

line & column = 5

	\emptyset	A	C	D	E	G	F
\emptyset	0	$\leftarrow 0$	$\leftarrow 0$	$\leftarrow 0$	$\leftarrow 0$	$\leftarrow 0$	$\leftarrow 0$
A	$\uparrow 0$	$\nwarrow 1$	$\leftarrow 1$	$\leftarrow 1$	$\leftarrow 1$	$\leftarrow 1$	$\leftarrow 1$
B	$\uparrow 0$	$\uparrow 1$	$\leftarrow 1 \uparrow$	$\leftarrow 1 \uparrow$	$\leftarrow 1 \uparrow$	$\leftarrow 1 \uparrow$	$\leftarrow 1 \uparrow$
C	$\uparrow 0$	$\uparrow 1$	$\nwarrow 2$	$\leftarrow 2$	$\leftarrow 2$	$\leftarrow 2$	$\leftarrow 2$
D	$\uparrow 0$	$\uparrow 1$	$\uparrow 2$	$\nwarrow 3$	$\leftarrow 3$	$\leftarrow 3$	$\leftarrow 3$
E	$\uparrow 0$	$\uparrow 1$	$\uparrow 2$	$\uparrow 3$	$\nwarrow 4$	$\leftarrow 4$	$\leftarrow 4$
F	$\uparrow 0$	$\uparrow 1$	$\uparrow 2$	$\uparrow 3$	$\uparrow 4$	$\leftarrow 4 \uparrow$	

Longest Common Subsequence

line & column = 6

	\emptyset	A	C	D	E	G	F
\emptyset	0	$\leftarrow 0$	$\leftarrow 0$	$\leftarrow 0$	$\leftarrow 0$	$\leftarrow 0$	$\leftarrow 0$
A	$\uparrow 0$	$\nwarrow 1$	$\leftarrow 1$	$\leftarrow 1$	$\leftarrow 1$	$\leftarrow 1$	$\leftarrow 1$
B	$\uparrow 0$	$\uparrow 1$	$\leftarrow 1 \uparrow$	$\leftarrow 1 \uparrow$	$\leftarrow 1 \uparrow$	$\leftarrow 1 \uparrow$	$\leftarrow 1 \uparrow$
C	$\uparrow 0$	$\uparrow 1$	$\nwarrow 2$	$\leftarrow 2$	$\leftarrow 2$	$\leftarrow 2$	$\leftarrow 2$
D	$\uparrow 0$	$\uparrow 1$	$\uparrow 2$	$\nwarrow 3$	$\leftarrow 3$	$\leftarrow 3$	$\leftarrow 3$
E	$\uparrow 0$	$\uparrow 1$	$\uparrow 2$	$\uparrow 3$	$\nwarrow 4$	$\leftarrow 4$	$\leftarrow 4$
F	$\uparrow 0$	$\uparrow 1$	$\uparrow 2$	$\uparrow 3$	$\uparrow 4$	$\leftarrow 4 \uparrow$	$\nwarrow 5$

Longest Common Subsequence

	\emptyset	A	C	D	E	G	F
\emptyset	0	$\leftarrow 0$	$\leftarrow 0$	$\leftarrow 0$	$\leftarrow 0$	$\leftarrow 0$	$\leftarrow 0$
A	$\uparrow 0$	$\nwarrow 1$	$\leftarrow 1$	$\leftarrow 1$	$\leftarrow 1$	$\leftarrow 1$	$\leftarrow 1$
B	$\uparrow 0$	$\uparrow 1$	$\leftarrow 1 \uparrow$	$\leftarrow 1 \uparrow$	$\leftarrow 1 \uparrow$	$\leftarrow 1 \uparrow$	$\leftarrow 1 \uparrow$
C	$\uparrow 0$	$\uparrow 1$	$\nwarrow 2$	$\leftarrow 2$	$\leftarrow 2$	$\leftarrow 2$	$\leftarrow 2$
D	$\uparrow 0$	$\uparrow 1$	$\uparrow 2$	$\nwarrow 3$	$\leftarrow 3$	$\leftarrow 3$	$\leftarrow 3$
E	$\uparrow 0$	$\uparrow 1$	$\uparrow 2$	$\uparrow 3$	$\nwarrow 4$	$\leftarrow 4$	$\leftarrow 4$
F	$\uparrow 0$	$\uparrow 1$	$\uparrow 2$	$\uparrow 3$	$\uparrow 4$	$\leftarrow 4 \uparrow$	$\nwarrow 5$

Improvements

- The first implementation of Unix Diff (Hunt & McIlroy, 1976) uses a variation of this LCS algorithm
 - $O(n)$ space complexity
 - $O(n^2 \times \log n)$ time complexity
- The current implementation of Unix Diff (Miller & Myers, 1985) does not fill the whole matrix
 - $O(n)$ space complexity
 - $O(n \times d)$ time complexity, where d is the edit distance

Miller & Myers algorithm

- bootstrap: Add zeros to the diagonal while the symbols match
- While the lowermost and rightmost cell is empty
 - Rule 1: For each filled cell, inserts its value added by one in the cell in the right
 - Rule 2: For each filled cell, inserts its value added by one in the cell in the bottom
 - Rule 3: For each filled cell, recursively inserts its value in the cell in the diagonal (bottom right) if the symbols in the diagonal match

Shortest Edit Distance

$d = 0$: bootstrap

	\emptyset	A	C	D	E	G	F
\emptyset	0						
A		$\nwarrow 0$					
B							
C							
D							
E							
F							

Shortest Edit Distance

$d = 1$: rule 1

	\emptyset	A	C	D	E	G	F
\emptyset	0	$\leftarrow 1$					
A		$\nwarrow 0$	$\leftarrow 1$				
B							
C							
D							
E							
F							

Shortest Edit Distance

$d = 1$: rule 2

	\emptyset	A	C	D	E	G	F
\emptyset	0	$\leftarrow 1$					
A	$\uparrow 1$	$\nwarrow 0$	$\leftarrow 1$				
B		$\uparrow 1$					
C							
D							
E							
F							

Shortest Edit Distance

$d = 1$: rule 3

	\emptyset	A	C	D	E	G	F
\emptyset	0	$\leftarrow 1$					
A	$\uparrow 1$	$\nwarrow 0$	$\leftarrow 1$				
B		$\uparrow 1$					
C			$\nwarrow 1$				
D				$\nwarrow 1$			
E					$\nwarrow 1$		
F							

Shortest Edit Distance

$d = 2$: rule 1

	\emptyset	A	C	D	E	G	F
\emptyset	0	$\leftarrow 1$	$\leftarrow 2$				
A	$\uparrow 1$	$\nwarrow 0$	$\leftarrow 1$	$\leftarrow 2$			
B		$\uparrow 1$	$\leftarrow 2$				
C			$\nwarrow 1$	$\leftarrow 2$			
D				$\nwarrow 1$	$\leftarrow 2$		
E					$\nwarrow 1$	$\leftarrow 2$	
F							

Shortest Edit Distance

$d = 2$: rule 2

	\emptyset	A	C	D	E	G	F
\emptyset	0	$\leftarrow 1$	$\leftarrow 2$				
A	$\uparrow 1$	$\nwarrow 0$	$\leftarrow 1$	$\leftarrow 2$			
B	$\uparrow 2$	$\uparrow 1$	$\leftarrow 2$				
C		$\uparrow 2$	$\nwarrow 1$	$\leftarrow 2$			
D			$\uparrow 2$	$\nwarrow 1$	$\leftarrow 2$		
E				$\uparrow 2$	$\nwarrow 1$	$\leftarrow 2$	
F					$\uparrow 2$		

Shortest Edit Distance

$d = 2$: rule 3

	\emptyset	A	C	D	E	G	F
\emptyset	0	$\leftarrow 1$	$\leftarrow 2$				
A	$\uparrow 1$	$\nwarrow 0$	$\leftarrow 1$	$\leftarrow 2$			
B	$\uparrow 2$	$\uparrow 1$	$\leftarrow 2$				
C		$\uparrow 2$	$\nwarrow 1$	$\leftarrow 2$			
D			$\uparrow 2$	$\nwarrow 1$	$\leftarrow 2$		
E				$\uparrow 2$	$\nwarrow 1$	$\leftarrow 2$	
F					$\uparrow 2$		$\nwarrow 2$

Shortest Edit Distance

	\emptyset	A	C	D	E	G	F
\emptyset	0	$\leftarrow 1$	$\leftarrow 2$				
A	$\uparrow 1$	$\nwarrow 0$	$\leftarrow 1$	$\leftarrow 2$			
B	$\uparrow 2$	$\uparrow 1$	$\leftarrow 2$				
C		$\uparrow 2$	$\nwarrow 1$	$\leftarrow 2$			
D			$\uparrow 2$	$\nwarrow 1$	$\leftarrow 2$		
E				$\uparrow 2$	$\nwarrow 1$	$\leftarrow 2$	
F					$\uparrow 2$		$\nwarrow 2$

Diff Algorithms in Git

- Myers
 - Diff algorithm proposed by Myers with speed optimizations that may lead to a non-minimal edit distance
- Minimal
 - Myers with a guarantee of minimal edit distance
- Patience
 - Just considers the unique lines in both files for computing the LCS, potentially leading to a more precise result
- Histogram
 - Extends the Patience algorithm to support low-occurrence common lines instead of just unique lines, potentially leading to faster executions

References

- Cormen, T. H., Leiserson, C. E., Rivest, R. L., Stein, C., 2001. Introduction to Algorithms, 2nd ed., MIT Press.
- Hunt, J., McIlroy, M., “An Algorithm for Differential File Comparison”, Bell Laboratories, 1976.
- Miller, W., Myers, E., “A File Comparison Program”, Software: Practice and Experience, v. 15, n. 11, p. 1025-1040, 1985.

Diff

Leonardo Gresta Paulino Murta

leomurta@ic.uff.br