Algorithm 2.30: Compute-AC-Fail

Input: AC trie: *root*, *child*() and *patterns*() Output: AC failure function *fail*() and updated *patterns*() (1) Create new node *fallback*

- (2) for $c \in \Sigma$ do child(fallback, c) \leftarrow root
- (3) $fail(root) \leftarrow fallback$

- (3) fan(160c) < fan(160c)(4) $queue \leftarrow \{root\}$ (5) while $queue \neq \emptyset$ do (6) $u \leftarrow popfront(queue)$ (7) for $c \in \Sigma$ such that $child(u, c) \neq \bot$ do
- (8) $v \leftarrow child(u, c)$
- (9)
- $w \leftarrow fail(u)$ while child(w, c) = \perp do $w \leftarrow fail(w)$ (10) $fail(v) \leftarrow Child(w, c)$ (11)
- (12)́ $patterns(v) \leftarrow patterns(v) \cup patterns(fail(v))$
- pushback(queue, v)(13)
- (14)return (fail(), patterns())

The algorithm does a breath first traversal of the trie. This ensures that correct values of fail() and patterns() are already computed when needed.

105

Assuming σ is constant:

- The search time is O(n).
- The space complexity is $\mathcal{O}(m)$, where $m = ||\mathcal{P}||$.
- Implementation of patterns() requires care (exercise).
- The preprocessing time is $\mathcal{O}(m)$, where $m = ||\mathcal{P}||$.
 - The only non-trivial issue is the while-loop on line (10).
 - Let $root, v_1, v_2, \ldots, v_\ell$ be the nodes on the path from root to a node representing a pattern P_i . Let $w_j = fail(v_j)$ for all j. Let depth(v) be the depth of a node v (depth(root) = 0).
 - When processing v_j and computing $w_j = fail(v_j)$, we have $depth(w_j) = depth(w_{j-1}) + 1$ before line (10) and $depth(w_j) \le depth(w_{j-1}) + 1 t_j$ after line (10), where t_j is the number of rounds in the while-loop.
 - Thus, the total number of rounds in the while-loop when processing the nodes v_1,v_2,\ldots,v_ℓ is at most $\ell=|P_i|,$ and thus over the whole algorithm at most $||\mathcal{P}||.$

The analysis when σ is not constant is left as an exercise.

107

109

3. Approximate String Matching

Often in applications we want to search a text for something that is similar to the pattern but not necessarily exactly the same.

To formalize this problem, we have to specify what does "similar" mean. This can be done by defining a similarity or a distance measure.

A natural and popular distance measure for strings is the edit distance, also known as the Levenshtein distance.

fail(v) is correctly computed on lines (8)–(11):

- The nodes that represent suffixes of S_v that are exactly $fail^{*}(v) = \{v, fail(v), fail(fail(v)), \dots, root\}.$
- Let u = parent(v) and child(u, c) = v. Then $S_v = S_u c$ and a string S is a suffix of S_u iff S_c is suffix of S_v . Thus for any node w
 - If $w \in fail^{*}(v)$, then parent(fail(v)) $\in fail^{*}(u)$.
 - If $w \in fail^*(u)$ and $child(w,c) \neq \bot$, then $child(w,c) \in fail^*(v)$.
- Therefore, fail(v) = child(w, c), where w is the first node in $fail^*(u)$ other than u such that $child(w,c) \neq \bot$.

patterns(v) is correctly computed on line (12):

 $patterns(v) = \{i \mid P_i \text{ is a suffix of } S_v\}$ $= \{i \mid P_i = S_w \text{ and } w \in fail^*(v)\}$ $= \{i \mid P_i = S_v\} \cup patterns(fail(v))$

106

Summary: Exact String Matching

Exact string matching is a fundamental problem in stringology. We have seen several different algorithms for solving the problem.

The properties of the algorithms vary with respect to worst case time complexity, average case time complexity, type of alphabet (ordered/integer) and even space complexity.

The algorithms use a wide range of completely different techniques:

- There exists numerous algorithms for exact string matching but almost all them are based on these techniques.
- Many of the techniques can be adapted to other problems. All of the techniques have some uses in practice too.
 - 108

Edit distance

The edit distance ed(A, B) of two strings A and B is the minimum number of edit operations needed to change A into B. The allowed edit operations are:

- S Substitution of a single character with another character.
- I Insertion of a single character.
- D Deletion of a single character.

Example 3.1: Let A = Lewensteinn and B = Levenshtein. Then ed(A, B) = 3.

The set of edit operations can be described

with an edit sequence:	NNSNNNINNND
or with an alignment:	Lewens-teinn
	Levenshtein-

In the edit sequence, N means No edit.

110

There are many variations and extension of the edit distance, for example: • Hamming distance allows only the subtitution operation.

- Damerau-Levenshtein distance adds an edit operation:
- T Transposition swaps two adjacent characters.
- With weighted edit distance, each operation has a cost or weight, which can be other than one.
- Allow insertions and deletions (indels) of factors at a cost that is lower than the sum of character indels.

We will focus on the basic Levenshtein distance.

Levenshtein distance has the following two useful properties, which are not shared by all variations (exercise):

- Levenshtein distance is a metric.
- If ed(A, B) = k, there exists an edit sequence and an alignment with k edit operations, but no edit sequence or alignment with less than k edit operations. An edit sequence and an alignment with ed(A, B) edit operations is called optimal.

Computing Edit Distance

Given two strings A[1..m] and B[1..n], define the values d_{ij} with the recurrence:

> $d_{00} = 0$, $d_{i0} = i, \ 1 \le i \le m,$ $d_{0j} = j, 1 \leq j \leq n$, and $d_{i-1,j-1} + \delta(A[i], B[j])$ $d_{ii} = \min$ $d_{i-1,j} + 1$ $1 \le i \le m, 1 \le j \le n,$ $d_{i,j-1} + 1$

where

$$\delta(A[i], B[j]) = \begin{cases} 1 & \text{if } A[i] \neq B[j] \\ 0 & \text{if } A[i] = B[j] \end{cases}$$

Theorem 3.2: $d_{ij} = ed(A[1..i], B[1..j])$ for all $0 \le i \le m, 0 \le j \le n$. In particular, $d_{mn} = ed(A, B)$.

Example 3.3: A = ballad, B = handball

d		h	a	n	d	b	a	1	1
	0	1	2	3	4	5	6	7	8
b	1	1	2	3	4	4	5	6	7
a	2	2	1	2	З	4	4	5	6
1	3	3	2	2	З	4	5	4	5
1	4	4	3	3	З	4	5	5	4
a	5	5	4	4	4	4	4	5	5
d	6	6	5	5	4	5	5	5	6

 $ed(A, B) = d_{mn} = d_{6,8} = 6.$

113

The recurrence gives directly a dynamic programming algorithm for computing the edit distance.

Algorithm 3.4: Edit distance Input: strings A[1..m] and B[1..n]Output: ed(A, B)

(1) for $i \leftarrow 0$ to $m \text{ do } d_{i0} \leftarrow i$ (2) for $j \leftarrow 1$ to n do $d_{0j} \leftarrow j$ (3) for $j \leftarrow 1$ to n do (4) for $i \leftarrow 1$ to m do (5) $d_{ij} \leftarrow \min\{d_{i-1,j-1} + \delta(A[i], B[j]), d_{i-1,j} + 1, d_{i,j-1} + 1\}$ (6)return d_{mn}

The time and space complexity is $\mathcal{O}(mn)$.

115

Algorithm 3.5: Edit distance in $\mathcal{O}(m)$ space Input: strings A[1..m] and B[1..n]Output: ed(A, B)(1) for $i \leftarrow 0$ to m do $C[i] \leftarrow i$ (2) for $j \leftarrow 1$ to n do (3) $c \leftarrow C[0]; C[0] \leftarrow j$

(4) for
$$i \leftarrow 1$$
 to m do
(5) $d \leftarrow \min\{c + \delta(A[i], B[j]), C[i-1] + 1, C[i] + 1\}$
(6) $c \leftarrow C[i]$
(7) $C[i] \leftarrow d$
(8) return $C[m]$

• Note that because ed(A, B) = ed(B, A) (exercise), we can assume that m < n.

Example 3.6: A = ballad, B = handball

d	ł	ı		a		n		d		b		a		1		1	
	0 =	> ∶	1 :	⇒	2	\Rightarrow	3	\Rightarrow	4	\rightarrow	5	\rightarrow	6	\rightarrow	7	\rightarrow	8
b	1	Ы		\searrow		\searrow		\searrow		$\underline{\mathbb{N}}$							
	1		1 ·	\rightarrow	2	\rightarrow	3	\rightarrow	4		4	\rightarrow	5	\rightarrow	6	\rightarrow	7
a	↓ `	. к	Ļ	$\not\bowtie$								\searrow					
	2	2	2		1	\Rightarrow	2	\rightarrow	3	\rightarrow	4		4	\rightarrow	5	\rightarrow	6
1	↓ `	. ĸ	Ļ		\downarrow	$\underline{\aleph}$		$\underline{\aleph}$		\searrow		\searrow	\downarrow	$\underline{\aleph}$		\searrow	
	3	:	3		2		2	\Rightarrow	3	\rightarrow	4	\rightarrow	5		4	\rightarrow	5
1	↓ `	. к	Ļ		\downarrow	\searrow	\downarrow	$\underline{\aleph}$		$\underline{\aleph}$		\searrow		\searrow	\downarrow	$\underline{\aleph}$	
	4		4		3		3		3	\Rightarrow	4	\rightarrow	5		5		4
a	↓ `	. к	Ļ	\searrow	\downarrow	\searrow	\downarrow	\searrow	\downarrow	\searrow		\not					₩
	5	!	5		4		4		4		4		4	\Rightarrow	5		5
d	↓ `	. к	Ļ		\downarrow	\searrow	\downarrow	\searrow		\searrow	\downarrow	\searrow	\downarrow	Ø		\mathcal{U}	₩
	6		6		5		5		4	\rightarrow	5		5		5	\Rightarrow	6

There are 7 paths from (0,0) to (6,8) corresponding to 7 different optimal edit sequences and alignments, including the following three:

IIIINNNDD	SNISSNIS	SNSSINSI
ballad	ba-lla-d	ball-ad-
handball	handball	handball

Proof of Theorem 3.2. We use induction with respect to i + j. For brevity, write $A_i = A[1.i]$ and $B_j = B[1..j]$. Basis:

$$\begin{aligned} d_{00} &= 0 = ed(\epsilon, \epsilon) \\ d_{i0} &= i = ed(A_i, \epsilon) \quad (i \text{ deletions}) \\ d_{0j} &= j = ed(\epsilon, B_j) \quad (j \text{ insertions}) \end{aligned}$$

Induction step: We show that the claim holds for d_{ij} , $1 \le i \le m, 1 \le j \le n$. By induction assumption, $d_{pq} = ed(A_p, B_q)$ when p + q < i + j.

(j insertions)

Let E_{ij} be an optimal edit sequence with the cost $ed(A_i, B_j)$. We have three cases depending on what the last operation symbol in E_{ij} is:

N or S:
$$E_{ij} = E_{i-1,j-1}$$
N or $E_{ij} = E_{i-1,j-1}$ S and
 $ed(A_i, B_j) = ed(A_{i-1}, B_{j-1}) + \delta(A[i], B[j]) = d_{i-1,j-1} + \delta(A[i], B[j]).$
I: $E_{ij} = E_{i,j-1}$ I and $ed(A_i, B_j) = ed(A_i, B_{j-1}) + 1 = d_{i,j-1} + 1.$
D: $E_{ij} = E_{i-1,j}$ D and $ed(A_i, B_j) = ed(A_{i-1}, B_j) + 1 = d_{i-1,j} + 1.$

One of the cases above is always true, and since the edit sequence is optimal, it must be one with the minimum cost, which agrees with the definition of d_{ij} .

114

П

The space complexity can be reduced by noticing that each column of the matrix (d_{ij}) depends only on the previous column. We do not need to store older columns.

A more careful look reveals that, when computing d_{ij} , we only need to store the bottom part of column j-1 and the already computed top part of column j. We store these in an array C[0..m] and variables c and d as shown below:



116

It is also possible to find optimal edit sequences and alignments from the matrix d_{ij} .

An edit graph is a directed graph, where the nodes are the cells of the edit distance matrix, and the edges are as follows:

- If A[i] = B[j] and $d_{ij} = d_{i-1,j-1}$, there is an edge $(i 1, j 1) \rightarrow (i, j)$ labelled with N.
- If $A[i] \neq B[j]$ and $d_{ij} = d_{i-1,j-1} + 1$, there is an edge $(i-1,j-1) \rightarrow (i,j)$ labelled with S.
- If $d_{ij} = d_{i,j-1} + 1$, there is an edge $(i, j 1) \rightarrow (i, j)$ labelled with I.
- If d_{ij} = d_{i-1,j} + 1, there is an edge (i − 1, j) → (i, j) labelled with D.

Any path from (0,0) to (m,n) is labelled with an optimal edit sequence.



117