

Methods of Analysis of Textual Data (MATD)

Jiří Dvorský

March 7, 2023

Department of Computer Science
VSB – TU Ostrava

1/59

Lectures Outline

1. Pattern Matching

Exact Pattern Matching

- Searching for One Patterns

- Searching for Finite Set of Patterns

- Searching for (Regular) Infinite Set of Patterns in Text

Approximate Pattern Matching

2/59

Pattern Matching

Jiří Dvorský

Department of Computer Science
VSB – TU Ostrava

3/59

Pattern Matching

Exact Pattern Matching

Pattern Matching Algorithms – common declarations

```
1 const size_t PatternNotFound = -1;
2
3 const int AlphabetSize = 256;
```

4/59

Brute Force Algorithm – source code

```
1 size_t BruteForce(const string& Pattern, const string&
   Text, const size_t StartPosition)
2 {
3   for (size_t i = StartPosition; i < Text.length() -
   Pattern.length() + 1; i++)
4   {
5     size_t j = 0;
6     while (j < Pattern.length())
7     {
8       if (Text[i + j] != Pattern[j])
9         break;
10      j += 1;
11    }
12    if (j == Pattern.length())
13      return i;
```

5/59

Brute Force Algorithm – source code (cont.)

```
14 }
15 return PatternNotFound;
16 }
```

6/59

Brute Force Algorithm – example

First attempt

G	C	A	T	C	G	C	A	G	A	G	A	G	T	A	T	A	C	A	G	T	A	C	G	
1	2	3	4																					
G	C	A	G	A	G	A	G																	

Shift by 1

Second attempt

G	C	A	T	C	G	C	A	G	A	G	A	G	T	A	T	A	C	A	G	T	A	C	G	
	1																							
	G	C	A	G	A	G	A	G																

Shift by 1

7/59

Brute Force Algorithm – example (cont.)

Third attempt

G	C	A	T	C	G	C	A	G	A	G	A	G	T	A	T	A	C	A	G	T	A	C	G
		1																					
		G	C	A	G	A	G	A	G														

Shift by 1

Fourth attempt

G	C	A	T	C	G	C	A	G	A	G	A	G	T	A	T	A	C	A	G	T	A	C	G	
			1																					
			G	C	A	G	A	G	A	G														

Shift by 1

8/59

Brute Force Algorithm – example (cont.)

Fifth attempt

G	C	A	T	C	G	C	A	G	A	G	A	G	T	A	T	A	C	A	G	T	A	C	G	
				1																				
				G	C	A	G	A	G	A	G													

Shift by 1

Sixth attempt

G	C	A	T	C	G	C	A	G	A	G	A	G	T	A	T	A	C	A	G	T	A	C	G	
					1	2	3	4	5	6	7	8												
					G	C	A	G	A	G	A	G												

Shift by 1

9/59

Brute Force Algorithm – example (cont.)

Seventh attempt

G	C	A	T	C	G	C	A	G	A	G	A	G	T	A	T	A	C	A	G	T	A	C	G	
						1																		
						G	C	A	G	A	G	A	G											

Shift by 1

Eighth attempt

G	C	A	T	C	G	C	A	G	A	G	A	G	T	A	T	A	C	A	G	T	A	C	G	
							1																	
							G	C	A	G	A	G	A	G										

Shift by 1

10/59

Brute Force Algorithm – example (cont.)

Nineth attempt

G	C	A	T	C	G	C	A	G	A	G	A	G	T	A	T	A	C	A	G	T	A	C	G	
								1	2															
								G	C	A	G	A	G	A	G									

Shift by 1

Tenth attempt

G	C	A	T	C	G	C	A	G	A	G	A	G	T	A	T	A	C	A	G	T	A	C	G	
									1															
									G	C	A	G	A	G	A	G								

Shift by 1

11/59

Brute Force Algorithm – example (cont.)

Eleventh attempt

G	C	A	T	C	G	C	A	G	A	G	A	G	T	A	T	A	C	A	G	T	A	C	G
										1	2												
										G	C	A	G	A	G	A	G						

Shift by 1

Twelfth attempt

G	C	A	T	C	G	C	A	G	A	G	A	G	A	G	T	A	T	A	C	A	G	T	A	C	G	
															1											
															G	C	A	G	A	G	A	G				

Shift by 1

12/59

Brute Force Algorithm – example (cont.)

Thirteenth attempt

G	C	A	T	C	G	C	A	G	A	G	A	G	A	G	A	G	T	A	T	A	C	A	G	T	A	C	G		
																		1	2										
																		G	T	A	T	A	C	A	G	T	A	C	G
																		G	C	A	G	A	G	A	G				

Shift by 1

Fourteenth attempt

G	C	A	T	C	G	C	A	G	A	G	A	G	A	G	A	G	T	A	T	A	C	A	G	T	A	C	G					
																					1											
																					G	T	A	T	A	C	A	G	T	A	C	G
																					G	C	A	G	A	G	A	G				

Shift by 1

13/59

Brute Force Algorithm – example (cont.)

Fifteenth attempt

G	C	A	T	C	G	C	A	G	A	G	A	G	A	G	T	A	T	A	C	A	G	T	A	C	G			
																					1							
																					G	C	A	G	A	G	A	G

Shift by 1

Sixteenth attempt

G	C	A	T	C	G	C	A	G	A	G	A	G	A	G	T	A	T	A	C	A	G	T	A	C	G						
																					G	T	A	C	A	G	T	A	C	G	
																					G	C	A	G	A	G	A	G			

Shift by 1

14/59

Brute Force Algorithm – example (cont.)

Seventeenth attempt

G	C	A	T	C	G	C	A	G	A	G	A	G	A	G	T	A	T	A	C	A	G	T	A	C	G						
																					G	T	A	C	A	G	T	A	C	G	
																					G	C	A	G	A	G	A	G			

Shift by 1

The algorithm performs 30 character comparisons.

15/59

Morris-Pratt Algorithm – preprocessing

```
1 void MorrisPrattPreprocessing(const string& Pattern,
2     vector<int>& Next)
3 {
4     int i = 0;
5     int j = Next[0] = -1;
6     while (i < Pattern.length())
7     {
8         while (j > -1 && Pattern[i] != Pattern[j])
9             j = Next[j];
10        Next[++i] = ++j;
11    }
12 }
13 }
```

16/59

Morris-Pratt Algorithm – source code

```
1 size_t MorrisPratt(const string& Pattern, const string&
2     Text, const size_t StartPosition)
3 {
4     vector<int> Next(Pattern.length() + 1, 0);
5     MorrisPrattPreprocessing(Pattern, Next);
6     int i = 0;
7     size_t j = StartPosition;
8     while (j < Text.length())
9     {
10        while (i > -1 && Pattern[i] != Text[j])
11            i = Next[i];
12        i += 1;
13        j += 1;
14    }
```

17/59

Morris-Pratt Algorithm – source code (cont.)

```
15     if (i >= Pattern.length())
16     {
17         return j - i;
18     }
19 }
20 return PatternNotFound;
21 }
```

18/59

Morris-Pratt Algorithm

First attempt

G	C	A	T	C	G	C	A	G	A	G	A	G	T	A	T	A	C	A	G	T	A	C	G
1	2	3	4																				
G	C	A	G	A	G	A	G																

Shift by $i - \text{Next}[i] = 3 - 0 = 3$

Second attempt

G	C	A	T	C	G	C	A	G	A	G	A	G	T	A	T	A	C	A	G	T	A	C	G
			1																				
			G	C	A	G	A	G	A	G													

Shift by $i - \text{Next}[i] = 0 - (-1) = 1$

19/59

Morris-Pratt Algorithm (cont.)

Third attempt

G	C	A	T	C	G	C	A	G	A	G	A	G	T	A	T	A	C	A	G	T	A	C	G
				1																			
				G	C	A	G	A	G	A	G												

Shift by $i - Next[i] = 0 - (-1) = 1$

Fourth attempt

G	C	A	T	C	G	C	A	G	A	G	A	G	T	A	T	A	C	A	G	T	A	C	G
					1	2	3	4	5	6	7	8											
					G	C	A	G	A	G	A	G											

Shift by $i - Next[i] = 8 - 1 = 7$

20/59

Morris-Pratt Algorithm (cont.)

Fifth attempt

G	C	A	T	C	G	C	A	G	A	G	A	G	T	A	T	A	C	A	G	T	A	C	G
										1	2												
										G	T												
										G	C												

Shift by $i - Next[i] = 1 - 0 = 1$

Sixth attempt

G	C	A	T	C	G	C	A	G	A	G	A	G	T	A	T	A	C	A	G	T	A	C	G
													1										
													T										
													G										

Shift by $i - Next[i] = 0 - (-1) = 1$

21/59

Morris-Pratt Algorithm (cont.)

Seventh attempt

G	C	A	T	C	G	C	A	G	A	G	A	G	T	A	T	A	C	A	G	T	A	C	G
														1									
														A									
														G									

Shift by $i - Next[i] = 0 - (-1) = 1$

Eighth attempt

G	C	A	T	C	G	C	A	G	A	G	A	G	T	A	T	A	C	A	G	T	A	C	G
															1								
															T								
															G								

Shift by $i - Next[i] = 0 - (-1) = 1$

22/59

Morris-Pratt Algorithm (cont.)

Ninth attempt

G	C	A	T	C	G	C	A	G	A	G	A	G	T	A	T	A	C	A	G	T	A	C	G
																	1						
																	A						
																	G						

Shift by $i - Next[i] = 0 - (-1) = 1$

The algorithm performs 19 character comparisons.

23/59

Knuth-Morris-Pratt Algorithm

First attempt

G	C	A	T	C	G	C	A	G	A	G	A	G	T	A	T	A	C	A	G	T	A	C	G
1	2	3	4																				
G	C	A	G	A	G	A	G																

Shift by $i - Next[i] = 3 - (-1) = 4$

Second attempt

G	C	A	T	C	G	C	A	G	A	G	A	G	T	A	T	A	C	A	G	T	A	C	G
				1																			
				G	C	A	G	A	G	A	G												

Shift by $i - Next[i] = 0 - (-1) = 1$

24/59

Knuth-Morris-Pratt Algorithm (cont.)

Third attempt

G	C	A	T	C	G	C	A	G	A	G	A	G	T	A	T	A	C	A	G	T	A	C	G
					1	2	3	4	5	6	7	8											
					G	C	A	G	A	G	A	G											

Shift by $i - Next[i] = 8 - 1 = 7$

Fourth attempt

G	C	A	T	C	G	C	A	G	A	G	A	G	T	A	T	A	C	A	G	T	A	C	G
													1	2									
													G	C	A	G	A	G	A	G			

Shift by $i - Next[i] = 1 - 0 = 1$

25/59

Knuth-Morris-Pratt Algorithm (cont.)

Fifth attempt

G	C	A	T	C	G	C	A	G	A	G	A	G	T	A	T	A	C	A	G	T	A	C	G
														1									
														G	C	A	G	A	G	A	G		

Shift by $i - Next[i] = 0 - (-1) = 1$

Sixth attempt

G	C	A	T	C	G	C	A	G	A	G	A	G	T	A	T	A	C	A	G	T	A	C	G
															1								
															G	C	A	G	A	G	A	G	

Shift by $i - Next[i] = 0 - (-1) = 1$

26/59

Knuth-Morris-Pratt Algorithm (cont.)

Seventh attempt

G	C	A	T	C	G	C	A	G	A	G	A	G	T	A	T	A	C	A	G	T	A	C	G
																1							
																G	C	A	G	A	G	A	G

Shift by $i - Next[i] = 0 - (-1) = 1$

Eighth attempt

G	C	A	T	C	G	C	A	G	A	G	A	G	T	A	T	A	C	A	G	T	A	C	G	
																	1							
																	G	C	A	G	A	G	A	G

Shift by $i - Next[i] = 0 - (-1) = 1$

The algorithm performs 18 character comparisons.

27/59

Knuth-Morris-Pratt Algorithm – preprocessing

```
1 void KnuthMorrisPrattPreprocessing(const string&
   Pattern, vector<int>& Next)
2 {
3     int i = 0;
4     int j = Next[0] = -1;
5     while (i < Pattern.length())
6     {
7         while (j > -1 && Pattern[i] != Pattern[j])
8         {
9             j = Next[j];
10        }
11        i += 1;
12        j += 1;
13        Next[i] = Pattern[i] == Pattern[j] ? Next[j] : j;
14    }
```

28/59

Knuth-Morris-Pratt Algorithm – preprocessing (cont.)

```
15 }
```

29/59

Knuth-Morris-Pratt Algorithm – source code

```
1 size_t KnuthMorrisPratt(const string& Pattern, const
   string& Text, const size_t StartPosition)
2 {
3     vector<int> Next(Pattern.length() + 1, 0);
4     KnuthMorrisPrattPreprocessing(Pattern, Next);
5     int i = 0;
6     size_t j = StartPosition;
7     while (j < Text.length())
8     {
9         while (i > -1 && Pattern[i] != Text[j])
10        {
11            i = Next[i];
12        }
13        i += 1;
14        j += 1;
```

30/59

Knuth-Morris-Pratt Algorithm – source code (cont.)

```
15     if (i >= Pattern.length())
16     {
17         return j - i;
18     }
19 }
20 return PatternNotFound;
21 }
```

31/59

Searching for (Regular) Infinite Set of Patterns in Text

1. How to describe infinite set of pattern i.e. string?

Regular Expressions

2. What shall we use to perform matching?

Finite Automata

32/59

Regular Expressions and Languages

Regular expression R	Value of expression $h(R)$
Atomic expressions	
\emptyset	\emptyset
ε	$\{\varepsilon\}$
$a, a \in \Sigma$	$\{a\}$
Operations	
$U \cdot V$	$\{uv \mid u \in h(U) \wedge v \in h(V)\}$
$U + V$	$h(U) \cup h(V)$
$V^k = \underbrace{V \cdot V \cdot \dots \cdot V}_{k \text{ times}}$	
$V^+ = V^1 + V^2 + V^3 + \dots$	
$V^* = V^0 + V^1 + V^2 + \dots$	

33/59

Regular Expression Features

$$U + (V + W) = (U + V) + W$$

$$U \cdot (V \cdot W) = (U \cdot V) \cdot W$$

$$U + V = V + U$$

$$(U + V) \cdot W = (U \cdot W) + (V \cdot W)$$

$$U \cdot (V + W) = (U \cdot V) + (U \cdot W)$$

$$U + U = U$$

$$\varepsilon \cdot U = U$$

$$\emptyset \cdot U = \emptyset$$

$$U + \emptyset = U$$

$$U^* = \varepsilon + U^+$$

34/59

Deterministic Finite Automaton

Definition

Deterministic Finite Automaton (DFA) is a quintuple $A = (Q, \Sigma, q_0, \delta, F)$, where

- Q is a finite set of states
- Σ is an alphabet
- $q_0 \in Q$ is an initial state
- $\delta : Q \times \Sigma \rightarrow Q$ is a transition function
- $F \subseteq Q$ is a set of final states

35/59

Deterministic Finite Automaton (cont.)

Configuration of Finite Automaton

$$(q, w) \in Q \times \Sigma^*$$

Transition of Finite Automaton is a relation

$$\mapsto: (Q \times \Sigma^*) \times (Q \times \Sigma^*)$$

such as

$$(q, aw) \mapsto (q', w) \iff \delta(q, a) = q'$$

Automaton accepts word w if

$$(q_0, w) \mapsto^* (q, \epsilon), q \in F$$

36/59

Nondeterministic Finite Automaton

Definition

Nondeterministic Finite Automaton (NFA) is a quintuple $A = (Q, \Sigma, q_0, \delta, F)$, where

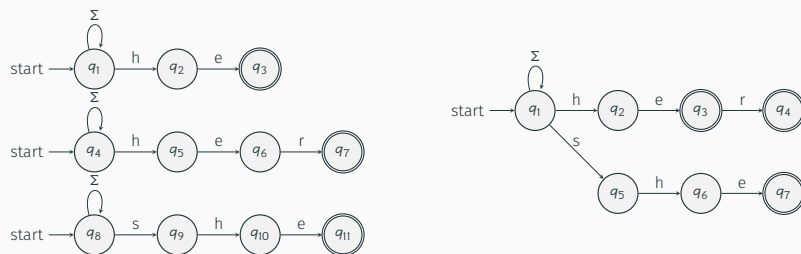
- Q is a finite set of states
- Σ is an alphabet
- $q_0 \in Q$ is an initial state
- $\delta: Q \times \Sigma \rightarrow P(Q)$ is a transition function
- $F \subseteq Q$ is a set of final states

- Alternatively NFA can be defined as $A = (Q, \Sigma, S, \delta, F)$, where $S \subseteq Q$ is a set of initial states.
- For each NFA, there is a DFA such that it recognizes the same formal language.

37/59

Nondeterministic Finite Automaton – example

Set of patterns $P = \{\text{he, her, she}\}$



38/59

NFA \rightarrow DFA Conversion

The DFA can be constructed using the **powerset construction**.

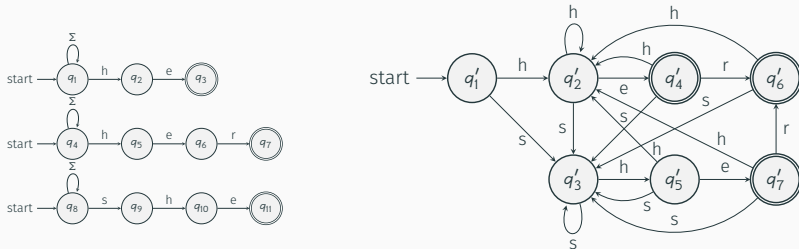
$$\text{NFA } A = (Q, \Sigma, S, \delta, F) \longrightarrow \text{DFA } A' = (Q', \Sigma', q'_0, \delta', F')$$

- $Q' \subseteq P(Q)$
- $\Sigma' = \Sigma$
- $q'_0 = S$
- $\delta'(q', x) = \cup \delta(q, x)$ for all $q \in q'$
- $F' = \{q' \in Q' \mid q' \cap F \neq \emptyset\}$

39/59

NFA → DFA Conversion I

State	Label	<i>e</i>	<i>h</i>	<i>r</i>	<i>s</i>	other
{1, 4, 8}	q'_1	{1, 4, 8}	{1, 2, 4, 5, 8}	{1, 4, 8}	{1, 4, 8, 9}	{1, 4, 8}
{1, 2, 4, 5, 8}	q'_2	{1, 3, 4, 6, 8}	{1, 2, 4, 5, 8}	{1, 4, 8}	{1, 4, 8, 9}	{1, 4, 8}
{1, 4, 8, 9}	q'_3	{1, 4, 8}	{1, 2, 4, 5, 8, 10}	{1, 4, 8}	{1, 4, 8, 9}	{1, 4, 8}
{1, 3, 4, 6, 8}	q'_4	{1, 4, 8}	{1, 2, 4, 5, 8}	{1, 4, 7, 8}	{1, 4, 8, 9}	{1, 4, 8}
{1, 2, 4, 5, 8, 10}	q'_5	{1, 3, 4, 6, 8, 11}	{1, 2, 4, 5, 8}	{1, 4, 8}	{1, 4, 8, 9}	{1, 4, 8}
{1, 4, 7, 8}	q'_6	{1, 4, 8}	{1, 2, 4, 5, 8}	{1, 4, 8}	{1, 4, 8, 9}	{1, 4, 8}
{1, 3, 4, 6, 8, 11}	q'_7	{1, 4, 8}	{1, 2, 4, 5, 8}	{1, 4, 7, 8}	{1, 4, 8, 9}	{1, 4, 8}

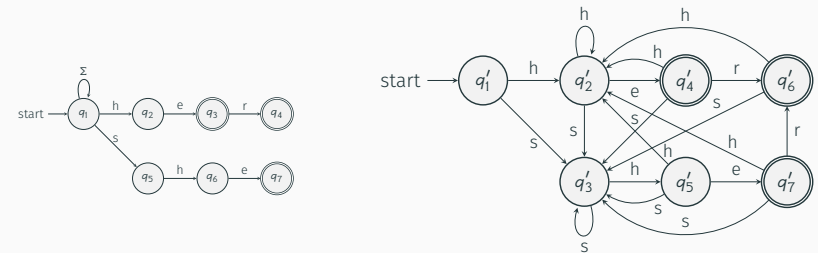


Only reachable states, transitions to state q_1 are not shown.

40/59

NFA → DFA Conversion II

State	Label	<i>e</i>	<i>h</i>	<i>r</i>	<i>s</i>	other
{}	q'_1	{}	{1, 2}	{}	{1, 5}	{}
{1, 2}	q'_2	{1, 3}	{1, 2}	{}	{1, 5}	{}
{1, 5}	q'_3	{}	{1, 2, 6}	{}	{1, 5}	{}
{1, 3}	q'_4	{}	{1, 2}	{1, 4}	{1, 5}	{}
{1, 2, 6}	q'_5	{1, 3, 7}	{1, 2}	{}	{1, 5}	{}
{1, 4}	q'_6	{}	{1, 2}	{}	{1, 5}	{}
{1, 3, 7}	q'_7	{}	{1, 2}	{1, 4}	{1, 5}	{}



41/59

Derivation of Regular Expression

For given regular expression R , derivation is defined as

$$h\left(\frac{dR}{dx}\right) = \{y \mid xy \in h(R)\}$$

Example

For $R = a + shell + stop + plot$ and its value $h(R) = \{a, shell, stop, plot\}$ derivations are

$$h\left(\frac{dR}{da}\right) = \{\varepsilon\}$$

$$h\left(\frac{dR}{ds}\right) = \{hell, top\}$$

$$h\left(\frac{dR}{dt}\right) = \emptyset$$

42/59

Derivation of Regular Expression – properties

$$\frac{d\emptyset}{da} = \emptyset, \forall a \in \Sigma$$

$$\frac{d\varepsilon}{da} = \emptyset, \forall a \in \Sigma$$

$$\frac{da}{da} = \varepsilon, \forall a \in \Sigma$$

$$\frac{db}{da} = \emptyset, \forall b \neq a$$

$$\frac{d(U + V)}{da} = \frac{dU}{da} + \frac{dV}{da}$$

$$\frac{d(U \cdot V)}{da} = \frac{dU}{da} \cdot V, \varepsilon \notin U$$

$$\frac{d(U \cdot V)}{da} = \frac{dU}{da} \cdot V + \frac{dV}{da}, \varepsilon \in U$$

$$\frac{dV^*}{da} = \frac{dV}{da} \cdot V^*$$

$$\frac{dV}{dx} = \frac{d}{da_n} \left(\frac{d}{da_{n-1}} \left(\dots \frac{d}{da_2} \left(\frac{dV}{da_1} \right) \right) \right), \text{ for } x = a_1 a_2 \dots a_n$$

43/59

Construction of DFA Derivations of RE

- Derivation of regular expressions allows directly and algorithmically build DFA for any regular expression.
 - Let V is given regular expression in alphabet Σ .
 - Each state of DFA defines a set of words, that move the DFA from this state to any of final states.
- So, every state can be associated with regular expression, defining this set of words

$$\begin{aligned}
 q_0 &= V \\
 \delta(q, x) &= \frac{dq}{dx} \\
 F &= \{q \in Q | \varepsilon \in h(q)\}
 \end{aligned}$$

44/59

Construction of DFA Derivations of RE – example

Lest's have $V = (0 + 1)^* \cdot 01$ over alphabet $\Sigma\{0, 1\}$.

Then $q_0 = (0 + 1)^* \cdot 01$

Example of derivations:

$$\begin{aligned}
 \frac{d((0 + 1)^* \cdot 01)}{d0} &= \frac{d((0 + 1)^*)}{d0} \cdot 01 + \frac{d01}{d0} \\
 &= \frac{d(0 + 1)}{d0} \cdot (0 + 1)^* \cdot 01 + 1 \\
 &= \left(\frac{d0}{d0} + \frac{d1}{d0} \right) \cdot (0 + 1)^* \cdot 01 + 1 \\
 &= (\varepsilon + \emptyset) \cdot (0 + 1)^* \cdot 01 + 1 \\
 &= (0 + 1)^* \cdot 01 + 1
 \end{aligned}$$

45/59

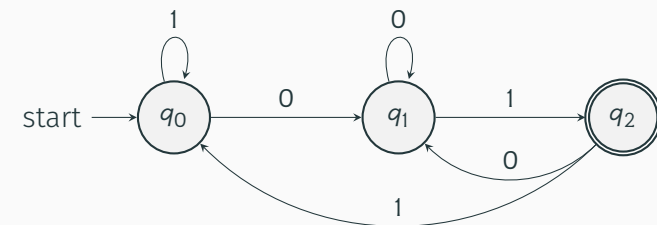
Construction of DFA Derivations of RE – example (cont.)

$$\begin{aligned}
 \frac{d((0 + 1)^* \cdot 01)}{d1} &= \frac{d((0 + 1)^*)}{d1} \cdot 01 + \frac{d01}{d1} \\
 &= \frac{d(0 + 1)}{d1} \cdot (0 + 1)^* \cdot 01 + \emptyset \\
 &= \left(\frac{d0}{d1} + \frac{d1}{d1} \right) \cdot (0 + 1)^* \cdot 01 \\
 &= (\emptyset + \varepsilon) \cdot (0 + 1)^* \cdot 01 \\
 &= (0 + 1)^* \cdot 01
 \end{aligned}$$

46/59

Construction of DFA Derivations of RE – example (cont.)

Regular Expression	State	0	1
$(0 + 1)^* \cdot 01$	q_0	$(0 + 1)^* \cdot 01 + 1$	$(0 + 1)^* \cdot 01$
$(0 + 1)^* \cdot 01 + 1$	q_1	$(0 + 1)^* \cdot 01 + 1$	$(0 + 1)^* \cdot 01 + \varepsilon$
$(0 + 1)^* \cdot 01 + \varepsilon$	q_2	$(0 + 1)^* \cdot 01 + 1$	$(0 + 1)^* \cdot 01$



47/59

Pattern Matching

Approximate Pattern Matching

Approximate Pattern Matching

- **String metric** (string distance function) is a metric that measures distance between two text strings for approximate string matching.
- String metric can be considered as “inverse similarity” – how two strings are dissimilar.
- There are two classic metrics
 1. Hamming distance
 2. Levenshtein distance
- Yes, string dissimilarity, distance can be measured. Both distances are metrics from mathematical point of view – non-negativity, identity, symmetry, and triangle inequality.

48/59

Hamming distance

Definition

Hamming distance between two strings of equal length is the number of positions at which the corresponding symbols are different.

In other words, it measures the minimum number of substitutions required to change one string into the other.

Example

Hamming distance of “karolin” and “kathrin” is 3.

k	a	r	o	l	i	n
k	a	t	h	r	i	n
0	0	1	1	1	0	0

49/59

Levenshtein distance

Definition

Levenshtein distance (1965) between two strings is the minimum number of single-character edits (insertions, deletions or substitutions) required to change one string into the other.

50/59

Levenshtein distance (cont.)

Example

Levenshtein distance between “kitten” and “sitting” is 3:

1. kitten → sitten (substitution of “s” for “k”)
2. sitten → sittin (substitution of “i” for “e”)
3. sittin → sitting (insertion of “g” at the end).

There is no way to do it with fewer than three edits.

51/59

Levenshtein distance (cont.)

Upper and lower bounds

The Levenshtein distance has several simple upper and lower bounds:

- It is at least the difference of the sizes of the two strings.
- It is at most the length of the longer string.
- It is zero if and only if the strings are equal.
- If the strings are the same size, the Hamming distance is an upper bound on the Levenshtein distance.
- The Levenshtein distance between two strings is no greater than the sum of their Levenshtein distances from a third string (triangle inequality).

52/59

Levenshtein distance (cont.)

$$d(i, j) = \begin{cases} i, & \text{if } j = 0 \\ j, & \text{if } i = 0 \\ \min \begin{pmatrix} d(i-1, j) + 1, \\ d(i, j-1) + 1, \\ d(i-1, j-1) + c(i, j) \end{pmatrix} \end{cases}$$

where

$$c(i, j) = \begin{cases} 0 & \text{if } a_i = b_j \\ 1 & \text{otherwise} \end{cases}$$

First element in the minimum corresponds to deletion (from a), the second to insertion (to b) and the third to match or mismatch.

53/59

Levenshtein distance (cont.)

```
1 int LevenshteinDistance(const char *s, int len_s, const
   char *t, int len_t)
2 {
3     int cost;
4
5     /* base case: empty strings */
6     if (len_s == 0) return len_t;
7     if (len_t == 0) return len_s;
8
9     /* test if last characters of the strings match */
10    if (s[len_s-1] == t[len_t-1])
11        cost = 0;
12    else
13        cost = 1;
14
```

54/59

Levenshtein distance (cont.)

```

15 /* return minimum of delete char from s, delete char
    from t, and delete char from both */
16 return minimum
17 (
18     LevenshteinDistance(s, len_s-1, t, len_t) + 1,
19     LevenshteinDistance(s, len_s, t, len_t-1) + 1,
20     LevenshteinDistance(s, len_s-1, t, len_t-1) + cost
21 );
22 }
    
```

55/59

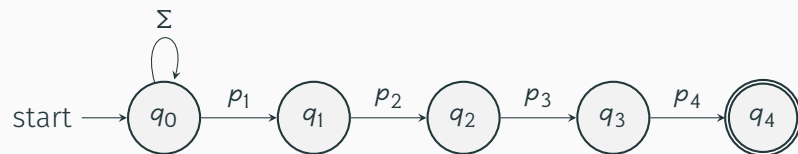
Levenshtein distance (cont.)

	0	k	i	t	t	e	n
s	1	1 ^a	2	3	4	5	6
i	2	2	1 ^b	2	3	4	5
t	3	3	2	1 ^c	2	3	4
t	4	4	3	2	1 ^d	2	3
i	5	5	4	3	2	2 ^e	3
n	6	6	5	4	3	3	2 ^f
g	7	7	6	5	4	4	3 ^g

^asubst. of k for s
^bi is equal i
^ct is equal t
^dt is equal t
^esubst. of e for i
^fn is equal n
^gdelete g

56/59

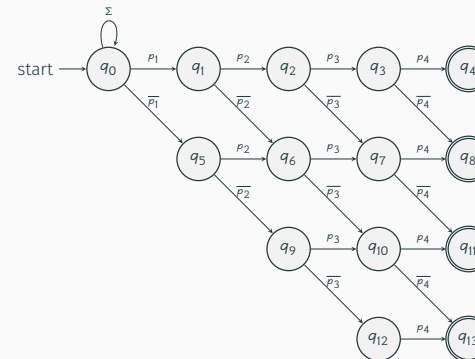
Approximate pattern matching using finite automata



NFA for the exact string matching ($m = 4$)

57/59

Approximate pattern matching using finite automata (cont.)

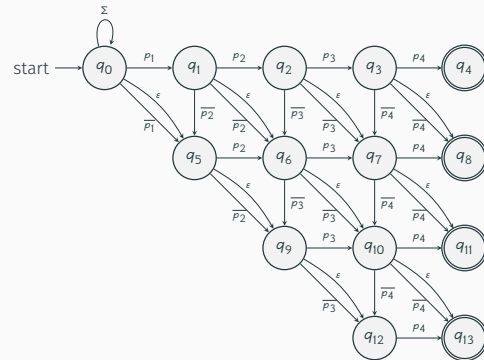


left to right – match
 diagonal – replace

NFA for the approximate string matching using the Hamming distance ($m = 4, k = 3$)

58/59

Approximate pattern matching using finite automata (cont.)



left to right – match
diagonal – replace
down – insert
diagonal ϵ – delete

NFA for the approximate string matching using the Levenshtein distance ($m = 4, k = 3$)