# Introduction to computational linguistics

## Exercise session 1: "Algorithms for matching"

*Thursday May 7 2009*

a. Show how to compute $Z_i$ stepwise for $i > 1$ (using the notion of Z-boxes) for the following strings:

   i. AABCAABXAAZ

   ii. ABCDXABCYABDXY

b. Apply the Knuth-Morris-Pratt algorithm to find occurrences of ABXYABXZ in XABXYABXYABXZABXZABXYABXZA

S = AABCAABXAAZ

**Step 0)**

Compute $Z_2(S)$ by comparing left-to-right S[2..ISI] and S[1..ISI] until a mismatch is found; $Z_2(S)$ is the length of that string. If $Z_2(S) > 0$ then r=r2=$Z_2(S)$+1 and l=2, else l=r=0

| S | A | A | B | C | A | A | B | X | A | A | Z |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 |
| $Z_i(S)$ | -- | 1 | | | | | | | | | |

$Z_2(S)$=1: { $A$ **A** B ... } so l=2, r=$Z_2(S)$+1=1+1=2

Exercise solutions: Matching  (Introduction to Comp.Ling)

Thursday, May 7, 2009

Thursday, May 7, 2009

**Step 1)**

Deutsches Forschungszentrum für Künstliche Intelligenz
German Research Center for Artificial Intelligence

## Step 1)

k > r: 3 > (r=2) so find $Z_3(S)$ by comparing S[3...|S|] to S[1..|S|] until a mismatch is found; if $Z_3(S) > 0$ then l=3, r=3+$Z_3(S)$-1

Thursday, May 7, 2009

## Step 1)

k > r: 3 > (r=2) so find $Z_3$(S) by comparing S[3...ISI] to S[1..ISI] until a mismatch is found; if $Z_3$(S) > 0 then l=3, r=3+$Z_3$(S)-1

S(3)='B' ≠ S(1)='A',  hence $Z_3$(S)=0, l and r remain as they are: l=r=2

Thursday, May 7, 2009

Deutsches Forschungszentrum für Künstliche Intelligenz
German Research Center for Artificial Intelligence

## Step 1)

k > r: 3 > (r=2) so find $Z_3(S)$ by comparing S[3...|S|] to S[1..|S|] until a mismatch is found; if $Z_3(S) > 0$ then l=3, r=3+$Z_3(S)$-1

S(3)='B' ≠ S(1)='A',  hence $Z_3(S)=0$, l and r remain as they are: l=r=2

| S | A | A | B | C | A | A | B | X | A | A | Z |
|---|---|---|---|---|---|---|---|---|---|---|---|
|   | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 |
| $Z_i(S)$ | -- | 1 | 0 |  |  |  |  |  |  |  |  |

$Z_3(S)=0$ so l=2, r=2

Exercise solutions: Matching  (Introduction to Comp.Ling)

Thursday, May 7, 2009

Exercise solutions: Matching  (Introduction to Comp.Ling)

Thursday, May 7, 2009

**Step 1)**

Thursday, May 7, 2009

## Step 1)

k > r: 4 > (r=2) so find $Z_4$(S) by comparing S[4...ISI] to S[1..ISI] until a mismatch is found; if $Z_4$(S) > 0 then l=4, r=4+$Z_4$(S)-1

Exercise solutions: Matching  (Introduction to Comp.Ling)

Thursday, May 7, 2009

## Step 1)

k > r: 4 > (r=2) so find $Z_4$(S) by comparing S[4...|S|] to S[1..|S|] until a mismatch is found; if $Z_4$(S) > 0 then l=4, r=4+$Z_4$(S)-1

S(4)='C' ≠ S(1)='A',  hence $Z_4$(S)=0, l and r remain as they are: l=r=2

Exercise solutions: Matching  (Introduction to Comp.Ling)

Thursday, May 7, 2009

**Step 1)**

$k > r$: $4 > (r=2)$ so find $Z_4(S)$ by comparing $S[4...|S|]$ to $S[1..|S|]$ until a mismatch is found; if $Z_4(S) > 0$ then $l=4$, $r=4+Z_4(S)-1$

$S(4)='C' \neq S(1)='A'$, hence $Z_4(S)=0$, l and r remain as they are: $l=r=2$

| S | A | A | B | C | A | A | B | X | A | A | Z |
|---|---|---|---|---|---|---|---|---|---|---|---|
|   | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 |
| $Z_i(S)$ | -- | 1 | 0 | 0 | | | | | | | |

$Z_4(S)=0$ so $l=2$, $r=2$

Exercise solutions: Matching  (Introduction to Comp.Ling)

Thursday, May 7, 2009

Deutsches Forschungszentrum für Künstliche Intelligenz
German Research Center for Artificial Intelligence

Thursday, May 7, 2009

Deutsches Forschungszentrum für Künstliche Intelligenz
German Research Center for Artificial Intelligence

**Step 1)**

Exercise solutions: Matching  (Introduction to Comp.Ling)

Thursday, May 7, 2009

**Step 1)**

k > r: 5 > (r=2) so find $Z_5(S)$ by comparing S[5...ISI] to S[1..ISI] until a mismatch is found; if $Z_5(S) > 0$ then l=5, r=5+$Z_5(S)$-1

Exercise solutions: Matching  (Introduction to Comp.Ling)

Thursday, May 7, 2009

## Step 1)

k > r: 5 > (r=2) so find $Z_5(S)$ by comparing S[5...ISI] to S[1..ISI] until a mismatch is found; if $Z_5(S) > 0$ then I=5, r=5+$Z_5(S)$-1

S[5..7]="A A B" matches S[1..3]="A A B", hence $Z_5(S)$=3, and I and r are set as follows: I=5, r=5+$Z_5(S)$-1=5+3-1=7

Exercise solutions: Matching  (Introduction to Comp.Ling)

Thursday, May 7, 2009

## Step 1)

k > r: 5 > (r=2) so find $Z_5$(S) by comparing S[5...ISI] to S[1..ISI] until a mismatch is found; if $Z_5$(S) > 0 then l=5, r=5+$Z_5$(S)-1

S[5..7]="A A B" matches S[1..3]="A A B", hence $Z_5$(S)=3, and l and r are set as follows: l=5, r=5+$Z_5$(S)-1=5+3-1=7

| S | A | A | B | C | A | A | B | X | A | A | Z |
|---|---|---|---|---|---|---|---|---|---|---|---|
|  | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 |
| $Z_i$(S) | -- | 1 | 0 | 0 | 3 |  |  |  |  |  |  |

$Z_5$(S)=3 so l=5, r=7

Exercise solutions: Matching  (Introduction to Comp.Ling)

Thursday, May 7, 2009

Exercise solutions: Matching  (Introduction to Comp.Ling)

Thursday, May 7, 2009

**Step 2)**

Thursday, May 7, 2009

## Step 2)

6 ≤ (r=7): position k=6 is contained in a Z-box (namely, "AAB"=S[5..7], with S(6)='A').

## Step 2)

$6 \leq (r=7)$: position k=6 is contained in a Z-box (namely, "AAB"=S[5..7], with S(6)='A').

Hence S(6) also appears in k'=k-l=6-5+1=2: S(6)=S(2)='A'

## Step 2)

$6 \leq (r=7)$: position k=6 is contained in a Z-box (namely, "AAB"=S[5..7], with S(6)='A').

Hence S(6) also appears in k'=k-l=6-5+1=2: S(6)=S(2)='A'

Therefore, S[6..7] must match S[2..3], which it does

Exercise solutions: Matching  (Introduction to Comp.Ling)

Thursday, May 7, 2009

## Step 2)

$6 \leq (r=7)$: position k=6 is contained in a Z-box (namely, "AAB"=S[5..7], with S(6)='A').

Hence S(6) also appears in k'=k-l=6-5+1=2: S(6)=S(2)='A'

Therefore, S[6..7] must match S[2..3], which it does

Furthermore, there must be a match to a prefix of S of length minimum [$Z_2$(S), lS[2..3]l], i.e. minimum [ 1,r-k+1=2] = 2

## Step 2)

6 ≤ (r=7): position k=6 is contained in a Z-box (namely, "AAB"=S[5..7], with S(6)='A').

Hence S(6) also appears in k'=k-l=6-5+1=2: S(6)=S(2)='A'

Therefore, S[6..7] must match S[2..3], which it does

Furthermore, there must be a match to a prefix of S of length minimum [$Z_2$(S), IS[2..3]I], i.e. minimum [ 1,r-k+1=2] = 2

## Step 2a)

Exercise solutions: Matching  (Introduction to Comp.Ling)

Thursday, May 7, 2009

## Step 2)

$6 \leq (r=7)$: position k=6 is contained in a Z-box (namely, "AAB"=S[5..7], with S(6)='A').

Hence S(6) also appears in k'=k-l=6-5+1=2: S(6)=S(2)='A'

Therefore, S[6..7] must match S[2..3], which it does

Furthermore, there must be a match to a prefix of S of length minimum [$Z_2$(S), lS[2..3]l], i.e. minimum [ 1,r-k+1=2] = 2

## Step 2a)

$Z_6$(S)=$Z_2$(S)=1 which is smaller than the length of S[2..3], hence l and r stay the same

Exercise solutions: Matching  (Introduction to Comp.Ling)

Thursday, May 7, 2009

## Step 2)

6 ≤ (r=7): position k=6 is contained in a Z-box (namely, "AAB"=S[5..7], with S(6)='A').

Hence S(6) also appears in k'=k-l=6-5+1=2: S(6)=S(2)='A'

Therefore, S[6..7] must match S[2..3], which it does

Furthermore, there must be a match to a prefix of S of length minimum [$Z_2$(S), IS[2..3]l], i.e. minimum [ 1,r-k+1=2] = 2

## Step 2a)

$Z_6$(S)=$Z_2$(S)=1 which is smaller than the length of S[2..3], hence l and r stay the same

| S | A | A | B | C | A | A | B | X | A | A | Z |
|---|---|---|---|---|---|---|---|---|---|---|---|
|  | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 |
| $Z_i$(S) | -- | 1 | 0 | 0 | 3 | 1 |  |  |  |  |  |

$Z_6$(S)=$Z_2$(S)=1 so l and r remain the same: l=5, r=7

Exercise solutions: Matching  (Introduction to Comp.Ling)

Thursday, May 7, 2009

Exercise solutions: Matching  (Introduction to Comp.Ling)

Thursday, May 7, 2009

**Step 2)**

Exercise solutions: Matching  (Introduction to Comp.Ling)

Thursday, May 7, 2009

**Step 2)**

7 ≤ (r=7): position k=7 is contained in S[5..7], with S(7)='B'.

Exercise solutions: Matching  (Introduction to Comp.Ling)

Thursday, May 7, 2009

**Step 2)**

$7 \leq (r=7)$: position $k=7$ is contained in $S[5..7]$, with $S(7)=$'B'.

Hence $S(7)$ also appears in $k'=k-l=7-5+1=3$: $S(7)=S(3)=$'B'

Exercise solutions: Matching  (Introduction to Comp.Ling)

Thursday, May 7, 2009

## Step 2)

7 ≤ (r=7): position k=7 is contained in S[5..7], with S(7)='B'.

Hence S(7) also appears in k'=k-l=7-5+1=3: S(7)=S(3)='B'

Therefore, S[7..7] must match S[3..3], i.e. S(7)=S(3), which it does

Thursday, May 7, 2009

## Step 2)

$7 \leq (r=7)$: position k=7 is contained in S[5..7], with S(7)='B'.

Hence S(7) also appears in k'=k-l=7-5+1=3: S(7)=S(3)='B'

Therefore, S[7..7] must match S[3..3], i.e. S(7)=S(3), which it does

Furthermore, there must be a match to a prefix of S of length minimum [$Z_3$(S), lS[3..3]l], i.e. minimum [ 0,r-k+1=1] = 1

## Step 2)

$7 \leq (r=7)$: position k=7 is contained in S[5..7], with S(7)='B'.

Hence S(7) also appears in k'=k-l=7-5+1=3: S(7)=S(3)='B'

Therefore, S[7..7] must match S[3..3], i.e. S(7)=S(3), which it does

Furthermore, there must be a match to a prefix of S of length minimum [$Z_3$(S), IS[3..3]I], i.e. minimum [ 0,r-k+1=1] = 1

## Step 2a)

## Step 2)

$7 \leq (r=7)$: position k=7 is contained in S[5..7], with S(7)='B'.

Hence S(7) also appears in k'=k-l=7-5+1=3: S(7)=S(3)='B'

Therefore, S[7..7] must match S[3..3], i.e. S(7)=S(3), which it does

Furthermore, there must be a match to a prefix of S of length minimum [$Z_3$(S), IS[3..3]I], i.e. minimum [ 0,r-k+1=1] = 1

## Step 2a)

$Z_7$(S)=$Z_3$(S)=0 which is smaller than the length of S[3..3], hence l and r stay the same

Thursday, May 7, 2009

## Step 2)

7 ≤ (r=7): position k=7 is contained in S[5..7], with S(7)='B'.

Hence S(7) also appears in k'=k-l=7-5+1=3: S(7)=S(3)='B'

Therefore, S[7..7] must match S[3..3], i.e. S(7)=S(3), which it does

Furthermore, there must be a match to a prefix of S of length minimum [$Z_3$(S), lS[3..3]l], i.e. minimum [ 0,r-k+1=1] = 1

## Step 2a)

$Z_7$(S)=$Z_3$(S)=0 which is smaller than the length of S[3..3], hence l and r stay the same

| S | A | A | B | C | A | A | B | X | A | A | Z |
|---|---|---|---|---|---|---|---|---|---|---|---|
|   | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 |
| $Z_i$(S) | -- | 1 | 0 | 0 | 3 | 1 | 0 |   |   |   |   |

$Z_7$(S)=$Z_3$(S)=0 so l and r remain the same: l=5, r=7

Exercise solutions: Matching  (Introduction to Comp.Ling)

Thursday, May 7, 2009

k=8 > (r=7) so step 1:

match S[8..|S|] to S[1..|S|]: mismatch, so $Z_8$(S)=0, l and r remain the same

| S | A | A | B | C | A | A | B | X | A | A | Z |
|---|---|---|---|---|---|---|---|---|---|---|---|
|   | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 |
| $Z_i$(S) | -- | 1 | 0 | 0 | 3 | 1 | 0 | 0 |   |   |   |

$Z_8$(S)=0 so l=5, r=7

k=9 > (r=7) so step 1:

match S[9..|S|] to S[1..|S|]: match S[9..10]=S[1..2], so $Z_9$(S)=2, l=9 and r=10

| S | A | A | B | C | A | A | B | X | A | A | Z |
|---|---|---|---|---|---|---|---|---|---|---|---|
|   | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 |
| $Z_i$(S) | -- | 1 | 0 | 0 | 3 | 1 | 0 | 0 | 2 |   |   |

$Z_9$(S)=2 so l=9, r=10

Exercise solutions: Matching  (Introduction to Comp.Ling)

Thursday, May 7, 2009

k=10 ≤ (r=10) so step 2:

S(10) contained in S[9..10]; S(10) matches S(10-9+1)=S(2)='A'; $Z_2$(S)=1 ≥ l

S[10..10]l=10-10+1=1, hence **Step 2b)** but mismatch

| S | A | A | B | C | A | A | B | X | A | A | Z |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | l | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | l0 | ll |
| $Z_i$(S) | -- | l | 0 | 0 | 3 | l | 0 | 0 | 2 | l | |

$Z_{10}$(S)=1

k=11 > (r=10) so step 1:

match S[11..ISI] to S[1..ISI]: mismatch so $Z_{11}$(S)=0

| S | A | A | B | C | A | A | B | X | A | A | Z |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | l | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | l0 | ll |
| $Z_i$(S) | -- | l | 0 | 0 | 3 | l | 0 | 0 | 2 | l | 0 |

$Z_{11}$(S)=0

Exercise solutions: Matching (Introduction to Comp.Ling)

Thursday, May 7, 2009

Thursday, May 7, 2009

$Z_2(S)$: $S(2) \neq S(1)$ so $Z_2(S)=0$, r=l=0

Thursday, May 7, 2009

$Z_2(S)$: $S(2) \neq S(1)$ so $Z_2(S)=0$, r=l=0

| S | A | B | C | D | X | A | B | C | Y | A | B | D | X | Y | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | | $Z_2(S)=0$ |
| $Z_i(S)$ | -- | 0 | | | | | | | | | | | | | | |

Exercise solutions: Matching  (Introduction to Comp.Ling)

Thursday, May 7, 2009

$Z_2(S)$: $S(2) \neq S(1)$ so $Z_2(S)=0$, r=l=0

| S | A | B | C | D | X | A | B | C | Y | A | B | D | X | Y |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 |
| $Z_i(S)$ | -- | 0 | | | | | | | | | | | | |

$Z_2(S)=0$

i=3..5: $Z_i(S)$: $S(i) \neq S(1)$ so $Z_i(S)=0$, r=l=0

Exercise solutions: Matching (Introduction to Comp.Ling)

Thursday, May 7, 2009

$Z_2(S)$: S(2) ≠ S(1) so $Z_2(S)$=0, r=l=0

| S | A | B | C | D | X | A | B | C | Y | A | B | D | X | Y |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 |
| $Z_i(S)$ | -- | 0 | | | | | | | | | | | | |

$Z_2(S)$=0

i=3..5: $Z_i(S)$: S(i) ≠ S(1) so $Z_i(S)$=0, r=l=0

| S | A | B | C | D | X | A | B | C | Y | A | B | D | X | Y |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 |
| $Z_i(S)$ | -- | 0 | 0 | 0 | 0 | | | | | | | | | |

$Z_{i=3..5}(S)$=0

$Z_2(S)$: $S(2) \neq S(1)$ so $Z_2(S)=0$, r=l=0

| S | A | B | C | D | X | A | B | C | Y | A | B | D | X | Y | $Z_2(S)=0$ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | |
| $Z_i(S)$ | -- | 0 | | | | | | | | | | | | | |

i=3..5: $Z_i(S)$: $S(i) \neq S(1)$ so $Z_i(S)=0$, r=l=0

| S | A | B | C | D | X | A | B | C | Y | A | B | D | X | Y | $Z_{i=3..5}(S)=0$ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | |
| $Z_i(S)$ | -- | 0 | 0 | 0 | 0 | | | | | | | | | | |

$Z_6(S)$: $S(6) = S(1)$: $S[6..8]$ matches $S[1..3]$, so $Z_6(S)=3$, l=6 and r=8

Exercise solutions: Matching  (Introduction to Comp.Ling)

Thursday, May 7, 2009

# a.ii) $Z_i$ for ABCDXABCYABDXY

$Z_2(S)$: $S(2) \neq S(1)$ so $Z_2(S)=0$, r=l=0

| S | A | B | C | D | X | A | B | C | Y | A | B | D | X | Y |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 |
| $Z_i(S)$ | -- | 0 | | | | | | | | | | | | |

$Z_2(S)=0$

i=3..5: $Z_i(S)$: $S(i) \neq S(1)$ so $Z_i(S)=0$, r=l=0

| S | A | B | C | D | X | A | B | C | Y | A | B | D | X | Y |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 |
| $Z_i(S)$ | -- | 0 | 0 | 0 | 0 | | | | | | | | | |

$Z_{i=3..5}(S)=0$

$Z_6(S)$: $S(6) = S(1)$: $S[6..8]$ matches $S[1..3]$, so $Z_6(S)=3$, l=6 and r=8

| S | A | B | C | D | X | A | B | C | Y | A | B | D | X | Y |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 |
| $Z_i(S)$ | -- | 0 | 0 | 0 | 0 | 3 | | | | | | | | |

$Z_6(S)=3$

Exercise solutions: Matching  (Introduction to Comp.Ling)

Thursday, May 7, 2009

Exercise solutions: Matching  (Introduction to Comp.Ling)

Thursday, May 7, 2009

$Z_7(S)$: $7 \leq (r=8)$ hence $S(7)=S(7-6+1)=S(2)=$'B', $Z_2(S)=0$ whereas $|S[7..8]|=2$,

hence $Z_7(S)=Z_2(S)=0$ and l and remain as they are: l=6 and r=8

Exercise solutions: Matching  (Introduction to Comp.Ling)

Thursday, May 7, 2009

$Z_7(S)$: $7 \leq (r=8)$ hence $S(7)=S(7-6+1)=S(2)='B'$, $Z_2(S)=0$ whereas $|S[7..8]|=2$,

hence $Z_7(S)=Z_2(S)=0$ and l and remain as they are: l=6 and r=8

| S | A | B | C | D | X | A | B | C | Y | A | B | D | X | Y |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 |
| $Z_i(S)$ | -- | 0 | 0 | 0 | 0 | 3 | 0 | | | | | | | |

$Z_7(S)=0$

Exercise solutions: Matching  (Introduction to Comp.Ling)

Thursday, May 7, 2009

$Z_7(S)$: $7 \leq (r=8)$ hence $S(7)=S(7-6+1)=S(2)=$'B', $Z_2(S)=0$ whereas $|S[7..8]|=2$,

hence $Z_7(S)=Z_2(S)=0$ and l and remain as they are: l=6 and r=8

| S | A | B | C | D | X | A | B | C | Y | A | B | D | X | Y |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 |
| $Z_i(S)$ | -- | 0 | 0 | 0 | 0 | 3 | 0 | | | | | | | |

$Z_7(S)=0$

$Z_8(S)$: $8 \leq (r=8)$ hence $S(8)=S(8-6+1)=S(3)=$'C', $Z_3(S)=0$ whereas $|S[8..8]|=1$,

hence $Z_8(S)=Z_3(S)=0$ and l and remain as they are: l=6 and r=8

Exercise solutions: Matching  (Introduction to Comp.Ling)

Thursday, May 7, 2009

$Z_7(S)$: $7 \le (r=8)$ hence $S(7)=S(7-6+1)=S(2)=$'B', $Z_2(S)=0$ whereas $|S[7..8]|=2$,

hence $Z_7(S)=Z_2(S)=0$ and l and remain as they are: l=6 and r=8

| S | A | B | C | D | X | A | B | C | Y | A | B | D | X | Y |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
|  | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 |
| $Z_i(S)$ | -- | 0 | 0 | 0 | 0 | 3 | 0 |  |  |  |  |  |  |  |

$Z_7(S)=0$

$Z_8(S)$: $8 \le (r=8)$ hence $S(8)=S(8-6+1)=S(3)=$'C', $Z_3(S)=0$ whereas $|S[8..8]|=1$,

hence $Z_8(S)=Z_3(S)=0$ and l and remain as they are: l=6 and r=8

| S | A | B | C | D | X | A | B | C | Y | A | B | D | X | Y |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
|  | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 |
| $Z_i(S)$ | -- | 0 | 0 | 0 | 0 | 3 | 0 | 0 |  |  |  |  |  |  |

$Z_8(S)=0$

Exercise solutions: Matching  (Introduction to Comp.Ling)

Thursday, May 7, 2009

$Z_7(S)$: $7 \leq (r=8)$ hence $S(7)=S(7-6+1)=S(2)=$'B', $Z_2(S)=0$ whereas $|S[7..8]|=2$,

hence $Z_7(S)=Z_2(S)=0$ and I and remain as they are: $l=6$ and $r=8$

| S | A | B | C | D | X | A | B | C | Y | A | B | D | X | Y |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 |
| $Z_i(S)$ | -- | 0 | 0 | 0 | 0 | 3 | 0 | | | | | | | |

$Z_7(S)=0$

$Z_8(S)$: $8 \leq (r=8)$ hence $S(8)=S(8-6+1)=S(3)=$'C', $Z_3(S)=0$ whereas $|S[8..8]|=1$,

hence $Z_8(S)=Z_3(S)=0$ and I and remain as they are: $l=6$ and $r=8$

| S | A | B | C | D | X | A | B | C | Y | A | B | D | X | Y |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 |
| $Z_i(S)$ | -- | 0 | 0 | 0 | 0 | 3 | 0 | 0 | | | | | | |

$Z_8(S)=0$

$Z_9(S)$: $9 > (r=8)$ but $S(9) \neq S(1)$ hence $Z_9(S)=0$ and I and remain as they are: $l=6$ and $r=8$

Exercise solutions: Matching  (Introduction to Comp.Ling)

Thursday, May 7, 2009

$Z_7(S)$: $7 \leq (r=8)$ hence $S(7)=S(7-6+1)=S(2)=$'B', $Z_2(S)=0$ whereas $|S[7..8]|=2$,

hence $Z_7(S)=Z_2(S)=0$ and l and remain as they are: l=6 and r=8

| S | A | B | C | D | X | A | B | C | Y | A | B | D | X | Y |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
|  | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 |
| $Z_i(S)$ | -- | 0 | 0 | 0 | 0 | 3 | 0 |  |  |  |  |  |  |  |

$Z_7(S)=0$

$Z_8(S)$: $8 \leq (r=8)$ hence $S(8)=S(8-6+1)=S(3)=$'C', $Z_3(S)=0$ whereas $|S[8..8]|=1$,

hence $Z_8(S)=Z_3(S)=0$ and l and remain as they are: l=6 and r=8

| S | A | B | C | D | X | A | B | C | Y | A | B | D | X | Y |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
|  | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 |
| $Z_i(S)$ | -- | 0 | 0 | 0 | 0 | 3 | 0 | 0 |  |  |  |  |  |  |

$Z_8(S)=0$

$Z_9(S)$: $9 > (r=8)$ but $S(9) \neq S(1)$ hence $Z_9(S)=0$ and l and remain as they are: l=6 and r=8

| S | A | B | C | D | X | A | B | C | Y | A | B | D | X | Y |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
|  | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 |
| $Z_i(S)$ | -- | 0 | 0 | 0 | 0 | 3 | 0 | 0 | 0 |  |  |  |  |  |

$Z_9(S)=0$

Thursday, May 7, 2009

Exercise solutions: Matching  (Introduction to Comp.Ling)

Thursday, May 7, 2009

$Z_{10}(S)$: 10 > (r=8), S(10)=S(1), match S[10..1] with S[1..2], hence $Z_{10}(S)$=2 and l=10 and r=11

Thursday, May 7, 2009

$Z_{10}(S)$: $10 > (r=8)$, $S(10)=S(1)$, match $S[10..1]$ with $S[1..2]$, hence $Z_{10}(S)=2$ and $l=10$ and $r=11$

| S | A | B | C | D | X | A | B | C | Y | A | B | D | X | Y | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | $Z_{10}(S)=0$ |
| $Z_i(S)$ | -- | 0 | 0 | 0 | 0 | 3 | 0 | 0 | 0 | 2 | | | | | |

Exercise solutions: Matching  (Introduction to Comp.Ling)

Thursday, May 7, 2009

Deutsches Forschungszentrum für Künstliche Intelligenz
German Research Center for Artificial Intelligence

$Z_{10}(S)$: 10 > (r=8), S(10)=S(1), match S[10..1] with S[1..2], hence $Z_{10}(S)$=2 and l=10 and r=11

| S | A | B | C | D | X | A | B | C | Y | A | B | D | X | Y |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
|  | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 |
| $Z_i(S)$ | -- | 0 | 0 | 0 | 0 | 3 | 0 | 0 | 0 | 2 |  |  |  |  |

$Z_{10}(S)$=0

$Z_{11}(S)$: 11 ≤ (r=11) hence S(11)=S(11-10+1)=S(2)='B', $Z_2(S)$=0 whereas |S[11..11]|=1,
hence $Z_{11}(S)$=$Z_2(S)$=0 and l and remain as they are: l=10 and r=11

Exercise solutions: Matching  (Introduction to Comp.Ling)

Thursday, May 7, 2009

$Z_{10}(S)$: 10 > (r=8), S(10)=S(1), match S[10..1] with S[1..2], hence $Z_{10}(S)$=2 and l=10 and r=11

| S | A | B | C | D | X | A | B | C | Y | A | B | D | X | Y |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 |
| $Z_i(S)$ | -- | 0 | 0 | 0 | 0 | 3 | 0 | 0 | 0 | 2 | | | | |

$Z_{10}(S)$=0

$Z_{11}(S)$: 11 ≤ (r=11) hence S(11)=S(11-10+1)=S(2)='B', $Z_2(S)$=0 whereas |S[11..11]|=1, hence $Z_{11}(S)$=$Z_2(S)$=0 and l and remain as they are: l=10 and r=11

| S | A | B | C | D | X | A | B | C | Y | A | B | D | X | Y |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 |
| $Z_i(S)$ | -- | 0 | 0 | 0 | 0 | 3 | 0 | 0 | 0 | 2 | 0 | | | |

$Z_{11}(S)$=0

Exercise solutions: Matching  (Introduction to Comp.Ling)

Thursday, May 7, 2009

$Z_{10}(S)$: 10 > (r=8), S(10)=S(1), match S[10..1] with S[1..2], hence $Z_{10}(S)$=2 and l=10 and r=11

| S | A | B | C | D | X | A | B | C | Y | A | B | D | X | Y | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | $Z_{10}(S)$=0 |
| $Z_i(S)$ | -- | 0 | 0 | 0 | 0 | 3 | 0 | 0 | 0 | 2 | | | | | |

$Z_{11}(S)$: 11 ≤ (r=11) hence S(11)=S(11-10+1)=S(2)='B', $Z_2(S)$=0 whereas |S[11..11]|=1,

hence $Z_{11}(S)$=$Z_2(S)$=0 and l and remain as they are: l=10 and r=11

| S | A | B | C | D | X | A | B | C | Y | A | B | D | X | Y | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | $Z_{11}(S)$=0 |
| $Z_i(S)$ | -- | 0 | 0 | 0 | 0 | 3 | 0 | 0 | 0 | 2 | 0 | | | | |

i=12..14: $Z_i(S)$=0

Thursday, May 7, 2009

$Z_{10}(S)$: 10 > (r=8), S(10)=S(1), match S[10..1] with S[1..2], hence $Z_{10}(S)$=2 and l=10 and r=11

| S | A | B | C | D | X | A | B | C | Y | A | B | D | X | Y |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 |
| $Z_i(S)$ | -- | 0 | 0 | 0 | 0 | 3 | 0 | 0 | 0 | 2 | | | | |

$Z_{10}(S)=0$

$Z_{11}(S)$: 11 ≤ (r=11) hence S(11)=S(11-10+1)=S(2)='B', $Z_2(S)$=0 whereas |S[11..11]|=1, hence $Z_{11}(S)$=$Z_2(S)$=0 and l and remain as they are: l=10 and r=11

| S | A | B | C | D | X | A | B | C | Y | A | B | D | X | Y |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 |
| $Z_i(S)$ | -- | 0 | 0 | 0 | 0 | 3 | 0 | 0 | 0 | 2 | 0 | | | |

$Z_{11}(S)=0$

i=12..14: $Z_i(S)$=0

| S | A | B | C | D | X | A | B | C | Y | A | B | D | X | Y |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 |
| $Z_i(S)$ | -- | 0 | 0 | 0 | 0 | 3 | 0 | 0 | 0 | 2 | 0 | 0 | 0 | 0 |

i=12..14: $Z_i(S)$=0

Exercise solutions: Matching  (Introduction to Comp.Ling)

Thursday, May 7, 2009

# b) Knuth-Morris-Pratt

- "Apply the Knuth-Morris-Pratt algorithm to find occurrences of ABXYABXZ in XABXYABXYABXZABXZABXYABXZA"

- Pre-processing

  - For each position i in the pattern we need to define $sp_i(P)$ to be the length of the longest proper suffix of P[1..i] that matches a prefix of P.

  - *Optimization*: let $sp'_i(P)$ be $sp_i(P)$ with the added condition that characters P(i+1) and P($sp'_i$+1) are unequal

  - Compute $sp'_i(P)$ on the basis of the Z-values for the pattern; compute the failure function on the basis of the $sp'_i(P)$ values

Exercise solutions: Matching  (Introduction to Comp.Ling)

Thursday, May 7, 2009

- **Basic idea**

  - Shift smarter than the naive method does

$$P = \overset{\text{1 2 3 4 5 6 7 8 9}}{A\ B\ C\ X\ A\ B\ C\ D\ E}$$

  - A mismatch with P(8) means we can shift 4 places

  - Deduction on P alone: no need to know T, or how P and T are aligned

- **Complexity of the algorithm**

  - The algorithm is linear, not -possibly- sublinear like Boyer-Moore

  - Extension: the Aho-Corasick algorithm for matching sets of patterns

Exercise solutions: Matching  (Introduction to Comp.Ling)

Thursday, May 7, 2009

- Basic idea

  - Shift smarter than the naive method does

$$\begin{matrix} & 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 \\ P= & A & B & C & X & A & B & C & D & E \end{matrix}$$

  ⬆

  - A mismatch with P(8) means we can shift 4 places

  - Deduction on P alone: no need to know T, or how P and T are aligned

- Complexity of the algorithm

  - The algorithm is linear, not -possibly- sublinear like Boyer-Moore

  - Extension: the Aho-Corasick algorithm for matching sets of patterns

Thursday, May 7, 2009

# The Knuth-Morris-Pratt algorithm

- Basic idea

  - Shift smarter than the naive method does

$$P = \overset{\text{1}\quad\text{2}\quad\text{3}\quad\text{4}\quad\text{5}\quad\text{6}\quad\text{7}\quad\text{8}\quad\text{9}}{A\ B\ C\ X\ A\ B\ C\ D\ E}$$

  - A mismatch with P(8) means we can shift 4 places

  - Deduction on P alone: no need to know T, or how P and T are aligned

- Complexity of the algorithm

  - The algorithm is linear, not -possibly- sublinear like Boyer-Moore

  - Extension: the Aho-Corasick algorithm for matching sets of patterns

Exercise solutions: Matching  (Introduction to Comp.Ling)

Thursday, May 7, 2009

- Basic idea

  - Shift smarter than the naive method does

  $$P = \overset{1}{A}\ \overset{2}{B}\ \overset{3}{C}\ \overset{4}{X}\ \overset{5}{A}\ \overset{6}{B}\ \overset{7}{C}\ \overset{8}{D}\ \overset{9}{E}$$

  ⬆

  - A mismatch with P(8) means we can shift 4 places

  - Deduction on P alone: no need to know T, or how P and T are aligned

- Complexity of the algorithm

  - The algorithm is linear, not -possibly- sublinear like Boyer-Moore

  - Extension: the Aho-Corasick algorithm for matching sets of patterns

Thursday, May 7, 2009

- **Basic idea**

  - Shift smarter than the naive method does

$$P= \quad \overset{1}{A} \ \overset{2}{B} \ \overset{3}{C} \ \overset{4}{X} \ \boxed{\overset{5}{A} \ \overset{6}{B} \ \overset{7}{C}} \ \overset{8}{D} \ \overset{9}{E}$$

  - A mismatch with P(8) means we can shift 4 places

  - Deduction on P alone: no need to know T, or how P and T are aligned

- **Complexity of the algorithm**

  - The algorithm is linear, not -possibly- sublinear like Boyer-Moore

  - Extension: the Aho-Corasick algorithm for matching sets of patterns

Exercise solutions: Matching  (Introduction to Comp.Ling)

Thursday, May 7, 2009

- Basic idea

  - Shift smarter than the naive method does

$$P= \boxed{\begin{matrix} \overset{1}{A} & \overset{2}{B} & \overset{3}{C} \end{matrix}} \overset{4}{X} \boxed{\begin{matrix} \overset{5}{A} & \overset{6}{B} & \overset{7}{C} \end{matrix}} \overset{8}{D}\ \overset{9}{E}$$

  - A mismatch with P(8) means we can shift 4 places

  - Deduction on P alone: no need to know T, or how P and T are aligned

- Complexity of the algorithm

  - The algorithm is linear, not -possibly- sublinear like Boyer-Moore

  - Extension: the Aho-Corasick algorithm for matching sets of patterns

Thursday, May 7, 2009

# The Knuth-Morris-Pratt algorithm

- Basic idea

  - Shift smarter than the naive method does

$$P= \begin{array}{|ccc|cccc|cc} 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 \\ A & B & C & X & A & B & C & D & E \end{array}$$

  - A mismatch with P(8) means we can shift 4 places *(like good suffix rule!)*

  - Deduction on P alone: no need to know T, or how P and T are aligned

- Complexity of the algorithm

  - The algorithm is linear, not -possibly- sublinear like Boyer-Moore

  - Extension: the Aho-Corasick algorithm for matching sets of patterns

Exercise solutions: Matching  (Introduction to Comp.Ling)

Thursday, May 7, 2009

- **Definition**

  For each position i in P, define $sp_i(P)$ to be the length of the longest proper suffix of P[1...i] that matches a prefix of P.

  $$
  \begin{array}{ccccccccccc}
  1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 & 10 & 11
  \end{array}
  $$

  P= A B C A E A B C A B D

- **Optimization**

  For each position i in P, define $sp'_i(P)$ to be the length of the longest proper suffix of P[1...i] that matches a prefix of P, with the added condition that characters P(i+1) and P($sp'_i$+1) are unequal.

16

Thursday, May 7, 2009

- ## Definition

  For each position i in P, define $sp_i(P)$ to be the length of the longest proper suffix of P[1...i] that matches a prefix of P.

  | | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 |
  |---|---|---|---|---|---|---|---|---|---|----|----|
  | P= | A | B | C | A | E | A | B | C | A | B | D |
  | sp | 0 | 0 | 0 | 1 | 0 | 1 | 2 | 3 | 4 | 2 | 0 |

- ## Optimization

  For each position i in P, define $sp'_i(P)$ to be the length of the longest proper suffix of P[1...i] that matches a prefix of P, with the added condition that characters P(i+1) and P($sp'_i$+1) are unequal.

Thursday, May 7, 2009

# The Knuth-Morris-Pratt shift rule

- Alignment of P and T, left-to-right matching

- The shift rule:

  For any alignment of P and T, if the first mismatch (comparing from left to right) occurs in position i+1 of P and position k of T, then shift P to the right (relative to T) so that $P[1..sp_i']$ aligns with $T[k-sp_i'..k-1]$. In other words, shift P exactly $i+1-(sp_i'+1)=i-sp_i'$ places to the right, so that character $sp_i'+1$ of P will align with character k in T. In the case that an occurrence of P has been found (no mismatch), shift P by $n-sp_i'$ places.

- Preprocessing using the Z values

  Position $j > 1$ maps to i if $i=j+Z_j(P)-1$. That is, j maps to i if i is the right end of a Z-box starting at j.

- Z-based Knuth-Morris-Pratt

  ```
  for i := 1 to n do

      sp_i' := 0;

  for j := n downto 2 do

      i := j + Z_j(P) -1;

      sp_i' := Z_j
  ```

Exercise solutions: Matching  (Introduction to Comp.Ling)

Thursday, May 7, 2009

Deutsches Forschungszentrum für Künstliche Intelligenz
German Research Center for Artificial Intelligence

- **Preprocessing using the Z values**

    Position j > 1 maps to i if i=j+Z$_j$(P)-1. That is, j maps to i if i is the right end of a Z-box starting at j.

S

I          $Z_{I_k}$                                              $I_k$        k        $r_k$

α                                                        α

- **Z-based Knuth-Morris-Pratt**

    for i := 1 to n do

        sp$_i$' := 0;

    for j := n downto 2 do

        i := j + Z$_j$(P) -1;

        sp$_i$' := Z$_j$

Thursday, May 7, 2009

# Computing the sp$_i$' values

- **Preprocessing using the Z values**

  Position $j > 1$ maps to i if $i=j+Z_j(P)-1$. That is, j maps to i if i is the right end of a Z-box starting at j.

  S $\qquad$ α $\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad$ α

  $\qquad$ I $\qquad\qquad\qquad$ $Z_{l_k}$ $\qquad\qquad\qquad\qquad\qquad$ $l_k$ =j $\qquad$ k $\qquad\qquad$ $r_k$

- **Z-based Knuth-Morris-Pratt**

  for i := 1 to n do

  $\qquad$ sp$_i$' := 0;

  for j := n downto 2 do

  $\qquad$ i := j + $Z_j(P)$ -1;

  $\qquad$ sp$_i$' := $Z_j$

Exercise solutions: Matching  (Introduction to Comp.Ling)

Thursday, May 7, 2009

- **Preprocessing using the Z values**

  Position $j > 1$ maps to i if $i = j + Z_j(P) - 1$. That is, j maps to i if i is the right end of a Z-box starting at j.

  $$S \quad \boxed{\alpha}_{\, l \qquad Z_{l_k}} \qquad\qquad\qquad \boxed{\alpha}_{\, l_k = j \quad k \quad r_k = i}$$

- **Z-based Knuth-Morris-Pratt**

  ```
  for i := 1 to n do

      sp_i' := 0;

  for j := n downto 2 do

      i := j + Z_j(P) -1;

      sp_i' := Z_j
  ```

Thursday, May 7, 2009

- **Preprocessing using the Z values**

  Position j > 1 maps to i if i=j+Z$_j$(P)-1. That is, j maps to i if i is the right end of a Z-box starting at j.

  S

  | α |          | α |
  
  l        Z$_{l_k}$                          l$_k$ =j    k    r$_k$ =i

- **Z-based Knuth-Morris-Pratt**

  for i := 1 to n do

      sp$_i$' := 0;

  for j := n downto 2 do

      i := j + Z$_j$(P) -1;

      sp$_i$' := Z$_j$

  sp'$_i$(P) is the length of the longest proper suffix of P[1...i] i.e. the length of the Z-box that starts at j (the suffix)

Deutsches Forschungszentrum für Künstliche Intelligenz
German Research Center for Artificial Intelligence

- **Preprocessing using the Z values**

  Position j > 1 maps to i if i=j+Z$_j$(P)-1. That is, j maps to i if i is the right end of a Z-box starting at j.



- **Z-based Knuth-Morris-Pratt**

  for i := 1 to n do

  sp$_i$' := 0;

  for j := n downto 2 do

  i := j + Z$_j$(P) -1;

  sp$_i$' := Z$_j$

sp'$_i$(P) is the length of the longest proper suffix of P[1...i] i.e. the length of the Z-box that starts at j (the suffix)

- ## Preliminaries

  - ### Shifts through pointers: p points into P, c points into T

  - ### For each position i from 1 to n+1, define the failure function F'(i) to be $sp'_{i-1} + 1$ (and define $F(i)=sp_{i-1} +1$); let $sp_0'$ and $sp_0$ be 0.

- ## The algorithm

  preprocess P to find $F'(k)=sp'_{k-1} + 1$ for k from 1 to n+1

      c := 1;

      p := 1;

      while c + (n-p) ≤ m do

          while P(p) = T(c) and p ≤ n

              p := p+1;

              c := c+1;

          if p = n+1 then

              report an occurrence of P starting at position c-n of T

          if p = 1 then c:=c+1

          p := F'(p)

The Z-values are as follows, we only have a Z-box starting at I=5: $Z_5(S)=3$

Exercise solutions: Matching  (Introduction to Comp.Ling)

Thursday, May 7, 2009

The Z-values are as follows, we only have a Z-box starting at I=5: $Z_5(S)=3$

| S | A | B | X | Y | A | B | X | Z |
|---|---|---|---|---|---|---|---|---|
| | I | 2 | 3 | 4 | 5 | 6 | 7 | 8 |
| $Z_i(S)$ | -- | 0 | 0 | 0 | 3 | 0 | 0 | 0 |

The Z-values are as follows, we only have a Z-box starting at I=5: $Z_5(S)=3$

| S | A | B | X | Y | A | B | X | Z |
|---|---|---|---|---|---|---|---|---|
| | I | 2 | 3 | 4 | 5 | 6 | 7 | 8 |
| $Z_i(S)$ | -- | 0 | 0 | 0 | 3 | 0 | 0 | 0 |

## Z-based Knuth-Morris-Pratt

```
for i := 1 to n do
    sp_i' := 0;
for j := n downto 2 do
    i := j + Z_j(P) -1;
    sp_i' := Z_j
```

Exercise solutions: Matching  (Introduction to Comp.Ling)

Thursday, May 7, 2009

The Z-values are as follows, we only have a Z-box starting at I=5: $Z_5(S)=3$

| S | A | B | X | Y | A | B | X | Z |
|---|---|---|---|---|---|---|---|---|
| | I | 2 | 3 | 4 | 5 | 6 | 7 | 8 |
| $Z_i(S)$ | -- | 0 | 0 | 0 | 3 | 0 | 0 | 0 |

## Z-based Knuth-Morris-Pratt

for i := 1 to n do
    sp$_i$' := 0;
for j := n downto 2 do
    i := j + Z$_j$(P) -1;
    sp$_i$' := Z$_j$

j

$Z_j(S)$

i

sp$_i$'

Exercise solutions: Matching  (Introduction to Comp.Ling)

Thursday, May 7, 2009

The Z-values are as follows, we only have a Z-box starting at I=5: $Z_5(S)=3$

| S | A | B | X | Y | A | B | X | Z |
|---|---|---|---|---|---|---|---|---|
| | I | 2 | 3 | 4 | 5 | 6 | 7 | 8 |
| $Z_i(S)$ | -- | 0 | 0 | 0 | 3 | 0 | 0 | 0 |

## Z-based Knuth-Morris-Pratt

```
for i := 1 to n do
    sp_i' := 0;
for j := n downto 2 do
    i := j + Z_j(P) -1;
    sp_i' := Z_j
```

| j | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| $Z_j(S)$ | | | | | | | | |
| i | | | | | | | | |
| $sp_i'$ | | | | | | | | |

The Z-values are as follows, we only have a Z-box starting at I=5: $Z_5(S)=3$

| S | A | B | X | Y | A | B | X | Z |
|---|---|---|---|---|---|---|---|---|
| | I | 2 | 3 | 4 | 5 | 6 | 7 | 8 |
| $Z_i(S)$ | -- | 0 | 0 | 0 | 3 | 0 | 0 | 0 |

**Z-based Knuth-Morris-Pratt**

for i := 1 to n do
    sp$_i$' := 0;

for j := n downto 2 do
    i := j + Z$_j$(P) -1;
    sp$_i$' := Z$_j$

| j | 8 | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| $Z_j(S)$ | 0 | | | | | | | |
| i | 8+0-I | | | | | | | |
| sp$_i$' | 0 | | | | | | | |

Exercise solutions: Matching  (Introduction to Comp.Ling)

Thursday, May 7, 2009

Deutsches Forschungszentrum für Künstliche Intelligenz
German Research Center for Artificial Intelligence

The Z-values are as follows, we only have a Z-box starting at I=5: $Z_5(S)=3$

| S | A | B | X | Y | A | B | X | Z |
|---|---|---|---|---|---|---|---|---|
| | I | 2 | 3 | 4 | 5 | 6 | 7 | 8 |
| $Z_i(S)$ | -- | 0 | 0 | 0 | 3 | 0 | 0 | 0 |

## Z-based Knuth-Morris-Pratt

for i := 1 to n do
    sp$_i$' := 0;
for j := n downto 2 do
    i := j + $Z_j$(P) -1;
    sp$_i$' := $Z_j$

| j | 8 | 7 | | | | | | |
|---|---|---|---|---|---|---|---|---|
| $Z_j(S)$ | 0 | 0 | | | | | | |
| i | 8+0-I | 7+0-I | | | | | | |
| sp$_i$' | 0 | 3 | | | | | | |

Exercise solutions: Matching  (Introduction to Comp.Ling)

Thursday, May 7, 2009

The Z-values are as follows, we only have a Z-box starting at I=5: $Z_5(S)=3$

| S | A | B | X | Y | A | B | X | Z |
|---|---|---|---|---|---|---|---|---|
| | I | 2 | 3 | 4 | 5 | 6 | 7 | 8 |
| $Z_i(S)$ | -- | 0 | 0 | 0 | 3 | 0 | 0 | 0 |

## Z-based Knuth-Morris-Pratt

for i := 1 to n do
    $sp_i$' := 0;
for j := n downto 2 do
    i := j + $Z_j$(P) -1;
    $sp_i$' := $Z_j$

| j | 8 | 7 | 6 | | | | | |
|---|---|---|---|---|---|---|---|---|
| $Z_j(S)$ | 0 | 0 | 0 | | | | | |
| i | 8+0-I | 7+0-I | 6+0-I | | | | | |
| $sp_i$' | 0 | 3 | 0 | | | | | |

Exercise solutions: Matching  (Introduction to Comp.Ling)

Thursday, May 7, 2009

Deutsches Forschungszentrum für Künstliche Intelligenz
German Research Center for Artificial Intelligence

The Z-values are as follows, we only have a Z-box starting at I=5: $Z_5(S)=3$

| S | A | B | X | Y | A | B | X | Z |
|---|---|---|---|---|---|---|---|---|
| | I | 2 | 3 | 4 | 5 | 6 | 7 | 8 |
| $Z_i(S)$ | -- | 0 | 0 | 0 | 3 | 0 | 0 | 0 |

**Z-based Knuth-Morris-Pratt**

for i := 1 to n do
    $sp_i$' := 0;
for j := n downto 2 do
    i := j + $Z_j$(P) -1;
    $sp_i$' := $Z_j$

| j | 8 | 7 | 6 | 5 | | | | |
|---|---|---|---|---|---|---|---|---|
| $Z_j(S)$ | 0 | 0 | 0 | 3 | | | | |
| i | 8+0-I | 7+0-I | 6+0-I | 5+3-I | | | | |
| $sp_i$' | 0 | 3 | 0 | 0 | | | | |

Exercise solutions: Matching (Introduction to Comp.Ling)

Thursday, May 7, 2009

The Z-values are as follows, we only have a Z-box starting at I=5: $Z_5(S)=3$

| S | A | B | X | Y | A | B | X | Z |
|---|---|---|---|---|---|---|---|---|
|   | I | 2 | 3 | 4 | 5 | 6 | 7 | 8 |
| $Z_i(S)$ | -- | 0 | 0 | 0 | 3 | 0 | 0 | 0 |

**Z-based Knuth-Morris-Pratt**

for i := 1 to n do
    sp$_i$' := 0;
for j := n downto 2 do
    i := j + $Z_j$(P) -1;
    sp$_i$' := $Z_j$

| j | 8 | 7 | 6 | 5 | 4 | | | |
|---|---|---|---|---|---|---|---|---|
| $Z_j(S)$ | 0 | 0 | 0 | 3 | 0 | | | |
| i | 8+0-I | 7+0-I | 6+0-I | 5+3-I | 4+0-I | | | |
| sp$_i$' | 0 | 3 | 0 | 0 | 0 | | | |

Exercise solutions: Matching  (Introduction to Comp.Ling)

Thursday, May 7, 2009

# Compute sp'$_i$(P) from Z for P

The Z-values are as follows, we only have a Z-box starting at I=5: $Z_5(S)=3$

| S | A | B | X | Y | A | B | X | Z |
|---|---|---|---|---|---|---|---|---|
| I | 2 | 3 | 4 | 5 | 6 | 7 | 8 |
| $Z_i(S)$ | -- | 0 | 0 | 0 | 3 | 0 | 0 | 0 |

**Z-based Knuth-Morris-Pratt**

for i := 1 to n do
　　sp$_i$' := 0;
for j := n downto 2 do
　　i := j + $Z_j$(P) -1;
　　sp$_i$' := $Z_j$

| j | 8 | 7 | 6 | 5 | 4 | 3 | | |
|---|---|---|---|---|---|---|---|---|
| $Z_j(S)$ | 0 | 0 | 0 | 3 | 0 | 0 | | |
| i | 8+0-I | 7+0-I | 6+0-I | 5+3-I | 4+0-I | 3+0-I | | |
| sp$_i$' | 0 | 3 | 0 | 0 | 0 | 0 | | |

Exercise solutions: Matching (Introduction to Comp.Ling)

Thursday, May 7, 2009

The Z-values are as follows, we only have a Z-box starting at I=5: $Z_5(S)=3$

| S | A | B | X | Y | A | B | X | Z |
|---|---|---|---|---|---|---|---|---|
| | I | 2 | 3 | 4 | 5 | 6 | 7 | 8 |
| $Z_i(S)$ | -- | 0 | 0 | 0 | 3 | 0 | 0 | 0 |

## Z-based Knuth-Morris-Pratt

for i := 1 to n do
$\quad$ sp$_i$' := 0;
for j := n downto 2 do
$\quad$ i := j + Z$_j$(P) -1;
$\quad$ sp$_i$' := Z$_j$

| j | 8 | 7 | 6 | 5 | 4 | 3 | 2 | |
|---|---|---|---|---|---|---|---|---|
| $Z_j(S)$ | 0 | 0 | 0 | 3 | 0 | 0 | 0 | |
| i | 8+0-I | 7+0-I | 6+0-I | 5+3-I | 4+0-I | 3+0-I | 2+0-I | |
| sp$_i$' | 0 | 3 | 0 | 0 | 0 | 0 | 0 | |

Exercise solutions: Matching  (Introduction to Comp.Ling)

Thursday, May 7, 2009

The Z-values are as follows, we only have a Z-box starting at I=5: $Z_5(S)=3$

| S | A | B | X | Y | A | B | X | Z |
|---|---|---|---|---|---|---|---|---|
|  | I | 2 | 3 | 4 | 5 | 6 | 7 | 8 |
| $Z_i(S)$ | -- | 0 | 0 | 0 | 3 | 0 | 0 | 0 |

## Z-based Knuth-Morris-Pratt

for i := 1 to n do
$\quad$ sp$_i$' := 0;

for j := n downto 2 do
$\quad$ i := j + $Z_j$(P) -1;
$\quad$ sp$_i$' := $Z_j$

| j | 8 | 7 | 6 | 5 | 4 | 3 | 2 | I |
|---|---|---|---|---|---|---|---|---|
| $Z_j(S)$ | 0 | 0 | 0 | 3 | 0 | 0 | 0 | 0 |
| i | 8+0-I | 7+0-I | 6+0-I | 5+3-I | 4+0-I | 3+0-I | 2+0-I | -- |
| sp$_i$' | 0 | 3 | 0 | 0 | 0 | 0 | 0 | 0 |

Exercise solutions: Matching (Introduction to Comp.Ling)

Thursday, May 7, 2009

The Z-values are as follows, we only have a Z-box starting at I=5: $Z_5(S)=3$

| S | A | B | X | Y | A | B | X | Z |
|---|---|---|---|---|---|---|---|---|
| | I | 2 | 3 | 4 | 5 | 6 | 7 | 8 |
| $Z_i(S)$ | -- | 0 | 0 | 0 | 3 | 0 | 0 | 0 |

## Z-based Knuth-Morris-Pratt

for i := 1 to n do
    sp$_i$' := 0;
for j := n downto 2 do
    i := j + Z$_j$(P) -1;
    sp$_i$' := Z$_j$

| j | 8 | 7 | 6 | 5 | 4 | 3 | 2 | I |
|---|---|---|---|---|---|---|---|---|
| $Z_j(S)$ | 0 | 0 | 0 | 3 | 0 | 0 | 0 | 0 |
| i | 8+0-I | 7+0-I | 6+0-I | 5+3-I | 4+0-I | 3+0-I | 2+0-I | -- |
| sp$_i$' | 0 | 3 | 0 | 0 | 0 | 0 | 0 | 0 |

sp$_7$'=Z$_8$(P)=0; sp$_6$'=Z$_7$(P)=0;
sp$_5$'=Z$_6$(P)=0;**sp$_7$'=Z$_5$(P)=3**;
sp$_3$'=Z$_4$(P)=0;...;sp$_1$'=Z$_2$(P)=0

Thursday, May 7, 2009

- Failure function F'(k)=sp'$_{k-1}$+1 for k from 1 to n+1

| k | -- | I | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|---|----|----|----|----|----|----|----|----|----|----|
| sp$_i$' | 0 | 0 | 0 | 0 | 0 | 2 | 2 | 0 | 0 | -- |
| F'(k) | -- | I | I | I | I | I | I | I | 4 | I |

Exercise solutions: Matching  (Introduction to Comp.Ling)

Thursday, May 7, 2009

preprocess P to find $F'(k)=sp'_{k-1} + 1$ for k from 1 to n+1

    c := 1;

    p := 1;

    while c + (n-p) ≤ m do

        while P(p) = T(c) and p ≤ n

            p := p+1;

            c := c+1;

        if p = n+1 then

            report an occurrence of P starting at position c-n of T

        if p = 1 then c:=c+1

        p := F'(p)

Thursday, May 7, 2009

preprocess P to find F'(k)=sp'$_{k-1}$ + 1 for k from 1 to n+1

    c := 1;

    p := 1;

    while c + (n-p) ≤ m do

```
while P(p) = T(c) and p ≤ n          MatchChar (MC)

    p := p+1;

    c := c+1;
```

      if p = n+1 then

        report an occurrence of P starting at position c-n of T

      if p = 1 then c:=c+1

      p := F'(p)

Exercise solutions: Matching  (Introduction to Comp.Ling)

Thursday, May 7, 2009

preprocess P to find $F'(k)=sp'_{k-1} + 1$ for k from 1 to n+1

    c := 1;

    p := 1;

    while c + (n-p) ≤ m do

        while P(p) = T(c) and p ≤ n           ***MatchChar (MC)***

           p := p+1;

           c := c+1;

        if p = n+1 then           ***MatchPattern (MP)***

           report an occurrence of P starting at position c-n of T

        if p = 1 then c:=c+1

        p := F'(p)

Exercise solutions: Matching  (Introduction to Comp.Ling)

Thursday, May 7, 2009

preprocess P to find F'(k)=sp'$_{k-1}$ + 1 for k from 1 to n+1

c := 1;

p := 1;

while c + (n-p) ≤ m do

| | |
|---|---|
| while P(p) = T(c) and p ≤ n | *MatchChar (MC)* |
|     p := p+1; | |
|     c := c+1; | |
| if p = n+1 then | *MatchPattern (MP)* |
|     report an occurrence of P starting at position c-n of T | |
| if p = 1 then c:=c+1 | *MismatchStart (F0)* |

p := F'(p)

Exercise solutions: Matching  (Introduction to Comp.Ling)

Thursday, May 7, 2009

preprocess P to find F'(k)=sp'$_{k-1}$ + 1 for k from 1 to n+1

c := 1;

p := 1;

while c + (n-p) ≤ m do

| | |
|---|---|
| while P(p) = T(c) and p ≤ n<br><br>p := p+1;<br><br>c := c+1; | **MatchChar (MC)** |
| if p = n+1 then<br><br>report an occurrence of P starting at position c-n of T | **MatchPattern (MP)** |
| if p = 1 then c:=c+1 | **MismatchStart (F0)** |
| p := F'(p) | **MismatchInternal (Fi)** |

Deutsches Forschungszentrum für Künstliche Intelligenz
German Research Center for Artificial Intelligence

c=1, p=1: F0

c=2, p=1: MC

c=3, p=2: MC

c=4, p=3: MC

c=5, p=4: MC

c=6, p=5: MC

c=7, p=6: MC

c=8, p=7: MC

c=9, p=8: Fi: F'(8)=4 $\Rightarrow$ p=4

c=9, p=4: MC

c=10, p=5: MC

c=11, p=6: MC

c=12, p=7: MC

c=13, p=8: MC

c=14, p=9: MP

Exercise solutions: Matching  (Introduction to Comp.Ling)

Thursday, May 7, 2009

c=1, p=1: F0

c=2, p=1: MC

c=3, p=2: MC

c=4, p=3: MC

c=5, p=4: MC

c=6, p=5: MC

c=7, p=6: MC

c=8, p=7: MC

c=9, p=8: Fi: F'(8)=4 $\Rightarrow$ p=4

c=9, p=4: MC

c=10, p=5: MC

c=11, p=6: MC

c=12, p=7: MC

c=13, p=8: MC

c=14, p=9: MP

Exercise solutions: Matching  (Introduction to Comp.Ling)

Thursday, May 7, 2009

| X | A | B | X | Y | A | B | X | Y | A | B | X | Z | A | B | X | Z | A | B | X | Y | A | B | X | Z |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 | 21 | 22 | 23 | 24 | 25 |

c=1, p=1: F0

c=2, p=1: MC

c=3, p=2: MC

c=4, p=3: MC

c=5, p=4: MC

c=6, p=5: MC

c=7, p=6: MC

c=8, p=7: MC

c=9, p=8: Fi: F'(8)=4 $\Rightarrow$ p=4

c=9, p=4: MC

c=10, p=5: MC

c=11, p=6: MC

c=12, p=7: MC

c=13, p=8: MC

c=14, p=9: MP

The matching grid rows:

| 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 | 21 | 22 | 23 | 24 | 25 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **A** | B | X | Y | A | B | X | Z | | | | | | | | | | | | | | | | | |
| | **A** | B | X | Y | A | B | X | Z | | | | | | | | | | | | | | | | |
| | A | **B** | X | Y | A | B | X | Z | | | | | | | | | | | | | | | | |
| | A | B | **X** | Y | A | B | X | Z | | | | | | | | | | | | | | | | |
| | A | B | X | **Y** | A | B | X | Z | | | | | | | | | | | | | | | | |
| | A | B | X | Y | **A** | B | X | Z | | | | | | | | | | | | | | | | |
| | A | B | X | Y | A | **B** | X | Z | | | | | | | | | | | | | | | | |
| | A | B | X | Y | A | B | **X** | Z | | | | | | | | | | | | | | | | |
| | A | B | X | Y | A | B | X | **Z** | | | | | | | | | | | | | | | | |
| | | | | | A | B | X | **Y** | A | B | X | Z | | | | | | | | | | | | |
| | | | | | A | B | X | Y | **A** | B | X | Z | | | | | | | | | | | | |
| | | | | | A | B | X | Y | A | **B** | X | Z | | | | | | | | | | | | |
| | | | | | A | B | X | Y | A | B | **X** | Z | | | | | | | | | | | | |
| | | | | | A | B | X | Y | A | B | X | **Z** | | | | | | | | | | | | |
| | | | | | **A** | **B** | **X** | **Y** | **A** | **B** | **X** | **Z** | | | | | | | | | | | | |
| | | | | | | | | | | | | | | | | | | | | | | | | |
| | | | | | | | | | | | | | | | | | | | | | | | | |
| | | | | | | | | | | | | | | | | | **A** | **B** | **X** | **Y** | **A** | **B** | **X** | **Z** |

Exercise solutions: Matching  (Introduction to Comp.Ling)

Thursday, May 7, 2009

| X | A | B | X | Y | A | B | X | Y | A | B | X | Z | A | B | X | Z | A | B | X | Y | A | B | X | Z |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 | 21 | 22 | 23 | 24 | 25 |

c=1, p=1: F0

c=2, p=1: MC

c=3, p=2: MC

c=4, p=3: MC

c=5, p=4: MC

c=6, p=5: MC

c=7, p=6: MC

c=8, p=7: MC

c=9, p=8: Fi: F'(8)=4 $\Rightarrow$ p=4

c=9, p=4: MC

c=10, p=5: MC

c=11, p=6: MC

c=12, p=7: MC

c=13, p=8: MC

c=14, p=9: MP

Exercise solutions: Matching  (Introduction to Comp.Ling)

Thursday, May 7, 2009

| X | A | B | X | Y | A | B | X | Y | A | B | X | Z | A | B | X | Z | A | B | X | Y | A | B | X | Z |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 | 21 | 22 | 23 | 24 | 25 |
| **A** | B | X | Y | A | B | X | Z | | | | | | | | | | | | | | | | | |
| | | | | | | | | | | | | | | | | | | | | | | | | |
| | | | | | | | | | | | | | | | | | | | | | | | | |
| | | | | | | | | | | | | | | | | | | | | | | | | |
| | | | | | | | | | | | | | | | | | | | | | | | | |
| | | | | | | | | | | | | | | | | | | | | | | | | |
| | | | | | | | | | | | | | | | | | | | | | | | | |
| | | | | | | | | | | | | | | | | | | | | | | | | |
| | | | | | | | | | | | | | | | | | | | | | | | | |
| | | | | | | | | | | | | | | | | | | | | | | | | |
| | | | | | | | | | | | | | | | | | | | | | | | | |
| | | | | | | | | | | | | | | | | | | | | | | | | |
| | | | | | | | | | | | | | | | | | | | | | | | | |
| | | | | | | | | | | | | | | | | | | | | | | | | |
| | | | | | | | | | | | | | | | | | | | | | | | | |
| | | | | | | | | | | | | | | | | | | | | | | | | |
| | | | | | | | | | | | | | | | | | | | | | | | | |
| | | | | | | | | | | | | | | | | | | | | | | | | |

c=1, p=1: F0

c=2, p=1: MC

c=3, p=2: MC

c=4, p=3: MC

c=5, p=4: MC

c=6, p=5: MC

c=7, p=6: MC

c=8, p=7: MC

c=9, p=8: Fi: F'(8)=4 $\Rightarrow$ p=4

c=9, p=4: MC

c=10, p=5: MC

c=11, p=6: MC

c=12, p=7: MC

c=13, p=8: MC

c=14, p=9: MP

Exercise solutions: Matching  (Introduction to Comp.Ling)

Thursday, May 7, 2009

Deutsches Forschungszentrum für Künstliche Intelligenz
German Research Center for Artificial Intelligence

| X | A | B | X | Y | A | B | X | Y | A | B | X | Z | A | B | X | Z | A | B | X | Y | A | B | X | Z |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 | 21 | 22 | 23 | 24 | 25 |
| **A** | B | X | Y | A | B | X | Z | | | | | | | | | | | | | | | | | |
| | **A** | B | X | Y | A | B | X | Z | | | | | | | | | | | | | | | | |

c=1, p=1: F0

c=2, p=1: MC

c=3, p=2: MC

c=4, p=3: MC

c=5, p=4: MC

c=6, p=5: MC

c=7, p=6: MC

c=8, p=7: MC

c=9, p=8: Fi: F'(8)=4 $\Rightarrow$ p=4

c=9, p=4: MC

c=10, p=5: MC

c=11, p=6: MC

c=12, p=7: MC

c=13, p=8: MC

c=14, p=9: MP

Exercise solutions: Matching  (Introduction to Comp.Ling)

Thursday, May 7, 2009

| X | A | B | X | Y | A | B | X | Y | A | B | X | Z | A | B | X | Z | A | B | X | Y | A | B | X | Z |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 | 21 | 22 | 23 | 24 | 25 |
| **A** | B | X | Y | A | B | X | Z | | | | | | | | | | | | | | | | | |
| | **A** | B | X | Y | A | B | X | Z | | | | | | | | | | | | | | | | |
| | A | **B** | X | Y | A | B | X | Z | | | | | | | | | | | | | | | | |
| | | | | | | | | | | | | | | | | | | | | | | | | |
| | | | | | | | | | | | | | | | | | | | | | | | | |
| | | | | | | | | | | | | | | | | | | | | | | | | |
| | | | | | | | | | | | | | | | | | | | | | | | | |
| | | | | | | | | | | | | | | | | | | | | | | | | |
| | | | | | | | | | | | | | | | | | | | | | | | | |
| | | | | | | | | | | | | | | | | | | | | | | | | |
| | | | | | | | | | | | | | | | | | | | | | | | | |
| | | | | | | | | | | | | | | | | | | | | | | | | |
| | | | | | | | | | | | | | | | | | | | | | | | | |
| | | | | | | | | | | | | | | | | | | | | | | | | |
| | | | | | | | | | | | | | | | | | | | | | | | | |
| | | | | | | | | | | | | | | | | | | | | | | | | |
| | | | | | | | | | | | | | | | | | | | | | | | | |
| | | | | | | | | | | | | | | | | | | | | | | | | |

c=1, p=1: F0

c=2, p=1: MC

c=3, p=2: MC

c=4, p=3: MC

c=5, p=4: MC

c=6, p=5: MC

c=7, p=6: MC

c=8, p=7: MC

c=9, p=8: Fi: F'(8)=4 $\Rightarrow$ p=4

c=9, p=4: MC

c=10, p=5: MC

c=11, p=6: MC

c=12, p=7: MC

c=13, p=8: MC

c=14, p=9: MP

Exercise solutions: Matching  (Introduction to Comp.Ling)

Thursday, May 7, 2009

# K-M-P matching P against T

c=1, p=1: F0

c=2, p=1: MC

c=3, p=2: MC

c=4, p=3: MC

c=5, p=4: MC

c=6, p=5: MC

c=7, p=6: MC

c=8, p=7: MC

c=9, p=8: Fi: F'(8)=4 $\Rightarrow$ p=4

c=9, p=4: MC

c=10, p=5: MC

c=11, p=6: MC

c=12, p=7: MC

c=13, p=8: MC

c=14, p=9: MP

| X | A | B | X | Y | A | B | X | Y | A | B | X | Z | A | B | X | Z | A | B | X | Y | A | B | X | Z |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 | 21 | 22 | 23 | 24 | 25 |
| **A** | B | X | Y | A | B | X | Z | | | | | | | | | | | | | | | | | |
| | **A** | B | X | Y | A | B | X | Z | | | | | | | | | | | | | | | | |
| | A | **B** | X | Y | A | B | X | Z | | | | | | | | | | | | | | | | |
| | A | B | **X** | Y | A | B | X | Z | | | | | | | | | | | | | | | | |

Exercise solutions: Matching (Introduction to Comp.Ling)

Thursday, May 7, 2009

# K-M-P matching P against T

c=1, p=1: F0

c=2, p=1: MC

c=3, p=2: MC

c=4, p=3: MC

c=5, p=4: MC

c=6, p=5: MC

c=7, p=6: MC

c=8, p=7: MC

c=9, p=8: Fi: F'(8)=4 $\Rightarrow$ p=4

c=9, p=4: MC

c=10, p=5: MC

c=11, p=6: MC

c=12, p=7: MC

c=13, p=8: MC

c=14, p=9: MP

| X | A | B | X | Y | A | B | X | Y | A | B | X | Z | A | B | X | Z | A | B | X | Y | A | B | X | Z |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 | 21 | 22 | 23 | 24 | 25 |
| **A** | B | X | Y | A | B | X | Z | | | | | | | | | | | | | | | | | |
| | **A** | B | X | Y | A | B | X | Z | | | | | | | | | | | | | | | | |
| | A | **B** | X | Y | A | B | X | Z | | | | | | | | | | | | | | | | |
| | A | B | **X** | Y | A | B | X | Z | | | | | | | | | | | | | | | | |
| | A | B | X | **Y** | A | B | X | Z | | | | | | | | | | | | | | | | |

Exercise solutions: Matching  (Introduction to Comp.Ling)

Thursday, May 7, 2009

| X | A | B | X | Y | A | B | X | Y | A | B | X | Z | A | B | X | Z | A | B | X | Y | A | B | X | Z |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 | 21 | 22 | 23 | 24 | 25 |

c=1, p=1: F0

c=2, p=1: MC

c=3, p=2: MC

c=4, p=3: MC

c=5, p=4: MC

c=6, p=5: MC

c=7, p=6: MC

c=8, p=7: MC

c=9, p=8: Fi: F'(8)=4 $\Rightarrow$ p=4

c=9, p=4: MC

c=10, p=5: MC

c=11, p=6: MC

c=12, p=7: MC

c=13, p=8: MC

c=14, p=9: MP

| | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| **A** | B | X | Y | A | B | X | Z | |
| | **A** | B | X | Y | A | B | X | Z |
| | A | **B** | X | Y | A | B | X | Z |
| | A | B | **X** | Y | A | B | X | Z |
| | A | B | X | **Y** | A | B | X | Z |
| | A | B | X | Y | **A** | B | X | Z |

Exercise solutions: Matching  (Introduction to Comp.Ling)

Thursday, May 7, 2009

# K-M-P matching P against T

| X | A | B | X | Y | A | B | X | Y | A | B | X | Z | A | B | X | Z | A | B | X | Y | A | B | X | Z |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 | 21 | 22 | 23 | 24 | 25 |

c=1, p=1: F0

c=2, p=1: MC

c=3, p=2: MC

c=4, p=3: MC

c=5, p=4: MC

c=6, p=5: MC

c=7, p=6: MC

c=8, p=7: MC

c=9, p=8: Fi: F'(8)=4 $\Rightarrow$ p=4

c=9, p=4: MC

c=10, p=5: MC

c=11, p=6: MC

c=12, p=7: MC

c=13, p=8: MC

c=14, p=9: MP

| A | B | X | Y | A | B | X | Z |
|---|---|---|---|---|---|---|---|

Exercise solutions: Matching (Introduction to Comp.Ling)

Thursday, May 7, 2009

# K-M-P matching P against T

c=1, p=1: F0

c=2, p=1: MC

c=3, p=2: MC

c=4, p=3: MC

c=5, p=4: MC

c=6, p=5: MC

c=7, p=6: MC

c=8, p=7: MC

c=9, p=8: Fi: F'(8)=4 $\Rightarrow$ p=4

c=9, p=4: MC

c=10, p=5: MC

c=11, p=6: MC

c=12, p=7: MC

c=13, p=8: MC

c=14, p=9: MP

| X | A | B | X | Y | A | B | X | Y | A | B | X | Z | A | B | X | Z | A | B | X | Y | A | B | X | Z |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 | 21 | 22 | 23 | 24 | 25 |
| **A** | B | X | Y | A | B | X | Z | | | | | | | | | | | | | | | | | |
| | **A** | B | X | Y | A | B | X | Z | | | | | | | | | | | | | | | | |
| | A | **B** | X | Y | A | B | X | Z | | | | | | | | | | | | | | | | |
| | A | B | **X** | Y | A | B | X | Z | | | | | | | | | | | | | | | | |
| | A | B | X | **Y** | A | B | X | Z | | | | | | | | | | | | | | | | |
| | A | B | X | Y | **A** | B | X | Z | | | | | | | | | | | | | | | | |
| | A | B | X | Y | A | **B** | X | Z | | | | | | | | | | | | | | | | |
| | A | B | X | Y | A | B | **X** | Z | | | | | | | | | | | | | | | | |

Exercise solutions: Matching  (Introduction to Comp.Ling)

Thursday, May 7, 2009

| X | A | B | X | Y | A | B | X | Y | A | B | X | Z | A | B | X | Z | A | B | X | Y | A | B | X | Z |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 | 21 | 22 | 23 | 24 | 25 |
| **A** | B | X | Y | A | B | X | Z | | | | | | | | | | | | | | | | | |
| | **A** | B | X | Y | A | B | X | Z | | | | | | | | | | | | | | | | |
| | A | **B** | X | Y | A | B | X | Z | | | | | | | | | | | | | | | | |
| | A | B | **X** | Y | A | B | X | Z | | | | | | | | | | | | | | | | |
| | A | B | X | **Y** | A | B | X | Z | | | | | | | | | | | | | | | | |
| | A | B | X | Y | **A** | B | X | Z | | | | | | | | | | | | | | | | |
| | A | B | X | Y | A | **B** | X | Z | | | | | | | | | | | | | | | | |
| | A | B | X | Y | A | B | **X** | Z | | | | | | | | | | | | | | | | |
| | A | B | X | Y | A | B | X | **Z** | | | | | | | | | | | | | | | | |

c=1, p=1: F0

c=2, p=1: MC

c=3, p=2: MC

c=4, p=3: MC

c=5, p=4: MC

c=6, p=5: MC

c=7, p=6: MC

c=8, p=7: MC

c=9, p=8: Fi: F'(8)=4 $\Rightarrow$ p=4

c=9, p=4: MC

c=10, p=5: MC

c=11, p=6: MC

c=12, p=7: MC

c=13, p=8: MC

c=14, p=9: MP

Exercise solutions: Matching  (Introduction to Comp.Ling)

Thursday, May 7, 2009

# K-M-P matching P against T

| X | A | B | X | Y | A | B | X | Y | A | B | X | Z | A | B | X | Z | A | B | X | Y | A | B | X | Z |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 | 21 | 22 | 23 | 24 | 25 |
| **A** | B | X | Y | A | B | X | Z | | | | | | | | | | | | | | | | | |
| | **A** | B | X | Y | A | B | X | Z | | | | | | | | | | | | | | | | |
| | A | **B** | X | Y | A | B | X | Z | | | | | | | | | | | | | | | | |
| | A | B | **X** | Y | A | B | X | Z | | | | | | | | | | | | | | | | |
| | A | B | X | **Y** | A | B | X | Z | | | | | | | | | | | | | | | | |
| | A | B | X | Y | **A** | B | X | Z | | | | | | | | | | | | | | | | |
| | A | B | X | Y | A | **B** | X | Z | | | | | | | | | | | | | | | | |
| | A | B | X | Y | A | B | **X** | Z | | | | | | | | | | | | | | | | |
| | A | B | X | Y | A | B | X | **Z** | | | | | | | | | | | | | | | | |
| | | | | | A | B | X | **Y** | A | B | X | Z | | | | | | | | | | | | |

c=1, p=1: F0

c=2, p=1: MC

c=3, p=2: MC

c=4, p=3: MC

c=5, p=4: MC

c=6, p=5: MC

c=7, p=6: MC

c=8, p=7: MC

c=9, p=8: Fi: F'(8)=4 ⇒ p=4

c=9, p=4: MC

c=10, p=5: MC

c=11, p=6: MC

c=12, p=7: MC

c=13, p=8: MC

c=14, p=9: MP

footer_navigation© 2005-2007 Geert-Jan M. Kruijff       23       Exercise solutions: Matching  (Introduction to Comp.Ling)

Thursday, May 7, 2009

| X | A | B | X | Y | A | B | X | Y | A | B | X | Z | A | B | X | Z | A | B | X | Y | A | B | X | Z |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 | 21 | 22 | 23 | 24 | 25 |

| c | row | | | | | | | | | | | | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| c=1, p=1: F0 | **A** | B | X | Y | A | B | X | Z | | | | | | | | | | | | | | | | | |
| c=2, p=1: MC | | **A** | B | X | Y | A | B | X | Z | | | | | | | | | | | | | | | | |
| c=3, p=2: MC | | A | **B** | X | Y | A | B | X | Z | | | | | | | | | | | | | | | | |
| c=4, p=3: MC | | A | B | **X** | Y | A | B | X | Z | | | | | | | | | | | | | | | | |
| c=5, p=4: MC | | A | B | X | **Y** | A | B | X | Z | | | | | | | | | | | | | | | | |
| c=6, p=5: MC | | A | B | X | Y | **A** | B | X | Z | | | | | | | | | | | | | | | | |
| c=7, p=6: MC | | A | B | X | Y | A | **B** | X | Z | | | | | | | | | | | | | | | | |
| c=8, p=7: MC | | A | B | X | Y | A | B | **X** | Z | | | | | | | | | | | | | | | | |
| c=9, p=8: Fi: F'(8)=4 ⇒ p=4 | | A | B | X | Y | A | B | X | **Z** | | | | | | | | | | | | | | | | |
| c=9, p=4: MC | | | | | | A | B | X | **Y** | A | B | X | Z | | | | | | | | | | | | |
| c=10, p=5: MC | | | | | | A | B | X | Y | **A** | B | X | Z | | | | | | | | | | | | |
| c=11, p=6: MC | | | | | | | | | | | | | | | | | | | | | | | | | |
| c=12, p=7: MC | | | | | | | | | | | | | | | | | | | | | | | | | |
| c=13, p=8: MC | | | | | | | | | | | | | | | | | | | | | | | | | |
| c=14, p=9: MP | | | | | | | | | | | | | | | | | | | | | | | | | |

Exercise solutions: Matching  (Introduction to Comp.Ling)

Thursday, May 7, 2009

# K-M-P matching P against T

| | X | A | B | X | Y | A | B | X | Y | A | B | X | Z | A | B | X | Z | A | B | X | Y | A | B | X | Z |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 | 21 | 22 | 23 | 24 | 25 |
| c=1, p=1: F0 | **A** | B | X | Y | A | B | X | Z | | | | | | | | | | | | | | | | | |
| c=2, p=1: MC | | **A** | B | X | Y | A | B | X | Z | | | | | | | | | | | | | | | | |
| c=3, p=2: MC | | A | **B** | X | Y | A | B | X | Z | | | | | | | | | | | | | | | | |
| c=4, p=3: MC | | A | B | **X** | Y | A | B | X | Z | | | | | | | | | | | | | | | | |
| c=5, p=4: MC | | A | B | X | **Y** | A | B | X | Z | | | | | | | | | | | | | | | | |
| c=6, p=5: MC | | A | B | X | Y | **A** | B | X | Z | | | | | | | | | | | | | | | | |
| c=7, p=6: MC | | A | B | X | Y | A | **B** | X | Z | | | | | | | | | | | | | | | | |
| c=8, p=7: MC | | A | B | X | Y | A | B | **X** | Z | | | | | | | | | | | | | | | | |
| c=9, p=8: Fi: F'(8)=4 ⇒ p=4 | | A | B | X | Y | A | B | X | **Z** | | | | | | | | | | | | | | | | |
| c=9, p=4: MC | | | | | | A | B | X | **Y** | A | B | X | Z | | | | | | | | | | | | |
| c=10, p=5: MC | | | | | | A | B | X | Y | **A** | B | X | Z | | | | | | | | | | | | |
| c=11, p=6: MC | | | | | | A | B | X | Y | A | **B** | X | Z | | | | | | | | | | | | |
| c=12, p=7: MC | | | | | | | | | | | | | | | | | | | | | | | | | |
| c=13, p=8: MC | | | | | | | | | | | | | | | | | | | | | | | | | |
| c=14, p=9: MP | | | | | | | | | | | | | | | | | | | | | | | | | |

Exercise solutions: Matching  (Introduction to Comp.Ling)

Thursday, May 7, 2009

| X | A | B | X | Y | A | B | X | Y | A | B | X | Z | A | B | X | Z | A | B | X | Y | A | B | X | Z |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 | 21 | 22 | 23 | 24 | 25 |

c=1, p=1: F0

c=2, p=1: MC

c=3, p=2: MC

c=4, p=3: MC

c=5, p=4: MC

c=6, p=5: MC

c=7, p=6: MC

c=8, p=7: MC

c=9, p=8: Fi: F'(8)=4 $\Rightarrow$ p=4

c=9, p=4: MC

c=10, p=5: MC

c=11, p=6: MC

c=12, p=7: MC

c=13, p=8: MC

c=14, p=9: MP

| | | | | | | | | | | | | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **A** | B | X | Y | A | B | X | Z | | | | | | | | | | | | | | | | | |
| | **A** | B | X | Y | A | B | X | Z | | | | | | | | | | | | | | | | |
| | A | **B** | X | Y | A | B | X | Z | | | | | | | | | | | | | | | | |
| | A | B | **X** | Y | A | B | X | Z | | | | | | | | | | | | | | | | |
| | A | B | X | **Y** | A | B | X | Z | | | | | | | | | | | | | | | | |
| | A | B | X | Y | **A** | B | X | Z | | | | | | | | | | | | | | | | |
| | A | B | X | Y | A | **B** | X | Z | | | | | | | | | | | | | | | | |
| | A | B | X | Y | A | B | **X** | Z | | | | | | | | | | | | | | | | |
| | A | B | X | Y | A | B | X | **Z** | | | | | | | | | | | | | | | | |
| | | | | | A | B | X | **Y** | A | B | X | Z | | | | | | | | | | | | |
| | | | | | A | B | X | Y | **A** | B | X | Z | | | | | | | | | | | | |
| | | | | | A | B | X | Y | A | **B** | X | Z | | | | | | | | | | | | |
| | | | | | A | B | X | Y | A | B | **X** | Z | | | | | | | | | | | | |

Thursday, May 7, 2009

# K-M-P matching P against T

| | X | A | B | X | Y | A | B | X | Y | A | B | X | Z | A | B | X | Z | A | B | X | Y | A | B | X | Z |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 | 21 | 22 | 23 | 24 | 25 |
| c=1, p=1: F0 | **A** | B | X | Y | A | B | X | Z | | | | | | | | | | | | | | | | | |
| c=2, p=1: MC | | **A** | B | X | Y | A | B | X | Z | | | | | | | | | | | | | | | | |
| c=3, p=2: MC | | A | **B** | X | Y | A | B | X | Z | | | | | | | | | | | | | | | | |
| c=4, p=3: MC | | A | B | **X** | Y | A | B | X | Z | | | | | | | | | | | | | | | | |
| c=5, p=4: MC | | A | B | X | **Y** | A | B | X | Z | | | | | | | | | | | | | | | | |
| c=6, p=5: MC | | A | B | X | Y | **A** | B | X | Z | | | | | | | | | | | | | | | | |
| c=7, p=6: MC | | A | B | X | Y | A | **B** | X | Z | | | | | | | | | | | | | | | | |
| c=8, p=7: MC | | A | B | X | Y | A | B | **X** | Z | | | | | | | | | | | | | | | | |
| c=9, p=8: Fi: F'(8)=4 ⇒ p=4 | | A | B | X | Y | A | B | X | **Z** | | | | | | | | | | | | | | | | |
| c=9, p=4: MC | | | | | | A | B | X | **Y** | A | B | X | Z | | | | | | | | | | | | |
| c=10, p=5: MC | | | | | | A | B | X | Y | **A** | B | X | Z | | | | | | | | | | | | |
| c=11, p=6: MC | | | | | | A | B | X | Y | A | **B** | X | Z | | | | | | | | | | | | |
| c=12, p=7: MC | | | | | | A | B | X | Y | A | B | **X** | Z | | | | | | | | | | | | |
| c=13, p=8: MC | | | | | | A | B | X | Y | A | B | X | **Z** | | | | | | | | | | | | |
| c=14, p=9: MP | | | | | | | | | | | | | | | | | | | | | | | | | |

Exercise solutions: Matching  (Introduction to Comp.Ling)

Thursday, May 7, 2009

| 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 | 21 | 22 | 23 | 24 | 25 |
|---|---|---|---|---|---|---|---|---|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|
| X | A | B | X | Y | A | B | X | Y | A | B | X | Z | A | B | X | Z | A | B | X | Y | A | B | X | Z |

**c=1, p=1: F0**

**c=2, p=1: MC**

**c=3, p=2: MC**

**c=4, p=3: MC**

**c=5, p=4: MC**

**c=6, p=5: MC**

**c=7, p=6: MC**

**c=8, p=7: MC**

**c=9, p=8: Fi: F'(8)=4 ⇒ p=4**

**c=9, p=4: MC**

**c=10, p=5: MC**

**c=11, p=6: MC**

**c=12, p=7: MC**

**c=13, p=8: MC**

**c=14, p=9: MP**

Pattern placements against T (columns 1–25):

| Step | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 |
|------|---|---|---|---|---|---|---|---|---|----|----|----|----|
| c=1  | **A** | B | X | Y | A | B | X | Z |   |    |    |    |    |
| c=2  |   | **A** | B | X | Y | A | B | X | Z |    |    |    |    |
| c=3  |   | A | **B** | X | Y | A | B | X | Z |    |    |    |    |
| c=4  |   | A | B | **X** | Y | A | B | X | Z |    |    |    |    |
| c=5  |   | A | B | X | **Y** | A | B | X | Z |    |    |    |    |
| c=6  |   | A | B | X | Y | **A** | B | X | Z |    |    |    |    |
| c=7  |   | A | B | X | Y | A | **B** | X | Z |    |    |    |    |
| c=8  |   | A | B | X | Y | A | B | **X** | Z |    |    |    |    |
| c=9 (p=8) |   | A | B | X | Y | A | B | X | **Z** |    |    |    |    |
| c=9 (p=4) |   |   |   |   |   | A | B | X | **Y** | A | B | X | Z |
| c=10 |   |   |   |   |   | A | B | X | Y | **A** | B | X | Z |
| c=11 |   |   |   |   |   | A | B | X | Y | A | **B** | X | Z |
| c=12 |   |   |   |   |   | A | B | X | Y | A | B | **X** | Z |
| c=13 |   |   |   |   |   | A | B | X | Y | A | B | X | **Z** |
| c=14 |   |   |   |   |   | **A** | **B** | **X** | **Y** | **A** | **B** | **X** | **Z** |

Exercise solutions: Matching  (Introduction to Comp.Ling)

Thursday, May 7, 2009

| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 | 21 | 22 | 23 | 24 | 25 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | X | A | B | X | Y | A | B | X | Y | A | B | X | Z | A | B | X | Z | A | B | X | Y | A | B | X | Z |
| c=1, p=1: F0 | **A** | B | X | Y | A | B | X | Z | | | | | | | | | | | | | | | | | |
| c=2, p=1: MC | | **A** | B | X | Y | A | B | X | Z | | | | | | | | | | | | | | | | |
| c=3, p=2: MC | | A | **B** | X | Y | A | B | X | Z | | | | | | | | | | | | | | | | |
| c=4, p=3: MC | | A | B | **X** | Y | A | B | X | Z | | | | | | | | | | | | | | | | |
| c=5, p=4: MC | | A | B | X | **Y** | A | B | X | Z | | | | | | | | | | | | | | | | |
| c=6, p=5: MC | | A | B | X | Y | **A** | B | X | Z | | | | | | | | | | | | | | | | |
| c=7, p=6: MC | | A | B | X | Y | A | **B** | X | Z | | | | | | | | | | | | | | | | |
| c=8, p=7: MC | | A | B | X | Y | A | B | **X** | Z | | | | | | | | | | | | | | | | |
| c=9, p=8: Fi: F'(8)=4 ⟹ p=4 | | A | B | X | Y | A | B | X | **Z** | | | | | | | | | | | | | | | | |
| c=9, p=4: MC | | | | | | A | B | X | **Y** | A | B | X | Z | | | | | | | | | | | | |
| c=10, p=5: MC | | | | | | A | B | X | Y | **A** | B | X | Z | | | | | | | | | | | | |
| c=11, p=6: MC | | | | | | A | B | X | Y | A | **B** | X | Z | | | | | | | | | | | | |
| c=12, p=7: MC | | | | | | A | B | X | Y | A | B | **X** | Z | | | | | | | | | | | | |
| c=13, p=8: MC | | | | | | A | B | X | Y | A | B | X | **Z** | | | | | | | | | | | | |
| c=14, p=9: MP | | | | | | **A** | **B** | **X** | **Y** | **A** | **B** | **X** | **Z** | | | | | | | | | | | | |

where $F'(8)=4 \Rightarrow p=4$

Exercise solutions: Matching  (Introduction to Comp.Ling)

Thursday, May 7, 2009

# K-M-P matching P against T

| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 | 21 | 22 | 23 | 24 | 25 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | X | A | B | X | Y | A | B | X | Y | A | B | X | Z | A | B | X | Z | A | B | X | Y | A | B | X | Z |
| c=1, p=1: F0 | **A** | B | X | Y | A | B | X | Z | | | | | | | | | | | | | | | | | |
| c=2, p=1: MC | | **A** | B | X | Y | A | B | X | Z | | | | | | | | | | | | | | | | |
| c=3, p=2: MC | | A | **B** | X | Y | A | B | X | Z | | | | | | | | | | | | | | | | |
| c=4, p=3: MC | | A | B | **X** | Y | A | B | X | Z | | | | | | | | | | | | | | | | |
| c=5, p=4: MC | | A | B | X | **Y** | A | B | X | Z | | | | | | | | | | | | | | | | |
| c=6, p=5: MC | | A | B | X | Y | **A** | B | X | Z | | | | | | | | | | | | | | | | |
| c=7, p=6: MC | | A | B | X | Y | A | **B** | X | Z | | | | | | | | | | | | | | | | |
| c=8, p=7: MC | | A | B | X | Y | A | B | **X** | Z | | | | | | | | | | | | | | | | |
| c=9, p=8: Fi: F'(8)=4 ⇒ p=4 | | A | B | X | Y | A | B | X | **Z** | | | | | | | | | | | | | | | | |
| c=9, p=4: MC | | | | | | A | B | X | **Y** | A | B | X | Z | | | | | | | | | | | | |
| c=10, p=5: MC | | | | | | A | B | X | Y | **A** | B | X | Z | | | | | | | | | | | | |
| c=11, p=6: MC | | | | | | A | B | X | Y | A | **B** | X | Z | | | | | | | | | | | | |
| c=12, p=7: MC | | | | | | A | B | X | Y | A | B | **X** | Z | | | | | | | | | | | | |
| c=13, p=8: MC | | | | | | A | B | X | Y | A | B | X | **Z** | | | | | | | | | | | | |
| c=14, p=9: MP | | | | | | **A** | **B** | **X** | **Y** | **A** | **B** | **X** | **Z** | | | | | | | | | | | | |

Exercise solutions: Matching  (Introduction to Comp.Ling)

Thursday, May 7, 2009

| X | A | B | X | Y | A | B | X | Y | A | B | X | Z | A | B | X | Z | A | B | X | Y | A | B | X | Z |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 | 21 | 22 | 23 | 24 | 25 |

c=1, p=1: F0

c=2, p=1: MC

c=3, p=2: MC

c=4, p=3: MC

c=5, p=4: MC

c=6, p=5: MC

c=7, p=6: MC

c=8, p=7: MC

c=9, p=8: Fi: F'(8)=4 $\Rightarrow$ p=4

c=9, p=4: MC

c=10, p=5: MC

c=11, p=6: MC

c=12, p=7: MC

c=13, p=8: MC

c=14, p=9: MP

| **A** | B | X | Y | A | B | X | Z | | | | | | | | | | | | | | | | | |
| | **A** | B | X | Y | A | B | X | Z | | | | | | | | | | | | | | | | |
| | A | **B** | X | Y | A | B | X | Z | | | | | | | | | | | | | | | | |
| | A | B | **X** | Y | A | B | X | Z | | | | | | | | | | | | | | | | |
| | A | B | X | **Y** | A | B | X | Z | | | | | | | | | | | | | | | | |
| | A | B | X | Y | **A** | B | X | Z | | | | | | | | | | | | | | | | |
| | A | B | X | Y | A | **B** | X | Z | | | | | | | | | | | | | | | | |
| | A | B | X | Y | A | B | **X** | Z | | | | | | | | | | | | | | | | |
| | A | B | X | Y | A | B | X | **Z** | | | | | | | | | | | | | | | | |
| | | | | | A | B | X | **Y** | A | B | X | Z | | | | | | | | | | | | |
| | | | | | A | B | X | Y | **A** | B | X | Z | | | | | | | | | | | | |
| | | | | | A | B | X | Y | A | **B** | X | Z | | | | | | | | | | | | |
| | | | | | A | B | X | Y | A | B | **X** | Z | | | | | | | | | | | | |
| | | | | | A | B | X | Y | A | B | X | **Z** | | | | | | | | | | | | |
| | | | | | **A** | **B** | **X** | **Y** | **A** | **B** | **X** | **Z** | | | | | | | | | | | | |
| | | | | | | | | | | | | | | | | | **A** | **B** | **X** | **Y** | **A** | **B** | **X** | **Z** |

Exercise solutions: Matching (Introduction to Comp.Ling)

Thursday, May 7, 2009

```java
public int match (String pattern, String text) {
    int p = 0;
    int s = 0;
    int t = 0;
    int matches = 0; // number of matches, return value
    while (t < text.length()) {
        if (p < pattern.length()) System.out.println("pattern("+p+") \t"+pattern.charAt(p));
        System.out.println("text("+t+") \t"+text.charAt(t));
        if (pattern.charAt(p) == text.charAt(t)) {
            // make sure to check against length-1 else OutOfRange exception!
            if (p < pattern.length()-1) {
                p = p+1;
                t = t+1;
            } else {
                System.out.println("Match found at position "+(s+1));
                p = 0;
                s = s+1;
                t = s;
                matches = matches+1;
            } // end if..else check for full occurrence of P
        } else {
            p = 1;
            s = s+1;
            t = s;
        } // end if..else check for character match
    } // end while over the text
    return matches;
} // end match
```

Thursday, May 7, 2009

```
public int match (String pattern, String text) {

    // ---------------------------------------
    // Initialization
    // ---------------------------------------
    int p = 0;
    int s = 0;
    int t = 0;
    int matches = 0; // number of matches, return value

    Vector patternVec = new Vector();
    Vector textVec = new Vector ();

    // Represent the pattern as a sequence of words

    StringTokenizer pst = new StringTokenizer(pattern);
    while (pst.hasMoreTokens()) {
       String word = (String)pst.nextToken();
       patternVec.addElement(word);
    } // end while

    // Represent the text as a sequence of words

    StringTokenizer tst = new StringTokenizer(text);
    while (tst.hasMoreTokens()) {
       String word = (String)tst.nextToken();
       textVec.addElement(word);
    } // end
```

Thursday, May 7, 2009

```
// ----------------------------------------
// Loop
// ----------------------------------------
// Note that the conditions now refer to the vectors, not to the original strings.

while (t < textVec.size()) {
    if (p < patternVec.size()) System.out.println("pattern("+p+") \t<"+(String)patternVec.elementAt(p)+">");
    System.out.println("text("+t+") \t<"+(String)textVec.elementAt(t)+">");
    if (((String)patternVec.elementAt(p)).equals((String)textVec.elementAt(t))) {
        // make sure to check against length-1 else OutOfRange exception!
        if (p < patternVec.size()-1) {
            p = p+1;
            t = t+1;
        } else {
            System.out.println("Match found at position "+(s+1));
            p = 0;
            s = s+1;
            t = s;
            matches = matches+1;
        } // end if..else check for full occurrence of P
    } else {
        p = 0;
        s = s+1;
        t = s;
    } // end if..else check for character match
} // end while over the text
return matches;
} // end match
```

- records:

    - the label on the edge to the vertex *v*

    - the path, i.e. the concatenation of the words on the path to *v*: L(v)

    - the parent of the vertex

    - the children of a vertex i

- basic accessor methods for adding, getting and setting

```java
public KeywordTreeNode buildBranch (Vector pvec) {
    KeywordTreeNode broot  = new KeywordTreeNode ();
    KeywordTreeNode parent = broot;
    boolean rootSet = false;
    String path = "";
    Iterator pvIter = pvec.iterator();
    while (pvIter.hasNext()) {
        String word = (String)pvIter.next();
        // check whether the root has been set; if not, initialize
        // the root, otherwise create a new node, and add it to
        // the current parent.
        if (!rootSet) {
            broot.setEdge(word);
            rootSet = true;
            broot.setPath(word);
            path = word;
        } else {
            KeywordTreeNode node = new KeywordTreeNode(word);
            path = path+" "+word;
            node.setPath(path);
            node.setParent(parent);
            parent.addChild(node);
            // Set the parent to be the current node
            parent = node;
        } // end if..else check for root or child
    } // end while
    return broot;
} // end
```

Thursday, May 7, 2009

```java
KeywordTreeNode branch = this.buildBranch(patternVec);

// Next, go down the tree as far as possible to find the lowest attachment point for this branch. From the root
// of the branch we go down tree, until we get to a point where none of the children on the branch would be
// matched; that is where we insert the (remainder of the) branch.

KeywordTreeNode branchnode = branch;   // the current branch node
KeywordTreeNode attachment = treeroot; // the node where to attach
boolean golower = (treeroot.isLeaf())?false:true;
while (golower) {
    // cycle over the children of the current attachment node
    boolean matchfound = false;
    Iterator chIter = attachment.getChildren();
    while (chIter.hasNext() && !matchfound) {
        KeywordTreeNode child = (KeywordTreeNode) chIter.next();
        // if this child has the same edge, and the same
        // path, as the current node in the branch, then
        // decend one node down the branch and set the
        // current node as the attachment point for the
        // remainder of the branch.
        if (child.getEdge().equals(branchnode.getEdge()) && child.getPath().equals(branchnode.getPath())) {
            matchfound = true;
            attachment = child;
            branchnode = branchnode.getFirstChild();
        } // end if.. check whether match found
    } // end while over children
    // continue if we found a matching child, and descended accordingly
    golower = matchfound;
} // end while
attachment.addChild(branchnode);
} // end for
```

Exercise solutions: Matching  (Introduction to Comp.Ling)

Thursday, May 7, 2009