

T-(538|725)-MALV, Málvinnsla Reglulegar segðir og Perl

Hrafn Loftsson¹ Hannes Högni Vilhjálmsson¹

¹Tölvunarfræðideild, Háskólinn í Reykjavík

Ágúst 2007

1 Strengir og mál

2 Reglulegar segðir

3 Forritunarmálið Perl

1 Strengir og mál

2 Reglulegar segðir

3 Forritunarmálið Perl

Stafróf

- Endanlegt mengi stafa.
- Dæmi: $\{0,1\}$ er tvíundarstafrófið.

Strengur

- Strengur s úr stafrófi Σ er endanleg runa stafa sem dregnir eru úr Σ .
- $|s|$ táknar lengdina á s .
- ϵ táknar tóma strenginn; lengd hans er 0.

Stafróf

- Endanlegt mengi stafa.
- Dæmi: $\{0,1\}$ er tvíundarstafrófið.

Strengur

- Strengur s úr stafrófi Σ er endanleg runa stafa sem dregnir eru úr Σ .
- $|s|$ tákna lengdina á s .
- ϵ tákna tóma strenginn; lengd hans er 0.

Skilgreining

- Mengi strengja.
- Dæmi: \emptyset , $\{\epsilon\}$, $\{ab, ba\}$, $\{011, 101, 111\}$.

Samtenging og margföldun

- Ef x og y eru strengir þá er samtenging þeirra, xy , strengur sem fæst með því að bæta y við x .
- $s\epsilon = \epsilon s = s$
- $s^0 = \epsilon$, $s^1 = s$, $s^2 = ss$,
- $s^i = ss^{i-1}$, $i > 0$

Skilgreining

- Mengi strengja.
- Dæmi: \emptyset , $\{\epsilon\}$, $\{ab, ba\}$, $\{011, 101, 111\}$.

Samtenging og margföldun

- Ef x og y eru strengir þá er samtenging þeirra, xy , strengur sem fæst með því að bæta y við x .
- $s\epsilon = \epsilon s = s$
- $s^0 = \epsilon$, $s^1 = s$, $s^2 = ss$,
- $s^i = ss^{i-1}$, $i > 0$

- $L \cup M = \{s \mid s \in L \text{ eða } s \in M\}$
- $LM = \{st \mid s \in L \text{ og } t \in M\}$
- Kleene closure: 0 eða fleiri samtengingar af L
 - $L^* = \bigcup_{i=0}^{\infty} L^i$
- Positive closure: 1 eða fleiri samtengingar af L
 - $L^+ = \bigcup_{i=1}^{\infty} L^i$

Dæmi um mál

$L = \{A, B, \dots, Z, a, b, \dots, z\}$ og $D = \{0, 1, \dots, 9\}$.

Hvaða mál (mengi strengja) eru þá:

- $L \cup D$
- LD
- L^4
- L^*
- $L(L \cup D)^*$
- D^+

Outline

1 Strengir og mál

2 Reglulegar segðir

3 Forritunarmálið Perl

Reglulegar segðir (e. regular expressions)

- Mál sem notað er til að lýsa mengi strengja.
- Sérstaklega öflugt til að lýsa mynstrum (e. patterns) sem leita skal að í texta.
- Sérhver regluleg segð r stendur fyrir mál $L(r)$.
- Samanstanda af venjulegum stöfum (t.d. abc) ásamt stöfum sem hafa sérstaka merkingu (þ.e. “metacharacters”) eins og “*” og “+”.
- Hægt að búa til flóknari reglulegar segðir úr einfaldri segðum með því að nota sérstakar reglur.

Fyrir stafróf Σ :

- 1 ϵ er regluleg segð (RS) sem stendur fyrir $\{\epsilon\}$.
- 2 Ef $a \in \Sigma$, þá er a RS sem stendur fyrir $\{a\}$.
- 3 G.r.f. að r og s séu RS sem standa fyrir málin $L(r)$ og $L(s)$. Þá gildir:
 - $(r)|(s)$ er RS sem táknar $L(r) \cup L(s)$.
 - $(r)(s)$ er RS sem táknar $L(r)L(s)$.
 - $(r)^*$ er RS sem táknar $(L(r))^*$.
 - (r) er RS sem táknar $L(r)$.

Forgangur:

- * hefur hæstan forgang.
- Samtenging hefur næst hæstan forgang.
- | hefur lægsta forgang.
- Því gildir t.d. að: $(a)|((b)*(c)) = a|b*c$

Dæmi um reglulegar segðir

Hvaða mál tákna reglulegu segðirnar:

- $a|b$
- $(a|b)(a|b)$
- a^*
- $a|b^*c$

Meira um reglulegar segðir

Aðrir stafir með sérstaka merkingu

Í mörgum tólum sem bjóða upp á reglulegar segðir hafa þessir stafir einnig sérstaka merkingu:

- `? + . {n}`
- Sjá skýringar í töflu 2.9 bls. 37

Meira um reglulegar segðir

“Character classes”

- Listi af stöfum innan í hornklofum stemmir við hvaða staf sem er í listanum.
- Reglulega segðin $[abc]$ merkir eitt tilvik af a , eða b eða c ($a|b|c$).

“Complement and range”

- $[\^a]$ merkir hvaða stafur sem er ekki a .
- $[a-zA-Z]$ merkir $a, b, \dots, z, A, B, \dots, Z$.

Meira um reglulegar segðir

“Character classes”

- Listi af stöfum innan í hornklofum stemmir við hvaða staf sem er í listanum.
- Reglulega segðin `[abc]` merkir eitt tilvik af `a`, eða `b` eða `c` (`a|b|c`).

“Complement and range”

- `[^a]` merkir hvaða stafur sem er ekki `a`.
- `[a-zA-Z]` merkir `a`, `b`, ..., `z`, `A`, `B`, ..., `Z`.

“Longest match”

Margræðni

- “String matching” getur verið margræð.
- T.d. strengurinn $s = \text{“aabbc”}$ og reglulega segðin a^+b^*
- Reglulega segðin passar við þessa hlutstrengi úr s : a , aa , ab , aab , abb , $aabb$

Einræðing - tvær reglur

Tól sem styðja reglulegar segðir:

- Byrja að “matcha” eins snemma og hægt er í strengnum.
- “Matcha” eins marga stafi í einu og mögulegt er.

“Longest match”

Margræðni

- “String matching” getur verið margræð.
- T.d. strengurinn $s = \text{“aabbcc”}$ og reglulega segðin a^+b^*
- Reglulega segðin passar við þessa hlutstrengi úr s : a , aa , ab , aab , abb , $aabb$

Einræðing - tvær reglur

Tól sem styðja reglulegar segðir:

- Byrja að “matcha” eins snemma og hægt er í strengnum.
- “Matcha” eins marga stafi í einu og mögulegt er.

Tengsl við endanlegar stöðuvélar

- Hægt að breyta reglulegri segð í endanlega stöðuvél á vélrænan hátt.
 - T.d. tekið fyrir í námskeiðinu *Þýðendur*.
- Endanleg stöðuvél getur þá borið kennsl á þá strengi sem tiltekin regluleg segð stendur fyrir.

Ýmis tól og forritunarmál

- `grep/egrep` (Unix/Linux-tól)
 - `grep 'ab*c' myFile`
 - Prentar út allar línur úr `myFile` sem innihalda strengina `ac`, `abc`, `abbc`, `abbbc`, o.s.frv.
 - Undir Windows getið þið sett upp *Cygwin*
<http://www.cygwin.com/> sem er “Linux-like environment for Windows”.
- Reglulegar segðir eru líka notaðar í Perl, Python, Java, C#.

Outline

1 Strengir og mál

2 Reglulegar segðir

3 Forritunarmálið Perl

Forritunarmálið Perl

Uppruni

- PERL = "Practical Extraction and Report Language"
- Búið til árið 1987 af Larry Wall
- Allrahandatöl fyrir UNIX notendur
- Öflugra en skeljaforrit og einfaldar í notkun en C
- Sérstaklega hröð textavinnsla

Halló Heimur

```
#Reading a string and printing it back out
print("What is your name? ");
$name = <STDIN>;
print "Pleased to meet you ", $name;
```

Útlit forrits

- Venjuleg textaskrá (*.pl) með röð Perl skipana
- Auðbil hunsað í upphafi línu
- Hver skipun endar á semikommu (";")
- Fjöldi auðbila utan strengja skiptir ekki máli
- Allt á eftir "#" er hunsað (athugasemd/skjölun)

Stutt forrit

```
#Prentar smá texta
print("Þetta er texti\n",
      "sem nær yfir tvær línur");
```


Útlit forrits

- Venjuleg textaskrá (*.pl) með röð Perl skipana
- Auðbil hunsað í upphafi línu
- Hver skipun endar á semikommu (";")
- Fjöldi auðbila utan strengja skiptir ekki máli
- Allt á eftir "#" er hunsað (athugasemd/skjölun)

Stutt forrit

```
#Prentar smá texta
print("Þetta er texti\n",
      "sem nær yfir tvær línur");
```

Kverður (Scalars)

- Allir strengir og tölur eru kverður
- Kverðubreytur byrja alltaf á "\$" (dollar)
- Perl breytir sjálfkrafa á milli strengja og talna

Dæmi

- `$a=2; $b=6; $c=$a.$b; $d=$c/2; print $d;`
- `$nafn='jón'; print "Ég heiti :\t$nafn\n";`

Kverður (Scalars)

- Allir strengir og tölur eru kverður
- Kverðubreytur byrja alltaf á "\$" (dollar)
- Perl breytir sjálfkrafa á milli strengja og talna

Dæmi

- `$a=2; $b=6; $c=$a.$b; $d=$c/2; print $d;`
- `$nafn='jón'; print "Ég heiti :\t$nafn\n";`

Fylki (Arrays)

- Fylki er safn kverða
- Fylkjabreytur byrja alltaf á "@" (at)
- Vísað er á stak með tölu í hornklofa "[númer]"

Dæmi

- `@g=('jón',25); print "$g[0] er $g[1] ára";`
- `@a=(1,2); @b=(3,4); @c=(@a,@b);`

Fylki (Arrays)

- Fylki er safn kverða
- Fylkjabreytur byrja alltaf á "@" (at)
- Vísað er á stak með tölu í hornklofa "[númer]"

Dæmi

- `@g=('jón',25); print "$g[0] er $g[1] ára";`
- `@a=(1,2); @b=(3,4); @c=(@a,@b);`

Tengifylki (Associative Array)

- Sama og tætitafla (hash table) eða kort (map)
- Fylki þar sem stök eru sótt með lyklastrengjum

Dæmi

- `%enska=('epli', 'apple', 'pera', 'pear');`
- `print enska{'epli'};`
- `print keys(%enska);`
- `print values(%enska);`

Skráarhald (File Handle)

- Strengur (upphafsstafastrengur) sem bendir á skrá
- Fyrst tengt við skrá og svo notað við lestur og skrif
- Perl er þegar búið að tengja **STDIN**, **STDOUT** og **STDERR**

Dæmi

```
open(INNSKRA, "text.txt");
open(UTSKRA, ">nidurstodur.txt");
while(<INNSKRA>) { #Ein lína í einu geymd í $_
    print UTSKRA $_, "\n——\n";
}
close(INNSKRA); close(UTSKRA);
```

Skráarhald (File Handle)

- Strengur (upphafsstafastrengur) sem bendir á skrá
- Fyrst tengt við skrá og svo notað við lestur og skrif
- Perl er þegar búið að tengja **STDIN**, **STDOUT** og **STDERR**

Dæmi

```
open(INNSKRA, "text.txt");
open(UTSKRA, ">nidurstodur.txt");
while(<INNSKRA>) { #Ein lína í einu geymd í $_
    print UTSKRA $_, "\n---\n";
}
close(INNSKRA); close(UTSKRA);
```


Sjálfgefið skráarhald

- Ef skráarhaldið er tómi strengurinn, þ.e. `<>` þá
 - Nota skrár tilgreindar á skipanalínu
 - Nota `STDIN` ef skrár ekki tilgreindar

Dæmi

```
while(<>) {  
    print $_; #Ath: má sleppa $_  
}
```

Sjálfgefið skráarhald

- Ef skráarhaldið er tómi strengurinn, þ.e. `<>` þá
 - Nota skrár tilgreindar á skipanalínu
 - Nota `STDIN` ef skrár ekki tilgreindar

Dæmi

```
while(<>) {  
    print $_; #Ath: má sleppa $_  
}
```

Núverandi lína

- Sérstaka breytan `$_` geymir oft þá línu sem verið er að lesa frá inntaki, þ.e. núverandi línu
- Ef þessa línu þarf að geyma, má nota aðra breytu:
`$sidast = $_;`
- Ýmsar aðgerðir gera ráð fyrir að unnið sé með þessa breytu, nema annað sé tekið fram

Skilyrtar skipanir

- `print "Tómt" if ! $fullt ;`
- `print "Tókst!" unless $villa >2;`
- `$n=100; print "$n\n" while $n-- > 0;`
- `$n=0; print "$n\n" until $n++ > 100;`

Skilyrtar blokkir

```
if ($n > 100) {  
    print "sigur!";  
} elsif ($s > 50) {  
    print "samt sigur!";  
} else {  
    print "tap";  
}
```

Skilyrtar skipanir

- `print "Tómt" if ! $fullt ;`
- `print "Tókst!" unless $villa >2;`
- `$n=100; print "$n\n" while $n-- > 0;`
- `$n=0; print "$n\n" until $n++ > 100;`

Skilyrtar blokkir

```
if ($n > 100) {  
    print "sigur!";  
} elsif ($s > 50) {  
    print "samt sigur!";  
} else {  
    print "tap";  
}
```

foreach lykkjan

- Mjög öflug lykjkuskipun til að vinna með stök fylkja
- `@tolur = ('einn', 'tveir', 'thrir');`
- `foreach $tala (@tolur) { print "Talan $tala ...\n"; }`

Notkun reglulegra segða í Perl

- Mátun (match): `m/regex/modifiers`
- Skipting (substitution): `s/regex/replacement/modifiers`
- Þýðing (translation):
`tr/search_list/replacement_list/modifiers`

Dæmi

```
while(<>) {  
    print if (m/ab*c/i);  
}
```

Notkun reglulegra segða í Perl

- Mátun (match): `m/regex/modifiers`
- Skipting (substitution): `s/regex/replacement/modifiers`
- Þýðing (translation):
`tr/search_list/replacement_list/modifiers`

Dæmi

```
while(<>) {  
    print if (m/ab*c/i);  
}
```


Annað dæmi

```
while($line = <>) {  
    if($line =~ m/ab*c/i) {  
        $line =~ s/ab*c/ABC/g;  
    } else {  
        $line =~ tr/A-Z/a-z/;  
    }  
}
```

Vísað í mátaðan streng

- Oft er þörf á að vísa aftur á mátaðan strenghluta, t.d. til að sjá hvort hann endurtekur sig eða til að prenta hann
- Strenghluti er settur í biðminni með svigum () og vísað á hann með `\númer` innan sömu segðar eða með `$númer` utan hennar

Dæmi

```
while($line = <>) {  
    while($line =~ m/\$ *([0-9]+)\.?([0-9]*)/g) {  
        print "Dollars: ", $1, " Cents: ", $2, "\n";  
    }  
}
```

Vísað í mátaðan streng

- Oft er þörf á að vísa aftur á mátaðan strenghluta, t.d. til að sjá hvort hann endurtekur sig eða til að prenta hann
- Strenghluti er settur í biðminni með svigum () og vísað á hann með `\númer` innan sömu segðar eða með `$númer` utan hennar

Dæmi

```
while($line = <>) {  
    while($line =~ m/\$( [0-9]+ )\.(?([0-9]*)/g) {  
        print "Dollars: ", $1, " Cents: ", $2, "\n";  
    }  
}
```

Dæmi

```
perl -i.old -p -e "s/foo/bar/g" wiffle .bat
```

- `-e` keyrir þær Perl skipanir sem eru í strengnum fyrir aftan
- Skipunin sem er keyrð hér er `s/foo/bar/g` sem skiptir út öllum strengjum `"foo"` í sjálfgefnu inntaki (hér er það hver lína í skránni `"wiffle.bat"`) fyrir strenginn `"bar"`
- `-p` prentar út niðurstöðuna, en vegna `-i` er það ekki gert á skjáinn, heldur í sömu skrá og inntakið er úr, en upphaflega skráin er vistuð með endingunni `".old"`