

R: analiza sekvenc

A. Blejec

December 14, 2015

Contents

1	Paket seqinr	1
2	ACNUC baze	4
2.1	Funkcije paketa seqinr	6
3	Primerjava zaporedij	9
3.1	Iskanje podobnosti zaporedij	9
4	dot-plot	10
5	Privzem podatkov iz baze Swissprot	22
6	Bioconductor in paket Biostrings	23
6.1	Paket Biostrings	23
6.2	Poravnanve DNA zaporedij z Needleman-Wunsch algoritmom . . .	23
6.3	Poravnanve zaporedij proteinov z Needleman-Wunsch algoritmom .	24
6.4	Daljša poravnava	26
	References	27

1 Paket seqinr

V paketu **seqinr** najdemo funkcije za manipulacijo sekvenc ([Charif and Lobry, 2007](#)). Paket je dobro dokumentiran, dokumentacija je dostopna na R-Forge v [SeqinR 2.0-7](#).

```
> library(seqinr)
> if(interactive()) help(package="seqinr")
>
```

Nekatere funkcije paketa **seqinr**

```
> tablecode()
>
```

Genetic code 1 : standard							
TTT	Phe	TCT	Ser	TAT	Tyr	TGT	Cys
TTC	Phe	TCC	Ser	TAC	Tyr	TGC	Cys
TTA	Leu	TCA	Ser	TAA	Stp	TGA	Stp
TTG	Leu	TCG	Ser	TAG	Stp	TGG	Trp
CTT	Leu	CCT	Pro	CAT	His	CGT	Arg
CTC	Leu	CCC	Pro	CAC	His	CGC	Arg
CTA	Leu	CCA	Pro	CAA	Gln	CGA	Arg
CTG	Leu	CCG	Pro	CAG	Gln	CGG	Arg
ATT	Ile	ACT	Thr	AAT	Asn	AGT	Ser
ATC	Ile	ACC	Thr	AAC	Asn	AGC	Ser
ATA	Ile	ACA	Thr	AAA	Lys	AGA	Arg
ATG	Met	ACG	Thr	AAG	Lys	AGG	Arg
GTT	Val	GCT	Ala	GAT	Asp	GGT	Gly
GTC	Val	GCC	Ala	GAC	Asp	GGC	Gly
GTA	Val	GCA	Ala	GAA	Glu	GGA	Gly
GTG	Val	GCG	Ala	GAG	Glu	GGG	Gly

> *tablecode(latexfile="tablecode.tex")*

TTT	Phe	TCT	Ser	TAT	Tyr	TGT	Cys
TTC	Phe	TCC	Ser	TAC	Tyr	TGC	Cys
TTA	Leu	TCA	Ser	TAA	Stp	TGA	Stp
TTG	Leu	TCG	Ser	TAG	Stp	TGG	Trp
CTT	Leu	CCT	Pro	CAT	His	CGT	Arg
CTC	Leu	CCC	Pro	CAC	His	CGC	Arg
CTA	Leu	CCA	Pro	CAA	Gln	CGA	Arg
CTG	Leu	CCG	Pro	CAG	Gln	CGG	Arg
ATT	Ile	ACT	Thr	AAT	Asn	AGT	Ser
ATC	Ile	ACC	Thr	AAC	Asn	AGC	Ser
ATA	Ile	ACA	Thr	AAA	Lys	AGA	Arg
ATG	Met	ACG	Thr	AAG	Lys	AGG	Arg
GTT	Val	GCT	Ala	GAT	Asp	GGT	Gly
GTC	Val	GCC	Ala	GAC	Asp	GGC	Gly
GTA	Val	GCA	Ala	GAA	Glu	GGA	Gly
GTG	Val	GCG	Ala	GAG	Glu	GGG	Gly

Table 1: Genetic code number 1: standard.

2 ACNUC baze

Opis s spletne strani o [ACNUC](#):

ACNUC is a retrieval system for the nucleotide and protein sequence databases GenBank, EMBL, UniProt/SWISS-PROT or NBRF-PIR, and for many other databases following the same formats.

Kratika ACNUC je izpeljana kot okrajšava francoskega izraza ACides NUCleiques (<http://pbil.univ-lyon1.fr/databases/acnuc/acnuc.html>)

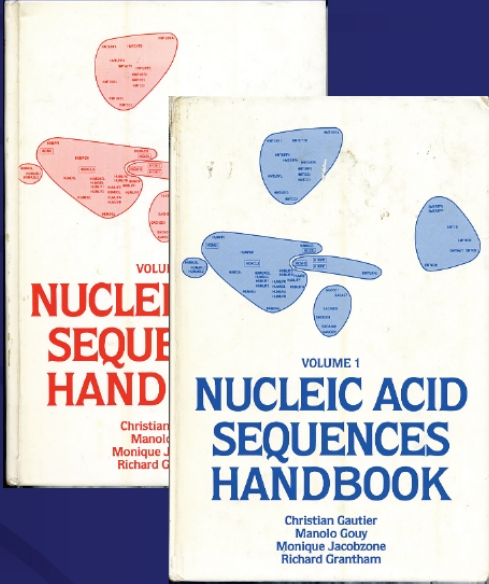
Interaktiven dostop do ACNUC baz omogoča tudi program [raa_query_win.exe](#)

Katere baze podatkov so dostopne?

```
> choosebank()
[1] "genbank"           "embl"           "emblwgs"
[4] "swissprot"        "ensembl"        "hogenom"
[7] "hogenomdna"       "hovergendna"    "hovergen"
[10] "hogenom5"         "hogenom5dna"    "hogenom4"
[13] "hogenom4dna"      "homolens"       "homolensdna"
[16] "hobacnucl"        "hobacprot"      "phever2"
[19] "phever2dna"       "refseq"         "greviews"
[22] "bacterial"        "archaeal"       "protozoan"
[25] "ensprotists"     "ensfungi"       "ensmetazoa"
[28] "ensplants"        "ensemblbacteria" "mito"
[31] "polymorphix"     "emglib"         "refseqViruses"
[34] "taxodb"
>
```

Prehistory and history

- First sequences library were available in books:
 - Atlas of Protein Sequences* de Dayhoff (1965-1978).
 - Nucleic Acid Sequences Hand-book* de Gautier *et al.* (1981):
 - 539 pages.
 - 1095 sequences.
 - 525 506 bp.
- First computer databases in the beginning of 80s:
 - GenBank (1979).
 - EMBL (1981).
 - PIR (1984).
 - SWISS-PROT (1986).



Pôle Bioinformatique Lyonnais – <http://pbil.univ-lyon1.fr> - Pôle Rhône-Alpes de Bioinformatique – <http://prabi.fr>

../clp/capture-ACNUC.jpg

http://www.genouest.org/documents/Evenements/200710_-ReNabiWorkshop/20071022_SimonPenel.pdf

Prvi dve knjigi sta debeli kakih 4.5cm in sta vsebovali 526506 baznih parov. Kako dolgo bi bilo danes v tiskani obliki vse kar je v bazi genebank? S pomočjo R in posegom v bazo genebank lahko ekstrahiramo podatke iz opisa.

```
> acnucbooksize <- 4.5 # cm
> acnucbp <- 526506 # baznih parov
> mybank <- choosebank("genbank")
> closebank()
> mybank$details

[1] "          ****      ACNUC Data Base Content      ****
[2] "          GenBank Release 210 (15 October 2015) Last Updated: Dec 13, 2015"
[3] "203,963,237,830 bases; 189,229,721 sequences; 28,010,697 subseqs; 855,553
[4] "Software M. Gouy, Lab. Biometrie et Biologie Evolutive, Universite Lyon I

> bpbk <- unlist(strsplit(mybank$details[3], split = " ")) [1]
> bpbk

[1] "203,963,237,830"

> bpbk <- as.numeric(paste(unlist(strsplit(bpbk, split = ",")),
+ collapse = ""))
> widthcm <- acnucbooksize * bpbk/acnucbp
> cat("Debelina:", round(widthkm <- widthcm/10^5), "km\n")

Debelina: 17 km
```

Če bi vse skupaj natisnili, bi bilo knjig za 17.4km!!
Leta 2014 jih je bilo za 15.6km, leta 2011 pa za 11.3km.

2.1 Funkcije paketa seqinr

```
> lseqinr()
[1] "a" "aaa"
[3] "AAstat" "acnucfclose"
[5] "acnucopen" "al2bp"
[7] "alllistranks" "alr"
[9] "amb" "as.alignment"
[11] "as.matrix.alignment" "as.SeqAcnucWeb"
[13] "as.SeqFastaAA" "as.SeqFastadna"
[15] "as.SeqFrag" "autosocket"
[17] "baselineabif" "bma"
[19] "c2s" "cai"
[21] "cfl" "choosebank"
[23] "circle" "clfcd"
[25] "clientid" "closebank"
[27] "col2alpha" "comp"
[29] "computePI" "con"
[31] "consensus" "count"
[33] "countfreelists" "countsubseqs"
[35] "crelistfromclientdata" "css"
[37] "dia.bactgensize" "dia.db.growth"
[39] "dist.alignment" "dotchart.uco"
[41] "dotPlot" "draw.oriloc"
[43] "draw.rearranged.oriloc" "draw.recstat"
[45] "exseq" "extract.breakpoints"
[47] "extractseqs" "fastacc"
[49] "gb2fasta" "gbk2g2"
[51] "gbk2g2.euk" "GC"
[53] "GC1" "GC2"
[55] "GC3" "GCpos"
[57] "get.db.growth" "get.ncbi"
[59] "getAnnot" "getAnnot.default"
[61] "getAnnot.list" "getAnnot.logical"
[63] "getAnnot.qaw" "getAnnot.SeqAcnucWeb"
[65] "getAnnot.SeqFastaAA" "getAnnot.SeqFastadna"
[67] "getAttributesocket" "getFrag"
[69] "getFrag.character" "getFrag.default"
[71] "getFrag.list" "getFrag.logical"
[73] "getFrag.qaw" "getFrag.SeqAcnucWeb"
[75] "getFrag.SeqFastaAA" "getFrag.SeqFastadna"
[77] "getFrag.SeqFrag" "getKeyword"
[79] "getKeyword.default" "getKeyword.list"
[81] "getKeyword.logical" "getKeyword.qaw"
[83] "getKeyword.SeqAcnucWeb" "getLength"
[85] "getLength.character" "getLength.default"
[87] "getLength.list" "getLength.logical"
[89] "getLength.qaw" "getLength.SeqAcnucWeb"
[91] "getLength.SeqFastaAA" "getLength.SeqFastadna"
[93] "getLength.SeqFrag" "getlistrank"
[95] "getliststate" "getLocation"
[97] "getLocation.default" "getLocation.list"
[99] "getLocation.logical" "getLocation.qaw"
[101] "getLocation.SeqAcnucWeb" "getName"
```

[103]	"getName.default"	"getName.list"
[105]	"getName.logical"	"getName.qaw"
[107]	"getName.SeqAcnucWeb"	"getName.SeqFastaAA"
[109]	"getName.SeqFastadna"	"getName.SeqFrag"
[111]	"getNumber.socket"	"getSequence"
[113]	"getSequence.character"	"getSequence.default"
[115]	"getSequence.list"	"getSequence.logical"
[117]	"getSequence.qaw"	"getSequence.SeqAcnucWeb"
[119]	"getSequence.SeqFastaAA"	"getSequence.SeqFastadna"
[121]	"getSequence.SeqFrag"	"getTrans"
[123]	"getTrans.character"	"getTrans.default"
[125]	"getTrans.list"	"getTrans.logical"
[127]	"getTrans.qaw"	"getTrans.SeqAcnucWeb"
[129]	"getTrans.SeqFastadna"	"getTrans.SeqFrag"
[131]	"getType"	"gfrag"
[133]	"ghelp"	"gln"
[135]	"glr"	"gls"
[137]	"is.SeqAcnucWeb"	"is.SeqFastaAA"
[139]	"is.SeqFastadna"	"is.SeqFrag"
[141]	"isenum"	"isn"
[143]	"kaks"	"kdb"
[145]	"knowndbs"	"lseqinr"
[147]	"modifylist"	"move"
[149]	"mv"	"n2s"
[151]	"ncbi.fna.url"	"ncbi.gbk.url"
[153]	"ncbi.ptt.url"	"ncbi.stats"
[155]	"oriloc"	"parser.socket"
[157]	"peakabif"	"permutation"
[159]	"pga"	"plot.SeqAcnucWeb"
[161]	"plotabif"	"plotladder"
[163]	"plotPanels"	"pmw"
[165]	"prepgetannots"	"prettyseq"
[167]	"print.qaw"	"print.SeqAcnucWeb"
[169]	"query"	"quitacnuc"
[171]	"read.abif"	"read.alignment"
[173]	"read.fasta"	"readBins"
[175]	"readfirstrec"	"readPanels"
[177]	"readsmj"	"rearranged.oriloc"
[179]	"recstat"	"residuecount"
[181]	"reverse.align"	"rho"
[183]	"rot13"	"s2c"
[185]	"s2n"	"savelist"
[187]	"SEQINR.UTIL"	"setlistname"
[189]	"splitseq"	"stresc"
[191]	"stutterabif"	"summary.SeqFastaAA"
[193]	"summary.SeqFastadna"	"swap"
[195]	"syncodons"	"synsequence"
[197]	"tablecode"	"test.co.recstat"
[199]	"test.li.recstat"	"translate"
[201]	"trimSpace"	"uco"
[203]	"ucoweight"	"where.is.this.acc"
[205]	"words"	"words.pos"
[207]	"write.fasta"	"zscore"

>

3 Primerjava zaporedij

3.1 Iskanje podobnosti zaporedij

Spremeba zaporedja v vektor znakov

```
> (seq1 <- "GGGATCACG")
[1] "GGGATCACG"
> (v <- s2c(seq1))
[1] "G" "G" "G" "A" "T" "C" "A" "C" "G"
```

Primerjava lege nukleotidov

```
> (P <- outer(v, v, "==")+0)
      [,1] [,2] [,3] [,4] [,5] [,6] [,7] [,8] [,9]
[1,]    1    1    1    0    0    0    0    0    1
[2,]    1    1    1    0    0    0    0    0    1
[3,]    1    1    1    0    0    0    0    0    1
[4,]    0    0    0    1    0    0    1    0    0
[5,]    0    0    0    0    1    0    0    0    0
[6,]    0    0    0    0    0    1    0    1    0
[7,]    0    0    0    1    0    0    1    0    0
[8,]    0    0    0    0    0    1    0    1    0
[9,]    1    1    1    0    0    0    0    0    1
```

Naredimo preglednejšo z dodatkom imen. Ničle so izpisane manj vidno.

```
> dimnames(P) <- list(v, v)
> print.table(P, zero.print=".")
  G G G A T C A C G
G 1 1 1 . . . . 1
G 1 1 1 . . . . 1
G 1 1 1 . . . . 1
A . . . 1 . . 1 . .
T . . . . 1 . . . .
C . . . . . 1 . 1 .
A . . . 1 . . 1 . .
C . . . . . 1 . 1 .
G 1 1 1 . . . . . 1
```

4 dot-plot

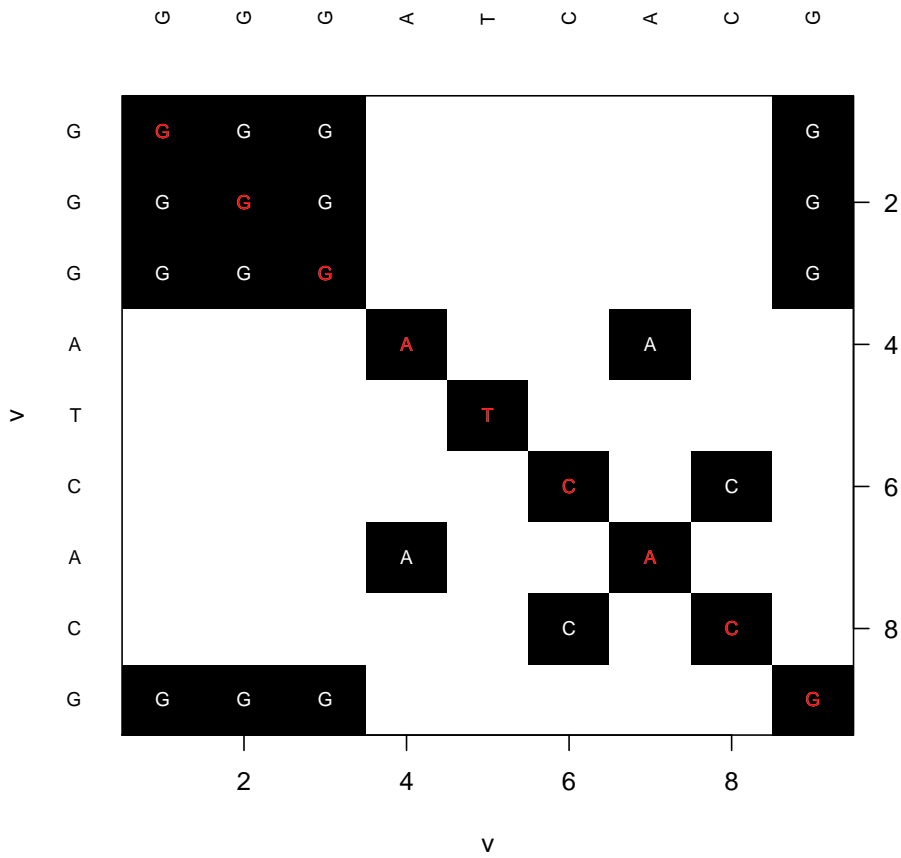
Dot plot je v analize sekvenc vpeljal Maizel (?).

```
> dotPlot <-
+ function (seq1, seq2, wsize = 1, wstep = 1, nmatch = 1, col = c("white",
+   "black"), xlab = deparse(substitute(seq1)), ylab = deparse(substitute(se
+   label=TRUE,
+   ...))
+ {
+   if (nchar(seq1[1]) > 1)
+     stop("seq1 should be provided as a vector of single chars")
+   if (nchar(seq2[1]) > 1)
+     stop("seq2 should be provided as a vector of single chars")
+   if (wsize < 1)
+     stop("non allowed value for wsize")
+   if (wstep < 1)
+     stop("non allowed value for wstep")
+   if (nmatch < 1)
+     stop("non allowed value for nmatch")
+   if (nmatch > wsize)
+     stop("nmatch > wsize is not allowed")
+   mkwin <- function(seq, wsize, wstep) {
+     sapply(seq(from = 1, to = length(seq) - wsize + 1, by = wstep),
+       function(i) c2s(seq[i:(i + wsize - 1)]))
+   }
+   wseq1 <- mkwin(seq1, wsize, wstep)
+   wseq2 <- mkwin(seq2, wsize, wstep)
+   if (nmatch == wsize) {
+     xy <- outer(wseq1, rev(wseq2), "==")
+   }
+   else {
+     "%==" <- function(x, y) colSums(sapply(x, s2c) == sapply(y,
+       s2c)) >= nmatch
+     xy <- outer(wseq1, rev(wseq2), "%==")
+   }
+   image(x = seq(from = 1, to = length(wseq1), length = length(wseq1)),
+     y = seq(from = 1, to = length(wseq2), length = length(wseq2)),
+     z = xy, col = col, xlab = xlab, ylab = ylab, axes=FALSE, ...)
+   axis(1)
+
+     axis(4, at=length(wseq2)-axTicks(4)+1, labels=axTicks(4), las=2)
+     n <- length(wseq1)
+     text(1:n, n+1.5, wseq1, cex=0.75, xpd=TRUE, srt=90, adj=0)
+     m <- length(wseq2)
+     text(0, m:1, wseq2, cex=0.75, xpd=TRUE, adj=1)
+
+     if(label) {text(row(xy), col(xy),
+       wseq1, cex=0.75/wsize, col="white")
+     text(1:n,
+       m:1,
+       wseq2, col=2, cex=0.75/wsize)
+   }
+   box()
+ }
```

```
+ invisible(list(xy=xy, seq1=wseq1, seq2=wseq2))  
+ }
```

Prikaz s funkcijo dotPlot()

```
> dotPlot(v, v, wsize=1)
```



Palindrom

```
> seq2 <- "pericarezeracirep"
```

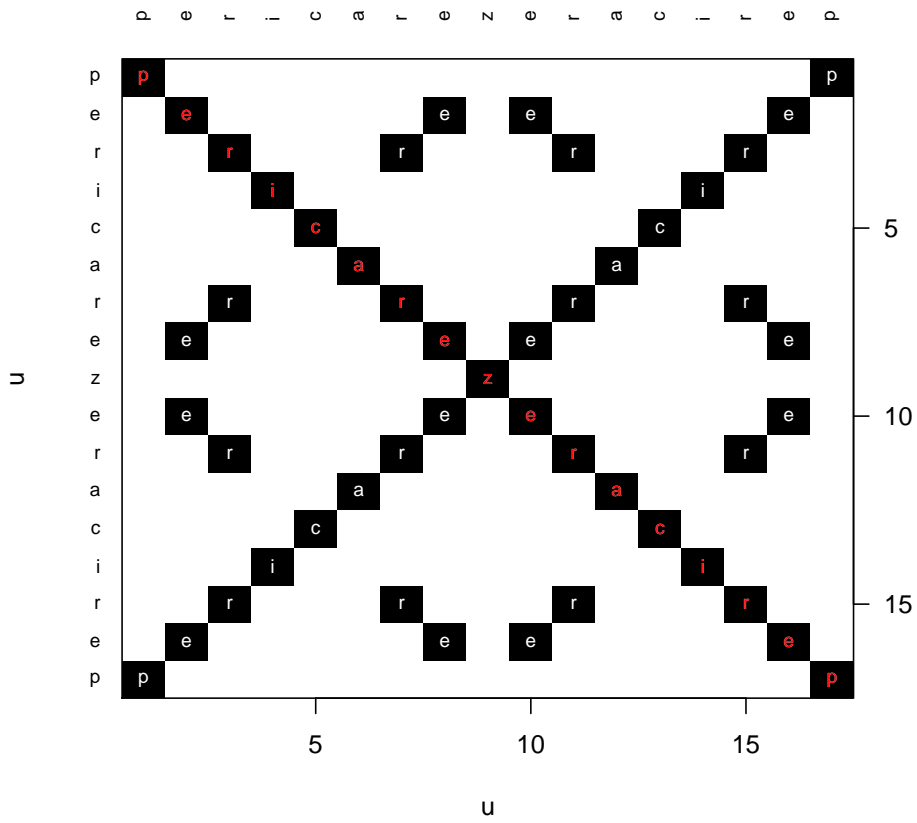
```
> (u <- s2c(seq2))
```

```
[1] "p" "e" "r" "i" "c" "a" "r" "e" "z" "e" "r" "a" "c" "i" "r" "e"
```

```
[17] "p"
```

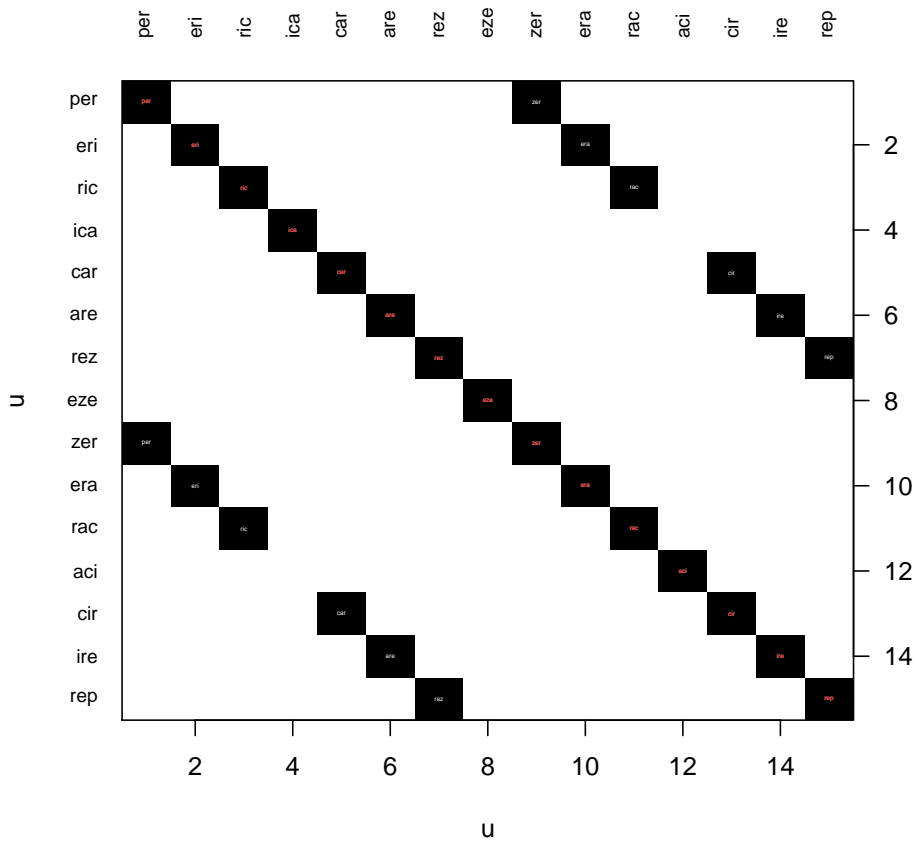
```
> dotPlot(u, u)
```

```
>
```



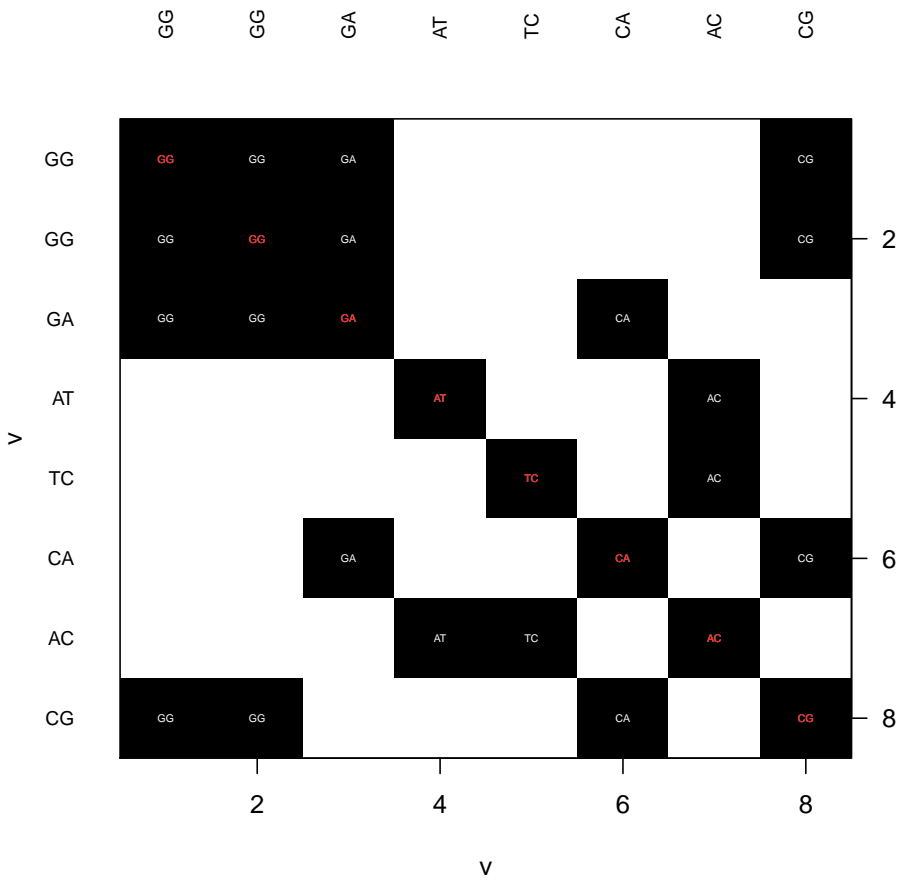
Primerjava z daljšim drsečim oknom

```
> dotPlot(u, u, wsize=3, nmatch=2)
>
```

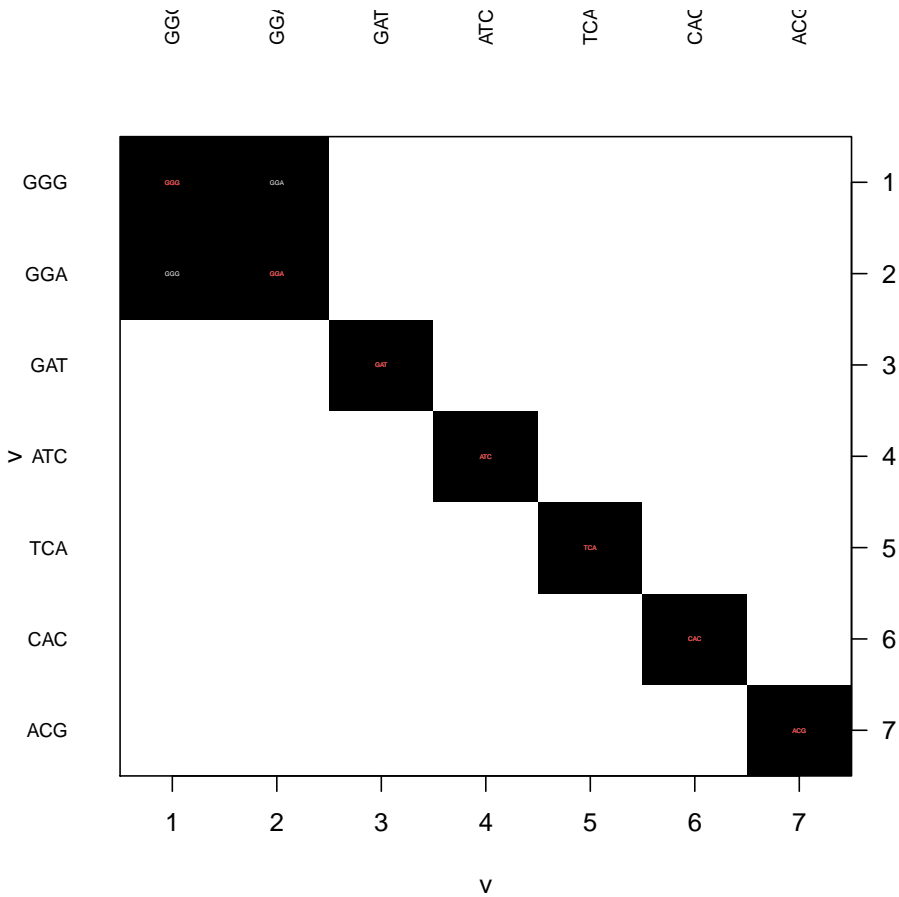


Daljšo okno, sprememba števila ujemanj

> dotPlot(v, v, wsize=2)



> dotPlot (v, v, wsize=3, nmatch=2)



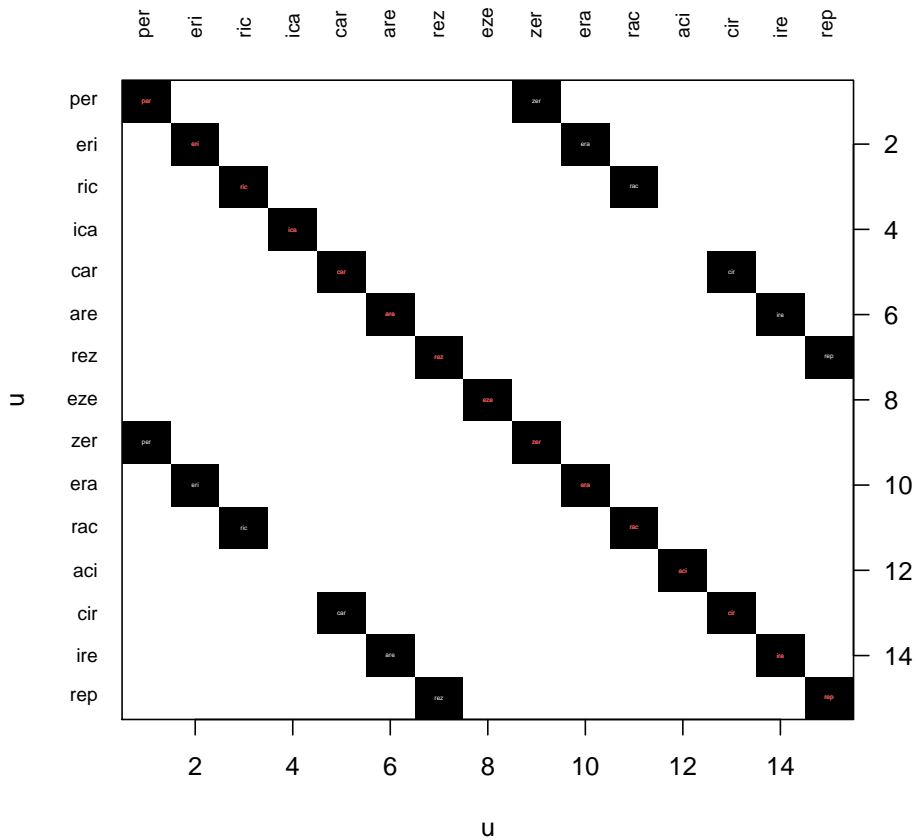
```

> seq2 <- "pericarezeracirep"
> (u<- s2c(seq2))

[1] "p" "e" "r" "i" "c" "a" "r" "e" "z" "e" "r" "a" "c" "i" "r" "e"
[17] "p"

> dotPlot (u, u, wsize=3, nmatch=2)
>

```



Daljša sekvenca

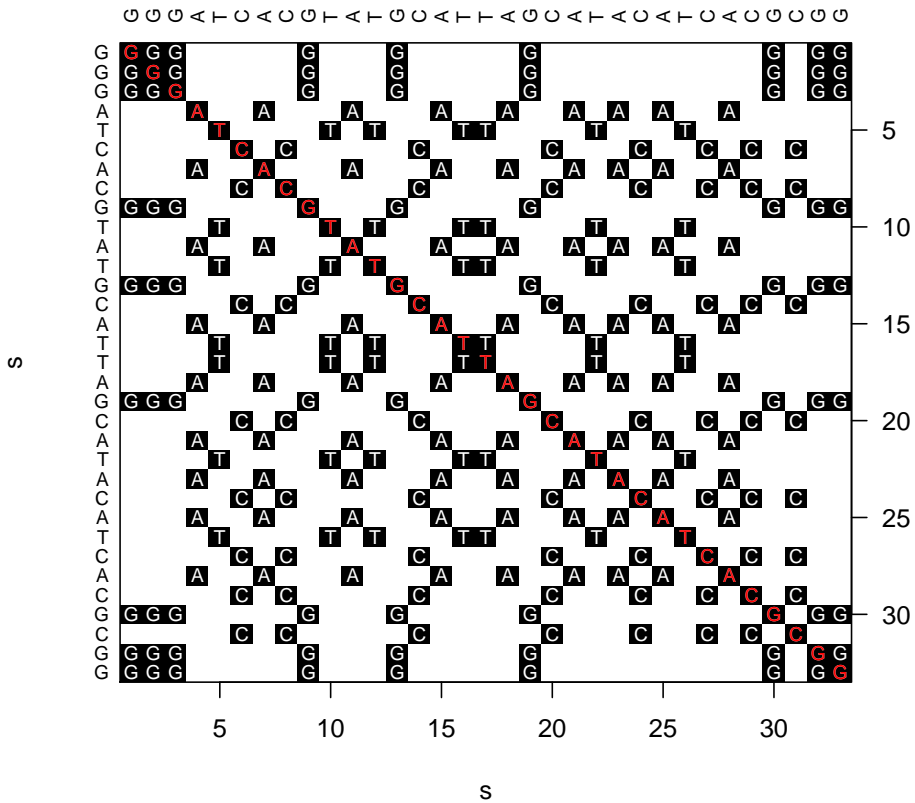
```

> seq1 <- "GGGATCACGTATGCATTAGCATAACATCACGCGG"
> wsize <- 1
> (s <- s2c(seq1))

[1] "G" "G" "G" "A" "T" "C" "A" "C" "G" "T" "A" "T" "G" "C" "A" "T"
[17] "T" "A" "G" "C" "A" "T" "A" "C" "A" "T" "C" "A" "C" "G" "C" "G"
[33] "G"

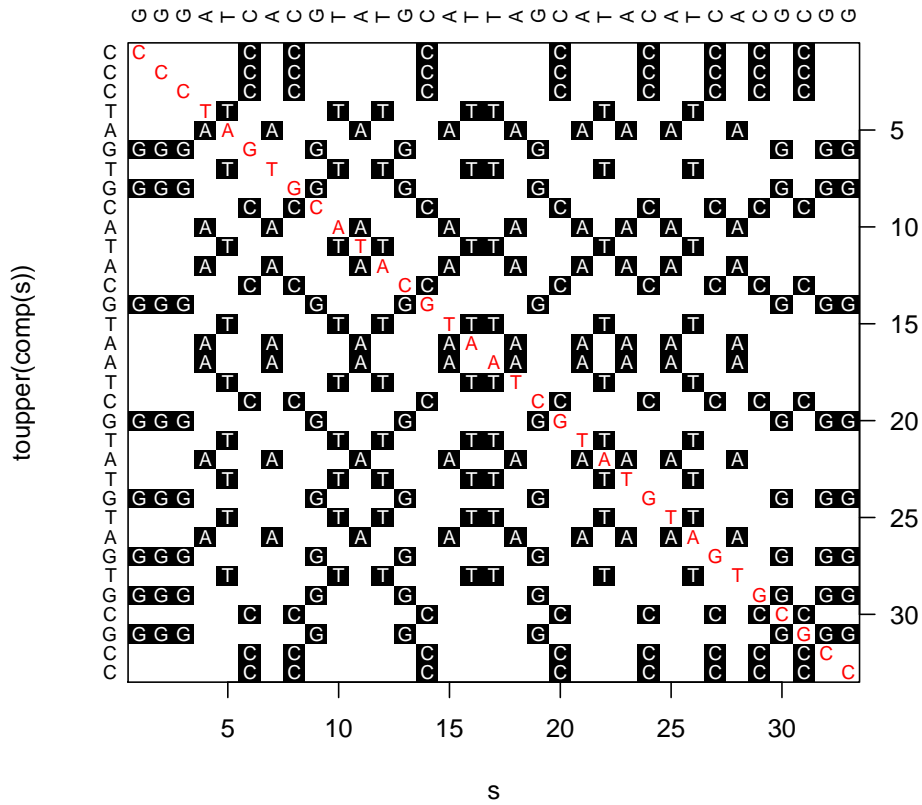
> dp <- dotPlot (s, s, wsize=wsize)

```

Funkcija `comp()` naredi komplement sekvence

```
> seq1 <- "GGGATCAGTATGCATTAGCATAACATCACGCGG"
> wsize <- 1
> (s <- s2c(seq1))
[1] "G" "G" "G" "A" "T" "C" "A" "C" "G" "T" "A" "T" "G" "C" "A" "T"
[17] "T" "A" "G" "C" "A" "T" "A" "C" "A" "T" "C" "A" "C" "G" "C" "G"
[33] "G"
> dp <- dotPlot(s,toupper(comp(s)),wsize=wsize)
```

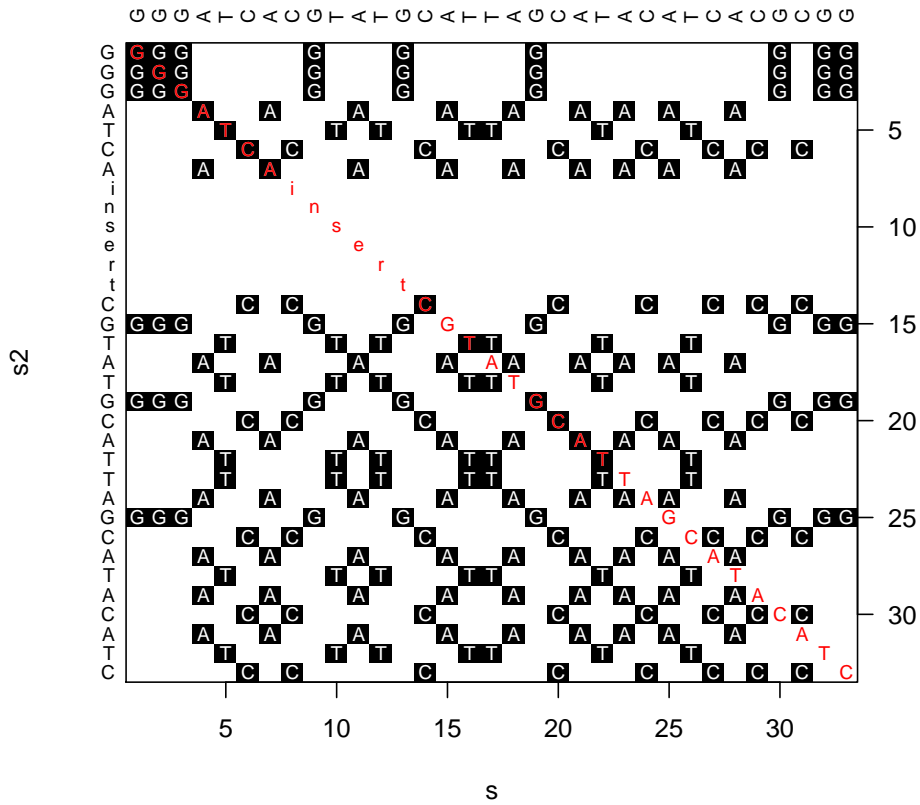


Vrinjena sekvenca

```

> seq1 <- "GGGATCACGTATGCATTAGCATAACATCACGCGG"
> s <- s2c(seq1)
> s2 <-c(s[1:7],s2c("insert"),s[8:(length(s)-6)])
> seq1
[1] "GGGATCACGTATGCATTAGCATAACATCACGCGG"
> c2s(s2)
[1] "GGGATCAinsertCGTATGCATTAGCATAACATC"
> dotPlot(s,s2)

```



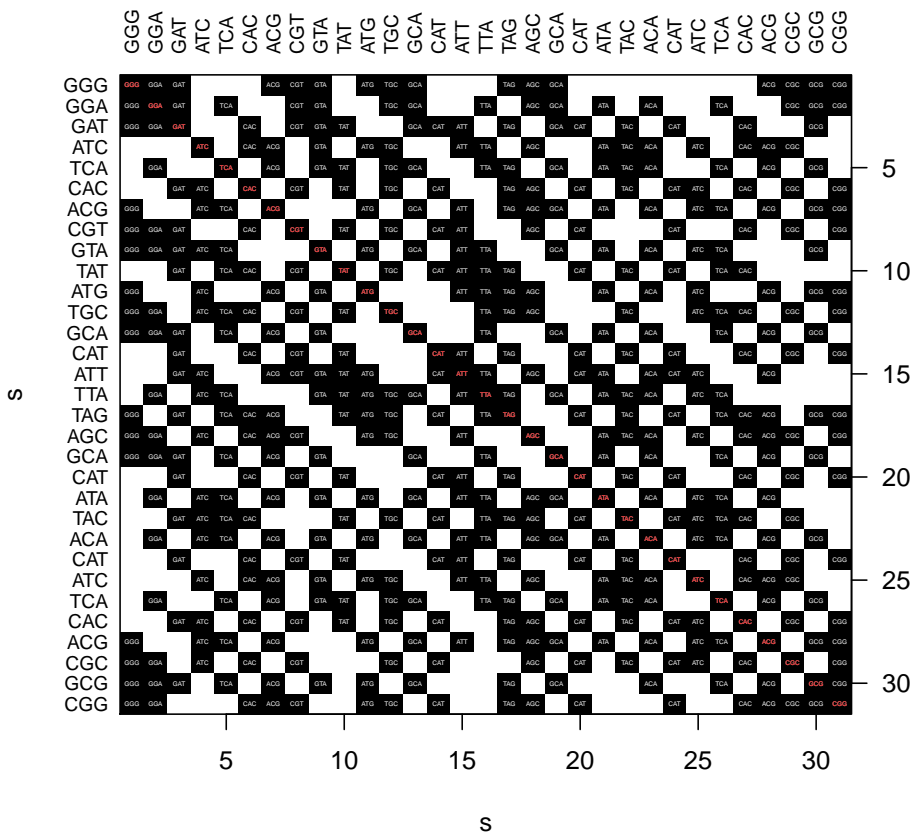
```

> seq1 <- "GGGATCACGTATGCATTAGCATAACATCACGCGG"
> wsize <- 3
> (s <- s2c(seq1))

[1] "G" "G" "G" "A" "T" "C" "A" "C" "G" "T" "A" "T" "G" "C" "A" "T"
[17] "T" "A" "G" "C" "A" "T" "A" "C" "A" "T" "C" "A" "C" "G" "C" "G"
[33] "G"

> dp <- dotPlot(s,s,wsize=wsize)
>

```



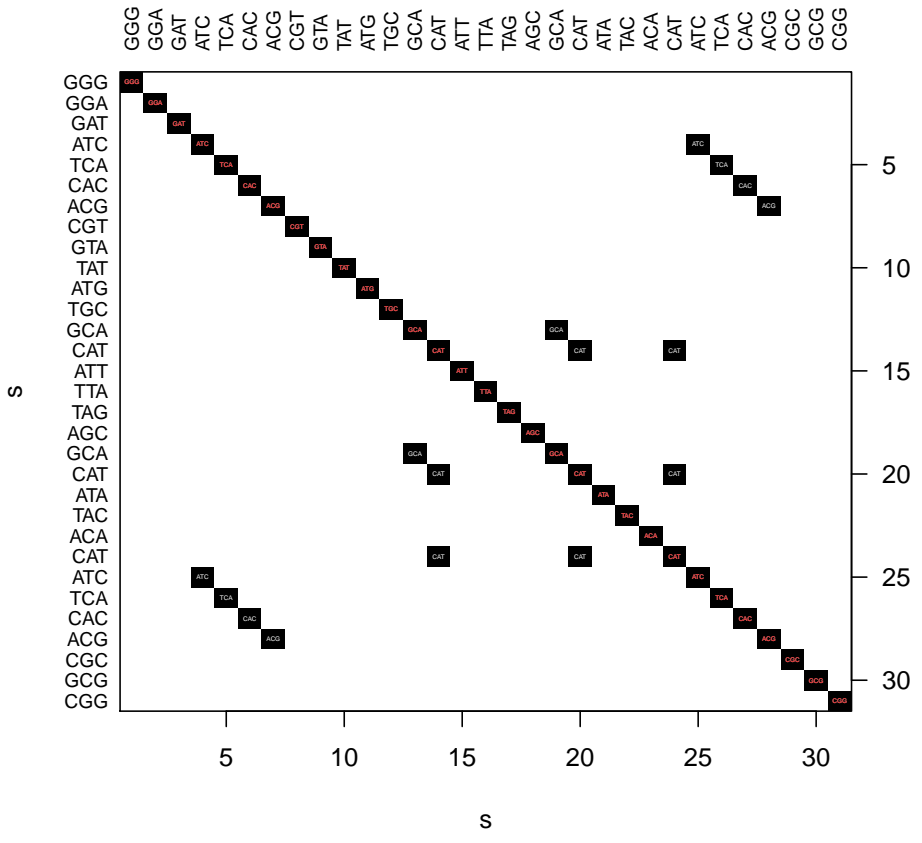
```

> seq1 <- "GGGATCACGTATGCATTAGCATAACATCACGCGG"
> wsize <- 3
> nmatch <- 3
> (s <- s2c(seq1))

[1] "G" "G" "G" "A" "T" "C" "A" "C" "G" "T" "A" "T" "G" "C" "A" "T"
[17] "T" "A" "G" "C" "A" "T" "A" "C" "A" "T" "C" "A" "C" "G" "C" "G"
[33] "G"

> dp <- dotPlot(s,s,wsize=wsize,nmatch=nmatch)

```



5 Privzem podatkov iz baze Swissprot

Izberemo dva proteina iz Swissprot. Accession Q9CD83 in A0PQ23.

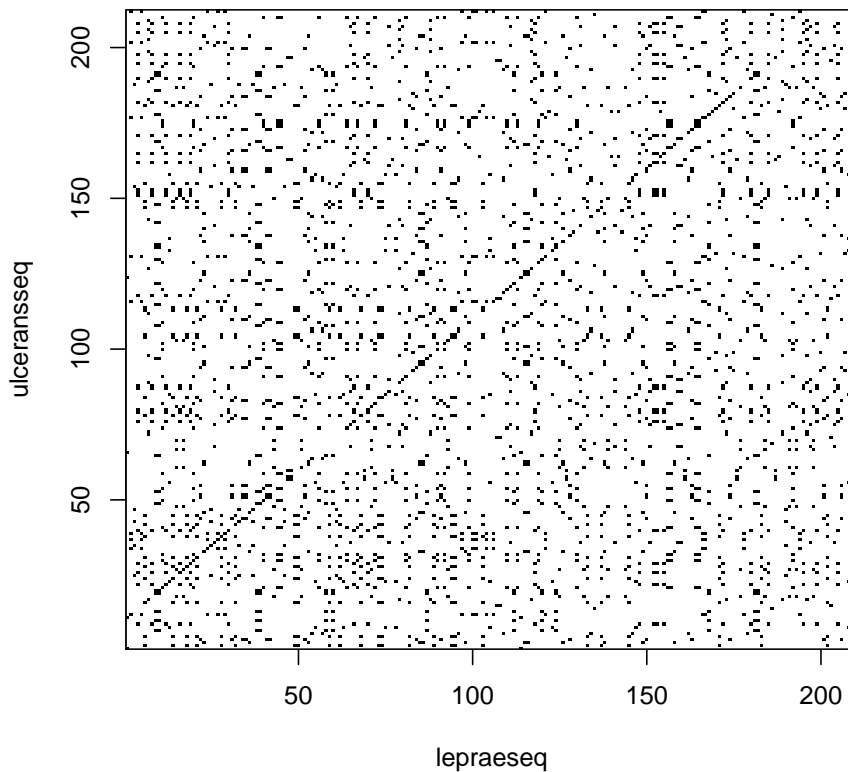
```
> library("seqinr")
> choosebank("swissprot")
> leprae <- query("leprae", "AC=Q9CD83")
> lepraeseq <- getSequence(leprae$req[[1]])
> ulcerans <- query("ulcerans", "AC=A0PQ23")
> ulceransseq <- getSequence(ulcerans$req[[1]])
> closebank()
> lepraeseq # Display the contents of "lepraeseq"

 [1] "M" "T" "N" "R" "T" "L" "S" "R" "E" "E" "I" "R" "K" "L" "D" "R"
[17] "D" "L" "R" "I" "L" "V" "A" "T" "N" "G" "T" "L" "T" "R" "V" "L"
[33] "N" "V" "V" "A" "N" "E" "E" "I" "V" "V" "D" "I" "I" "N" "Q" "Q"
[49] "L" "L" "D" "V" "A" "P" "K" "I" "P" "E" "L" "E" "N" "L" "K" "I"
[65] "G" "R" "I" "L" "Q" "R" "D" "I" "L" "L" "K" "G" "Q" "K" "S" "G"
[81] "I" "L" "F" "V" "A" "A" "E" "S" "L" "I" "V" "I" "D" "L" "L" "P"
[97] "T" "A" "I" "T" "T" "Y" "L" "T" "K" "T" "H" "H" "P" "I" "G" "E"
[113] "I" "M" "A" "A" "S" "R" "I" "E" "T" "Y" "K" "E" "D" "A" "Q" "V"
[129] "W" "I" "G" "D" "L" "P" "C" "W" "L" "A" "D" "Y" "G" "Y" "W" "D"
[145] "L" "P" "K" "R" "A" "V" "G" "R" "R" "Y" "R" "I" "I" "A" "G" "G"
[161] "Q" "P" "V" "I" "I" "T" "T" "E" "Y" "F" "L" "R" "S" "V" "F" "Q"
[177] "D" "T" "P" "R" "E" "E" "L" "D" "R" "C" "Q" "Y" "S" "N" "D" "I"
[193] "D" "T" "R" "S" "G" "D" "R" "F" "V" "L" "H" "G" "R" "V" "F" "K"
[209] "N" "L"

>
```

Dotplot - diagonala kaže enake aminokislino na podobnih ali enakih mestih v obeh proteinih.

```
> rm(dotPlot) # odstranim svoj dotPlot
> dotPlot(lepraeseq, ulceransseq)
>
```



6 Bioconductor in paket Biostrings

Za poravnave sekvenc lahko uporabite paket **Biostrings**, ki je del obširnega in dobro dokumentiranega sistema *Bioconductor*.

6.1 Paket Biostrings

Nekaj o tem si lahko preberete v poglavju [Pairwise alignment](#) v spletni knjigi [Little Book of R for Bioinformatics](#) (Coghlan, 2012).

6.2 Poravnanve DNA zaporedij z Needleman-Wunsch algoritmom

```
> library(Biostrings)
> sigma <- nucleotideSubstitutionMatrix(
+       match = 2, mismatch = -1, baseOnly = TRUE)
> sigma # Print out the matrix
```

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

Optimalna poravnava

```

> s1 <- "GAATTC"
> s2 <- "GATTA"
> globalAligns1s2 <- pairwiseAlignment(s1, s2,
+ substitutionMatrix = sigma,
+ gapOpening = -2,
+ gapExtension = -8,
+ scoreOnly = FALSE)
> #
> globalAligns1s2 # Print out the optimal alignment and its score
Global PairwiseAlignmentsSingleSubject (1 of 1)
pattern: [1] GAATTC
subject: [1] GA-TTA
score: -3

```

6.3 Poravnanve zaporedij proteinov z Needleman-Wunsch algoritmom

```

> data(BLOSUM50)
> BLOSUM50 # Print out the data

```

	A	R	N	D	C	Q	E	G	H	I	L	K	M	F	P	S	T	W	Y	V	B	Z
A	5	-2	-1	-2	-1	-1	-1	0	-2	-1	-2	-1	-1	-3	-1	1	0	-3	-2	0	-2	-1
R	-2	7	-1	-2	-4	1	0	-3	0	-4	-3	3	-2	-3	-3	-1	-1	-3	-1	-3	-1	0
N	-1	-1	7	2	-2	0	0	0	1	-3	-4	0	-2	-4	-2	1	0	-4	-2	-3	4	0
D	-2	-2	2	8	-4	0	2	-1	-1	-4	-4	-1	-4	-5	-1	0	-1	-5	-3	-4	5	1
C	-1	-4	-2	-4	13	-3	-3	-3	-3	-2	-2	-3	-2	-2	-4	-1	-1	-5	-3	-1	-3	-3
Q	-1	1	0	0	-3	7	2	-2	1	-3	-2	2	0	-4	-1	0	-1	-1	-1	-3	0	4
E	-1	0	0	2	-3	2	6	-3	0	-4	-3	1	-2	-3	-1	-1	-1	-3	-2	-3	1	5
G	0	-3	0	-1	-3	-2	-3	8	-2	-4	-4	-2	-3	-4	-2	0	-2	-3	-3	-4	-1	-2
H	-2	0	1	-1	-3	1	0	-2	10	-4	-3	0	-1	-1	-2	-1	-2	-3	2	-4	0	0
I	-1	-4	-3	-4	-2	-3	-4	-4	-4	5	2	-3	2	0	-3	-3	-1	-3	-1	4	-4	-3
L	-2	-3	-4	-4	-2	-2	-3	-4	-3	2	5	-3	3	1	-4	-3	-1	-2	-1	1	-4	-3
K	-1	3	0	-1	-3	2	1	-2	0	-3	-3	6	-2	-4	-1	0	-1	-3	-2	-3	0	1
M	-1	-2	-2	-4	-2	0	-2	-3	-1	2	3	-2	7	0	-3	-2	-1	-1	0	1	-3	-1
F	-3	-3	-4	-5	-2	-4	-3	-4	-1	0	1	-4	0	8	-4	-3	-2	1	4	-1	-4	-4
P	-1	-3	-2	-1	-4	-1	-1	-2	-2	-3	-4	-1	-3	-4	10	-1	-1	-4	-3	-3	-2	-1
S	1	-1	1	0	-1	0	-1	0	-1	-3	-3	0	-2	-3	-1	5	2	-4	-2	-2	0	0
T	0	-1	0	-1	-1	-1	-1	-2	-2	-1	-1	-1	-1	-2	-1	2	5	-3	-2	0	0	-1
W	-3	-3	-4	-5	-5	-1	-3	-3	-3	-3	-2	-3	-1	1	-4	-4	-3	15	2	-3	-5	-2
Y	-2	-1	-2	-3	-3	-1	-2	-3	2	-1	-1	-2	0	4	-3	-2	-2	2	8	-1	-3	-2
V	0	-3	-3	-4	-1	-3	-3	-4	-4	4	1	-3	1	-1	-3	-2	0	-3	-1	5	-4	-3
B	-2	-1	4	5	-3	0	1	-1	0	-4	-4	0	-3	-4	-2	0	0	-5	-3	-4	5	2
Z	-1	0	0	1	-3	4	5	-2	0	-3	-3	1	-1	-4	-1	0	-1	-2	-2	-3	2	5
X	-1	-1	-1	-1	-2	-1	-1	-2	-1	-1	-1	-1	-1	-2	-2	-1	0	-3	-1	-1	-1	-1
*	-5	-5	-5	-5	-5	-5	-5	-5	-5	-5	-5	-5	-5	-5	-5	-5	-5	-5	-5	-5	-5	-5
X	*																					
A	-1	-5																				
R	-1	-5																				
N	-1	-5																				
D	-1	-5																				
C	-2	-5																				
Q	-1	-5																				
E	-1	-5																				


```

G -2 -5
H -1 -5
I -1 -5
L -1 -5
K -1 -5
M -1 -5
F -2 -5
P -2 -5
S -1 -5
T  0 -5
W -3 -5
Y -1 -5
V -1 -5
B -1 -5
Z -1 -5
X -1 -5
* -5  1

```

Matrike vrednotenja zamenjav

```
> data(package="Biostrings")
```

```
> data(BLOSUM50)
```

```
> s3 <- "PAWHEAE"
```

```
> s4 <- "HEAGAWGHEE"
```

```
> globalAligns3s4 <- pairwiseAlignment(s3, s4,
```

```
+ substitutionMatrix = "BLOSUM50",
```

```
+ gapOpening = -2,
```

```
+ gapExtension = -8,
```

```
+ scoreOnly = FALSE)
```

```
> #
```

```
> globalAligns3s4 # Print out the optimal global alignment and its score
```

```
Global PairwiseAlignmentsSingleSubject (1 of 1)
```

```
pattern: [1] P---AWHEAE
```

```
subject: [1] HEAGAWGHEE
```

```
score: -5
```

6.4 Daljša poravnava

Pretvorba v nize znakov

```
> lepraeseqstring <- c2s(lepraeseq)      # Make a string that contains the sequ
> ulceransseqstring <- c2s(ulceransseq) # Make a string that contains the sequ
> ulceransseq[1:15]
[1] "M" "L" "A" "V" "L" "P" "E" "K" "R" "E" "M" "T" "E" "C" "H"
> ulceransseqstring
[1] "MLAVLPEKREMTECHLSDEEIRKLNRLRILVIATNGTLTRILNVLANDEIVVEIVKQIQDAAPEMDGDHSS
```

Če je potrebno, spremenimo v velike črke

```
> lepraeseqstring <- toupper(lepraeseqstring)
> ulceransseqstring <- toupper(ulceransseqstring)

> c(lepraeseqstring,ulceransseqstring)
[1] "MTNRTLSREEIRKLDRLRILVATNGTLTRVLNVVANEEIVVDIINQQLLDVAPKIPLENLKIQRILQDIL
[2] "MLAVLPEKREMTECHLSDEEIRKLNRLRILVIATNGTLTRILNVLANDEIVVEIVKQIQDAAPEMDGDHSS
```

Poravnava

```
> globalAlignLepraeUlcerans <- pairwiseAlignment(
+ lepraeseqstring, ulceransseqstring,
+ substitutionMatrix = BLOSUM50,
+ gapOpening = -2, gapExtension = -8, scoreOnly = FALSE)
> globalAlignLepraeUlcerans # Print out the optimal global alignment and its s
Global PairwiseAlignmentsSingleSubject (1 of 1)
pattern: [1] MT-----NR--T---LSREEIRKLDRLRI...EELDRCQYSNDIDTRSGDRFVLHGRVFKN
subject: [1] MLAVLPEKREMTECHLSDEEIRKLNRLRI...EPIRHQRS--VGT-SA-R---SGRSICT
score: 627
```

References

Charif, D. and J. Lobry (2007). SeqinR 1.0-2: a contributed package to the R project for statistical computing devoted to biological sequences retrieval and analysis. In H. R. U. Bastolla, M. Porto and M. Vendruscolo (Eds.), *Structural approaches to sequence evolution: Molecules, networks, populations*, Biological and Medical Physics, Biomedical Engineering, pp. 207–232. New York: Springer Verlag. ISBN : 978-3-540-35305-8. 1

Coghlan, A. (2012). Little book of R for bioinformatics.

<http://a-little-book-of-r-for-bioinformatics.readthedocs.org/en/latest/>.

23

SessionInfo

Windows 7 x64 (build 7601) Service Pack 1

- R version 3.2.0 (2015-04-16), x86_64-w64-mingw32
- Locale: LC_COLLATE=Slovenian_Slovenia.1250, LC_CTYPE=Slovenian_Slovenia.1250, LC_MONETARY=Slovenian_Slovenia.1250, LC_NUMERIC=C, LC_TIME=Slovenian_Slovenia.1250
- Base packages: base, datasets, graphics, grDevices, methods, parallel, stats, stats4, utils
- Other packages: ade4 1.7-2, BiocGenerics 0.14.0, Biostrings 2.36.0, IRanges 2.2.0, patchDVI 1.9.1601, S4Vectors 0.6.0, seqinr 3.1-3, XVector 0.8.0
- Loaded via a namespace (and not attached): tools 3.2.0, zlibbioc 1.14.0

Project path: D:/_Y/R/Bioinformatika

View as vignette

Project files can be viewed by pasting this code to R console:

```
> projectName <-"Bioinformatika"; mainFile <-"dotplot"

> commandArgs ()
> library(tkWidgets)
> # getrootpath <- function() {
> # fp <- (strsplit(getwd(), "/"))[[1]]
> # file <- file.path(paste(fp[-length(fp)], collapse = "/"))
> # return(file)
> # }
> # fileName <- function(name="bla", ext="PDF") paste(name, ext, sep=".")
> openPDF(file.path(dirname(getwd()), "doc", paste(mainFile, "PDF", sep=".")))
> viewVignette("viewVignette", projectName, file.path("../doc", paste(mainFile
>
```