

# Package ‘prozor’

April 20, 2017

**Type** Package

**Title** Minimal Protein Set Explaining Peptide Spectrum Matches

**Version** 0.2.3

**Author** Witold Wolski

**Maintainer** Witold Wolski <[wewolski@gmail.com](mailto:wewolski@gmail.com)>

**Description** Determine minimal protein set explaining peptide spectrum matches. Utility functions for creating libraries with decoys. Planned is peptide FDR estimation for search results.

**License** GPL-3

**LazyData** TRUE

**Imports** AhoCorasickTrie, Matrix, doParallel, foreach, plyr, seqinr, stringr

**URL** <https://github.com/protviz/prozor>

**BugReports** <https://github.com/protviz/prozor/issues>

**Repository** CRAN

**RoxygenNote** 6.0.1

**NeedsCompilation** no

**Date/Publication** 2017-04-20 11:43:20 UTC

## R topics documented:

annotateAHO . . . . .	2
annotatePeptides . . . . .	3
annotateVec . . . . .	3
createDecoyDB . . . . .	4
filterSequences . . . . .	5
greedy . . . . .	5
greedyRes2Matrix . . . . .	6
loadContaminantsFasta . . . . .	6
loadContaminantsNoHumanFasta . . . . .	7

loadHomoSapiensSignalPeptides . . . . .	7
loadMusMusculusSignalPeptides . . . . .	8
makeID . . . . .	8
makeIDUnip . . . . .	9
pepprot . . . . .	9
prepareMatrix . . . . .	10
protpepmeta . . . . .	11
protpepmetashort . . . . .	11
prozor . . . . .	11
readPeptideFasta . . . . .	11
removeSignalPeptide . . . . .	12
reverseSeq . . . . .	12
writeFasta . . . . .	13

<b>Index</b>	<b>14</b>
--------------	-----------

---

annotateAHO	<i>annotate peptides using AhoCorasickTrie</i>
-------------	--

---

## Description

peptides which do not have protein assignment drop out

## Usage

```
annotateAHO(pepseq, fasta)
```

## Arguments

pepseq           - list of peptides - sequence, optional modified sequence, charge state.  
 fasta            - object as created by readPeptideFasta

## Examples

```
library(prozor)
file = file.path(path.package("prozor"), "extdata/shortfasta.fasta" )
fasta = readPeptideFasta(file = file)
#res = annotateVec(pepprot[1:20,"peptideSeq"],fasta)
system.time(res2 <- annotateAHO(pepprot[1:20,"peptideSeq"],fasta))
colnames(res2)
```

annotatePeptides      *Annotate peptides with protein ids*

**Description**

peptides which do not have protein assignment drop out

**Usage**

```
annotatePeptides(pepinfo, fasta, prefix = "([RK]|(^)|(M))", suffix = "")
```

**Arguments**

pepinfo      - list of peptides - sequence, optional modified sequence, charge state.  
 fasta      - object as created by readPeptideFasta  
 prefix      - default "([RK]|(^)|(M))"  
 suffix      - default ""

**Examples**

```
library(prozor)
data(pepprot)

file = file.path(path.package("prozor"), "extdata/shortfasta.fasta" )

fasta = readPeptideFasta(file = file)
res = annotatePeptides(pepprot[1:20,], fasta)
head(res)
res = annotatePeptides(pepprot[1:20,"peptideSeq"], fasta)
length(res)
```

annotateVec      *annotate vector of pepptide sequences against fasta file (Deprecated)*

**Description**

annotate vector of pepptide sequences against fasta file (Deprecated)

**Usage**

```
annotateVec(pepseq, fasta, digestPattern = "([RK]|(^)|(M))",
    mcCores = NULL)
```

**Arguments**

pepseq	peptide sequences
fasta	fasta file
digestPattern	digest pattern as regex
mcCores	nr of cores to use

**Examples**

```
library(prozor)
file = file.path(path.package("prozor"), "extdata/shortfasta.fasta" )
fasta = readPeptideFasta(file = file)

res = annotateVec(pepprot[1:20, "peptideSeq"], fasta)
head(res)
```

---

createDecoyDB	<i>create db with decoys and contaminants</i>
---------------	---

---

**Description**

create db with decoys and contaminants

**Usage**

```
createDecoyDB(dbs, useContaminants = TRUE, revLab = "REV_",
  annot = "zz|sourceOf|database")
```

**Arguments**

dbs	a path to a fasta file or an array of files
useContaminants	add fgcz contaminants
revLab	label for reversed peptides (if NULL do not generate decoys)
annot	source of database

**Examples**

```
file = file.path(path.package("prozor"), "extdata/uniprot_taxonomy_Oryctolagus_cuniculus.fasta.gz")
cont <- loadContaminantsFasta()
rabbit <- readPeptideFasta(file)
tmp <- 2*(2*length(rabbit)+length(cont)) + 1

res <- createDecoyDB(c(file, file))
length(res)
tmp
stopifnot(length(res) == tmp)
```

```

res <- createDecoyDB(c(file,file), revLab=NULL)
stopifnot(length(res) == (2*length(rabbit)+length(cont) + 1))
res <- createDecoyDB(c(file,file), revLab=NULL, useContaminants = FALSE)
stopifnot(length(res) == (2*length(rabbit) + 1) )

```

---

filterSequences	<i>Filter for specific residues</i>
-----------------	-------------------------------------

---

### Description

Will check if AA at Offset is a valid cleavage site

### Usage

```
filterSequences(matches, prefix = "([RK]|(^)|(^M))", suffix = "")
```

### Arguments

matches	must have 2 columns proteinSequunce and Offset
prefix	- regular expression describing the prefix of the peptide sequence e.g. (([RK] (^) (^M))
suffix	- regular expression describing the suffix of the peptide sequence

---

greedy	<i>given matrix (columns protein rows peptides), compute minimal protein set using greedy algorithm</i>
--------	---

---

### Description

given matrix (columns protein rows peptides), compute minimal protein set using greedy algorithm

### Usage

```
greedy(pepprot)
```

### Arguments

pepprot	matrix as returned by prepareMatrix
---------	-------------------------------------

### Value

list of peptide protein assignment

**Examples**

```

library(prozor)

data(protpepmetashort)
colnames(protpepmetashort)
dim(unique(protpepmetashort[,4:5]))
xx = prepareMatrix(protpepmetashort, weight= "one")
dim(xx)
stopifnot(dim(xx)[1] == dim(unique(protpepmetashort[,4:5]))[1])
es = greedy(as.matrix(xx))
stopifnot(length(unique(names(es))) == dim(unique(protpepmetashort[,4:5]))[1])

```

---

greedyRes2Matrix	<i>converts result of greedy function to a matrix with 3 columns - peptide - charge and protein</i>
------------------	---

---

**Description**

converts result of greedy function to a matrix with 3 columns - peptide - charge and protein

**Usage**

```
greedyRes2Matrix(res)
```

**Arguments**

res                    result of function prozor::greedy

**Value**

matrix of peptide protein assignments

---

loadContaminantsFasta	<i>load list of contaminant sequences</i>
-----------------------	---

---

**Description**

load list of contaminant sequences

**Usage**

```
loadContaminantsFasta()
```

### Examples

```
library(prozor)
cont <- loadContaminantsFasta()
cont[[1]]
#example how to create a protein db with decoy sequences
```

---

```
loadContaminantsNoHumanFasta
load list of contaminant without human sequences
```

---

### Description

load list of contaminant without human sequences

### Usage

```
loadContaminantsNoHumanFasta()
```

### Examples

```
library(prozor)
cont <- loadContaminantsNoHumanFasta()
cont[[1]]
#example how to create a protein db with decoy sequences
```

---

```
loadHomoSapiensSignalPeptides
load human signal peptides
```

---

### Description

load human signal peptides

### Usage

```
loadHomoSapiensSignalPeptides()
```

### Examples

```
library(prozor)
signal <- loadHomoSapiensSignalPeptides()
```

---

```
loadMusMusculusSignalPeptides
      load mus musculus signal peptides
```

---

**Description**

load mus musculus signal peptides

**Usage**

```
loadMusMusculusSignalPeptides()
```

**Examples**

```
library(prozor)
signal <- loadMusMusculusSignalPeptides()
head(signal)
```

---

```
makeID      make id for chain in format sp|P30443|1A01_HUMANs25
```

---

**Description**

make id for chain in format sp|P30443|1A01\_HUMANs25

**Usage**

```
makeID(sequence, id, sp)
```

**Arguments**

sequence	- aa sequence as string
id	uniprot id id: sp P30443 1A01_HUMAN
sp	start position of chain numeric

**Examples**

```
seq <- "MAVMAPRTL L L L L L S G A L A L T Q T W A G S H S M R Y F F T S V S R P G R \
G E P R F I A V G Y V D D T Q F V R F D S D A A S Q K M E P R A P W I E Q E G P E Y W D Q E T R N \
M K A H S Q T D R A N L G T L R G Y Y N Q S E D G S H T I Q I M Y G C D V G P D G R F L R G Y R Q \
D A Y D G K D Y I A L N E D L R S W T A A D M A A Q I T K R K W E A V H A A E Q R R V Y L E G R C \
V D G L R R Y L E N G K E T L Q R T D P P K T H M T H H P I S D H E A T L R C W A L G F Y P A E I \
T L T W Q R D G E D Q T Q D E L V E T R P A G D G T F Q K W A A V V V P S G E E Q R Y T C H V Q \
H E G L P K P L T L R W E L S S Q P T I P I V G I I A G L V L L G A V I T G A V V A A V M W R R K \
S S D R K G G S Y T Q A A S S D S A Q G S D V S L T A C K V "
```

```
nam <- "sp|P30443|1A01_HUMAN"
sp <- 24
makeID(seq, nam, sp)
```



---

makeIDUnip	<i>make id for chain compatible with uniprot</i>
------------	--

---

**Description**

make id for chain compatible with uniprot

**Usage**

```
makeIDUnip(sequence, id, sp)
```

**Arguments**

sequence	- aa sequence as string
id	uniprot id id: sp P30443 1A01_HUMAN
sp	start position of chain numeric

**Examples**

```
seq <- "MAVMAPRTL L L L L L S G A L A L T Q T W A G S H S M R Y F F T S V S R P G R \
G E P R F I A V G Y V D D T Q F V R F D S D A A S Q K M E P R A P W I E Q E G P E Y W D Q E T R N \
M K A H S Q T D R A N L G T L R G Y Y N Q S E D G S H T I Q I M Y G C D V G P D G R F L R G Y R Q \
D A Y D G K D Y I A L N E D L R S W T A A D M A A Q I T K R K W E A V H A A E Q R R V Y L E G R C \
V D G L R R Y L E N G K E T L Q R T D P P K T H M T H H P I S D H E A T L R C W A L G F Y P A E I \
T L T W Q R D G E D Q T Q D E L V E T R P A G D G T F Q K W A A V V V P S G E E Q R Y T C H V Q \
H E G L P K P L T L R W E L S S Q P T I P I V G I I A G L V L L G A V I T G A V V A A V M W R R K \
S S D R K G G S Y T Q A A S S D S A Q G S D V S L T A C K V"
nam <- "sp|P30443|1A01_HUMAN"
sp <- 24
makeIDUnip(seq, nam, sp)
```

---

pepprot	<i>Table containing peptide information</i>
---------	---

---

**Description**

Table containing peptide information

---

prepareMatrix                    *given table of peptide protein assignments generate matrix*

---

### Description

given table of peptide protein assignments generate matrix

### Usage

```
prepareMatrix(data, weighting = "one", sep = "|")
```

### Arguments

data	generated by annotatePeptides
weighting	weight type to use. Options are "one", "AA" - amino acids, "coverage" - coverage, "inverse" - inverse peptide frequencies
sep	separator for precursor (rownames)

### Value

sparse matrix

### Examples

```
library(prozor)
data(protpepmetashort)
library(Matrix)
colnames(protpepmetashort)
head(protpepmetashort)
dim(protpepmetashort)
count = prepareMatrix( protpepmetashort)
inverse = prepareMatrix( protpepmetashort, weight = "inverse")
aa = prepareMatrix(protpepmetashort, weight = "AA")
#xx = prepareMatrix(protpepmetashort, weight = "coverage")
image( as.matrix(count) )

corProt = cor( as.matrix(count) )
par(mfrow =c(1,2))
image(corProt)

#penalise peptides matching many proteins
corProtn = cor( as.matrix(aa) )
image(corProtn)
```

---

protpepmeta	<i>Generated from pepprot using method annotatePeptides</i>
-------------	---

---

**Description**

Generated from pepprot using method annotatePeptides

---

protpepmetashort	<i>Small version of pepprot dataset to speed up computation</i>
------------------	---

---

**Description**

Small version of pepprot dataset to speed up computation

---

prozor	<i>Minimal Protein set Explaining Peptides</i>
--------	--

---

**Description**

Minimal Protein set Explaining Peptides

---

readPeptideFasta	<i>wrapper setting the correct parameters</i>
------------------	---

---

**Description**

peptides which do not have protein assignment drop out

**Usage**

```
readPeptideFasta(file)
```

**Arguments**

file            - fasta file

**Examples**

```
library(seqinr)
library(prozor)
file = file.path(path.package("prozor"), "extdata/fgcz_contaminants_20150123.fasta")
fasta = readPeptideFasta(file)
```

removeSignalPeptide     *remove signal peptides from main chain*

---

**Description**

remove signal peptides from main chain

**Usage**

```
removeSignalPeptide(db, signal, idfun = makeID)
```

**Arguments**

db	uniprot fasta database as list
signal	tab delimited file with signals
idfun	function to generate id's

**Examples**

```
## Not run:  
library(prozor)  
  
hsfasta <- loadHomoSapiensFasta()  
hssignal <- loadHomoSapiensSignalPeptides()  
xx <- removeSignalPeptide(hsfasta, hssignal)  
db <- hsfasta  
signal <- hssignal  
  
## End(Not run)
```

---

reverseSeq     *create rev sequences to fasta list*

---

**Description**

peptides which do not have protein assignment drop out

**Usage**

```
reverseSeq(fasta, revLab = "REV_")
```

**Arguments**

fasta	- an r list with SeqFastaAA
revLab	- how to label reverse sequences, default = REV_

**Examples**

```

library(seqinr)
library(prozor)

file = file.path(path.package("prozor"), "extdata/fgcz_contaminants_20150123.fasta")
file
fasta = readPeptideFasta(file = file)
x <- reverseSeq(fasta)

revseq <- reverseSeq(fasta ,revLab = "REV_")
stopifnot(length(revseq) == length(fasta))
stopifnot(grep("^REV_", "REV_zz|ZZ_FGCZCont0000|")==1)

tmp <- list(as.SeqFastaAA(("DYKDDDDK"), name="Flag|FLAG|p2079", Annot=""))

reverseSeq(tmp)

```

---

writeFasta

*write fasta lists into file*


---

**Description**

peptides which do not have protein assignment drop out

**Usage**

```
writeFasta(file, ...)
```

**Arguments**

file	where to write
...	fasta list or single file

**Examples**

```

#example how to create a protein db with decoy sequences
library(seqinr)
library(prozor)
file = file.path(path.package("prozor"), "extdata/fgcz_contaminants_20150123.fasta")
fasta = readPeptideFasta(file = file)
revfasta <- reverseSeq(fasta)
decoyDB <- c(fasta, revfasta)
stopifnot(length(decoyDB) == 2 * length(fasta))
writeFasta(decoyDB, file="test.fasta")

```

# Index

## \*Topic **data**

- pepprot, [9](#)
- protpepmeta, [11](#)
- protpepmetashort, [11](#)

annotateAHO, [2](#)  
annotatePeptides, [3](#)  
annotateVec, [3](#)

createDecoyDB, [4](#)

filterSequences, [5](#)

greedy, [5](#)  
greedyRes2Matrix, [6](#)

loadContaminantsFasta, [6](#)  
loadContaminantsNoHumanFasta, [7](#)  
loadHomoSapiensSignalPeptides, [7](#)  
loadMusMusculusSignalPeptides, [8](#)

makeID, [8](#)  
makeIDUnip, [9](#)

pepprot, [9](#)  
prepareMatrix, [10](#)  
protpepmeta, [11](#)  
protpepmetashort, [11](#)  
prozor, [11](#)  
prozor-package (prozor), [11](#)

readPeptideFasta, [11](#)  
removeSignalPeptide, [12](#)  
reverseSeq, [12](#)

writeFasta, [13](#)