

# Package ‘canprot’

October 19, 2020

**Date** 2020-10-19

**Version** 1.1.0

**Title** Compositional Analysis of Differentially Expressed Proteins in Cancer

**Maintainer** Jeffrey Dick <j3ffdick@gmail.com>

**Depends** R (>= 3.1.0)

**Imports** xtable, MASS, rmarkdown

**Suggests** CHNOSZ (>= 1.3.0), knitr, testthat

**Description** Compositional analysis of differentially expressed proteins in cancer and cell culture proteomics experiments. The data include lists of up- and down-regulated proteins in different cancer types (breast, colorectal, liver, lung, pancreatic, prostate) and laboratory conditions (hypoxia, hyperosmotic stress, high glucose, 3D cell culture, and proteins secreted in hypoxia), together with amino acid compositions computed for protein sequences obtained from UniProt. Functions are provided to calculate compositional metrics including protein length, carbon oxidation state, and stoichiometric hydration state. In addition, phylostrata (evolutionary ages) of protein-coding genes are compiled using data from Liebeskind et al. (2016) <doi:10.1093/gbe/evw113> or Trigos et al. (2017) <doi:10.1073/pnas.1617743114>. The vignettes contain plots of compositional differences, phylostrata for human proteins, and references for all datasets.

**Encoding** UTF-8

**License** GPL (>= 2)

**BuildResaveData** no

**VignetteBuilder** knitr

**URL** <https://github.com/jedick/canprot>

**NeedsCompilation** no

**Author** Jeffrey Dick [aut, cre] (<<https://orcid.org/0000-0002-0687-5890>>),  
Ben Bolker [ctb] (<<https://orcid.org/0000-0002-2127-0443>>)

**Repository** CRAN

**Date/Publication** 2020-10-19 04:30:02 UTC

## R topics documented:

canprot-package . . . . .	2
check_IDs . . . . .	3
cleanup . . . . .	4
CLES . . . . .	6
diffplot . . . . .	8
get_comptab . . . . .	10
human . . . . .	12
metrics . . . . .	13
mkvig . . . . .	16
pdat_ . . . . .	17
protcomp . . . . .	19
PS . . . . .	20
qdist . . . . .	21
xsummary . . . . .	22
<b>Index</b>	<b>24</b>

---

canprot-package	<i>Compositional Analysis of Differentially Expressed Proteins</i>
-----------------	--

---

## Description

**canprot** is a package for analysis of the chemical compositions of proteins from their amino acid compositions. The package compiles datasets for differentially expressed proteins in cancer and cell culture conditions from over 250 studies.

## Overview

This package includes datasets for differential expression of proteins in six cancer types (breast, colorectal, liver, lung, pancreatic, prostate), and four cell culture conditions (hypoxia, hyperosmotic stress, secreted proteins in hypoxia, and 3D compared to 2D growth conditions). The hyperosmotic stress data are divided into bacteria, archaea (both high- and low-salt experiments) and eukaryotes; the latter are further divided into salt and glucose experiments. Nearly all datasets use UniProt IDs; if not given in the original publications they have been added using the UniProt mapping tool (<https://www.uniprot.org/mapping/>).

The vignettes have plots for each cancer type and cell culture condition and references for all data sources used. Because of their size, pre-built vignette HTML files are not included with the package; use `mkvig` to compile and view any of the vignettes.

The functions in this package were originally derived from code used for the papers of Dick (2016 and 2017). Updated data compilations and revised metrics, plots, and vignettes were developed for the papers of Dick et al. (2020) and Dick (2020).

## References

- Dick, J. M. (2016) Proteomic indicators of oxidation and hydration state in colorectal cancer. *PeerJ* **4**, e2238. doi: [10.7717/peerj.2238](https://doi.org/10.7717/peerj.2238)
- Dick, J. M. (2017) Chemical composition and the potential for proteomic transformation in cancer, hypoxia, and hyperosmotic stress. *PeerJ* **5**, e3421. doi: [10.7717/peerj.3421](https://doi.org/10.7717/peerj.3421)
- Dick, J. M., Yu, M. and Tan, J. (2020) Uncovering chemical signatures of salinity gradients through compositional analysis of protein sequences. *Biogeosciences Discussions*. doi: [10.5194/bg2020146](https://doi.org/10.5194/bg2020146)
- Dick, J. M. (2020) Water as a reactant in the differential expression of proteins in cancer. *bioRxiv*. doi: [10.1101/2020.04.09.035022](https://doi.org/10.1101/2020.04.09.035022)

## See Also

The **JMDplots** package on GitHub (<https://github.com/jedick/JMDplots>) has vignettes showing analysis of data from The Cancer Genome Atlas and The Human Protein Atlas, which are not included here because of the large size of the data files.

## Examples

```
# List the data files for all studies
# (one study can have more than one dataset)
exprdata <- system.file("extdata/expression", package="canprot")
datafiles <- dir(exprdata, recursive=TRUE)
print(datafiles)
# Show the number of data files for each condition
table(dirname(datafiles))
```

---

check\_IDs

*Check UniProt IDs*

---

## Description

Find the first ID for each protein that matches a known UniProt ID.

## Usage

```
check_IDs(dat, IDcol, aa_file = NULL, updates_file = NULL)
```

## Arguments

dat	data frame, protein expression data
IDcol	character, name of column that has the UniProt IDs
aa_file	character, name of file with additional amino acid compositions
updates_file	character, name of file with old to new ID mappings

## Details

check\_IDs is used to check for known UniProt IDs and to update obsolete IDs. The source IDs should be provided in the IDcol column of dat; multiple IDs for one protein can be separated by a semicolon.

The function keeps the first “known” ID for each protein, which must be present in one of these groups:

- The [human\\_aa](#) dataset of amino acid compositions.
- Old UniProt IDs that are mapped to new UniProt IDs in [uniprot\\_updates](#) or in updates\_file if specified.
- IDs of proteins in aa\_file, which lists amino acid compositions in the format described for [human\\_aa](#) (see extdata/protein/human\_extra.csv for an example and [thermo\\$protein](#) for more details).

## Value

dat is returned with possibly changed values in the column designated by IDcol; old IDs are replaced with new ones, the first known ID for each protein is kept, then proteins with no known IDs are assigned NA.

## See Also

This function is used by the [pdat\\_](#) functions, where it is called before [cleanup](#).

## Examples

```
# Make up some data for this example
ID <- c("P61247;PXXXXX", "PYYYYY;P46777;P60174", "PZZZZZ")
dat <- data.frame(ID = ID, stringsAsFactors = FALSE)
# Get the first known ID for each protein; the third one is NA
check_IDs(dat, "ID")

# Update an old ID
dat <- data.frame(Entry = "P50224", stringsAsFactors = FALSE)
check_IDs(dat, "Entry")
```

---

cleanup

*Clean Up Data*

---

## Description

Remove proteins with unavailable IDs, ambiguous expression ratios, and duplicated IDs.

## Usage

```
cleanup(dat, IDcol, up2 = NULL)
```

**Arguments**

dat	data frame, protein expression data
IDcol	character, name of column that has the UniProt IDs
up2	logical, TRUE for up-regulated proteins, FALSE for down-regulated proteins

**Details**

cleanup is used in the `pdat_` functions to clean up the dataset given in `dat`. `IDcol` is the name of the column that has the UniProt IDs, and `up2` indicates the expression change for each protein. The function removes proteins with unavailable (NA or "") or duplicated IDs. If `up2` is provided, the function also removes unquantified proteins (those that have NA values of `up2`) and those with ambiguous expression ratios (up and down for the same ID). For each operation, a message is printed describing the number of proteins that are 'unavailable', 'unquantified', 'ambiguous', or 'duplicated'.

Alternatively, if `IDcol` is a logical value, it selects proteins to be unconditionally removed.

**See Also**

This function is used extensively by the `pdat_` functions, where it is called after `check_IDs` (if needed).

**Examples**

```
# Set up a simple workflow
extdatadir <- system.file("extdata", package="canprot")
datadir <- paste0(extdatadir, "/expression/pancreatic/")
dataset <- "CYD+05"
dat <- read.csv(paste0(datadir, dataset, ".csv.xz"), as.is = TRUE)
up2 <- dat$Ratio..cancer.normal. > 1
# Remove two unavailable and one duplicated proteins
dat <- cleanup(dat, "Entry", up2)
# Now we can retrieve the amino acid compositions
pcomp <- protcomp(dat$Entry)

# Read another data file
datadir <- paste0(system.file("extdata", package="canprot"), "/expression/colorectal/")
dataset <- "STK+15"
dat <- read.csv(paste0(datadir, "STK+15.csv.xz"), as.is = TRUE)
# Remove unavailable proteins
dat <- cleanup(dat, "uniprot")
# Remove proteins that have less than 2-fold expression ratio
dat <- cleanup(dat, abs(log2(dat$invratio)) < 1)
```

---

CLES

*Common Language Effect Size*

---

### Description

Calculate the common language effect size.

### Usage

```
CLES(x, y, distribution = "normal")
```

### Arguments

x	numeric, data
y	numeric, data
distribution	'normal' to use probabilities calculated for a normal distribution, or NA for empirical probabilities

### Details

The common language statistic is defined for continuous data as “the probability that a score sampled at random from one distribution will be greater than a score sampled from some other distribution” (McGraw and Wong, 1992).

Given the default value of `distribution` ('normal'), this function uses `pnorm` to calculate the probability that a random sample from the unit normal distribution is greater than the Z score (i.e. the mean of 'y' minus the mean of 'x') / square root of (variance of 'x' plus variance of 'y').

If `distribution` is NA, this function calculates the empirical probability that the difference is positive, that is, the fraction of all possible pairings between elements of x and y where the difference ('y' value - 'x' value) is positive. It may not be possible to calculate the empirical probability for very large samples because of memory limits.

The examples use *simulated data for normal distributions*, given the sample size, mean, and standard deviation of datasets cited by McGraw and Wong, 1992. Therefore, the empirical probability in the examples approaches the normal curve probability. However, the empirical probability for *nonnormal* distributions is distinct from the normal curve probability, as discussed on p. 364-365 of McGraw and Wong, 1992.

### References

McGraw, Kenneth O. and Wong, S. P. (1992) A common language effect size statistic. *Psychological Bulletin* **11**, 361–365. doi: [10.1037/00332909.111.2.361](https://doi.org/10.1037/00332909.111.2.361)

National Center for Health Statistics (1987) *Anthropometric Reference Data and Prevalence of Overweight: United States, 1976-1980*. Data from the National Health Survey, Series 11, No. 238. DHHS Publication (PHS) No. 87-1688. U.S. Government Printing Office, Washington, DC. [https://www.cdc.gov/nchs/data/series/sr\\_11/sr11\\_238.pdf](https://www.cdc.gov/nchs/data/series/sr_11/sr11_238.pdf)

**Examples**

```

# Example 1: Height differences between males and females
# a) Use statistics quoted by McGraw and Wong, 1992 from NCHS, 1987
# for heights in inches of 18-24 year-old males and females
# Table 14: number, mean height, and standard deviation of height of females
n1 <- 1066
M1 <- 64.3
SD1 <- 2.8
# Table 13: number, mean height, and standard deviation of height of males
n2 <- 988
M2 <- 69.7
SD2 <- 2.6
# b) Simulate data from a normal distribution with exact mean and SD
# use rnorm2 function from Ben Bolker's answer to
# https://stackoverflow.com/questions/18919091/generate-random-numbers-with-fixed-mean-and-sd
rnorm2 <- function(n, mean, sd) { mean + sd * scale(rnorm(n)) }
set.seed(1234)
height_female <- rnorm2(n1, M1, SD1)
height_male <- rnorm2(n2, M2, SD2)
# c) Calculate the CLES using the normal distribution and empirical probability
CLES_normal <- CLES(height_female, height_male)
CLES_empirical <- CLES(height_female, height_male, distribution = NA)
# d) Test numerical equivalence of the results
# The CLES is approximately 0.92 (McGraw and Wong, 1992)
# (note: because we used rnorm2, this doesn't depend on the seed)
stopifnot(all.equal(CLES_normal, 0.92, tol = 0.01))
# With this seed, the difference between the normal curve probability
# and empirical probability is less than 1%
stopifnot(all.equal(CLES_normal, CLES_empirical, tol = 0.01))

# Example 1.5: Use multiple simulated datasets to show approach of
# of empirical probability to normal curve probability
CLES_empirical_n <- sapply(1:100, function(x) {
  height_female <- rnorm2(n1, M1, SD1)
  height_male <- rnorm2(n2, M2, SD2)
  CLES(height_female, height_male, distribution = NA)
})
CLES_empirical <- mean(CLES_empirical_n)
# now we're even closer to the normal curve probability
stopifnot(all.equal(CLES_normal, CLES_empirical, tol = 0.0001))

# Example 2: Multiple datasets in Table 2 of McGraw and Wong, 1992
# Sample statistics for females
n1 <- c(638, 672, 3139, 420740, 19274, 104263, 207, 394, 1066, 982, 108, 108)
M1 <- c(103, 15, 103, 18.9, 30, 16.1, 6.9, 13.3, 64.3, 134, 45, 94)
SD1 <- sqrt(c(908, 74, 219, 27, 110, 59, 15, 164, 6.8, 688, 310, 1971))
# Sample statistics for males
n2 <- c(354, 359, 3028, 356704, 21768, 133882, 199, 469, 988, 988, 443, 443)
M2 <- c(112, 23, 100, 17.9, 33, 18.6, 9.3, 21.8, 69.7, 163, 86, 212)
SD2 <- sqrt(c(1096, 96, 202, 29, 110, 61, 15, 133, 7.8, 784, 818, 5852))
# A function to calculate the effect size using simulated data
CLESfun <- function(n1, M1, SD1, n2, M2, SD2, distribution) {

```

```

rnorm2 <- function(n, mean, sd) { mean + sd * scale(rnorm(n)) }
set.seed(1234)
x <- rnorm2(n1, M1, SD1)
y <- rnorm2(n2, M2, SD2)
CLES(x, y, distribution)
}
# Calculate 100 * CL for the normal curve probabilities
CLnorm <- sapply(1:12, function(i) {
  CL <- CLESfun(n1[i], M1[i], SD1[i], n2[i], M2[i], SD2[i], "normal")
  round(100 * CL)
})
# Calculate 100 * CL for empirical probabilities
CLEmp <- sapply(1:12, function(i) {
  # skip very large samples: not enough memory
  if(n1[i] > 5000 | n2[i] > 5000) NA else {
    CL <- CLESfun(n1[i], M1[i], SD1[i], n2[i], M2[i], SD2[i], NA)
    round(100 * CL)
  }
})
# The difference between the empirical and normal curve
# probabilities is not more than 1 percent
stopifnot(max(abs(CLEmp - CLnorm), na.rm = TRUE) <= 1)
# TODO: Why are some of the calculated values different from
# Table 2 of McGraw and Wong, 1992?
CLref <- c(54, 74, 44, 45, 56, 63, 67, 65, 92, 78, 89, 91)
# Differences range from -4 to 4
range(CLnorm - CLref)
#stopifnot(max(abs(CLnorm - CLref)) == 0)

```

---

diffplot

*Plot Compositional Differences*


---

### Description

Make a plot showing differences of selected compositional metrics.

### Usage

```

diffplot(comptab, vars = c("ZC", "nH2O"), col = "black", plot.rect = FALSE,
  pt.text = c(letters, LETTERS), cex.text = 0.85, oldstyle = FALSE,
  pch = 1, cex = 2.1, contour = TRUE, col.contour = par("fg"),
  probs = 0.5, add = FALSE, labtext = NULL, ...)
cplab

```

### Arguments

comptab	list or data frame, compositional differences generated by <a href="#">get_comptab</a>
vars	character, which variables to plot
col	character or numeric, color(s) for the points



<code>plot.rect</code>	logical, plot a reference rectangle?
<code>pt.text</code>	character, text labels for the points
<code>cex.text</code>	numeric, size of text labels
<code>oldstyle</code>	logical, use old style plot?
<code>pch</code>	numeric, point symbol
<code>cex</code>	numeric, point size
<code>contour</code>	logical, add contour lines?
<code>col.contour</code>	character or numeric, color of contour lines
<code>probs</code>	numeric, probability level(s) for contours
<code>add</code>	logical, add to an existing plot?
<code>labtext</code>	character, text to add to axis labels
<code>...</code>	other arguments passed to <code>plot</code>

## Details

A plot is created with points showing the differences between up- and down-regulated proteins for two compositional metrics, as calculated by `get_comptab`. The default setting of `vars` refers to average oxidation state of carbon ( $Z_C$ ) as the x-variable and stoichiometric hydration state ( $n_{H_2O}$ ) as the y-variable.

The colors of the points are controlled by `col`, which is recycled to be equal to the number of comparisons in `comptab`.

If `plot.rect` is TRUE, a shaded `rectangle` is drawn with coordinates -0.01, -0.01, 0.01, 0.01. This is useful for visualizing the different scales of multi-panel plots.

If `pt.text` is not NA or FALSE, `text` labels are added with size controlled by `cex.text`. The default value produces labels that are taken sequentially from the 26 lowercase Roman letters in alphabetical order (`letters`), followed by the set of uppercase letters (`LETTERS`).

For `labtext = NULL`, descriptive text (“median difference” or “mean difference”) is added to the axis labels in parentheses. This text can be changed by giving a value in `labtext` (for both axes), two values (for each axis), or NA to suppress the text.

`cp1ab` is a list of formatted labels used by `diffplot`. It is an exported object, available to the user and other packages.

## Plot style

The overall style of the plot is controlled by `oldstyle`.

`oldstyle = FALSE` This is the current default style. Use `pch` and `cex` to control the point symbol and size. Contours are added for confidence regions of highest probability density, computed using a 2-D kernel density estimate (`kde2d`). `probs` gives the probability level(s) and `col.contour` sets the color(s) of the contour lines. `contour` can be a logical vector, indicating which points to include; set it to FALSE to omit the contour lines.

The code to calculate the contour levels is modified from `HPDregionplot` in the **emdbook** package by Ben Bolker (<https://cran.r-project.org/package=emdbook>).

oldstyle = TRUE This style was used for the historical (2017) vignettes, which have been moved to the 'extdata/cpcp' directory in **JMDplots** (<https://github.com/jedick/JMDplots>). For each dataset, the point symbol is a filled square if the  $p$ -values of both the x-variable and y-variable are less than 0.05, a filled circle if the  $p$ -value of one of the x- or y-variables is less than 0.05, and an open circle otherwise. A solid line is drawn from the point to the corresponding axis if the rounded, absolute value of (CLEs in percent - 50) of the x- or y-variable is greater than or equal 10. Otherwise, a dashed line is drawn from the point to the corresponding axis if the  $p$ -value of the x- or y-variable is less than 0.05. Otherwise, no line is drawn.

### See Also

[qdist](#) to plot quantile distributions for a single dataset.

### Examples

```
library(CHNOSZ)
# Make an old-style plot for two datasets
comptab <- lapply(c("JKMF10", "WDO+15_C.N"), function(dataset) {
  pdat <- pdat_colorectal(dataset)
  get_comptab(pdat, oldstyle = TRUE)
})
diffplot(comptab, oldstyle = TRUE)
```

---

get\_comptab

*Calculate Compositional Differences*

---

### Description

Compute differences of carbon oxidation state, stoichiometric hydration state and other compositional metrics between groups of up- and down-regulated proteins.

### Usage

```
get_comptab(pdat, var1 = "ZC", var2 = "nH2O", plot.it = FALSE,
  mfun = "median", oldstyle = FALSE, basis = getOption("basis"))
```

### Arguments

pdat	list, data object generated by a <a href="#">pdat_</a> function
var1	character, the first variable
var2	character, the second variable
plot.it	logical, make a scatterplot?
mfun	character, either 'median' or 'mean'
oldstyle	logical, also calculate CLEs and $p$ -values?
basis	character, keyword for basis species to use

## Details

The available variables are:

'ZC'	average oxidation state of carbon ( $Z_C$ ; see <a href="#">ZCAA</a> )
'nH2O'	stoichiometric hydration state per residue ( $n_{H_2O}$ ; see <a href="#">H2OAA</a> )
'nO2'	stoichiometric oxidation state per residue ( $n_{O_2}$ ; see <a href="#">O2AA</a> )
'V0'	standard molal volume per residue
'nAA'	protein length (number of amino acids)
'GRAVY'	grand average of hydropathicity (see <a href="#">GRAVY</a> )
'pI'	isoelectric point (see <a href="#">pI</a> )
'PS_TPPG17'	phylostratum (see <a href="#">PS</a> )
'PS_LMM16'	phylostratum (see <a href="#">PS</a> )
'MW'	molecular weight per residue

Differentially expressed proteins are identified by the value of `pdat$up2` (TRUE for up-regulated proteins and FALSE for down-regulated proteins). The differences are calculated as (median for up-regulated proteins) - (median for down-regulated proteins); if `mfun` is 'mean', means of the groups are used instead. If `oldstyle` is TRUE, the function also calculates the common language effect size ([CLES](#), in percent) and *p*-value for each variable.

The `basis` argument is used to select the basis species, which are used for the calculation of  $n_{H_2O}$  and  $n_{O_2}$ . The default for `getOption("basis")` is to use the 'QEC' basis species (see [metrics](#)).

Volume is calculated using amino acid group additivity as described by Dick et al. (2006).

Phylostrata are not compositional metrics, but are retrieved by matching UniProt accession numbers in a data file (see [PS](#)). Because phylostratum numbers are discrete values, mean values are calculated regardless of the value of `mfun`.

Set `plot.it` to TRUE to make a scatterplot. Open red squares and filled blue circles stand for up-regulated and down-regulated proteins, respectively.

## Value

A data frame is returned invisibly containing the columns 'dataset', 'description', 'n1' (number of down-regulated proteins), 'n2' (number of up-regulated proteins), followed two sets of columns for the variables. These are denoted generically as ('var.mfun1', 'var.mfun2', 'var.diff', 'var.CLES', 'var.p.value'), where 'var' is replaced by the name of var1 or var2, and 'mfun' is replaced by the value of `mfun`. For example, 'ZC.median1' and 'ZC.median2' are the median  $Z_C$  of the down- and up-regulated proteins, respectively.

## References

Dick, J. M., LaRowe, D. E. and Helgeson, H. C. (2006) Temperature, pressure, and electrochemical constraints on protein speciation: Group additivity calculation of the standard molal thermodynamic properties of ionized unfolded proteins. *Biogeosciences* **3**, 311–336. <https://doi.org/10.5194/bg-3-311-2006>

## Examples

```
pd <- pdat_colorectal("JKMF10")
# default variables: ZC and nH2O
get_comptab(pd, plot.it = TRUE)
# protein length and per-residue volume
get_comptab(pd, "nAA", "V0", plot.it = TRUE)
```

---

human

*Amino Acid Compositions of Human Proteins*

---

## Description

Data for amino acid compositions of proteins and conversion from old to new UniProt IDs.

## Format

human\_aa is a data frame with 25 columns in the format used for amino acid compositions in **CHNOSZ** (see [thermo](#)):

protein	character	Identification of protein
organism	character	Identification of organism
ref	character	Reference key for source of sequence data
abbrv	character	Abbreviation or other ID for protein (e.g. gene name)
chains	numeric	Number of polypeptide chains in the protein
Ala...Tyr	numeric	Number of each amino acid in the protein

The protein column contains UniProt IDs in the format database|accession-isoform, where database is most often 'sp' (Swiss-Prot) or 'tr' (TrEMBL), and isoform is an optional suffix indicating the isoform of the protein (particularly in the human\_additional file).

## Details

The amino acid compositions of human proteins are stored in three files under extdata/protein.

- human\_base.rds contains amino acid compositions of canonical isoforms of manually reviewed proteins in the **UniProt** reference human proteome (computed from sequences in UP000005640\_9606.fasta.gz, dated 2016-04-03).
- human\_additional.rds contains amino acid compositions of additional proteins (UP000005640\_9606\_additional.fasta.gz) including isoforms and unreviewed sequences. In version 0.1.5, this file was trimmed to include only those proteins that are used in any of the datasets in the package.
- human\_extra.csv contains amino acid compositions of other ("extra") proteins used in a dataset but not listed in one of the files above. These proteins may include obsolete, unreviewed, or newer additions to the UniProt database. Most, but not all, sequences here are HUMAN (see the organism column and the ref column for the reference keys).

On loading the package, the individual data files are read and combined, and the result is assigned to the `human_aa` object in the human environment.

As an aid for processing datasets that list old (obsolete) UniProt IDs, the corresponding new (current) IDs are stored in `uniprot_updates`. These ID mappings have been manually added as needed for individual datasets, and include proteins from humans as well as other organisms. `check_IDs` performs the conversion of old to new IDs.

### See Also

Amino acid compositions of non-human proteins are stored under `extdata/aa` in directories `archaea`, `bacteria`, `cow`, `dog`, `mouse`, `rat`, and `yeast`. These files can be loaded in `protcomp` via the `aa_file` argument, which is used e.g. in `pdat_osmotic_bact`.

### Examples

```
# The number of proteins
nrow(get("human_aa", human))
# The number of old to new ID mappings
nrow(get("uniprot_updates", human))
```

---

metrics

*Calculate Compositional Metrics for Proteins*

---

### Description

These functions calculate compositional metrics of proteins given a data frame of amino acid compositions.

### Usage

```
ZCAA(AAcomp, nothing = NULL)
H2OAA(AAcomp, basis = getOption("basis"))
O2AA(AAcomp, basis = getOption("basis"))
GRAVY(AAcomp)
pI(AAcomp)
MWA(AAcomp)
basis.text(basis)
```

### Arguments

<code>AAcomp</code>	data frame, amino acid compositions
<code>nothing</code>	dummy argument
<code>basis</code>	character, basis species

## Details

Columns in AAcomp should be named with the three-letter abbreviations for the amino acids (e.g. 'Ala'). The metrics are described below:

ZCAA Average oxidation state of carbon ( $Z_C$ ) (Dick, 2014). nothing is an extra argument that does nothing. It is provided so that `do.call` can be used to run ZCAA or H2OAA with the same number of arguments.

This metric is independent of the choice of basis species.

H2OAA Stoichiometric hydration state ( $n_{H_2O}$ ) per residue. The available basis species are:

- 'QEC' - glutamine, glutamic acid, cysteine,  $H_2O$ ,  $O_2$  (Dick et al., 2020) (this is the default for `getOption("basis")`)
- 'QCa' - glutamine, cysteine, acetic acid,  $H_2O$ ,  $O_2$
- Any other valid basis specification for `basis`, such as 'CHNOS' for  $CO_2$ ,  $NH_3$ ,  $H_2S$ ,  $H_2O$ , and  $O_2$

O2AA Stoichiometric oxidation state ( $n_{O_2}$ ) per residue. The basis species also affect this calculation.

GRAVY Grand average of hydrophobicity. Values of the hydrophobicity index for individual amino acids are from Kyte and Doolittle (1982).

pI Isoelectric point. The net charge for each ionizable group was pre-calculated from pH 0 to 14 at intervals of 0.01. The isoelectric point is found as the pH where the sum of charges of all groups in the protein is closest to zero. The  $pK$  values for the terminal groups and sidechains are taken from Bjellqvist et al. (1993) and Bjellqvist et al. (1994); note that the calculation does not implement position-specific adjustments described in the latter paper. The number of N- and C-terminal groups is taken to be one, unless a value for `chains` (number of polypeptide chains) is given in AAcomp.

MWAA Molecular weight per residue.

Note that  $Z_C$  is a per-carbon average, but  $n_{H_2O}$  is a per-residue average. The contribution of  $H_2O$  from the terminal groups of proteins is counted, so shorter proteins have slightly greater  $n_{H_2O}$ .

Tests for a few proteins (see examples) indicate that GRAVY and pI are equal those calculated with the ProtParam tool (<https://web.expasy.org/protparam/>; Gasteiger et al., 2005).

`basis.text` is used in the vignettes to generate a textual description of the names of the basis species, except  $H_2O$  and  $O_2$ , for one of the keywords 'QEC' or 'QCa'.

## References

- Bjellqvist, B., Hughes, G. J., Pasquali, C., Paquet, N., Ravier, F., Sanchez, J.-C., Frutiger, S. and Hochstrasser, D. (1993) The focusing positions of polypeptides in immobilized pH gradients can be predicted from their amino acid sequences. *Electrophoresis* **14**, 1023–1031. doi: [10.1002/elps.11501401163](https://doi.org/10.1002/elps.11501401163)
- Bjellqvist, B. and Basse, B. and Olsen, E. and Celis, J. E. (1994) Reference points for comparisons of two-dimensional maps of proteins from different human cell types defined in a pH scale where isoelectric points correlate with polypeptide compositions. *Electrophoresis* **15**, 529–539. doi: [10.1002/elps.1150150171](https://doi.org/10.1002/elps.1150150171)
- Dick, J. M. (2014) Average oxidation state of carbon in proteins. *J. R. Soc. Interface* **11**, 20131095. doi: [10.1098/rsif.2013.1095](https://doi.org/10.1098/rsif.2013.1095)

Dick, J. M., Yu, M. and Tan, J. (2020) Uncovering chemical signatures of salinity gradients through compositional analysis of protein sequences. *Biogeosciences Discussions*. doi: [10.5194/bg2020146](https://doi.org/10.5194/bg2020146)

Gasteiger, E., Hoogland, C., Gattiker, A., Duvaud, S., Wilkins, M. R., Appel, R. D. and Bairoch, A. (2005) Protein identification and analysis tools on the ExPASy server. In J. M. Walker (Ed.), *The Proteomics Protocols Handbook* (pp. 571–607). Totowa, NJ: Humana Press Inc. doi: [10.1385/1592598900:571](https://doi.org/10.1385/1592598900:571)

Kyte, J. and Doolittle, R. F. (1982) A simple method for displaying the hydropathic character of a protein. *J. Mol. Biol.* **157**, 105–132. doi: [10.1016/00222836\(82\)905150](https://doi.org/10.1016/00222836(82)905150)

### See Also

Use [PS](#) to get phylostrata of proteins (based on matching IDs in a data file, not amino acid composition). For calculation of  $Z_C$  from a chemical formula instead of amino acid composition, see the [ZC](#) function in **CHNOSZ**.

### Examples

```
# we need CHNOSZ for these examples
require(CHNOSZ)

# for reference, compute ZC of alanine and glycine "by hand"
ZC.Gly <- ZC("C2H5NO2")
ZC.Ala <- ZC("C3H7NO2")
# define the composition of a Gly-Ala-Gly tripeptide
AAcomp <- data.frame(Gly = 2, Ala = 1)
# calculate the ZC of the tripeptide (value: 0.571)
ZC.GAG <- ZCAA(AAcomp)
# this is equal to the carbon-number-weighted average of the amino acids
nC.Gly <- 2 * 2
nC.Ala <- 1 * 3
ZC.average <- (nC.Gly * ZC.Gly + nC.Ala * ZC.Ala) / (nC.Ala + nC.Gly)
stopifnot(all.equal(ZC.GAG, ZC.average))

# compute the per-residue nH2O of Gly-Ala-Gly
basis("QEC")
nH2O.GAG <- species("Gly-Ala-Gly")$H2O
# divide by the length to get residue average (we keep the terminal H-OH)
nH2O.residue <- nH2O.GAG / 3
# compare with the value calculated by H2OAA() (-0.2)
nH2O.H2OAA <- H2OAA(AAcomp, "QEC")
stopifnot(all.equal(nH2O.residue, nH2O.H2OAA))

# calculate GRAVY for a few proteins
# first get the protein index in CHNOSZ's list of proteins
iprotein <- pinfo(c("LYSC_CHICK", "RNAS1_BOVIN", "AMYA_PYRFU"))
# then get the amino acid compositions
AAcomp <- pinfo(iprotein)
# then calculate GRAVY
Gcalc <- as.numeric(GRAVY(AAcomp))
# these are equal to values obtained with ProtParam on uniprot.org
# https://web.expasy.org/cgi-bin/protparam/protparam1?P00698@19-147@
```

```

# https://web.expasy.org/cgi-bin/protparam/protparam1?P61823@27-150@
# https://web.expasy.org/cgi-bin/protparam/protparam1?P49067@2-649@
Gref <- c(-0.472, -0.663, -0.325)
stopifnot(all.equal(round(Gcalc, 3), Gref))
# also calculate molecular weight of the proteins
MWcalc <- as.numeric(MWAA(AAcomp)) * protein.length(iprotein)
MWref <- c(14313.14, 13690.29, 76178.25)
stopifnot(all.equal(round(MWcalc, 2), MWref))

# calculate pI for a few proteins
iprotein <- pinfo(c("LYSC_CHICK", "RNAS1_BOVIN", "AMYA_PYRFU", "CSG_HALJP"))
AAcomp <- pinfo(iprotein)
pI_calc <- pI(AAcomp)
# reference values calculated with ProtParam on uniprot.org
# LYSC_CHICK: residues 19-147 (sequence v1)
# RNAS1_BOVIN: residues 27-150 (sequence v1)
# AMYA_PYRFU: residues 2-649 (sequence v2)
# CSG_HALJP: residues 35-862 (sequence v1)
pI_ref <- c(9.32, 8.64, 5.46, 3.37)
stopifnot(all.equal(as.numeric(pI_calc), pI_ref))

```

---

mkvig

*Compile and View Vignettes from the Command Line*


---

## Description

Compile the indicated vignette and open it in the browser.

## Usage

```
mkvig(vig = NULL)
```

## Arguments

`vig` character, name of a vignette without `‘.Rmd’` extension

## Details

In order to reduce package space and check time, pre-built vignettes are not included in the package. `mkvig` is a convenience function to compile the vignettes on demand and view them in a browser.

The available vignettes for `mkvig` are listed here:

- *Cell culture* – ‘hypoxia’, ‘secreted’, ‘osmotic\_bact’, ‘osmotic\_euk’, ‘osmotic\_halo’, ‘glucose’, ‘3D’
- *Cancer* – ‘breast’, ‘colorectal’, ‘liver’, ‘lung’, ‘pancreatic’, ‘prostate’

Note that `pandoc` (including `pandoc-citeproc`), as a system dependency of **rmarkdown**, is required to build the vignettes.



## See Also

The vignettes can also be run using e.g. `demo("glucose")`, and through the interactive help system (`help.start > Packages > canprot > Code demos`).

## Examples

```
## Not run:
mkvig("colorectal")

## End(Not run)
```

---

pdat\_ *Get Protein Expression Data*

---

## Description

Get differentially expressed proteins and amino acid compositions.

## Usage

```
pdat_breast(dataset = 2020)
pdat_colorectal(dataset = 2020)
pdat_liver(dataset = 2020)
pdat_lung(dataset = 2020)
pdat_pancreatic(dataset = 2020)
pdat_prostate(dataset = 2020)
pdat_hypoxia(dataset = 2020)
pdat_secreted(dataset = 2020)
pdat_3D(dataset = 2020)
pdat_glucose(dataset = 2020)
pdat_osmotic_bact(dataset = 2020)
pdat_osmotic_euk(dataset = 2020)
pdat_osmotic_halo(dataset = 2020)
.pdat_multi(dataset = 2020)
.pdat_osmotic(dataset = 2017)
```

## Arguments

dataset            character, dataset name

## Details

The `pdat_` functions assemble lists of up- and down-regulated proteins and retrieve their amino acid compositions using `protcomp`. The result can be used with `get_comptab` to make a table of compositional metrics that can then be plotted with `diffplot`.

If `dataset` is '2020' (the default) or '2017', the function returns the names of all datasets in the compilation for the respective year.

Each dataset name starts with a reference key indicating the study (i.e. paper or other publication) where the data were reported. The reference keys are made by combining the first characters of the authors' family names with the 2-digit year of publication.

If a study has more than one dataset, the reference key is followed by an underscore and an identifier for the particular dataset. This identifier is saved in the variable named `stage` in the functions, but can be any descriptive text.

To retrieve the data, provide a single dataset name in the `dataset` argument. Protein expression data is read from the CSV files stored in `extdata/expression/`, under the subdirectory corresponding to the name of the `pdat_` function. Some of the functions also read amino acid compositions (e.g. for non-human proteins) from the files in `extdata/aa/`.

Descriptions for each function:

- `pdat_colorectal`, `pdat_pancreatic`, `pdat_breast`, `pdat_lung`, `pdat_prostate`, and `pdat_liver` retrieve data for protein expression in different cancer types.
- `pdat_hypoxia` gets data for cellular extracts in hypoxia and `pdat_secreted` gets data for secreted proteins (e.g. exosomes) in hypoxia.
- `pdat_3D` retrieves data for 3D (e.g. tumor spheroids and aggregates) compared to 2D (monolayer) cell culture.
- `.pdat_osmotic` retrieves data for hyperosmotic stress, for the 2017 compilation only. In 2020, this compilation was expanded and split into `pdat_osmotic_bact` (bacteria), `pdat_osmotic_euk` (eukaryotic cells) and `pdat_osmotic_halo` (halophilic bacteria and archaea).
- `pdat_glucose` gets data for high-glucose experiments in eukaryotic cells.
- `.pdat_multi` retrieves data for studies that have multiple types of datasets (e.g. both cellular and secreted proteins in hypoxia), and is used internally by the specific functions (e.g. `pdat_hypoxia` and `pdat_secreted`).

## Value

A list consisting of:

`dataset` Name of the dataset

`description` Descriptive text for the dataset, used for making the tables in the vignettes (see [mkvig](#))

`pcomp` UniProt IDs together with amino acid compositions obtained using [protcomp](#)

`up2` Logical vector with length equal to the number of proteins; TRUE for up-regulated proteins and FALSE for down-regulated proteins

## Examples

```
# List datasets in the 2017 compilation for colorectal cancer
pdat_colorectal(2017)
# Get proteins and amino acid compositions for one dataset
pdat_colorectal("JKMF10")
```

---

protcomp                      *Protein Compositions*

---

## Description

Get amino acid compositions of proteins.

## Usage

```
protcomp(uniprot = NULL, aa = NULL, aa_file = NULL)
```

## Arguments

uniprot	character, UniProt IDs of proteins
aa	data frame, amino acid compositions
aa_file	character, file name

## Details

This function retrieves the amino acid compositions of one or more proteins specified by uniprot.

This function depends on the amino acid compositions of human proteins, which are stored in the [human](#) environment when the package is attached. If aa\_file is specified, additional amino acid compositions to be considered are read from this file, which should be in the same format as e.g. [human\\_extra.csv](#) (see also `thermo$protein`). Alternatively, the amino acid compositions can be given in aa, bypassing the search step.

## Value

The function returns a list with elements uniprot (UniProt IDs as given in the arguments) and aa (amino acid compositions of the proteins).

## See Also

[cleanup](#)

## Examples

```
protcomp("P24298")
```

PS

*Retrieve Phylostrata for Given UniProt IDs***Description**

Retrieves the phylostrata for protein-coding genes according to Liebeskind et al. (2016) or Trigos et al. (2017).

**Usage**

```
PS(uniprot, source = "TPPG17")
```

**Arguments**

uniprot	character, UniProt accession numbers
source	character, 'TPPG17' or 'LMM16'

**Details**

The phylostratum for each protein is found by matching the UniProt ID in one of these data files:

extdata/phylostrata/TPPG17.csv.xz This file has columns 'GeneID' (gene name), 'Entrez', 'Entry', and 'Phylostrata'. Except for 'Entry', the values are from Dataset S1 of Trigos et al. (2017). UniProt accession numbers in 'Entry' were generated using the UniProt mapping tool first for 'Entrez', followed by 'GeneID' for the unmatched genes. 'Entry' is NA for genes that remain unmatched to any proteins after both mapping steps.

extdata/phylostrata/LMM16.csv.xz This file has columns 'UniProt', 'modeAge', and 'PS'. The data are from file main\_HUMAN.csv in Gene-Ages v1.0 (<https://zenodo.org/record/51708>; Liebeskind et al. (2016)). The modeAges were converted to phylostrata values 1-8 ('PS' column) in this order: Cellular\_organisms, Euk\_Archaea, Euk+Bac, Eukaryota, Opisthokonta, Eumetazoa, Vertebrata, Mammalia.

**References**

- Trigos, A. S. and Pearson, R. B. and Papenfuss, A. T. and Goode, D. L. (2017) Altered interactions between unicellular and multicellular genes drive hallmarks of transformation in a diverse range of solid tumors. *Proc. Natl. Acad. Sci.* **114**, 6406–6411. doi: [10.1073/pnas.1617743114](https://doi.org/10.1073/pnas.1617743114)
- Liebeskind, B. J. and McWhite, C. D. and Marcotte, E. M. (2016) Towards consensus gene ages. *Genome Biol. Evol.* **8**, 1812–1823. doi: [10.1093/gbe/evw113](https://doi.org/10.1093/gbe/evw113)

**See Also**

Call `get_comptab` with 'PS\_TPPG17' or 'PS\_LMM16' as a variable name to calculate mean differences of phylostrata for differentially expressed proteins.

## Examples

```
# Get protein expression data for one dataset
pd <- pdat_colorectal("JKMF10")
opar <- par(mfrow = c(2, 1))
# Plot nH2O vs phylostrata with Trigos et al. data
get_comptab(pd, "PS_TPPG17", plot.it = TRUE)
# Plot nH2O vs phylostrata with Liebeskind et al. data
get_comptab(pd, "PS_LMM16", plot.it = TRUE)
par(opar)

# compare the two sources
PSdir <- system.file("extdata/phylostrata", package = "canprot")
TPPG17 <- read.csv(file.path(PSdir, "TPPG17.csv.xz"))
LMM16 <- read.csv(file.path(PSdir, "LMM16.csv.xz"))
IDs <- intersect(TPPG17$Entry, LMM16$UniProt)
PS_TPPG17 <- TPPG17$Phylostrata[match(IDs, TPPG17$Entry)]
PS_LMM16 <- LMM16$PS[match(IDs, LMM16$UniProt)]
plot(jitter(PS_TPPG17), jitter(PS_LMM16), pch = ".")
```

qdist

*Quantile Distributions for One Dataset*

## Description

Make a plot showing quantile distributions for up- and down-regulated proteins.

## Usage

```
qdist(pdat, vars = c("ZC", "nH2O"), show.steps = FALSE)
```

## Arguments

pdat	list, output of a <code>pdat_</code> function for a single dataset
vars	character, which variables to plot
show.steps	logical, show the steps using <code>plot.ecdf</code> ?

## Details

This function makes a quantile distribution plot with lines for both up- and down-regulated proteins. The variable (var) can be 'ZC', 'H2O', or both (two plots are made for the latter). The horizontal axis is the variable and the vertical axis is the quantile point. A solid black line is drawn for the down-regulated proteins, and a dashed red line for the up-regulated proteins. The median difference is shown by a gray horizontal line drawn between the distributions at the 0.5 quantile point.

## References

Jimenez, C. R. and Knol, J. C. and Meijer, G. A. and Fijneman, R. J. A. (2010) Proteomics of colorectal cancer: Overview of discovery studies and identification of commonly identified cancer-associated proteins and candidate CRC serum markers. *J. Proteomics* **73**, 1873–1895. doi: [10.1016/j.jprot.2010.06.004](https://doi.org/10.1016/j.jprot.2010.06.004)

**See Also**

[diffplot](#) to plot median differences for multiple datasets.

**Examples**

```
# Plot the data of Jimenez et al., 2020 for colorectal cancer
pdat <- pdat_colorectal("JKMF10")
qdist()
```

---

xsummary

*Summarize Compositional Differences*


---

**Description**

Make an HTML table summarizing compositional differences.

**Usage**

```
xsummary(comptab, vars=c("ZC", "nH2O"))
xsummary2(comptab1, comptab2, comptab3)
xsummary3(comptab1, comptab2, comptab3)
```

**Arguments**

comptab	list or data frame, summary of comparisons generated by <a href="#">get_comptab</a>
vars	character, two variables to tabulate
comptab1	list, output of <a href="#">get_comptab</a>
comptab2	list, output of <a href="#">get_comptab</a>
comptab3	list, output of <a href="#">get_comptab</a>

**Details**

xsummary makes an HTML table (using [xtable](#)) and adds bold and underline formatting to highlight significant compositional differences. The  $p$ -value is bolded if it is less than 0.05, and the percent common language effect size (CLES) is bolded if it is  $\leq 40$  or  $\geq 60$ . The mean (or median) difference is [underlined / bolded] if [only one of / both] the  $p$ -value and CLES pass these cutoffs.

The generated table is written to the console, and can be used in a vignette using the `results = "asis"` chunk option. The function also returns (invisibly) the data frame used to make the table; this data frame differs from `comptab` by having row names added (alphabetical one-letter IDs for the datasets).

`xsummary2` is an updated version that is used in the current vignettes in the package. It shows negative numbers in bold ( $p$ -value and CLES are not shown). `xsummary3` is a further revision that shows GRAVY and pI instead of phylostrata; it is used in the ‘osmotic\_bact’ and ‘osmotic\_halo’ vignettes.

**Examples**

```
comptab <- lapply(c("JKMF10", "WDO+15_C.N"), function(dataset) {  
  pdat <- pdat_colorectal(dataset)  
  get_comptab(pdat, oldstyle = TRUE)  
})  
xsummary(comptab)
```

# Index

## \* Amino acid composition

human, 12

## \* Chemical composition

get\_comptab, 10

metrics, 13

protcomp, 19

## \* Compositional plots

diffplot, 8

qdist, 21

## \* Protein data

check\_IDs, 3

cleanup, 4

pdat\_, 17

PS, 20

## \* Statistics

CLES, 6

## \* Vignette utilities

mkvig, 16

xsummary, 22

## \* package

canprot-package, 2

.pdat\_multi (pdat\_), 17

.pdat\_osmotic (pdat\_), 17

basis, 14

basis.text (metrics), 13

canprot-package, 2

check\_IDs, 3, 5, 13

cleanup, 4, 4, 19

CLES, 6, 10, 11, 22

cplab (diffplot), 8

demo, 17

diffplot, 8, 17, 22

do.call, 14

get\_comptab, 8, 9, 10, 17, 20, 22

GRAVY, 11

GRAVY (metrics), 13

H2OAA, 11

H2OAA (metrics), 13

help.start, 17

human, 12, 19

human\_aa, 4

human\_aa (human), 12

human\_additional (human), 12

human\_base (human), 12

human\_extra, 19

human\_extra (human), 12

kde2d, 9

LETTERS, 9

letters, 9

metrics, 11, 13

mkvig, 2, 16, 18

MWAA (metrics), 13

O2AA, 11

O2AA (metrics), 13

pdat\_, 4, 5, 10, 17, 21

pdat\_3D (pdat\_), 17

pdat\_breast (pdat\_), 17

pdat\_colorectal (pdat\_), 17

pdat\_glucose (pdat\_), 17

pdat\_hypoxia (pdat\_), 17

pdat\_liver (pdat\_), 17

pdat\_lung (pdat\_), 17

pdat\_osmotic\_bact, 13

pdat\_osmotic\_bact (pdat\_), 17

pdat\_osmotic\_euk (pdat\_), 17

pdat\_osmotic\_halo (pdat\_), 17

pdat\_pancreatic (pdat\_), 17

pdat\_prostate (pdat\_), 17

pdat\_secreted (pdat\_), 17

pI, 11

pI (metrics), 13

plot, 9



plot.ecdf, [21](#)  
pnorm, [6](#)  
protcomp, [13](#), [17](#), [18](#), [19](#)  
PS, [11](#), [15](#), [20](#)  
  
qdist, [10](#), [21](#)  
  
rect, [9](#)  
  
text, [9](#)  
thermo, [4](#), [12](#), [19](#)  
  
uniprot\_updates, [4](#)  
uniprot\_updates (human), [12](#)  
  
xsummary, [22](#)  
xsummary2 (xsummary), [22](#)  
xsummary3 (xsummary), [22](#)  
xtable, [22](#)  
  
ZC, [15](#)  
ZCAA, [11](#)  
ZCAA (metrics), [13](#)