

Package ‘PCAWG7’

January 25, 2022

Title Repository of Data from paper 'The repertoire of Mutational Signatures in Human Cancer'

Version 0.1.3

Description Contains data from the paper by Alexandrov, Kim, Haradhvala, Huang et al., 'The repertoire of Mutational Signatures in Human Cancer' <[doi:10.1038/s41586-020-1943-3](https://doi.org/10.1038/s41586-020-1943-3)>. Please see ?PCAWG7. The funny name comes from the fact that this paper was generated by Working Group 7 of the Pan Cancer Analysis of Whole Genomes (PCAWG) consortium. The signature profiles were later placed on the COSMIC web site and have been subsequently updated.

License GPL-3

Language en-US

Encoding UTF-8

LazyData true

LazyDataCompression bzip2

Depends R (>= 3.5),

RoxygenNote 7.1.2

URL <https://github.com/steverozen/PCAWG7>

BugReports <https://github.com/steverozen/PCAWG7/issues>

Imports ICAMS

Suggests usethis, testthat (>= 3.0.0)

Config/testthat/edition 3

NeedsCompilation no

Author Steven G. Rozen [aut, cre] (<<https://orcid.org/0000-0002-4288-0056>>),
Nanhai Jiang [aut] (<<https://orcid.org/0000-0003-4974-2753>>)

Maintainer Steven G. Rozen <steverozen@pm.me>

R topics documented:

CancerTypes	2
exposure	2
exposure.stats	3
map_aliquot_ID_to_SP_ID	4

map_SP_ID_to_tumor_type	4
PCAWG.sample.id	5
PCAWG.sample.sheet	6
PCAWG7	6
SampleIDToCancerType	7
spectra	8
SplitMatrixBySampleType	8
SplitPCAWGMatrixByTumorType	9
TCGA_ID_to_ICGC_ID	10
Index	11

CancerTypes	<i>Return a character vector of some common cancer types</i>
-------------	--

Description

Return a character vector of some common cancer types

Usage

CancerTypes()

Examples

CancerTypes()[1:5]

exposure	<i>PCAWG7 SigProfiler signature assignments (numbers of mutations due to each signature in each tumor)</i>
----------	--

Description

PCAWG7 SigProfiler signature assignments (numbers of mutations due to each signature in each tumor)

Usage

exposure

Format

A list with the elements:

PCAWG A list with the elements:

- SBS96** Strand-agnostic single-base substitutions in trinucleotide context.
- DBS78** Strand-agnostic doublet-base substitutions.
- ID** Strand-agnostic indels. These are signature assignments for the PCAWG platinum genomes.

TCGA A list with the elements:

- SBS96** As above.

ID As above. These are signature assignments for the TCGA exomes.

other.genome A list with the element:

SBS96 As above. This contains signature assignments for non-TCGA genomes.

other.exome A list with the element:

SBS96 As above. This contains signature assignments for non-TCGA exomes.

Source

Files of <https://www.synapse.org/#!Synapse:syn12009743>, 2019 Oct 09,
populated by data-raw/sig.profiler.sures/load.package.variable.exposure.R.

Examples

```
SBS96.exposure <- exposure$PCAWG$SBS96
```

exposure.stats	<i>Exposure statistics from the PCAWG7 paper</i>
----------------	--

Description

Exposure statistics from the PCAWG7 paper

Usage

```
exposure.stats
```

Format

A list with one element, PCAWG, which has the sub-elements SBS96, DBS78, ID with statistics for the corresponding mutation types by cancer type. I.e. each element has a sub-element for each cancer type, and this element is a data.frame with one row for each signature and columns mean.of.those.present (the mean number of mutations for those tumors that have the mutation) and proportion.present (the proportion of tumors in which the signature is present).

Source

Computed from other package variables using GatherPCAWG7ExposureStatsSBS96.

Examples

```
exposure.stats$PCAWG$SBS96$`Biliary-AdenoCA`[1:3, ]
```

```
map_aliquot_ID_to_SP_ID
```

Translate aliquot IDs (e.g. e0fccaf5-925a-41f9-b87c-cd5ee4aecb59) to "SP" IDs (e.g. SP1682)

Description

Translate aliquot IDs (e.g. e0fccaf5-925a-41f9-b87c-cd5ee4aecb59) to "SP" IDs (e.g. SP1682)

Usage

```
map_aliquot_ID_to_SP_ID(aliquot.ids)
```

Arguments

`aliquot.ids` Character vector of aliquot IDs.

Details

If there are aliquot IDs that cannot be matched to any "SP" IDs, return NA with a warning.

Value

Character vector of corresponding "SP" IDs. If a corresponding aliquot ID cannot be found, then return NA with a warning.

Note

This function is mainly designed to translate the file names of PCAWG consensus callsets for SNV/Indel (https://dcc.icgc.org/api/v1/download?fn=/PCAWG/consensus_snv_indel/final_consensus_snv_indel_passon)

Examples

```
aliquot.ids <- c("e0fccaf5-925a-41f9-b87c-cd5ee4aecb59", "foo")
SP.ids <- map_aliquot_ID_to_SP_ID(aliquot.ids)
```

```
map_SP_ID_to_tumor_type
```

Given PCAWG "SP" IDs (e.g. SP123958) return either the "full" IDs (Kidney-ChRCC::SP123958) or just the tumor type (Kidney-ChRCC)

Description

Given PCAWG "SP" IDs (e.g. SP123958) return either the "full" IDs (Kidney-ChRCC::SP123958) or just the tumor type (Kidney-ChRCC)

Usage

```
map_SP_ID_to_tumor_type(SP.IDs, merge = TRUE)
```

Arguments

SP.IDs	A character vector of PCAWG "SP" IDs.
merge	If TRUE return a parallel vector of <tumor_type>::<SP_ID>; otherwise just <tumor_type>.

Details

Fails with an "subscript out of bounds" error if any of the elements of SP.IDs is unknown.

Examples

```
map_SP_ID_to_tumor_type(c("SP123958", "SP43633"))
map_SP_ID_to_tumor_type(c("SP123958", "SP43633"), merge = FALSE)
```

PCAWG.sample.id	<i>Vectors of the PCAWG tumor_wgs_icgc_specimen_ids</i>
-----------------	---

Description

Note that the PCAWG7 spectra catalogs have 2 sample ids that were blacklisted after the mutational signature analysis was underway. The blacklisted samples are SP116419 and SP116883, which are in PCAWG.sample.id\$black.

Usage

```
PCAWG.sample.id
```

Format

A list with the elements:

white Whitelisted IDs

grey Greylisted IDs

black Blacklisted IDs

Source

https://dcc.icgc.org/api/v1/download?fn=/PCAWG/data_releases/latest/release_may2016.v1.4.with_consensus_calls.tsv, 2019 Oct 09

Examples

```
PCAWG.white.ids <- PCAWG.sample.id$white
```

PCAWG.sample.sheet	<i>PCAWG sample sheet which contains various sample information</i>
--------------------	---

Description

PCAWG sample sheet which contains various sample information

Usage

```
PCAWG.sample.sheet
```

Format

A data table with the following columns:

- donor_unique_id
- donor_wgs_exclusion_white_gray
- submitter_donor_id
- icgc_donor_id
- dcc_project_code
- aliquot_id
- submitter_specimen_id
- icgc_sample_id
- dcc_specimen_type
- library_strategy

Source

https://dcc.icgc.org/api/v1/download?fn=/PCAWG/data_releases/latest/pcawg_sample_sheet.v1.4.2016-09-14.tsv, 2019 Oct 15

Examples

```
aliquot.ids <- PCAWG.sample.sheet$aliquot_id
```

PCAWG7	<i>PCAWG7: A package of data from paper 'Repertoire of Mutational Signatures in Human Cancer'</i>
--------	---

Description

This is a data package with 2 main package variables: `exposure` and `spectra`.

Details

There are also PDF plots of the signatures in `data-raw/plots/`.

There are also several functions for handling PCAWG identifiers:

- * `map_SP_ID_to_tumor_type`

- * `map_aliquot_ID_to_SP_ID`

- * `SampleIDToCancerType`

- * `SplitPCAWGMatrixByTumorType`

- * `SplitMatrixBySampleType`

The reference for the data is

Alexandrov, L.B., Kim, J., Haradhvala, N.J. et al. The repertoire of mutational signatures in human cancer. *Nature* 578, 94-101 (2020). doi: [10.1038/s4158602019433](https://doi.org/10.1038/s4158602019433).

SampleIDToCancerType	<i>Split out the cancer type from the sample ID for PCAWG IDs</i>
----------------------	---

Description

Split out the cancer type from the sample ID for PCAWG IDs

Usage

```
SampleIDToCancerType(PCAWGID)
```

Arguments

PCAWGID	A character vector of PCAWG IDs of the form <code><cancer.type>::<sample.id></code> .
---------	---

Value

A character vector parallel to `PCAWGID` containing only the `<cancer.type>` strings.

Examples

```
cancer.type <- SampleIDToCancerType("Biliary-AdenoCA::SP117655")
```

spectra	<i>PCAWG7 mutational spectra (catalogs)</i>
---------	---

Description

PCAWG7 mutational spectra (catalogs)

Usage

spectra

Format

A list with the elements:

SBS96 Deprecated.

DBS78 Deprecated.

PCAWG A list with the elements:

SBS96 Strand-agnostic single-base substitutions in trinucleotide context.

SBS192 Single-base substitutions in transcripts based on the sense strand.

SBS1536 Strand-agnostic single-base substitutions in pentanucleotide context.

DBS78 Strand-agnostic doublet-base substitutions.

ID Strand-agnostic indels.

TCGA A list with the same elements as the PCAWG element.

other.genome A list with the same elements as the PCAWG element but with ID omitted.

other.exome A list with the same elements as the PCAWG element but with ID omitted.

Source

Files below <https://www.synapse.org/#!Synapse:syn11801889>, 2019 Oct 09. Populated by data-raw/spectra/load.package.variable.specra.R.

Examples

```
SBS96.spectra <- spectra$PCAWG$SBS96
```

SplitMatrixBySampleType

Split an exposure matrix or spectrum matrix into a list of matrices, each for a single sample type

Description

Split an exposure matrix or spectrum matrix into a list of matrices, each for a single sample type

Usage

```
SplitMatrixBySampleType(M, sample.type)
```


Arguments

<code>M</code>	A numerical matrix or data frame or ICAMS catalog in which columns are samples (e.g. tumors) and rows are either mutational signatures (for exposures) or mutation types (for spectra), and, each element is the number of mutations due to a given mutational signature or mutation type in a single sample
<code>sample.type</code>	A character or numeric vector, each element of which indicates a particular sample type.

Value

Invisibly, the list of exposure or spectrum matrices created by splitting `M` by `sample.type`.

Examples

```
ff <- matrix(1, nrow=3, ncol = 2)
colnames(ff) <- c("sample1", "sample2")
xx <- SplitMatrixBySampleType(ff, c("sample.type.x", "sample.type.y"))
xx
```

SplitPCAWGMatrixByTumorType

Extract tumor type from column names and return the input matrix split by tumor type based on the PCAWG <tumor_type>::<sample_id> convention

Description

Extract tumor type from column names and return the input matrix split by tumor type based on the PCAWG <tumor_type>::<sample_id> convention

Usage

```
SplitPCAWGMatrixByTumorType(M)
```

Arguments

<code>M</code>	A numerical matrix or data frame or ICAMS catalog in which columns are samples (e.g. tumors) and rows are either mutational signatures (for exposures) or mutation types (for spectra), and each element is the number of mutations due to a given mutational signature or mutation type in a single sample. The column names must be of the the form <cancer.type>::<sample.ID>.
----------------	---

Value

Invisibly, the list of exposure matrices or [ICAMS](#) catalogs created by splitting `matrix` by the tumor type encoded in the column names.

Examples

```
mm <- SplitPCAWGMatrixByTumorType(spectra$PCAWG$DBS78)
```

TCGA_ID_to_ICGC_ID	<i>Translate TCGA (The Cancer Genome Atlas) IDs to ICGC (International Cancer Genome Consortium) IDs</i>
--------------------	--

Description

Translate TCGA (The Cancer Genome Atlas) IDs to ICGC (International Cancer Genome Consortium) IDs

Usage

```
TCGA_ID_to_ICGC_ID(tcga.ids)
```

Arguments

<code>tcga.ids</code>	Character vector of TCGA IDs.
-----------------------	-------------------------------

Details

If there are TCGA IDs that cannot be matched to any ICGC IDs, return NA with a warning.

Value

Character vector of corresponding ICGC IDs. If a corresponding ICGC ID cannot be found, then return NA with a warning.

Examples

```
tcga.ids <- c("TCGA-AA-A01V", "foo", "TCGA-CA-6717", "bar")
icgc.ids <- TCGA_ID_to_ICGC_ID(tcga.ids)
icgc.ids <- icgc.ids[nzchar(icgc.ids)]
```

Index

* datasets

- exposure, [2](#)
- exposure.stats, [3](#)
- PCAWG.sample.id, [5](#)
- PCAWG.sample.sheet, [6](#)
- spectra, [8](#)

CancerTypes, [2](#)

exposure, [2](#), [6](#)
exposure.stats, [3](#)

ICAMS, [9](#)

map_aliquot_ID_to_SP_ID, [4](#), [7](#)
map_SP_ID_to_tumor_type, [4](#), [7](#)

PCAWG.sample.id, [5](#)
PCAWG.sample.sheet, [6](#)
PCAWG7, [6](#)

SampleIDToCancerType, [7](#), [7](#)
spectra, [6](#), [8](#)
SplitMatrixBySampleType, [7](#), [8](#)
SplitPCAWGMatrixByTumorType, [7](#), [9](#)

TCGA_ID_to_ICGC_ID, [10](#)