

Package: ggmanh (via r-universe)

July 21, 2024

Title Visualization Tool for GWAS Result

Version 1.9.7

Description Manhattan plot and QQ Plot are commonly used to visualize the end result of Genome Wide Association Study. The ``ggmanh" package aims to keep the generation of these plots simple while maintaining customizability. Main functions include `manhattan_plot`, `qqunif`, and `thinPoints`.

biocViews Visualization, GenomeWideAssociation, Genetics

License MIT + file LICENSE

Encoding UTF-8

Roxygen list(markdown = TRUE)

RoxygenNote 7.3.1

Imports gdsfmt, ggrepel, grDevices, RColorBrewer, rlang, scales, SeqArray (>= 1.32.0), stats, tidyr

Depends methods, ggplot2

Suggests BiocStyle, rmarkdown, knitr, testthat (>= 3.0.0), markdown, GenomicRanges, magick

Config/testthat/edition 3

VignetteBuilder knitr

Repository <https://bioc.r-universe.dev>

RemoteUrl <https://github.com/bioc/ggmanh>

RemoteRef HEAD

RemoteSha 5114ae2995dcb9e776eac4f7f02ba15c761bc9cc

Contents

<code>calc_new_pos</code>	2
<code>default_gds_path</code>	3
<code>gds_annotate</code>	3
<code>ggmanh</code>	4

ggmanh_annotation_gds	5
manhattan_data_preprocess	6
manhattan_plot	8
qqunif	13
thinPoints	14

Index	16
--------------	-----------

calc_new_pos	<i>Calculate new x-position of each point</i>
--------------	---

Description

Calculate the actual x-positions of each point used for the manhattan plot. MPdata object contains the unscaled positions that has not been positioned according to the relative position and width of each chromosome.

Usage

```
calc_new_pos(mpdata)
```

Arguments

mpdata an MPdata object.

Details

This is used calculate the actual positions used for the inside manhattan_plot function. It was designed this way should the scaling and relative positioning of each chromosome be changed (e.g. gap between the)

Value

a numeric vector containing the scaled x-positions.

Examples

```
gwasdat <- data.frame(
  "chromosome" = rep(1:5, each = 30),
  "position" = c(replicate(5, sample(1:300, 30))),
  "pvalue" = rbeta(150, 1, 1)^5
)

mpdata <- manhattan_data_preprocess(
  gwasdat, pval.colname = "pvalue", chr.colname = "chromosome", pos.colname = "position",
  chr.order = as.character(1:5)
)

calc_new_pos(mpdata)
```

default_gds_path	<i>Path to Default GDS File</i>
------------------	---------------------------------

Description

Find path to the default gds file.

Usage

```
default_gds_path()
```

Value

A character vector.

Examples

```
default_gds_path()
```

gds_annotate	<i>Annotation with GDS File</i>
--------------	---------------------------------

Description

Retrieve variant annotation stored in a GDS file with chromosome location or rs.id.

Usage

```
gds_annotate(  
  x,  
  gdsfile = NULL,  
  annot.method = "position",  
  chr = NULL,  
  pos = NULL,  
  ref = NULL,  
  alt = NULL,  
  rs.id = NULL,  
  concat_char = "/",  
  verbose = TRUE,  
  annotation_names = c("annotation/info/symbol", "annotation/info/consequence",  
    "annotation/info/LoF")  
)
```

Arguments

<code>x</code>	a <code>data.frame</code> object to be annotated.
<code>gdsfile</code>	a character for GDS filename. If <code>NULL</code> , the default GDS file included with the package is used.
<code>annot.method</code>	a method for searching variants. "position" requires <code>chr</code> , <code>pos</code> , <code>ref</code> , and <code>alt</code> . "rs.id" requires <code>rs.id</code> .
<code>chr, pos, ref, alt, rs.id</code>	column names of <code>x</code> that contain chromosome, position, reference allele, alternate allele, and <code>rs.id</code> , respectively.
<code>concat_char</code>	a character used to separate multiple annotations returned from the gds file.
<code>verbose</code>	output messages.
<code>annotation_names</code>	a character vector of nodes of the <code>gdsfile</code> that are to be extracted.

Value

A character vector the length of `nrow(x)` if `concat_char` is a character. A data frame with `nrow(x)` rows and `length(annotation_names)` if `concat_char` is null.

Examples

```
vardata <- data.frame(
  chr = c(11,20,14),
  pos = c(12261002, 10033792, 23875025),
  ref = c("G", "G", "CG"),
  alt = c("A", "A", "C")
)

annotations <- gds_annotate(
  x = vardata, annot.method = "position",
  chr = "chr", pos = "pos", ref = "ref", alt = "alt"
)

print(annotations)
```

Description

ggmanh provides flexible tools for visualizing GWAS result for downstream analysis.

Details

Manhattan plot is commonly used to display significant Single Nucleotide Polymorphisms (SNPs) in Genome Wide Association Study (GWAS). This package comes with features useful for Manhattan plot creation, including annotation with [ggrepel](#), truncating data for faster plot generation, and manual rescaling of the y-axis. The Manhattan plot is generated in two steps: data preprocessing and plotting. This allows the user to iteratively customize the plot without having to process the GWAS summary data over and over again. Currently, `data.frame` and `GRanges` from `GenomicRanges` are supported.

A vignette detailing the usage of the package is accessible by `vignette("ggmanh")`.

Author(s)

Maintainer: John Lee <swannyy.stat@gmail.com>

Authors:

- John Lee <john.lee@abbvie.com> (AbbVie)

Other contributors:

- Xiuwen Zheng <xiuwen.zheng@abbvie.com> [contributor, data contributor]

ggmanh_annotation_gds *gnomAD Variant Annotation in SeqArray Format*

Description

ggmanh provides a GDS file whose path is accessible by `default_gds_path`. The original annotation file is from the gnomAD browser v2.1.1 release, available in this link: <https://gnomad.broadinstitute.org/downloads>. This gds file contains variants in the exome with the global minor allele frequency ≥ 0.0002 , and has been manually curated to fit the file size requirement for R Bioconductor packages.

Format

A GDS file with 1015430 variants with chromosome, position, allele, gene symbol, Ensembl VEP Consequence, and predicted LoF.

`manhattan_data_preprocess`*Preprocess GWAS Result*

Description

Preprocesses a result from Genome Wide Association Study before making a manhattan plot. It accepts a `data.frame`, which at bare minimum should contain a chromosome, position, and p-value. Additional options, such as chromosome color, label column names, and colors for specific variants, are provided here.

Usage

```
manhattan_data_preprocess(x, ...)

## Default S3 method:
manhattan_data_preprocess(x, ...)

## S3 method for class 'data.frame'
manhattan_data_preprocess(
  x,
  chromosome = NULL,
  signif = c(5e-08, 1e-05),
  pval.colname = "pval",
  chr.colname = "chr",
  pos.colname = "pos",
  highlight.colname = NULL,
  chr.order = NULL,
  signif.col = NULL,
  chr.col = NULL,
  highlight.col = NULL,
  preserve.position = FALSE,
  thin = NULL,
  thin.n = 1000,
  thin.bins = 200,
  pval.log.transform = TRUE,
  chr.gap.scaling = 1,
  ...
)

## S4 method for signature 'GRanges'
manhattan_data_preprocess(
  x,
  chromosome = NULL,
  signif = c(5e-08, 1e-05),
  pval.colname = "pval",
  highlight.colname = NULL,
```

```

    chr.order = NULL,
    signif.col = NULL,
    chr.col = NULL,
    highlight.col = NULL,
    preserve.position = FALSE,
    thin = NULL,
    thin.n = 100,
    thin.bins = 200,
    pval.log.transform = TRUE,
    chr.gap.scaling = 1,
    ...
  )

```

Arguments

<code>x</code>	a data frame or any other extension of data frame (e.g. a tibble). At bare minimum, it should contain chromosome, position, and p-value.
<code>...</code>	Additional arguments for <code>manhattan_data_preprocess</code> .
<code>chromosome</code>	a character. This is supplied if a manhattan plot of a single chromosome is desired. If <code>NULL</code> , then all the chromosomes in the data will be plotted.
<code>signif</code>	a numeric vector. Significant p-value thresholds to be drawn for manhattan plot. At least one value should be provided. Default value is <code>c(5e-08, 1e-5)</code>
<code>pval.colname</code>	a character. Column name of <code>x</code> containing p.value.
<code>chr.colname</code>	a character. Column name of <code>x</code> containing chromosome number.
<code>pos.colname</code>	a character. Column name of <code>x</code> containing position.
<code>highlight.colname</code>	a character. If you desire to color certain points (e.g. significant variants) rather than color by chromosome, you can specify the category in this column, and provide the color mapping in <code>highlight.col</code> . Ignored if <code>NULL</code> .
<code>chr.order</code>	a character vector. Order of chromosomes presented in manhattan plot.
<code>signif.col</code>	a character vector of equal length as <code>signif</code> . It contains colors for the lines drawn at <code>signif</code> . If <code>NULL</code> , the smallest value is colored black while others are grey.
<code>chr.col</code>	a character vector of equal length as <code>chr.order</code> . It contains colors for the chromosomes. Name of the vector should match <code>chr.order</code> . If <code>NULL</code> , default colors are applied using <code>RColorBrewer</code> .
<code>highlight.col</code>	a character vector. It contains color mapping for the values from <code>highlight.colname</code> .
<code>preserve.position</code>	a logical. If <code>TRUE</code> , the width of each chromosome reflect the number of variants and the position of each variant is correctly scaled? If <code>FALSE</code> , the width of each chromosome is equal and the variants are equally spaced.
<code>thin</code>	a logical. If <code>TRUE</code> , <code>thinPoints</code> will be applied. Defaults to <code>TRUE</code> if chromosome is <code>NULL</code> . Defaults to <code>FALSE</code> if chromosome is supplied.
<code>thin.n</code>	an integer. Number of max points per horizontal partitions of the plot. Defaults to 1000.

`thin.bins` an integer. Number of bins to partition the data. Defaults to 200.

`pval.log.transform` a logical. If TRUE, the p-value will be transformed to $-\log_{10}(\text{p-value})$.

`chr.gap.scaling` scaling factor for gap between chromosome if you desire to change it. This can also be set in `manhattan_plot`

Details

`manhattan_data_preprocess` gathers information needed to plot a manhattan plot and organizes the information as MPdata S3 object.

New positions for each points are calculated, and stored in the data.frame as "new_pos". By default, all chromosomes will have the same width, with each point being equally spaced. This behavior is changed when `preserve.position = TRUE`. The width of each chromosome will scale to the number of points and the points will reflect the original positions.

`chr.col` and `highlight.col`, maps the data values to colors. If they are an unnamed vector, then the function will try its best to match the values of `chr.colname` or `highlight.colname` to the colors. If they are a named vector, then they are expected to map all values to a color. If `highlight.colname` is supplied, then `chr.col` is ignored.

While feeding a data.frame directly into `manhattan_plot` does preprocessing & plotting in one step. If you plan on making multiple plots with different graphic options, you have the choice to preprocess separately and then generate plots.

Value

a MPdata object. This object contains all the necessary info for constructing a manhattan plot.

Examples

```
gwasdat <- data.frame(
  "chromosome" = rep(1:5, each = 30),
  "position" = c(replicate(5, sample(1:300, 30))),
  "pvalue" = rbeta(150, 1, 1)^5
)

manhattan_data_preprocess(
  gwasdat, pval.colname = "pvalue", chr.colname = "chromosome", pos.colname = "position",
  chr.order = as.character(1:5)
)
```

manhattan_plot

Manhattan Plotting

Description

A generic function for manhattan plot.

Usage

```
manhattan_plot(x, ...)

manhattan_plot.default(x, ...)

## S3 method for class 'data.frame'
manhattan_plot(
  x,
  chromosome = NULL,
  outfn = NULL,
  signif = c(5e-08, 1e-05),
  pval.colname = "pval",
  chr.colname = "chr",
  pos.colname = "pos",
  label.colname = NULL,
  highlight.colname = NULL,
  chr.order = NULL,
  signif.col = NULL,
  chr.col = NULL,
  highlight.col = NULL,
  rescale = TRUE,
  rescale.ratio.threshold = 5,
  signif.rel.pos = 0.2,
  chr.gap.scaling = 1,
  color.by.highlight = FALSE,
  preserve.position = FALSE,
  thin = NULL,
  thin.n = 1000,
  thin.bins = 200,
  pval.log.transform = TRUE,
  plot.title = ggplot2::waiver(),
  plot.subtitle = ggplot2::waiver(),
  plot.width = 10,
  plot.height = 5,
  point.size = 0.75,
  label.font.size = 2,
  max.overlaps = 20,
  x.label = "Chromosome",
  y.label = expression(-log[10](p)),
  ...
)

## S3 method for class 'MPdata'
manhattan_plot(
  x,
  chromosome = NULL,
  outfn = NULL,
  signif = NULL,
```

```

    signif.col = NULL,
    rescale = TRUE,
    rescale.ratio.threshold = 5,
    signif.rel.pos = 0.2,
    chr.gap.scaling = NULL,
    color.by.highlight = FALSE,
    label.colname = NULL,
    x.label = "Chromosome",
    y.label = expression(-log[10](p)),
    point.size = 0.75,
    label.font.size = 2,
    max.overlaps = 20,
    plot.title = ggplot2::waiver(),
    plot.subtitle = ggplot2::waiver(),
    plot.width = 10,
    plot.height = 5,
    ...
)

## S4 method for signature 'GRanges'
manhattan_plot(
  x,
  chromosome = NULL,
  outfn = NULL,
  signif = c(5e-08, 1e-05),
  pval.colname = "pval",
  label.colname = NULL,
  highlight.colname = NULL,
  chr.order = NULL,
  signif.col = NULL,
  chr.col = NULL,
  highlight.col = NULL,
  rescale = TRUE,
  rescale.ratio.threshold = 5,
  signif.rel.pos = 0.2,
  chr.gap.scaling = 1,
  color.by.highlight = FALSE,
  preserve.position = FALSE,
  thin = NULL,
  thin.n = 1000,
  thin.bins = 200,
  pval.log.transform = TRUE,
  plot.title = ggplot2::waiver(),
  plot.subtitle = ggplot2::waiver(),
  plot.width = 10,
  plot.height = 5,
  point.size = 0.75,
  label.font.size = 2,

```

```

    max.overlaps = 20,
    x.label = "Chromosome",
    y.label = expression(-log[10](p)),
    ...
)

```

Arguments

<code>x</code>	a <code>data.frame</code> , an extension of <code>data.frame</code> object (e.g. <code>tibble</code>), or an <code>MPdata</code> object.
<code>...</code>	additional arguments to be passed onto <code>geom_label_repel</code>
<code>chromosome</code>	a character. This is supplied if a manhattan plot of a single chromosome is desired. If <code>NULL</code> , then all the chromosomes in the data will be plotted.
<code>outfn</code>	a character. File name to save the Manhattan Plot. If <code>outfn</code> is supplied (i.e. <code>!is.null(outfn)</code>), then the plot is not drawn in the graphics window.
<code>signif</code>	a numeric vector. Significant p-value thresholds to be drawn for manhattan plot. At least one value should be provided. Default value is <code>c(5e-08, 1e-5)</code> . If <code>signif</code> is not <code>NULL</code> and <code>x</code> is an <code>MPdata</code> object, <code>signif</code> argument overrides the value inside <code>MPdata</code> .
<code>pval.colname</code>	a character. Column name of <code>x</code> containing p.value.
<code>chr.colname</code>	a character. Column name of <code>x</code> containing chromosome number.
<code>pos.colname</code>	a character. Column name of <code>x</code> containing position.
<code>label.colname</code>	a character. Name of the column in <code>MPdata\$data</code> to be used for labelling.
<code>highlight.colname</code>	a character. If you desire to color certain points (e.g. significant variants) rather than color by chromosome, you can specify the category in this column, and provide the color mapping in <code>highlight.col</code> . Ignored if <code>NULL</code> .
<code>chr.order</code>	a character vector. Order of chromosomes presented in manhattan plot.
<code>signif.col</code>	a character vector of equal length as <code>signif</code> . It contains colors for the lines drawn at <code>signif</code> . If <code>NULL</code> , the smallest value is colored black while others are grey. If <code>x</code> is an <code>MPdata</code> object, behaves similarly to <code>signif</code> .
<code>chr.col</code>	a character vector of equal length as <code>chr.order</code> . It contains colors for the chromosomes. Name of the vector should match <code>chr.order</code> . If <code>NULL</code> , default colors are applied using <code>RColorBrewer</code> .
<code>highlight.col</code>	a character vector. It contains color mapping for the values from <code>highlight.colname</code> .
<code>rescale</code>	a logical. If <code>TRUE</code> , the plot will rescale itself depending on the data. More on this in details.
<code>rescale.ratio.threshold</code>	a numeric. Threshold of that triggers the rescale.
<code>signif.rel.pos</code>	a numeric between 0.1 and 0.9. If the plot is rescaled, where should the significance threshold be positioned?
<code>chr.gap.scaling</code>	a numeric. scaling factor for gap between chromosome if you desire to change it if <code>x</code> is an <code>MPdata</code> object, then the gap will scale relative to the gap in the object.

<code>color.by.highlight</code>	a logical. Should the points be colored based on a highlight column?
<code>preserve.position</code>	a logical. If TRUE, the width of each chromosome reflect the number of variants and the position of each variant is correctly scaled? If FALSE, the width of each chromosome is equal and the variants are equally spaced.
<code>thin</code>	a logical. If TRUE, <code>thinPoints</code> will be applied. Defaults to TRUE if chromosome is NULL. Defaults to FALSE if chromosome is supplied.
<code>thin.n</code>	an integer. Number of max points per horizontal partitions of the plot. Defaults to 1000.
<code>thin.bins</code>	an integer. Number of bins to partition the data. Defaults to 200.
<code>pval.log.transform</code>	a logical. If TRUE, the p-value will be transformed to $-\log_{10}(\text{p-value})$.
<code>plot.title</code>	a character. Plot title
<code>plot.subtitle</code>	a character. Plot subtitle
<code>plot.width</code>	a numeric. Plot width in inches.
<code>plot.height</code>	a numeric. Plot height in inches.
<code>point.size</code>	a numeric. Size of the points.
<code>label.font.size</code>	a numeric. Size of the labels.
<code>max.overlaps</code>	an integer. Exclude text labels that overlaps too many things.
<code>x.label</code>	a character. x-axis label
<code>y.label</code>	a character. y-axis label

Details

This generic function accepts a result of a GWAS in the form of `data.frame` or a `MPdata` object produced by `manhattan_data_preprocess`. The function will throw an error if another type of object is passed.

Having `rescale = TRUE` is useful when there are points with very high $-\log_{10}(\text{p.value})$. In this case, the function attempts to split the plot into two different scales, with the split happening near the strictest significance threshold. More precisely, the plot is rescaled when

$$-\log_{10}(\text{pvalue}) / (\text{strictestsignificancethreshold}) \geq \text{rescale.ratio.threshold}$$

If you wish to add annotation to the points, provide the name of the column to `label.colname`. The labels are added with `ggrepel`.

Be careful though: if the annotation column contains a large number of variants, then the plotting could take a long time, and the labels will clutter up the plot. For those points with no annotation, you have the choice to set them as NA or "".

Value

gg object if `is.null(outfn)`, NULL if `!is.null(outf)`

Examples

```

gwasdat <- data.frame(
  "chromosome" = rep(1:5, each = 30),
  "position" = c(replicate(5, sample(1:300, 30))),
  "pvalue" = rbeta(150, 1, 1)^5
)

manhattan_plot(
  gwasdat, pval.colname = "pvalue", chr.colname = "chromosome", pos.colname = "position",
  chr.order = as.character(1:5)
)

mpdata <- manhattan_data_preprocess(
  gwasdat, pval.colname = "pvalue", chr.colname = "chromosome", pos.colname = "position",
  chr.order = as.character(1:5)
)

manhattan_plot(mpdata)

```

qqunif

*Plot Quantile-Quantile Plot of p-values against uniform distribution.***Description**

Plot Quantile-Quantile Plot of p-values against uniform distribution.

Usage

```

qqunif(
  x,
  outfn = NULL,
  conf.int = 0.95,
  plot.width = 5,
  plot.height = 5,
  thin = TRUE,
  thin.n = 500,
  zero.pval = "replace"
)

```

Arguments

<code>x</code>	a numeric vector of p-values. All values should be between 0 and 1.
<code>outfn</code>	a character. File name to save the QQ Plot. If <code>outfn</code> is supplied (i.e. <code>!is.null(outfn)</code>), then the plot is not drawn in the graphics window.
<code>conf.int</code>	a numeric between 0 and 1. Confidence band to draw around reference line. Set to NA to leave it out.

plot.width	a numeric. Plot width in inches.
plot.height	a numeric. Plot height in inches.
thin	a logical. Reduce number of data points when they are cluttered?
thin.n	an integer. Number of max points per horizontal partitions of the plot. Defaults to 500.
zero.pval	a character. Determine how to treat 0 pvals. "replace" will replace the p-value of zero with the non-zero minimum. "remove" will remove the p-value of zero.

Value

a ggplot object

Examples

```
x <- rbeta(1000, 1, 1)
qqunif(x)
```

thinPoints	<i>Thin Data Points</i>
------------	-------------------------

Description

Reduce the number of cluttered data points.

Usage

```
thinPoints(dat, value, n = 1000, nbins = 200, groupBy = NULL)
```

Arguments

dat	a data frame
value	column name of dat to be used for partitioning (see details)
n	number of points to sample for each partition
nbins	number of partitions
groupBy	column name of dat to group by before partitioning (e.g. chromosome)

Details

The result of Genome Wide Association Study can be very large, with the majority of points being being clustered below significance threshold. This unnecessarily increases the time to plot while making almost no difference. This function reduces the number of points by partitioning the points by a numeric column value into nbins and sampling n points.

Value

a data.frame

Examples

```
dat <- data.frame(  
  A1 = c(1:20, 20, 20),  
  A2 = c(rep(1, 12), rep(1,5), rep(20, 3), 20, 20) ,  
  B = rep(c("a", "b", "c", "d"), times = c(5, 7, 8, 2))  
)  
# partition "A1" into 2 bins and then sample 6 data points  
thinPoints(dat, value = "A1", n = 6, nbins = 2)  
# partition "A2" into 2 bins and then sample 6 data points  
thinPoints(dat, value = "A2", n = 6, nbins = 2)  
# group by "B", partition "A2" into 2 bins and then sample 3 data points  
thinPoints(dat, value = "A2", n = 3, nbins = 2, groupBy = "B")
```

Index

`calc_new_pos`, [2](#)

`default_gds_path`, [3](#)

`gds_annotate`, [3](#)

`ggmanh`, [4](#)

`ggmanh-package` (`ggmanh`), [4](#)

`ggmanh_annotation_gds`, [5](#)

`ggrepel`, [5](#), [12](#)

`manhattan_data_preprocess`, [6](#)

`manhattan_data_preprocess`, `GRanges`-method

 (`manhattan_data_preprocess`), [6](#)

`manhattan_data_preprocess.data.frame`

 (`manhattan_data_preprocess`), [6](#)

`manhattan_data_preprocess.default`

 (`manhattan_data_preprocess`), [6](#)

`manhattan_plot`, [8](#)

`manhattan_plot`, `GRanges`-method

 (`manhattan_plot`), [8](#)

`manhattan_plot.data.frame`

 (`manhattan_plot`), [8](#)

`manhattan_plot.default`

 (`manhattan_plot`), [8](#)

`manhattan_plot.MPdata` (`manhattan_plot`),
 [8](#)

`qqunif`, [13](#)

`thinPoints`, [14](#)