

# Package ‘MegaCollapsing’

January 23, 2020

**Type** Package

**Title** Collapsing Analysis at the MegaGene Level

**Version** 0.1.0

**Date** 2020-01-31

**Author** Quanli Wang and Slave Petrovski

**Maintainer** Quanli Wang<quanliwang20@gmail.com>

**Description** Taking results from a gene/sample level variant count matrix and a set of gene sets, the package collapses gene/sample level variant counts into gene-set (MegaGene)/sample level counts and then allow users to run Fisher Exact Test efficiently.

**License** GPL(>=3)

**LazyData** TRUE

## R topics documented:

as.mega.matrix . . . . .	2
burden.test . . . . .	2
burden.test.with.permutation . . . . .	3
exclude.genes . . . . .	4
get.pvalues . . . . .	4
random.gene.sets . . . . .	5
read.collapsing.data . . . . .	5
read.collapsing.matrix . . . . .	6
read.data . . . . .	7
read.gene.sets . . . . .	8
read.sample.list . . . . .	8
read.seed.list . . . . .	9

<b>Index</b>	<b>10</b>
--------------	-----------

---

<code>as.mega.matrix</code>	<i>Generate mega collapsing matrix/matrices for a list of gene sets.</i>
-----------------------------	--

---

### Description

Taking collapsing matrix/matrices and a list of gene sets, this function "collapses" gene/sample level variant information to gene set level variant information, which can further be used for (mega) collapsing analysis.

### Usage

```
as.mega.matrix(gene.sets, input.data)
```

### Arguments

<code>gene.sets</code>	The gene sets object that hold list of gene sets. It is typically obtained by calling function <code>read.gene.sets</code> .
<code>input.data</code>	The objects that holds the collapsing matrix/matrices. It is typically obtained by calling function <code>read.collapsing.data</code> .

### Details

A lot of explanation here if we wanted. Or provide a link to the method.

### Value

Return a list of objects that represent the mega collapsing matrix/matrices.

---

<code>burden.test</code>	<i>Run regression/test using collapsing matrix with covariates.</i>
--------------------------	---

---

### Description

Add a lot of explanations here to justify why we set up the model this way.

### Usage

```
burden.test(sample.list, gene.sets, syn.mat, non.syn.mat, X = NULL, Y = NULL)
```

### Arguments

<code>sample.list</code>	An object that hold information/covariates about samples.
<code>gene.sets</code>	The list of gene sets to be analyzed.
<code>syn.mat</code>	The syn collapsing matrix to be used. We should have an option not to use this one.
<code>non.syn.mat</code>	The non syn collapsing matrix to be used.
<code>X</code>	a character list of covariates to be used in regression models.
<code>Y</code>	The response variable. 1 for cases and 0 for controls.

**Details**

Explain.

**Value**

Return a list of test results and summary statistics. Will document this further once we decide what to keep.

---

burden.test.with.permutation

*Permutation based Burden Test*


---

**Description**

Don't do this unless we really want to do it. This will run permutation based on Burden Test for given gene list and input mega/collapsing matrix.

**Usage**

```
burden.test.with.permutation(sample.list, gene.set, syn.mat, non.syn.mat,
                             X = NULL, Y = NULL, n.permutations = 100)
```

**Arguments**

sample.list	An object that hold information/covaraite about samples.[Copied from burden test]
gene.set	The list of gene sets to be anallize.
syn.mat	The syn collsping matrix to be used. We should have an option not to use this one.
non.syn.mat	The non syn collsping matrix to be used.
X	a charactor list of covariates to be used in regression models.
Y	The response varaible. 1 for cases and 0 for controls.
n.permutations	Number of permutations to use in this test. Default to 100.

**Value**

A list that holds the actual burden test results and permuted results. Add more if we are to keep this function.

---

<code>exclude.genes</code>	<i>Exclude a set of genes from input collapsing matrices before running MegaCollapsing Analysis.</i>
----------------------------	--

---

### Description

It might be desire to remove genes that are causing trouble before running MegaCollapsing and this function can be used for that purpose.

### Usage

```
exclude.genes(input.data, exclude.file)
```

### Arguments

<code>input.data</code>	The output from function call to <code>read.collapsing.data</code> .
<code>exclude.file</code>	The file contains the list of genes to be excluded before analysis.

### Details

Add more about how it was implemented and the reason why we might want to do this.

### Value

A list that has the same structure as `input.data`.

---

<code>get.pvalues</code>	<i>Run permutation based Fisher Test on (Mega)collapsing matrix.</i>
--------------------------	--

---

### Description

This function allows one to run permutation based FET tests. Not sure if we need this one in this package though.

### Usage

```
get.pvalues(matrix, is.case, n.permutations = 1000)
```

### Arguments

<code>matrix</code>	A (mega) collapsing matrix.
<code>is.case</code>	A boolean vector indicating case/control status.
<code>n.permutations</code>	Number of permutations to run the tests.

### Details

Running permutation based FET test from (mega)Collapsing matrix. Will explain more if we will keep this function.

**Value**

A list that holds both FET p values and permutated p-values. These p-values can be used to do permutation based QQ plot

---

random.gene.sets	<i>Generate random gene sets from a list of genes that share similar structure with an input gene set</i>
------------------	---

---

**Description**

This function is used to generate randomized sets of genes that can be used to evaluate test results. This was mainly introduced to calibrate our test results for the giving input data. Need to investigate/explore more about its usage. Might want to remove it from our package if we are not going to use it.

**Usage**

```
random.gene.sets(gene.sets, genes)
```

**Arguments**

gene.sets	The actual gene sets that would be used in a MegaCollapsing analysis.
genes	A list of genes that would be used to generate randomized gene sets.

**Details**

Explain/explore why we had this function.

**Value**

A randomized sets of genes that share the same set length with the input gene sets but with each gene set populated with random genes.

---

read.collapsing.data	<i>Read all data needed for a "paired burden test."</i>
----------------------	---

---

**Description**

This function is currently used by CGR to run MegaCollapsing analysis that uses both synonymous and "non-synonymous" collapsing matrices. We need to come up with better names and descriptions for this kind of analysis.

**Usage**

```
read.collapsing.data(samples, syn, non.syn, sample.column = "IID")
```

**Arguments**

<code>samples</code>	The input PED file to define the list of samples to be used in analysis.
<code>syn</code>	The collapsing matrix for syn model.
<code>non.syn</code>	The collapsing matrix for matrix to be acutally tested.
<code>sample.column</code>	The column name that is used to identify the sample ID from input sample file.

**Details**

Need a lot of description here to describe we are trying to do. Obvious bug: the passed in "sample.column" was not used.

**Value**

Return a lost of objects that are parsed and matched for "paired burden test"

<code>sample.list</code>	The final sample list used for analysis
<code>syn</code>	The final matching collapsing matrix for syn model
<code>non.syn</code>	The final matching collapsing matrix for non-syn model

---

`read.collapsing.matrix`

*Read gene/sample varaint count (collapsing) matrix.*

---

**Description**

This is a helper function that parses the input gene/sample varaint count (collapsing) matrix into an R matrix.

**Usage**

```
read.collapsing.matrix(matrix_file)
```

**Arguments**

<code>matrix_file</code>	The input collapsing matrix.
--------------------------	------------------------------

**Details**

Describe what a collapsing matrix looks like.

**Value**

Returns a matrix that represents the collapsing matrix, with columns for samples and rows for genes.

---

read.data	<i>Read gene/sample variant count matrix.</i>
-----------	---

---

## Description

This is a helper function to extract gene/sample variant count matrix based on an input sample list and also optionally a list of genes to be used as filter.

## Usage

```
read.data(sample.file, matrix.file, filter.list = NULL, sample.column = 2,
          case.control.column = 6)
```

## Arguments

sample.file	An input sample list file. The default format is a standard PED file that is tab delimited, without column header and has second column indicating the sample ID/name and the 6 column indicating the control/case status, with 1 for control sample and 2 for case sample. User can also provide a tab delimited text file and indicating the column index for sample ID and case/control status.
matrix.file	The gene/sample count (collapsing) matrix to be read. A gene/sample count matrix is a tab delimited spreadsheet like text file with each column represents a sample, and each row represents a gene. Gene names will be stored in the first column and the sample names will be stored in the first row.
filter.list	Optional. A flat text file with each row gives a gene that is to be extracted from the input matrix file. By default, all genes found in the matrix file will be used.
sample.column	Optional. If the sample ID column from the sample file is not the second column, user will need to provide the index(one based) for actual sample ID column.
case.control.column	Optional. If the case/control status column from the sample file is not the second column, user will need to provide the index(one based) for actual sample ID column.

## Value

Return a list containing the extracted case/control status and the selected gene/sample count matrix.

data	A matrix for the extracted gene/sample count
is.case	A boolean vector indicating the case/control status

---

read.gene.sets	<i>Read the gene set definition from input text file</i>
----------------	--

---

### Description

This function read the gene set to be used for mega collapsing analysis.

### Usage

```
read.gene.sets(meta_gene_file)
```

### Arguments

meta\_gene\_file The input file that defines gene sets.

### Details

Need to add more details about the gene set format we are using here and the minimal information we will need to create the geneset information.

### Value

Return a list of gene sets.

---

read.sample.list	<i>Read and format a sample list file.</i>
------------------	--

---

### Description

This is a helper function to parse and standardize the sample names from sample list file. It assumes that the sample list file is tab delimited, has a column header and the sample ID column is labeled by column header "IID". To do: make this more generic and allow user to change the format.

### Usage

```
read.sample.list(sample_file)
```

### Arguments

sample\_file The input sample list file. Tab delimited, with column header and the sample ID column is indicated by column name IID.

### Details

This is a helper function used only when the sample file has column headers and the some sample names are just integers. In this case, we will add an "X" in front the integer to make it a proper sample name. We might want to remove this function later.

### Value

Retrun a data frame represent the sample list.



---

read.seed.list	<i>A helper function to read "seed genes"</i>
----------------	---

---

**Description**

Forget what this is supposed to be. Will need to investigate if we need it or not, or if we need to re-implement it to fit our purpose.

**Usage**

```
read.seed.list(seed_set)
```

**Arguments**

seed_set	An input for "seed set" of genes.
----------	-----------------------------------

**Details**

Need to figure how this is used (or not).

**Value**

Return a list of "seed set"

# Index

`as.mega.matrix`, [2](#)  
`burden.test`, [2](#)  
`burden.test.with.permutation`, [3](#)  
`exclude.genes`, [4](#)  
`get.pvalues`, [4](#)  
`random.gene.sets`, [5](#)  
`read.collapsing.data`, [5](#)  
`read.collapsing.matrix`, [6](#)  
`read.data`, [7](#)  
`read.gene.sets`, [8](#)  
`read.sample.list`, [8](#)  
`read.seed.list`, [9](#)