

GSEApot - R package

Civelek Lab, University of Virginia

Sarah Innis, Kelsie Reinaltt, Mete Civelek, Warren Anderson

May 23, 2020

Contents

Introduction to GSEApot	2
Load the Package	2
Geneset Database Options	2
Identifying and Loading Existing Databases	2
Setting up Data for Gene Set Enrichment Analysis using GSEApot	3
Create a Phenotype Input Function	4
GSEA Analysis	5
GSEApots: Function and Parameters	5
Function Outputs	6
Running GSEA on a Novel Geneset	10
Add to an Existing Gene Set Database	11
Ligand and Receptor Databases	11
Introduction of Transcription Factor Databases	11
Introduction of Kadoki Ligand and Receptor Database	12
References	14

Introduction to GSEApLOT

This vignette provides an introduction to the GSEApLOT package. This package builds on the Gene Set Enrichment Analysis (GSEA), a computational method published by the Broad Institute. This package uses the enrichment data from the Broad Institute's analysis to create plots that are easier to understand and modify, and expands the amount of gene sets analyzed, thus increasing the amount of plots output by the analysis. In addition to improved plotting functionality, this package also introduces the capability to create or customize geneset databases. The package includes all geneset databases published by the Broad Institute as well as 6 different transcription factor databases and two databases containing ligands and receptor pairs.

Three data files are needed in order to perform the gene set enrichment analysis. The **GSEApLOTS** function executes the analysis producing a summary results file for each phenotype in the working directory. Within the current working directory a folder will be created with the name assigned to doc.string in the **GSEApLOTS** function call. The results for each gene set within the database will be saved into this folder. The **plot.ES** function saves a pdf of all generated enrichment graphs into the working directory.

Load the Package

The following script will load the package into the user's library:

```
remotes::install_github("kelsiereinaltt/GSEApLOT")
library(devtools)
library(GSEApLOT)
```

Geneset Database Options

To improve on the GSEA code published by the Broad Institute, additional geneset database options were included within this package. The GSEApLOT function was made to be compatible with other data sets, whether the data comes from another source or is modified by the user. Detailed below.

Identifying and Loading Existing Databases

Calling the key included within the package shows which databases are currently saved in the package.

```
data(key)
head(key)
```

##	File	Description
## 1	C1.gs	positional all
## 2	C2.cgp.gs	curated chemical and genetic perturbations
## 3	C2.cp.biocarta.gs	curated canonical pathways biocarta
## 4	C2.cp.gs	curated canonical pathways
## 5	C2.cp.kegg.gs	curated canonical pathways kegg
## 6	C2.cp.reactome.gs	curated canonical pathways reactome
##	Source	
## 1	http://software.broadinstitute.org/gsea/msigdb	
## 2	http://software.broadinstitute.org/gsea/msigdb	
## 3	http://software.broadinstitute.org/gsea/msigdb	
## 4	http://software.broadinstitute.org/gsea/msigdb	
## 5	http://software.broadinstitute.org/gsea/msigdb	
## 6	http://software.broadinstitute.org/gsea/msigdb	

The key has three columns: file, description, and source. The description corresponding to a file makes the databases more accessible by detailing what a database contains. Users can identify a dataset of interest by reading its description then loading the database using its file name.

database_key may also return the descriptions for all databases when the function input is "all". Given the descriptions of all databases, the user may re-use the function with the description of interest to find the file of interest.

```
descriptive_names=database_key("all")
head(descriptive_names)

## [1] "positional all"
## [2] "curated chemical and genetic perturbations"
## [3] "curated canonical pathways biocarta"
## [4] "curated canonical pathways"
## [5] "curated canonical pathways kegg"
## [6] "curated canonical pathways reactome"
```

The function **database_key** returns the file name of a database given the description of a database.

```
GO_mf_filename=database_key("GO molecular function")
GO_mf_filename

## [1] "C5.mf.gs"

GO_mf_filename = database_key(descriptive_names[20]) #kelsie -test this
```

Once the file name is identified, the database of interest may be loaded into the workspace.

```
data(hallmark.gs)
data(C1.gs) #make this hhallmark
```

To check if a geneset database includes data of interest, all sets within a database can be viewed using the **get_genesets** function.

```
sets=get_genesets(hallmark.gs)
head(sets)

## [1] "HALLMARK_TNFA_SIGNALING_VIA_NFKB" "HALLMARK_HYPOXIA"
## [3] "HALLMARK_CHOLESTEROL_HOMEOSTASIS" "HALLMARK_MITOTIC_SPINDLE"
## [5] "HALLMARK_WNT_BETA_CATENIN_SIGNALING" "HALLMARK_TGF_BETA_SIGNALING"
```

Setting up Data for Gene Set Enrichment Analysis using GSEAplot

Gene Set Enrichment Analysis requires the user to provide three dataframes: an expression file, a phenotype file, and a geneset database. The user may use dataframes provided in this package, or create his/her own dataframes. To load data included in the package, call the desired R dataset within the **data()** function. If you have not saved data to the package, simply assign your data to the corresponding function parameter making sure that it follows the correct data formatting. Instructions for how to prepare each file are included on the Broad Institute page for their GSEA analysis function.

Formatting Data

The function parameter `input.ds.name` requires *differential expression* data. The following is an example of how gene expression data included this package can be loaded:

```
data(aagmex_expr)
expr.input=aagmex_expr
expr.input[1:4,1:6]
```

```
##      GSM2520983 GSM2520984 GSM2520985 GSM2520986 GSM2520991 GSM2520992
## 7A5      6.416175  6.303225  6.292974  6.288139  6.354787  6.365153
## A1BG      6.442692  6.474838  6.538247  6.453630  6.498097  6.404469
## A1CF      6.429408  6.354601  6.380667  6.369513  6.386864  6.357330
## A26C3     6.431504  6.370737  6.325280  6.377389  6.376220  6.388600
```

The parameter `input.cls.name` takes the *phenotype* data, which maps the differential expression data to one of two phenotypes. The following is an example of how phenotype data included this package can be loaded:

```
data(aagmex_pheno)
pheno.input=aagmex_pheno
pheno.input$phen
```

```
## [1] "Female" "Male"
```

```
head(pheno.input)
```

```
## $phen
## [1] "Female" "Male"
##
## $class.v
## [1] 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
## [38] 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
## [75] 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
## [112] 0 0 0 0 0 0 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
## [149] 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
## [186] 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
## [223] 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
```

The parameter `input.gs.name` takes a *geneset* or a geneset database (collection of genesets). The following is an example of how geneset data included this package can be loaded:

```
data(hallmark.gs)
gene.set.input=hallmark.gs
```

Create a Phenotype Input Function

The `pheno.input` parameter for **GSEA.plot** requires a list with a `phen` element and a `class.v` element.

```
data(aagmex_pheno)
aagmex_pheno$phen
```

```
## [1] "Female" "Male"
```

```
head(aagmex_pheno$class.v)
```

```
## [1] 0 0 0 0 0 0
```

The `create_phenoinput` function takes a character column with labels for the phenotype of each sample and makes it into a working `pheno.input` for the `GSEA.plot` function.

```
data(gtex_ann)
head(gtex_ann)
```

```
## [1] "Female" "Female" "Female" "Female" "Female" "Female"
```

```
pheno.input=create_phenoinput(gtex_ann)
pheno.input$phen
```

```
## [1] "Female" "Male"
```

```
head(pheno.input$class.v)
```

```
## [1] 0 0 0 0 0 0
```

```
data(gtex_expr)
expr.input=gtex_expr
```

GSEA Analysis

GSEAPlots: Function and Parameters

In addition to the three data frames, the GSEA analysis requires a number of other parameters.

AnalysisandInputParamaters :

- *input.ds.name* - references differential expression data.
- *input.cls.name* - references phenotype input.
- *gene.set.input* - references geneset database.
- *doc.string* - name of the folder where individual geneset results are saved to. The folder is created if it does not exist.
- *nperm* - determines the number of permutations. Permutations are used to analyze the significance of the association between gene expression and phenotype, and determine whether the associations are more significant than random ones.

- *abs.val* - When set to true (T), genes are ranked according to the absolute value of their signal to noise ratio. This value is set to false (F) by default.

PlottingParameters :

- *bar_percent*: the size of the enrichment tick mark relative to the size of the enrichment score plot.
- *gap_percent*: the size of the gap between the running enrichment score graph and the tick marks relative to the size of the enrichment score plot.
- *under_percent*: the size of the space below the tick marks relative to the size of the enrichment score plot.
- *upper_percent*: the size of the space above the enrichment plot relative to the size of the enrichment score plot.
- *color_line*: color of the line in the enrichment score plot.
- *color_tick*: color of the tick marks that indicate a gene symbol in the expression data that is within the given geneset.

The following script is an example of how to properly use the **GSEAplots** function. Please note that on a computer with a MacOS 10.14 (Mojave) operating system, 2.2 GHz Intel Core i7 processor, and 16 GB of memory, the run time for this function with its current parameters is approximately 3.244399 minutes. kelsie NOTE.

```
start_time <- Sys.time()

pp = GSEAplots(input.ds.name=expr.input,
               input.cls.name=pheno.input,
               gene.set.input=gene.set.input,
               doc.string="GSEA_plots",
               nperm=1000,
               abs.val=F,
               bar_percent=0.1,
               gap_percent=0.1,
               under_percent=0.1,
               upper_percent=0.1,
               color_line="black",
               color_tick="black")

end_time <- Sys.time()
end_time-start_time
```

Function Outputs

Outputs of **GSEAplots** include reports, enrichment scores, plots, and the variables `gene.set.reference.matrix` and `gene.set.leading`. All outputs will be saved to the working directory.

Reports

A report is generated for each phenotype (output tables include *report1* and *report2*). The table shows all of the information available in the report file. An example of the data is provided below. In this case, the first row contains the data for the geneset HALLMARK_PANCREAS_BETA_CELLS. The third column is the source which is omitted for easy viewing.

```
pp$report1[1,-3]
```

```
##                               GS SIZE      ES      NES NOM p-val FDR q-val
## 1 HALLMARK_OXIDATIVE_PHOSPHORYLATION 187 0.70426 2.3103      0      0
##   FWER p-val Tag % Gene % Signal FDR (median) glob.p.val
## 1      0 0.561 0.135 0.489      0      0
```

Enrichment Scores

The list *ES* provides access to the enrichment score data used to generate the enrichment plots along with its corresponding gene symbol.

```
data1=pp$ES[[50]] #kelsie
head(data1)
```

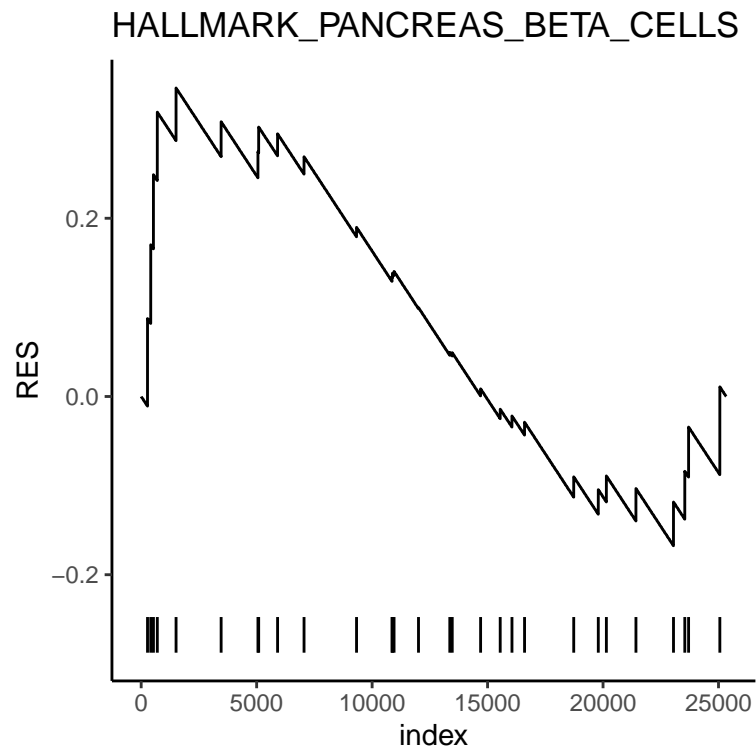
```
##   index gene.symbol      RES EStag
## 1     1   EIF1AXP1 -3.952101e-05    0
## 2     2   EIF2S3L -7.904201e-05    0
## 3     3   ANKRD30A -1.185630e-04    0
## 4     4     SAA1 -1.580840e-04    0
## 5     5   ATP8A2P1 -1.976050e-04    0
## 6     6  LINC00993 -2.371260e-04    0
```

Access to enrichment score data allows for deeper analysis of the gene set enrichment. Users of the package may create their own plotting functions using the data.

The following describes the procedure for producing plots in the package as well as an example:

Enrichment scores are plotted for each position within the gene set. A positive enrichment score indicates correlation with the first phenotype which in the case of the aagmex pheno data is Female and a negative enrichment score indicates correlation with the second phenotype, Male. The gene symbols within your gene set that are differentially expressed are indicated with a tick mark under this plot. In the enrichment score data these symbols are indicated with an EStag of 1. To specify the dimensions and positions of the tick marks a data frame is created to create line segments that have a regular height and are positioned at each position with an EStag of 1. The segments should be vertical so the vx or change in x over the segment should be 0. The y position in the data frame dictates the starting y position of the mark. It is currently set to 0.12 below the minimum of of the enrichment plot. vy dictates the end of each mark to be 0.04 higher than that starting position. In the plotting function, 0.12 and 0.04 are variables that can be modified by the bar_percent and under_percent parameters. Using ggplot, the enrichment score plot and tick marks are combined into one plot called p.

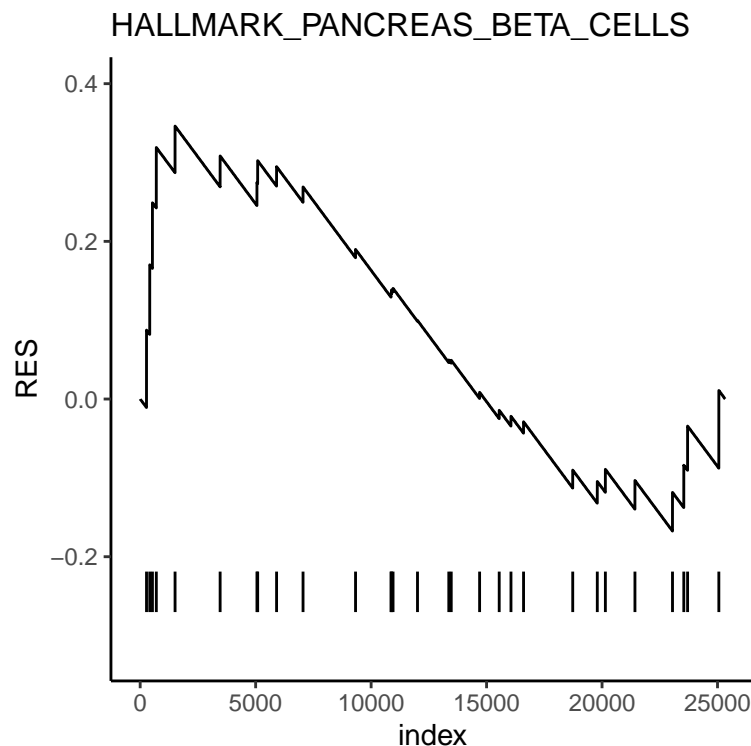
```
enrich_ind=which(data1$EStag==1)
d=data.frame(x=enrich_ind, y=matrix(min(data1$RES)-0.12,length(enrich_ind),1),
            vx=matrix(0,length(enrich_ind),1), vy=matrix(0.04,length(enrich_ind),1))
p <- ggplot(data1, aes(index,RES))+geom_line(col="black")
p <- p+geom_segment(data=d, mapping=aes(x=x, y=y, xend=x+vx, yend=y+vy),col="black")
p <- p+theme_classic()
p <- p+ggtitle(names(pp$gene.set.reference.matrix)[[50]])
print(p)
```



Plots

The following is an example of an enrichment plot from the first geneset:

```
pp$plots[[50]]
```

Use the function `plot.ES` to save a pdf of enrichment plots for all sets in the gene set database. The pdf will save to the working directory.

```
plot.ES(list.of.plots=pp$plots,plotname="GSEA_plots")
```

gene.set.reference.matrix and gene.set.leading

gene.set.reference.matrix is a table of the gene sets being studied.

```
names(pp$gene.set.reference.matrix)[[50]]
```

```
## [1] "HALLMARK_PANCREAS_BETA_CELLS"
```

```
head(pp$gene.set.reference.matrix[[50]])
```

```
## [1] "PAX6"      "NEUROD1" "ISL1"      "NKX2-2"   "PCSK1"    "NKX6-1"
```

gene.set.leading is a list containing the leading edge set (set of genes whose differential expression is most related to a change in phenotype) for each gene set. If differential expression is more strongly related to a change in the first phenotype, then the maximum absolute value of the running enrichment score (RES) will be positive. In this case the leading edge set is all of the enriched genes (have a tick mark) that are to the left of and including the max if it is enriched. If differential expression is more related to the second phenotype then the largest running enrichment score is negative. In this case the leading edge set is all of the enriched genes that are to the right of and including the max if it is enriched.

The following example displays the leading gene symbols for one of gene sets:

```
names(pp$gene.set.reference.matrix)[[50]]
```

```
## [1] "HALLMARK_PANCREAS_BETA_CELLS"
```

```
head(pp$gene.set.leading[[50]])
```

```
## [1] "VDR" "SRPRB" "AKT3" "FOXO1"
```

Running GSEA on a Novel Geneset

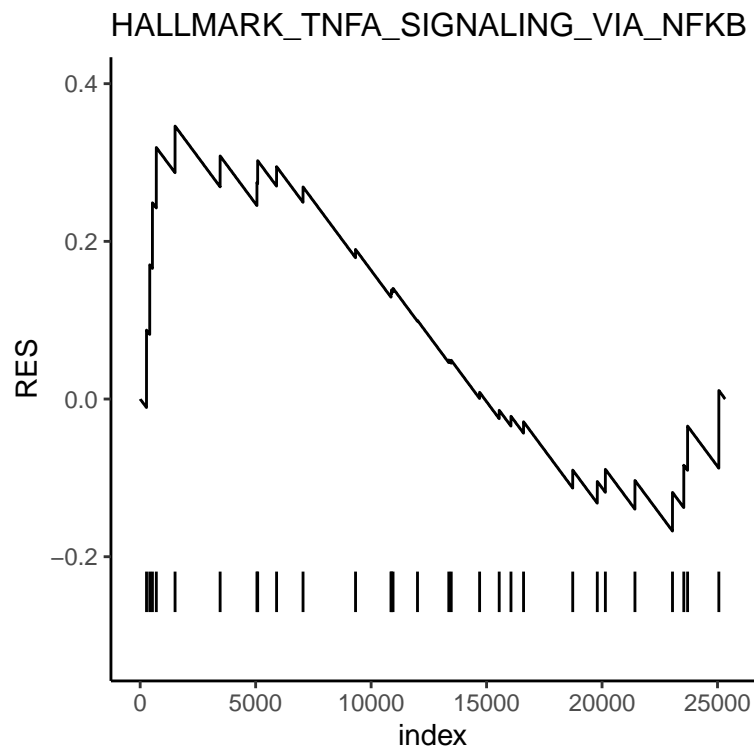
This package introduces the ability to create your own database. This requires the name of the list to be the names of the gene sets. The items in each element of the list need to include the description/source and gene symbols. For example, the names of a potential list would be the transcription factors and the elements would be the genes affected by that specific transcription factor. The `create_geneset_db()` function then converts this list into a format compatible with the **GSEAplots** function.

The following script creates a custom geneset database (containing the first three genesets in the Hallmark database) that is compatible with the **GSEAplots** function:

```
sets=get_genesets(hallmark.gs)
symbols=get_genesymbols(hallmark.gs)
entry1=c("source_1",symbols$HALLMARK_PANCREAS_BETA_CELLS)
entry2=c("source_2",symbols$HALLMARK_HYPOXIA)
entry3=c("source_3",symbols$HALLMARK_CHOLESTEROL_HOMEOSTASIS)
custom_db=list(entry1,entry2,entry3) #kelsie - if vignette doesn't run check here
names(custom_db)=c(sets[1],sets[2],sets[3]) #and here
gene.set.input=create_geneset_db(custom_db) #and here
```

```
custom_results= GSEAplots(input.ds.name=expr.input,
                           input.cls.name=pheno.input,
                           gene.set.input=gene.set.input,
                           doc.string="Aagmex_hall13",
                           nperm=1000,
                           bar_percent=0.1,
                           gap_percent=0.1,
                           under_percent=0.1,
                           upper_percent=0.1,
                           color_line="black",
                           color_tick="black",
                           abs.val=F)
```

```
custom_results$plots[[1]]
```



Add to an Existing Gene Set Database

This new function allows you to add genesets into an existing database. The database input needs to be already in the tab separated format that is used in the GSEAplots function. The genesets to add have a separate line for each geneset with the geneset name, source, and gene symbols. These will be converted into the tab separated format in the `add_to_database()` function.

To add the geneset HALLMARK_KLF14 create a row vector with the name of the set then the source and then the gene symbols.

```
data(transf)
data(transm)
tx = t(c("KLF14_targets", "custom", transf[,1], transm[,1]))
```

The `add_to_database` function adds this set to the hallmark sets and formats it for use in the function. An example of it being used is in the section below.

```
gene.set.input=add_to_database(database=hallmark.gs, addition=tx)
```

Ligand and Receptor Databases

This package offers new databases: transcription factor databases, and Kadoyama ligand and receptor databases.

Introduction of Transcription Factor Databases

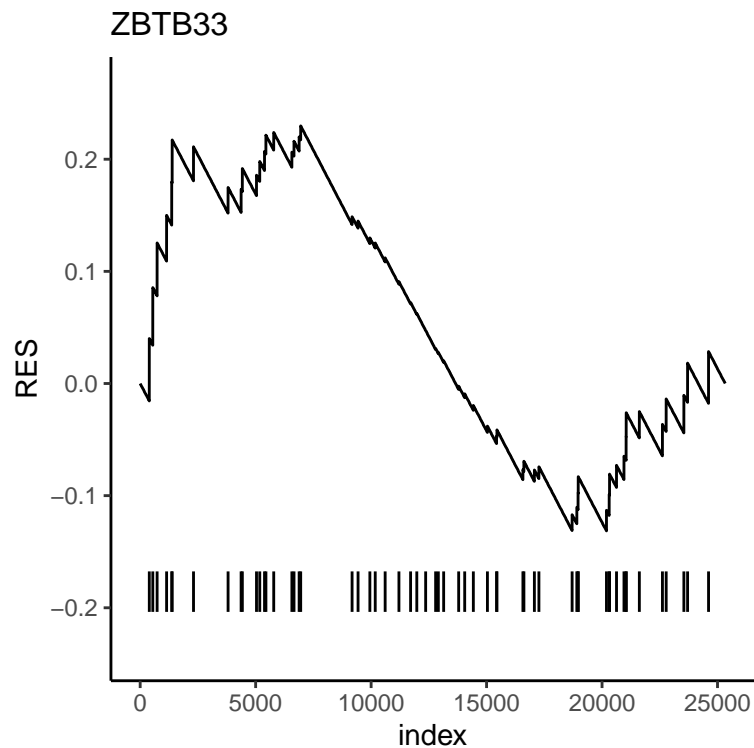
The following transcription factor databases were introduced: Neph2012, ENCODE, ITFP, Marbach2016, TRRUST, and TRED (<https://github.com/slowkow/tftargets>). Kelsie - reference the papers in the github

```
data(ENCODE.db)
gene.set.input=ENCODE.db

enrichment_TF= GSEAplots(input.ds.name=expr.input,
                          input.cls.name=pheno.input,
                          gene.set.input=gene.set.input,
                          doc.string="Aagmex_TF",
                          nperm=1000,
                          bar_percent=0.1,
                          gap_percent=0.1,
                          under_percent=0.1,
                          upper_percent=0.1,
                          color_line="black",
                          color_tick="black",
                          abs.val=F)
```

```
enrichment_TF$plots[1] #kelsie - plot the most significant one
```

```
## [[1]]
```



Introduction of Kadoki Ligand and Receptor Database

NOTE: Sarah does not specify where she got this data from - is it possible to find out?? [https://www.cell.com/cell/pdf/S0092-8674\(17\)30949-2.pdf](https://www.cell.com/cell/pdf/S0092-8674(17)30949-2.pdf)

```

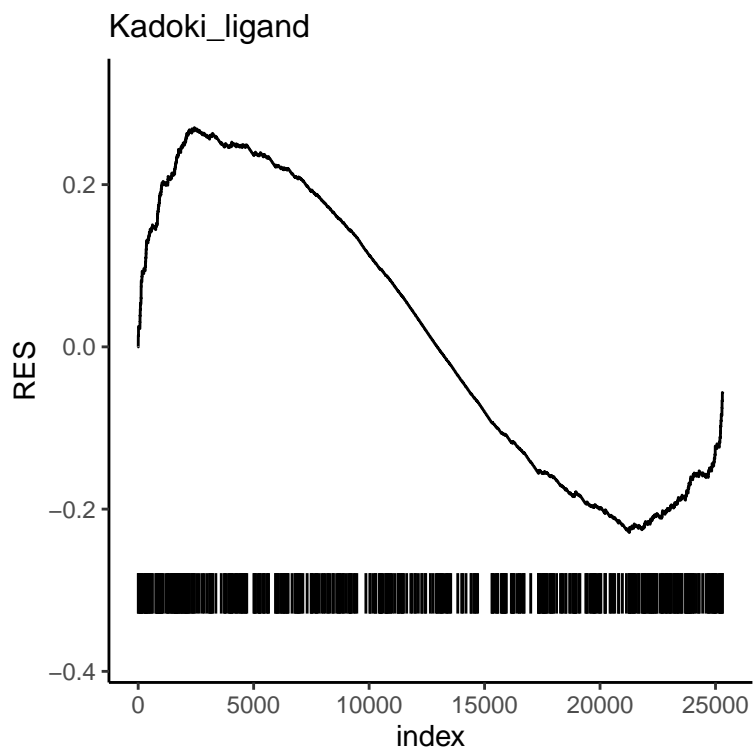
data(Kadoki_ligands.db)
gene.set.input=Kadoki_ligands.db

ligand_results=GSEAplots(input.ds.name=expr.input,
                          input.cls.name=pheno.input,
                          gene.set.input=gene.set.input,
                          doc.string="Aagmex_ligands",
                          nperm=1000,
                          bar_percent=0.1,
                          gap_percent=0.1,
                          under_percent=0.1,
                          upper_percent=0.1,
                          color_line="black",
                          color_tick="black",
                          abs.val=F)

```

```
ligand_results$plots[1]
```

```
## [[1]]
```



```

data(Kadoki_receptors.db)
gene.set.input=Kadoki_receptors.db

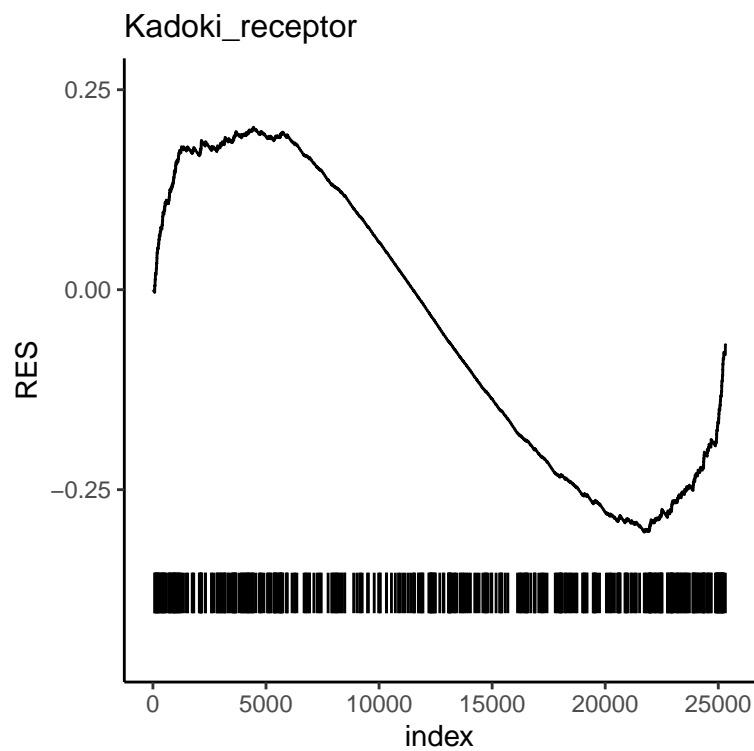
receptor_results=GSEAplots(input.ds.name=expr.input,
                            input.cls.name=pheno.input,
                            gene.set.input=gene.set.input,
                            doc.string="Aagmex_receptors",

```

```
nperm=1000,  
bar_percent=0.1,  
gap_percent=0.1,  
under_percent=0.1,  
upper_percent=0.1,  
color_line="black",  
color_tick="black",  
abs.val=F)
```

```
receptor_results$plots[1] #kelsie - plot the most significant one
```

```
## [[1]]
```



References

- Subramanian, Tamayo, et al. (2005), PNAS 102, 15545-15550, <http://www.broad.mit.edu/gsea/>
- <https://github.com/slowkow/tftargets>
- <https://www.ncbi.nlm.nih.gov/pubmed/28942919>
- [https://www.cell.com/cell/pdf/S0092-8674\(17\)30949-2.pdf](https://www.cell.com/cell/pdf/S0092-8674(17)30949-2.pdf)