

GSEA分析

软件和基因集下载

(<http://software.broadinstitute.org/gsea/downloads.jsp>)

Downloads

Software

There are several options for GSEA software. All options implement exactly the same algorithm. Usage recommendations and installation instructions are listed below. Current Java implementations of GSEA require Java 8.

See the [license terms page](#) for details about the license for the GSEA software and source code. Please note that the license terms vary for different versions of the software.

javaGSEA Desktop Application

- ▶ Easy-to-use graphical user interface.
- ▶ Runs on any desktop computer (Windows, macOS, Linux etc.) that supports Java 8. **Oracle Java is recommended as there are known issues when running with OpenJDK.**
- ▶ Produces richly annotated reports of enrichment results.
- ▶ This release is open source under a [BSD-style license](#). The source is available on our [GitHub repository](#). The changes are noted in the [Release Notes](#).
- ▶ **We recommend using a memory configuration smaller than your computer's total memory.**

Launch with

1GB (for 32 or 64-bit Java) memory:

Launch

根据内存大小选择

MSigDB

Use the following links to download individual gene set collections or the complete Molecular Signatures Database (MSigDB). For details on the MSigDB gene set collections refer to the [Molecular Signatures Database](#) page.

See the [license terms page](#) for details about the license for MSigDB. Please note that the license terms vary for different versions of MSigDB, and that certain gene sets have special access terms.

All gene sets	Current MSigDB gene sets, gene symbols	msigdb.v6.1.symbols.gmt
	Current MSigDB gene sets, Entrez IDs	msigdb.v6.1.entrez.gmt
	Current MSigDB xml file	msigdb_v6.1.xml
h: hallmark gene sets	hallmark gene sets, gene symbols	h.all.v6.1.symbols.gmt
	hallmark gene sets, Entrez IDs	h.all.v6.1.entrez.gmt
c1: positional gene sets	positional gene sets, gene symbols	c1.all.v6.1.symbols.gmt
	positional gene sets, Entrez IDs	c1.all.v6.1.entrez.gmt
c2: curated gene sets	all curated gene sets, gene symbols	c2.all.v6.1.symbols.gmt
	all curated gene sets, Entrez IDs	c2.all.v6.1.entrez.gmt
	chemical and genetic perturbations, gene symbols	c2.cgp.v6.1.symbols.gmt
	chemical and genetic perturbations, Entrez IDs	c2.cgp.v6.1.entrez.gmt
	all canonical pathways, gene symbols	c2.cp.v6.1.symbols.gmt
	all canonical pathways, Entrez IDs	c2.cp.v6.1.entrez.gmt
	BioCarta gene sets, gene symbols	c2.cp.biocarta.v6.1.symbols.gmt
	BioCarta gene sets, Entrez IDs	c2.cp.biocarta.v6.1.entrez.gmt
	KEGG gene sets, gene symbols	c2.cp.kegg.v6.1.symbols.gmt
	KEGG gene sets, Entrez IDs	c2.cp.kegg.v6.1.entrez.gmt
	Reactome gene sets, gene symbols	c2.cp.reactome.v6.1.symbols.gmt
	Reactome gene sets, Entrez IDs	c2.cp.reactome.v6.1.entrez.gmt

输入数据准备

1. 表达矩阵。常见表达矩阵格式，**tab**键分割，**txt**格式，第一列为基因名字（名字与注释数据库一致，同为GeneSymbol或EntrezID或其它自定义

名字)，第一行为标题行，含样品信息。也可为gct文件，具体见

<http://blog.genesino.com/2014/08/GSEA-usages/>

[illegible]

2. 样品分组信息

```
6 2 1
# fresh old
fresh fresh fresh old old old
```

分组信息示例

基因集信息示例

Each row represents one gene set

If editing in excel, watch out for its tendency to auto-format gene sets (SEP8 becomes 8-Sep)

First column are gene set names. Duplicates are not allowed

Second column contains a brief description. Its optional – you can fill in a dummy field (e.g. “na”)

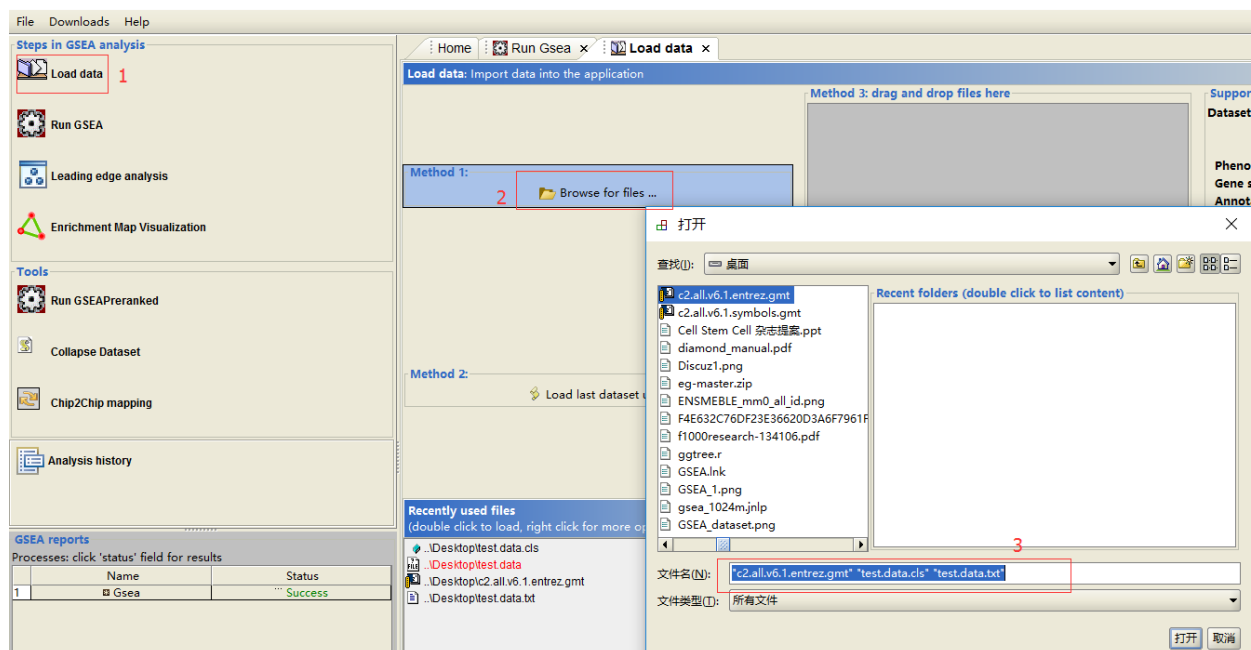
Unequal lengths (i.e # of genes) is allowed

	A	B	C	D	E	F	G
1	chr10q24	Cytogenetic band	PITX3	SPFH1	NEURL	C10orf12	NDUFB8
2	chr5q23	Cytogenetic band	ALDH7A1	IL13	8-Sep	ACSL6	
3	chr8q24	Cytogenetic band	HAS2	LRRC14	TSTA3	DGAT1	RECQL4
4	chr16q24	Cytogenetic band	RPL13	GALNS	FANCA	CPNE7	COTL1
5	chr13q14	Cytogenetic band	AKAP11	ARL11	ATP7B	C13orf1	C13orf9
6	chr7p21	Cytogenetic band	ARL4A	SCIN	GLCCI1	SP8	SOSTDC1
7	chr10q23	Cytogenetic band	SNCG	FER1L3	C10orf116	HHEX	TNKS2
8	chr14q12	Cytogenetic band	C14orf125	FOXG1C	HECTD1	SCFD1	AP4S1
9	chr13q13	Cytogenetic band	ALG5	RFXAP	DCAMKL1	MAB21L1	STOML3
10	chr1p34	Cytogenetic band	JMJD2A	MRPS15	HIVEP3	GJB3	CDCA8
11	chr10q21	Cytogenetic band	MBL2	C10orf70	DNAJC12	BICC1	CXXC6

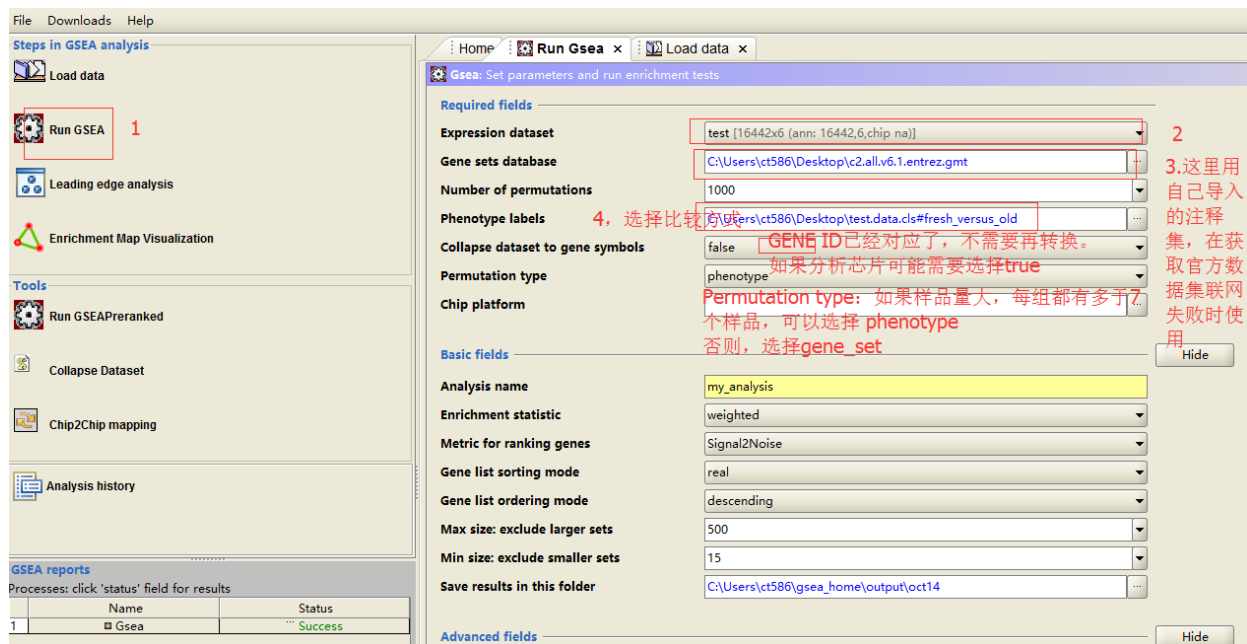
GMT format is convenient to store large databases of gene sets. For a handful of sets (<256) the gmx format offers greater excel-editability

软件运行 (每一步的步骤如有不明确的参考文后第一个链接)

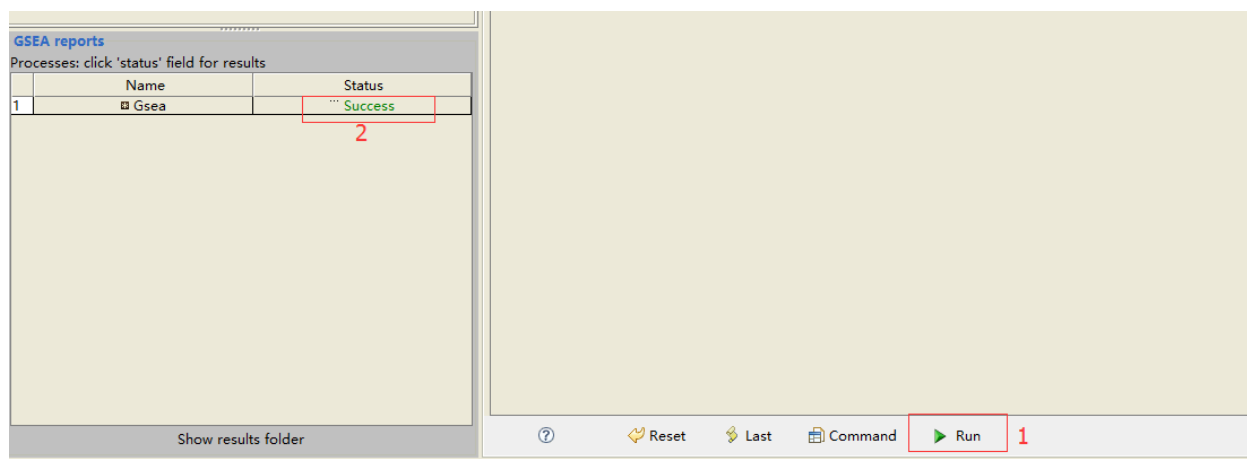
1. 导入数据



2. 运行GSEA (若每组样品都有多于7个样品, 则Permutation type选择 **phenotype**, 结果理论上更好; 否则选择 **gene_set**)



3. 设置好参数后，点击正下方的run，等待运行结束，左侧出现success



4. 点击success，查看结果

Enrichment in phenotype: fresh (3 samples)

- 3013 / 3449 gene sets are upregulated in phenotype **fresh**
- 2188 gene sets are significant at FDR < 25%
- 1723 gene sets are significantly enriched at nominal pvalue < 1%
- 1723 gene sets are significantly enriched at nominal pvalue < 5%
- [Snapshot](#) of enrichment results
- Detailed [enrichment results in html](#) format
- Detailed [enrichment results in excel](#) format (tab delimited text)
- [Guide to](#) interpret results

Enrichment in phenotype: old (3 samples)

- 436 / 3449 gene sets are upregulated in phenotype **old**
- 0 gene sets are significant at FDR < 25%
- 73 gene sets are significantly enriched at nominal pvalue < 1%
- 73 gene sets are significantly enriched at nominal pvalue < 5%
- [Snapshot](#) of enrichment results
- Detailed [enrichment results in html](#) format
- Detailed [enrichment results in excel](#) format (tab delimited text)
- [Guide to](#) interpret results

Dataset details

- The dataset has 16398 features (genes)
- No probe set => gene symbol collapsing was requested, so all 16398 features were used

Gene set details

- Gene set size filters (min=15, max=500) resulted in filtering out 1289 / 4738 gene sets
- The remaining 3449 gene sets were used in the analysis
- List of [gene sets used and their sizes](#) (restricted to features in the specified dataset)

顺着网页的导航一步步去查看结果，有耐心就好。主要解释下，最常见的这种图。

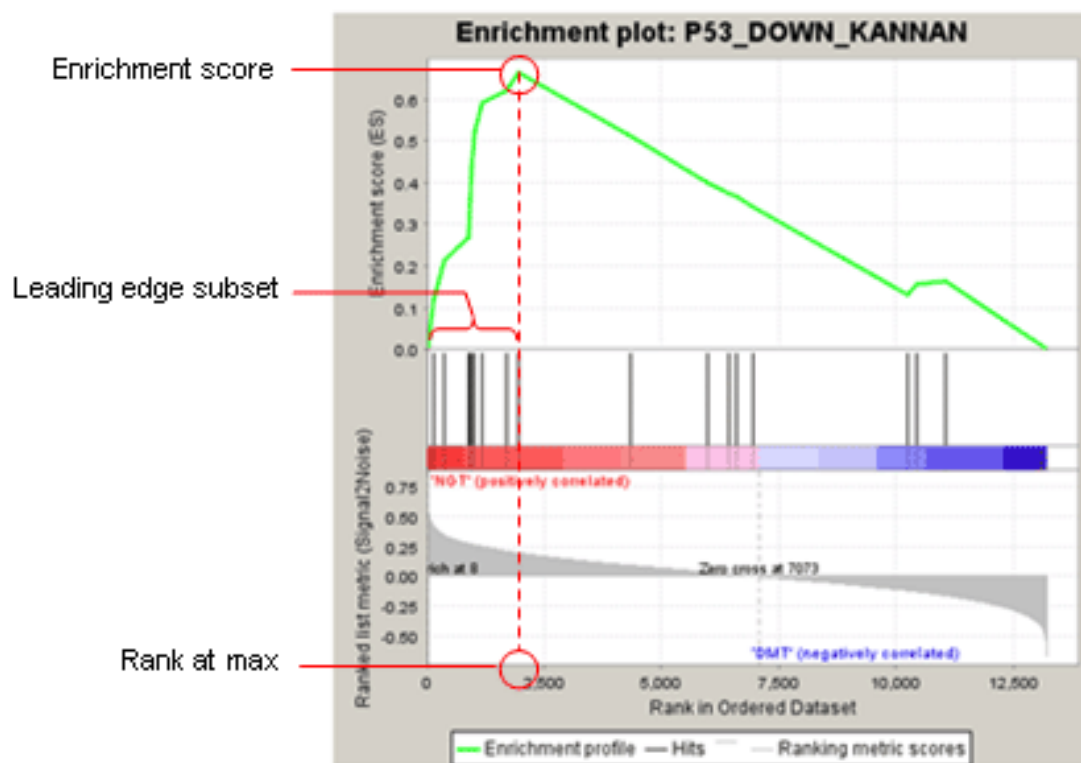


Fig 1: Enrichment plot: P53_DOWN_KANNAN
Profile of the Running ES Score & Positions of GeneSet Members on the Rank Ordered List

1. 图最上面部分展示的是ES的值计算过程，从左至右每到一个基因，计算出一个ES值，连成线。最高峰为富集得分(ES)。在最左侧或最右侧有一个特别明显的峰的基因集通常是感兴趣的基因集。
2. 图中间部分每一条线代表基因集中的一个基因，及其在基因列表中的排序位置。
3. 最下面部分展示的是基因与表型关联的矩阵，红色为与第一个表型(MUT)正相关，在MUT中表达高，蓝色与第二个表型(WT)正相关，在WT中表达高。
4. **Leading-edge subset** 对富集得分贡献最大的基因成员。若富集得分为正值，则是峰左侧的基因；若富集得分为负值，则是峰右侧的基因。
5. **FDR GSEA**默认提供所有的分析结果，并且设定**FDR<0.25**为可信的富集，最可能获得有功能研究价值的结果。但如果样品数目少，而且选择了**gene_set**作为**Permumation type**则需要使用更为严格的标准，比如**FDR<0.05**。

Leading-edge分析

主要对筛选感兴趣的基因有意义；选择一个或多个显著富集的基因集，查看其内 **Leading-edge** 基因的表达和重叠状态。

1. Leading edge analysis

2. Select a GSEA result from the application cache

3. Load GSEA Results

positive phenotype: na pos negative phenotype: old

Filter Gene Sets

3449 out of 3449 gene sets

Gene Set	Size	ES	NES	NOM p-val	FDR q-val	FWER p-val	Rank at Max	Leading Edge
PID_ER_NONGENOMIC...	35	0.586	2.279	0	0.053	0	3,325 tags=40%, list=20%	
PID_ANGIOPOIETIN_R...	43	0.579	2.254	0	0.053	0	3,544 tags=42%, list=23%	
KAMMINGA_SENESCENCE...	35	0.566	2.16	0	0.053	0.053	1,883 tags=37%, list=11%	
KEGG_PATHOGENIC_E...	48	0.488	2.159	0	0.053	0.053	1,507 tags=25%, list=9%	
SCIAH_INVERSED_TAR...	27	0.545	2.152	0	0.053	0.053	1,659 tags=26%, list=10%	
ALCALA_APOPTOSIS...	71	0.517	2.107	0	0.053	0.053	2,741 tags=55%, list=17%	
REACTOME_THROMB...	23	0.533	2.105	0	0.053	0.053	2,684 tags=30%, list=16%	
TOMLINSON_PROSTATE...	48	0.487	2.072	0	0.053	0.053	3,723 tags=44%, list=23%	
TARTE_PLASMA_CELL...	27	0.544	2.072	0	0.053	0.053	3,903 tags=48%, list=24%	
PID_IGF1_PATHWAY...	29	0.572	2.07	0	0.053	0.053	3,562 tags=45%, list=22%	
STEIN_ESRRA_TARGETS...	498	0.371	2.07	0	0.053	0.053	4,439 tags=34%, list=27%	
YAO_TEMPORAL_RESP...	137	0.419	2.069	0	0.053	0.053	3,764 tags=33%, list=23%	
BANDRES_RESPONSE...	24	0.628	2.054	0	0.053	0.053	2,343 tags=42%, list=14%	
LU_TUMOR_ANGIOGE...	22	0.713	2.054	0	0.053	0.053	2,236 tags=45%, list=14%	
PROVENZANI_METAS...	172	0.391	2.038	0	0.053	0.053	3,888 tags=28%, list=24%	
STEIN_ESRRA_TARGET...	25	0.488	2.03	0	0.053	0.053	4,499 tags=50%, list=27%	
REACTOME_G0_AND...	23	0.501	2.013	0	0.053	0.053	2,779 tags=30%, list=17%	
HU_ANGIOGENESIS_D...	34	0.481	2.012	0	0.053	0.053	2,956 tags=35%, list=18%	
REACTOME_THROMB...	17	0.539	2.01	0	0.053	0.053	2,684 tags=35%, list=16%	
WANG_METHYLATED...	29	0.629	2.005	0	0.053	0.053	1,284 tags=28%, list=8%	
BIOCARTA_BIOPEPTID...	39	0.555	2.002	0	0.053	0.053	1,530 tags=26%, list=9%	
OZEN_MIR125B1_TAR...	24	0.564	2.001	0	0.053	0.053	2,574 tags=29%, list=16%	
PUIFFE_INVASION_IN...	75	0.442	1.998	0	0.053	0.053	2,636 tags=29%, list=16%	
PID_VEGF1_2_PATH...	62	0.45	1.995	0	0.053	0.053	3,562 tags=35%, list=22%	
KEGG_PEROXISOME...	66	0.563	1.994	0	0.053	0.053	3,146 tags=39%, list=19%	
SA_TRKA_RECEPTOR...	15	0.674	1.989	0	0.053	0.053	2,088 tags=47%, list=13%	
CROMER_METASTAS...	69	0.447	1.987	0	0.053	0.053	3,143 tags=29%, list=19%	
REACTOME_AMINE_C...	18	0.662	1.984	0	0.053	0.053	1,871 tags=33%, list=11%	
REACTOME_CTNNB1...	15	0.585	1.981	0	0.053	0.053	1,329 tags=27%, list=8%	
PID_IL2_STATS_PATH...	22	0.55	1.978	0	0.053	0.053	2,982 tags=36%, list=18%	
REACTOME_PURINE...	30	0.544	1.974	0	0.053	0.053	2,772 tags=37%, list=17%	
CHANDRAN_METAS...	36	0.493	1.974	0	0.053	0.053	1,834 tags=25%, list=11%	
REACTOME_SHC1_EV...	19	0.576	1.972	0	0.053	0.053	1,530 tags=26%, list=9%	
ZHANG_BREAST_CAN...	126	0.406	1.967	0	0.053	0.053	3,591 tags=31%, list=22%	
ACOSTA_PROLIFERAT...	81	0.434	1.966	0	0.053	0.053	4,663 tags=43%, list=28%	
SHEPARD_CRUSH_AN...	180	0.411	1.96	0	0.053	0.053	3,035 tags=26%, list=19%	

4. Run leading edge analysis

5. Build HTML Report

6. For 6 selected gene sets

基因集中leading-edge基因的重合度，颜色越深重合度越高

基因集内Leading-edge基因的表达值，红到白到蓝代表高到底

Number Of Gene Sets

Gene

Number of Occurrences

Jaccard

Bin Width: 0.02 Update

MSigDB

GSEA团队整理好的基因集，可用于注释，也可下载下来搜寻自己感兴趣的方向的基因作为一个补充。每个注释都提供了基于**Gene Symbol**和**Entrez ID**的索引表格。

参考

1. 较早记录的一篇GSEA的使用，有脚本可以转换表达矩阵为gct, cls文件作为GSEA的输入。文档为英文，但软件操作步骤还算详细，可配合着看。

<http://blog.genesino.com/2014/08/GSEA-usages/>

2. 最开始学习的教程，每一步操作都比较详细。

<http://www.baderlab.org/Software/EnrichmentMap/Tutorial>

3. GSEA软件和数据集下载

<http://software.broadinstitute.org/gsea/downloads.jsp>

4. 原文对GSEA原理的讲解是很清晰的，可以读下，关键的内容也都摘录在第一个链接里。 <https://www.ncbi.nlm.nih.gov/pubmed/16199517>