

通路分析pipeline（GO+GSEA+KEGG）

卢建璋

一、总览：作用及数据背景

通路分析是常见的二代测序数据分析的下游分析方式，诸如RNA-seq，ChIP-seq，蛋白质组学等数据均可通过通路分析来获得更准确的结果。本文将以实验室RNA-seq数据（小鼠Pten/LKB1 KO：实验组 VS 小鼠Pten KO：对照组）为例演示通路分析的做法和如何生成有意义的图。

在我们进行基因差异表达分析的时候，我们会得到很多的DEG。我们很难找到一个规律去探究基因之间的联系。而通路分析则可以帮助我们获得差异通路，对RNA-seq数据进行进一步探索。

1. 数据准备：通过RNA-seq分析得到的基因表达矩阵。（如何从RNA-seq原数据得到表达矩阵请参考RNA-seq pipeline）。

注意：表达矩阵必须有以下三列：

- a. gene symbol或者gene id
- b. padj或者q.value
- c. Fold Change或者log(Fold Change)

gene_name	padj	FC
Aagab	0.0067925	1.612667548
AB124611	0.000154605	0.628886504
Abca1	0.000972927	0.432502743
Abca4	0.005908344	0.366970389
Abca4	4.01E-05	0.335148652
Abca5	0.000515771	0.631147085
Abcb10	0.001770376	0.70883929
Abcc12	0.008967195	0.569545216
Abcc5	0.00025144	1.633197706

2. R语言资源包准备

```
1 library(DOSE)
2 library(org.Hs.eg.db) # For Human data
3 library(org.Mm.eg.db) # For Mouse data
4 library(topGO)
5 library(clusterProfiler)
6 library(pathview)
7 library(tidyverse)
8 library(RColorBrewer)
9 library(enrichplot)
10 library(R.utils)
11 library(DO.db)
12 library(msigdf)
13 library(GseaVis)
```

3. 数据导入 (如果原数据是Fold Change，最好转成logFC)

```
1 # 3. 导入数据
2 data <- read.csv('Pten_LKB1 vs Pten.csv')
3 data$FC <- log2(data$FC) # 如果数据不是logFC，转化成logFC
4 colnames(data) <- c('gene_name', 'padj', 'log2FoldChange') # 规范格式，方便余下的分
```

	gene_name	padj	log2FoldChange
1	Aagab	0.006792500	0.6894491
2	AB124611	0.000154605	-0.6691284
3	Abca1	0.000972927	-1.2092188
4	Abca4	0.005908344	-1.4462644
5	Abca4	0.000040100	-1.5771270

注意：如果是microarray数据，极大概率会出现gene_name列有相同值的情况，这个时候选择保留logFC最显著的那一行数据。

```
1 # 2. 去重
2 data$sign <- case_when(
3   data$log2FoldChange >= 0 ~ "+",
4   data$log2FoldChange < 0 ~ "-",
5 )
6 data$log2FoldChange <- abs(data$log2FoldChange)
7 data <- data %>% arrange(data[,1], data[,3])
8 data <- data[-which(duplicated(data[, 1], fromLast = T)), ]
9 data[which(data$sign == '-'), 'log2FoldChange'] <- -data[which(data$sign == '-')]
10 data$sign <- NULL
11
12 # 3. 修改格式
13 rownames(data) <- data$gene_name
14 data$gene_name <- NULL
15 data <- na.omit(data)
```

最后的结果应该如图所示：

	padj	log2FoldChange
Aagab	0.006792500	0.6894491
AB124611	0.000154605	-0.6691284
Abca1	0.000972927	-1.2092188
Abca4	0.000040100	-1.5771270
Abca5	0.000515771	-0.6639518

二、Gene Ontology (GO)

GO分析要求输入所有的DEG，所以第一步要找出符合条件的DEG（筛选标准自定）

4. 找到DEG

```

1 # 4. 找到DEG
2 data$group <- case_when(
3   data$log2FoldChange > 1 & data$padj < 0.05 ~ "up",
4   data$log2FoldChange < -1 & data$padj < 0.05 ~ "down",
5   abs(data$log2FoldChange) <= 1 ~ "none",
6   data$padj >= 0.05 ~ "none"
7 )
8
9 up <- rownames(data)[data$group == "up"]
10 down <- rownames(data)[data$group == "down"]
11 diff <- c(up,down)

```

5. 基因Symbol转化基因ID

对于GO分析，最好先将gene_symbol转化成gene_id。

```

1 # 5. ID转换 (使用clusterProfiler包自带ID转换函数bitr)
2 ## 如果是人则用org.Hs.eg.db
3 up_entrez <- bitr(up,
4                   fromType = "SYMBOL",
5                   toType = "ENTREZID",
6                   OrgDb = "org.Mm.eg.db")
7
8 down_entrez <- bitr(down,
9                    fromType = "SYMBOL",
10                   toType = "ENTREZID",
11                   OrgDb = "org.Mm.eg.db")
12
13 diff_entrez <- bitr(diff,
14                    fromType = "SYMBOL",
15                   toType = "ENTREZID",
16                   OrgDb = "org.Mm.eg.db")

```

基因map失败是正常现象，但是如果失败比例太高则需要注意。

```

'select()' returned 1:1 mapping between keys and columns
Warning message:
In bitr(diff, fromType = "SYMBOL", toType = "ENTREZID", OrgDb = "org.Mm.eg.db") :
  10.06% of input gene IDs are fail to map...

```

6. 开始GO分析

GO一共分为三个类别：Molecular Function (MF)，Cellular Component (CC)，Biological Process (BP)。本文只以其中一类作为例子，其他类别的操作同理。

```

1 # 6. GO分析
2 GO_MF_diff <- enrichGO(gene = diff_entrez$ENTREZID, ## 分析所有DEG (up+down)
3                        OrgDb = org.Mm.eg.db, ## human: org.Hs.eg.db
4                        ont = "MF", ## 可选'CC'或者'BP'
5                        pAdjustMethod = "BH",
6                        pvalueCutoff = 0.05,
7                        qvalueCutoff = 0.05,
8                        readable = TRUE) ## 是否将gene ID映射到gene name
9 GO_MF_result <- GO_MF_diff@result

```

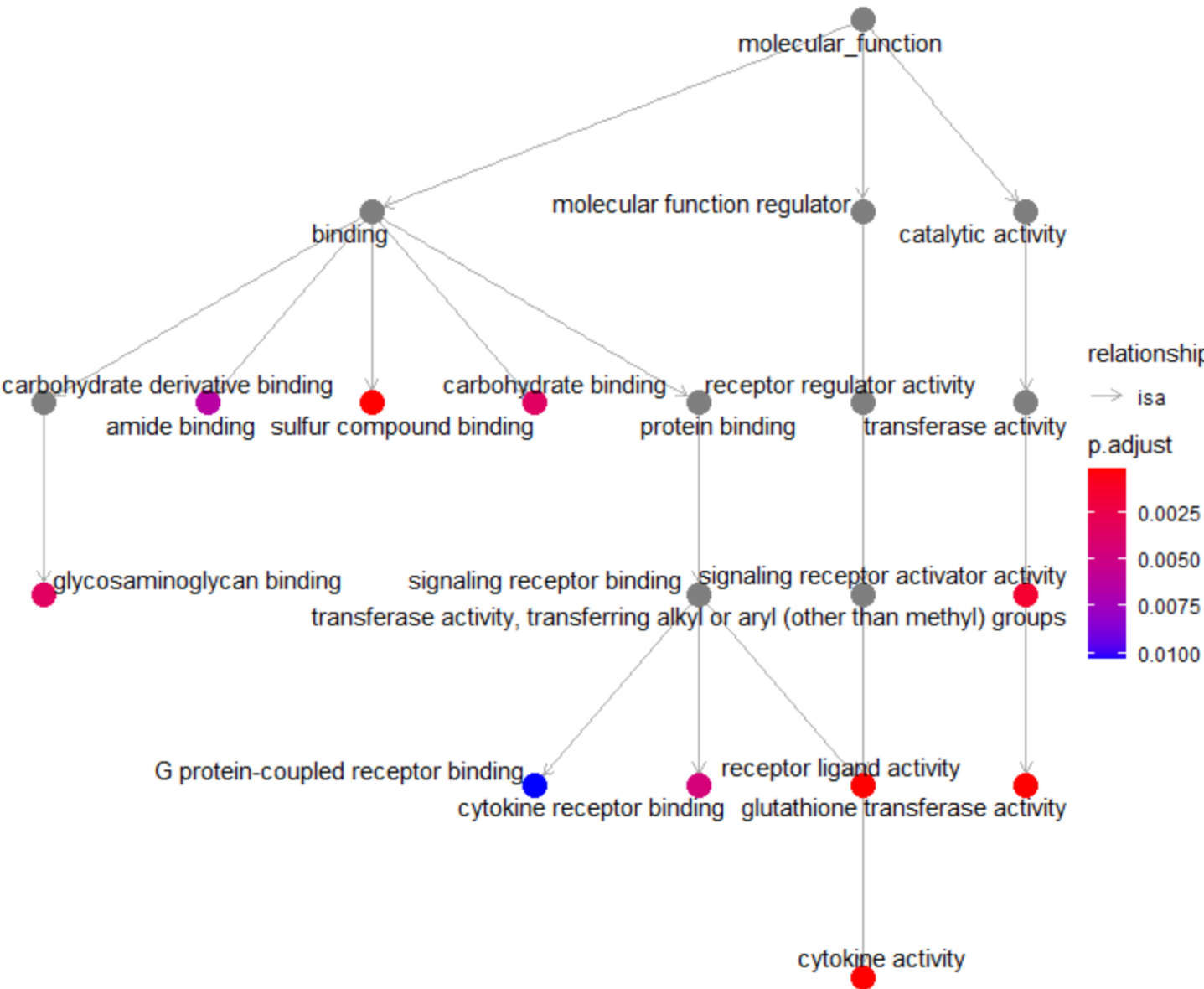
	ID	Description	GeneRatio	BgRatio	pvalue	p.adjust
GO:0048018	GO:0048018	receptor ligand activity	27/405	499/22669	3.862462e-07	0.000118
GO:0004364	GO:0004364	glutathione transferase activity	7/405	30/22669	7.876469e-07	0.000235
GO:0005125	GO:0005125	cytokine activity	17/405	231/22669	9.869832e-07	0.000297
GO:1901681	GO:1901681	sulfur compound binding	18/405	260/22669	1.152661e-06	0.000349
GO:0016765	GO:0016765	transferase activity, transferring alkyl or aryl (other than met...	8/405	60/22669	1.100499e-05	0.000350

7. 网络关系图

```

1 # 7. GO plot
2 goplot(GO_MF_diff)

```



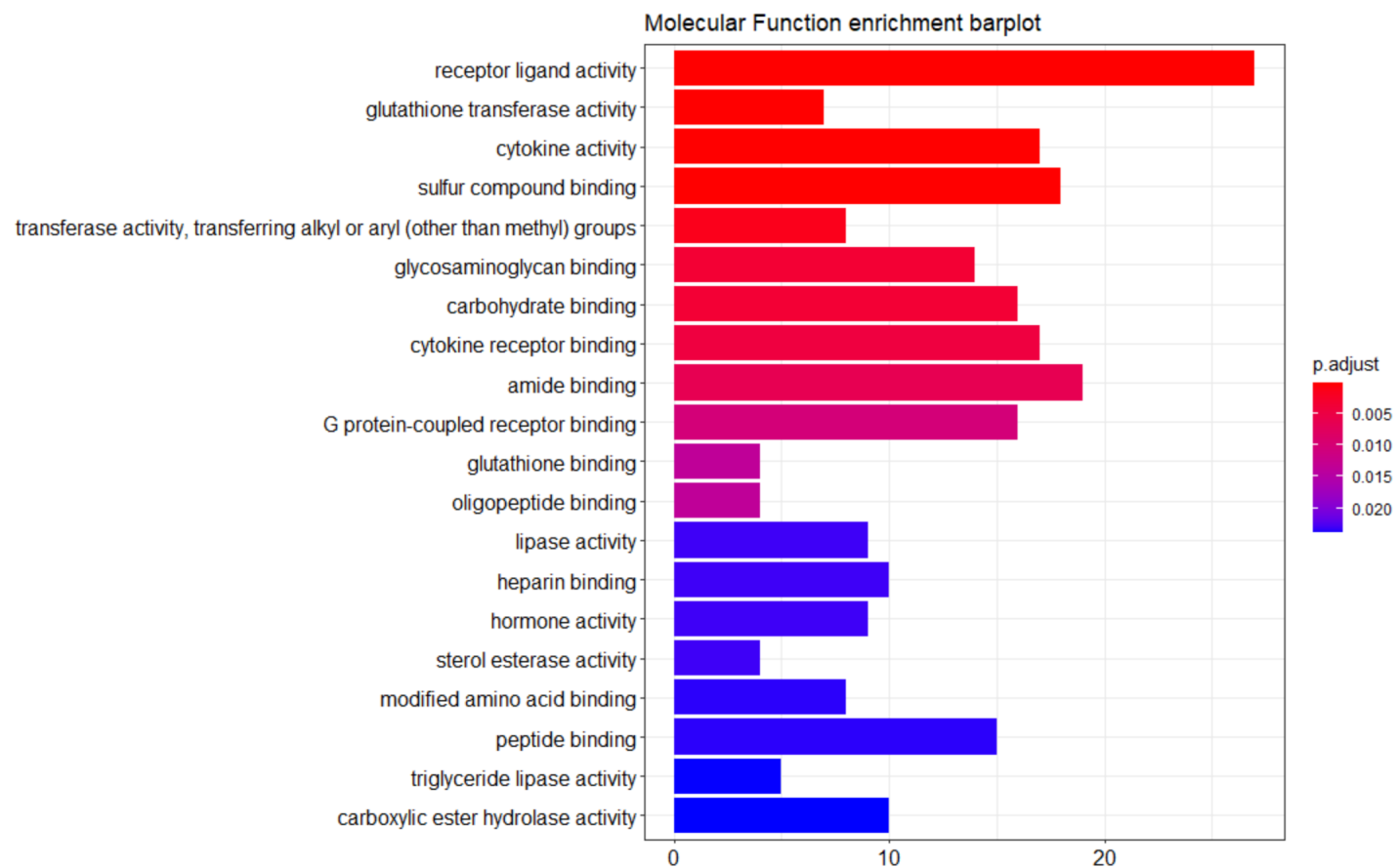
8. Barplot

接下来将以两种方式演示barplot的画法，第一种是现成的函数画图（便捷但是不灵活），第二种是ggplot自行画图（可根据需求自行调整）

```

1 # 8. Barplot — 现成函数
2 barplot(
3   GO_MF_diff,
4   x = "Count", # 或者"GeneRatio"
5   color = "p.adjust", # "pvalue" and "qvalue"
6   showCategory = 20, # 显示前20(enrichResult按照pvalue排序)
7   font.size = 12, # 字体大小
8   title = "Molecular Function enrichment barplot",
9   label_format = 30 # 超过30个字符串换行

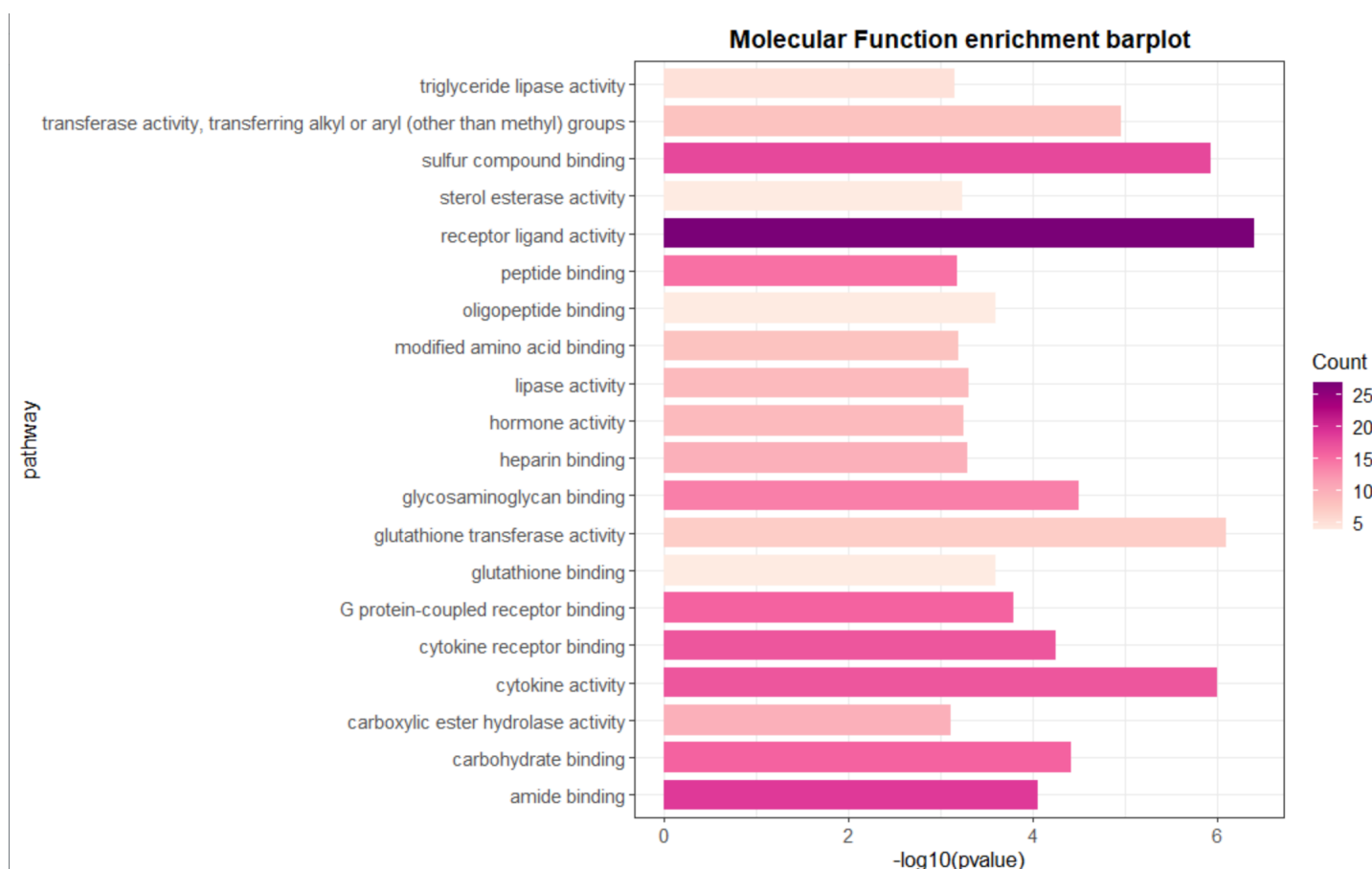
```



```

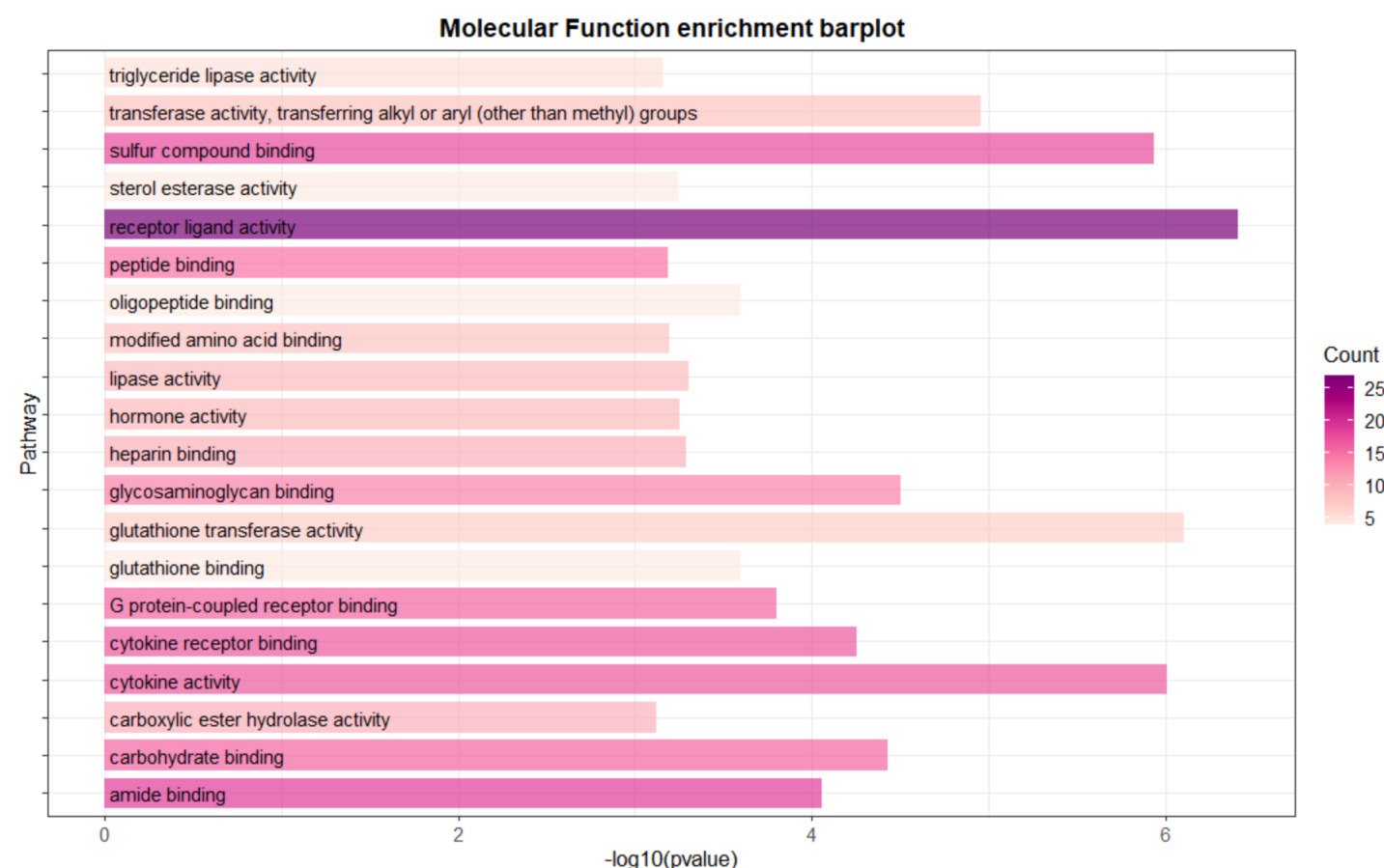
1 # 8. Barplot — ggplot
2 mytheme <- theme(axis.title = element_text(size = 13),
3                 axis.text = element_text(size = 11),
4                 plot.title = element_text(size = 14,
5                                           hjust = 0.5,
6                                           face = "bold"),
7                 legend.title = element_text(size = 13),
8                 legend.text = element_text(size = 11))
9
10 MF_top20 <- GO_MF_result[1:20, ]
11 ggplot(MF_top20,
12       aes(x = -log10(pvalue), y = Description, fill = Count)) +
13   scale_fill_distiller(palette = "RdPu", direction = 1) +
14   geom_bar(stat = "identity", width = 0.8) +
15   labs(x = "log10(pvalue)", y = "pathway",
16        title = "Molecular Function enrichment barplot") +
17   theme_bw() +
18   mytheme

```



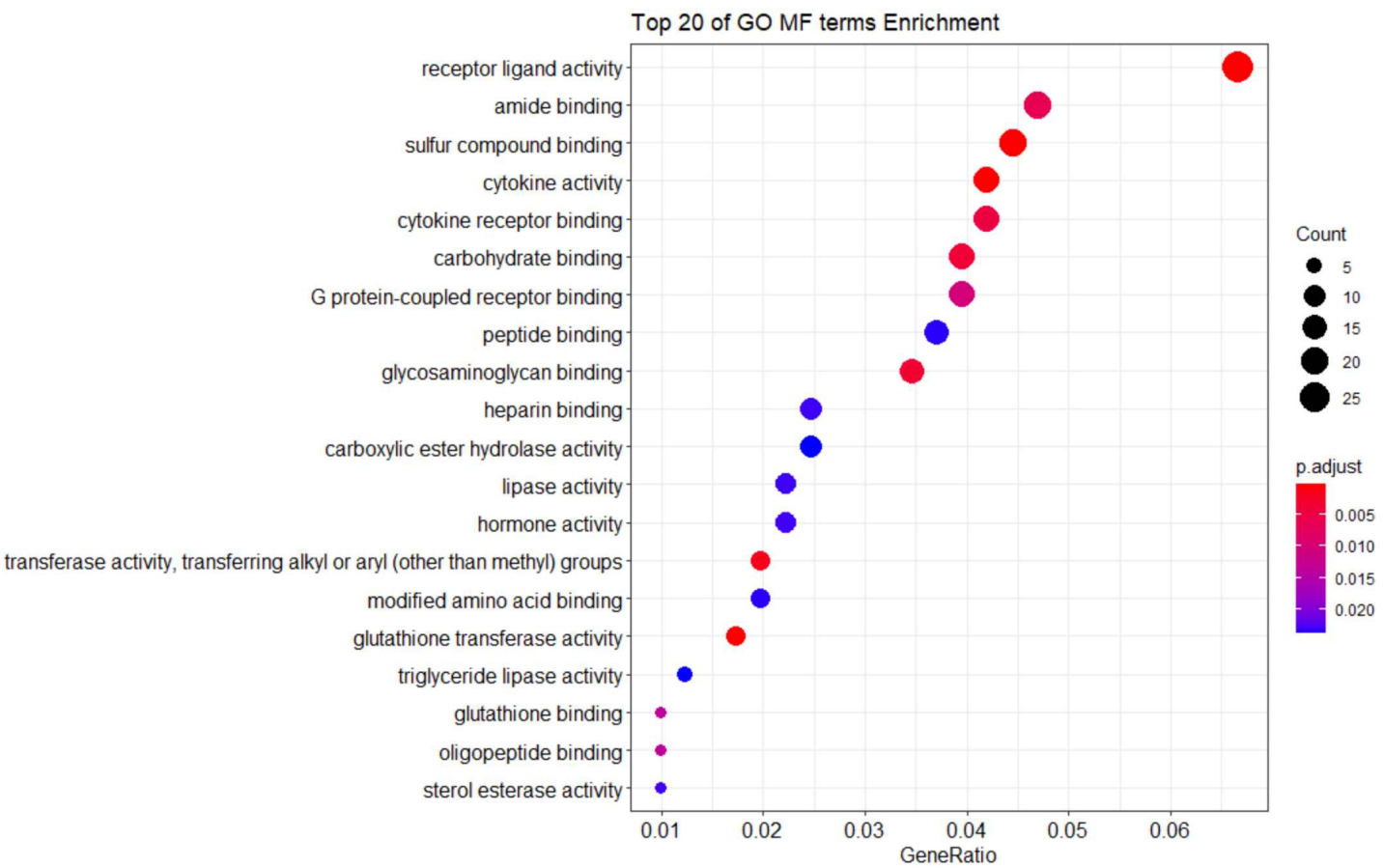
个人比较喜欢的一种方式：将pathway移到柱状图内部：

```
1 mytheme2 <- mytheme + theme(axis.text.y = element_blank())
2 MF_top20$text_x <- rep(0.03,20)
3 ggplot(data = MF_top20,
4       aes(x = -log10(pvalue), y = Description)) +
5   geom_bar(aes(fill = Count), stat = "identity", width = 0.8, alpha = 0.7) +
6   scale_fill_distiller(palette = "RdPu", direction = 1) +
7   labs(x = "-log10(pvalue)", y = "Pathway",
8       title = "Molecular Function enrichment barplot") +
9   geom_text(aes(x = text_x, label = Description),
10          hjust = 0) +
11   theme_bw() +
12   mytheme2
```



9. Dotplot

```
1 # 9. Dotplot
2 dotplot(
3   GO_MF_diff,
4   x = "GeneRatio",
5   color = "p.adjust",
6   title = "Top 20 of GO MF terms Enrichment",
7   showCategory = 20,
8   label_format = 30
9 )
```



三、Gene Set Enrichment Analysis (GSEA)

10. 总览

不同于GO，GSEA的输入则要求是排序好的gene_symbol，排序规则建议是按照logFC从大到小。

GSEA相比于GO的优势如下：

1. 分析所有基因，避免了GO中主动筛选DEG的影响。
2. GSEA分类众多（甚至可以分析GO数据集），可以根据不同需求进行不同类别的分析。

GSEA分类如下（想要查看具体的基因集可以去官网：<https://www.gsea-msigdb.org/gsea/msigdb>）

Human Collections

H **hallmark gene sets** are coherently expressed signatures derived by aggregating many MSigDB gene sets to represent well-defined biological states or processes.

C5 **ontology gene sets** consist of genes annotated by the same ontology term.

C1 **positional gene sets** corresponding to human chromosome cytogenetic bands.

C6 **oncogenic signature gene sets** defined directly from microarray gene expression data from cancer gene perturbations.

C2 **curated gene sets** from online pathway databases, publications in PubMed, and knowledge of domain experts.

C7 **immunologic signature gene sets** represent cell states and perturbations within the immune system.

C3 **regulatory target gene sets** based on gene target predictions for microRNA seed sequences and predicted transcription factor binding sites.

C8 **cell type signature gene sets** curated from cluster markers identified in single-cell sequencing studies of human tissue.

C4 **computational gene sets** defined by mining large collections of cancer-oriented microarray data.

Mouse Collections

MH **mouse-ortholog hallmark gene sets** are versions of gene sets in the MSigDB Hallmarks collection mapped to their mouse orthologs.

M3 **regulatory target gene sets** based on gene target predictions for microRNA seed sequences and predicted transcription factor binding sites.

M1 **positional gene sets** corresponding to mouse chromosome cytogenetic bands.

M5 **ontology gene sets** consist of genes annotated by the same ontology term.

M2 **curated gene sets** from online pathway databases, publications in PubMed, and knowledge of domain experts.

M8 **cell type signature gene sets** curated from cluster markers identified in single-cell sequencing studies of mouse tissue.

11. 处理数据

```
1 genelist <- data$log2FoldChange
2 # names(genelist) <- toupper(rownames(data))
3 names(genelist) <- rownames(data)
4 genelist <- sort(genelist, decreasing = T)
```

```
> genelist
      DMBT1      OOSP1      FBP1      THEM5      GM13880
7.4338866  5.9688772  5.6502578  5.4405683  4.2910346
      PLAC1L      LTF      HAL      MFSD2A      FXYD2
3.9255050  3.8595259  3.6917530  3.6683282  3.6010390
      HSD3B2      GM8909      CAP2      SLC22A19      GPRC5B
3.4470780  3.3992282  3.2912731  3.2760984  3.2591520
      PTPRN      KRT85      NPHS2      SOX9      PAX7
3.2390031  3.2339350  3.1809946  3.1609471  3.0756974
      IL6      IDI2      DPCR1      PMVK      GM15698
3.0200912  2.9006113  2.8748999  2.8666734  2.8475421
```


正常来说，我们需要从官网下载pathway注释文件（.gmt）并进行处理（比较麻烦）。感兴趣的同学可以研究研究以下处理的代码：

```
1 x <- readLines("c6.all.v7.5.1.symbols.gmt")
2 res <- strsplit(x, "\t")
3 names(res) <- vapply(res, function(y) y[1], character(1))
4 res <- lapply(res, "[", -c(1:2))
```

好在R语言中GSEA数据库被打包进了msigdf这个包，因此我们无需再下载gmt文件，只需要对msigdf中数据库进行提取即可。**注意：该数据库中没有gene_id列，只有gene_symbol列，所以无需像GO一样将symbol转化成ID。且在human数据库中symbol列中的值均为大写，因此处理human数据时需要进行大小写转化（mouse数据则不需要）。**

```
> msigdf.human
# A tibble: 3,850,745 x 4
  category_code category_subcode geneset symbol
  <chr>         <chr>          <chr>   <chr>
1 c1           all          chr1p12 VTCN1
2 c1           all          chr1p12 LINC01525
3 c1           all          chr1p12 MAN1A2
4 c1           all          chr1p12 VPS25P1
5 c1           all          chr1p12 TENT5C
6 c1           all          chr1p12 VDAC2P3
7 c1           all          chr1p12 PNRC2P1
8 c1           all          chr1p12 GDAP2
9 c1           all          chr1p12 WDR3
10 c1          all          chr1p12 SPAG17
```

```
> msigdf.mouse
# A tibble: 3,387,971 x 4
  category_code category_subcode geneset mouse.symbol
  <chr>         <chr>          <chr>   <chr>
1 c1           all          chr1p12 Vtcn1
2 c1           all          chr1p12 Man1a2
3 c1           all          chr1p12 Tent5c
4 c1           all          chr1p12 Gdap2
5 c1           all          chr1p12 Wdr3
6 c1           all          chr1p12 Spag17
7 c1           all          chr1p12 Tbx15
8 c1           all          chr1p12 Wars2
9 c1           all          chr1p12 Hao2
10 c1          all          chr1p12 Hsd3b1
# ... with 3,387,961 more rows
# i Use `print(n = ...)` to see more rows
```

12. 开始GSEA

为了方便分析，可以根据自己需求写一个小函数（以下仅是我的写法）。其中返回的是一个列表，列表第一个值是通路汇总矩阵，方便人为查看，列表第二个值是GSEA分析出的大类，方便软件包后续画图。

```

1 myGSEA <- function(c){
2   class <- msigdf.mouse %>% # msigdf.human
3     filter(category_code == c) %>%
4     select(geneset, mouse.symbol) %>%
5     as.data.frame
6   Class_GSEA <- GSEA(genelist,
7     TERM2GENE = class,
8     minGSSize = 10,
9     maxGSSize = 500,
10    pvalueCutoff = 0.05,
11    pAdjustMethod = "BH",
12    verbose = FALSE,
13    eps = 0)
14
15   result <- Class_GSEA@result
16   # Class_GSEA@result$Description <-
17   #   str_trunc(Class_GSEA@result$Description, width = 30, side = "right")
18   return(list(result, Class_GSEA))
19 }
20
21 H_result_list <- myGSEA("h")
22 c5_result_list <- myGSEA("c5")
23
24 H_result <- H_result_list[[1]]
25 c5_result <- c5_result_list[[1]]

```

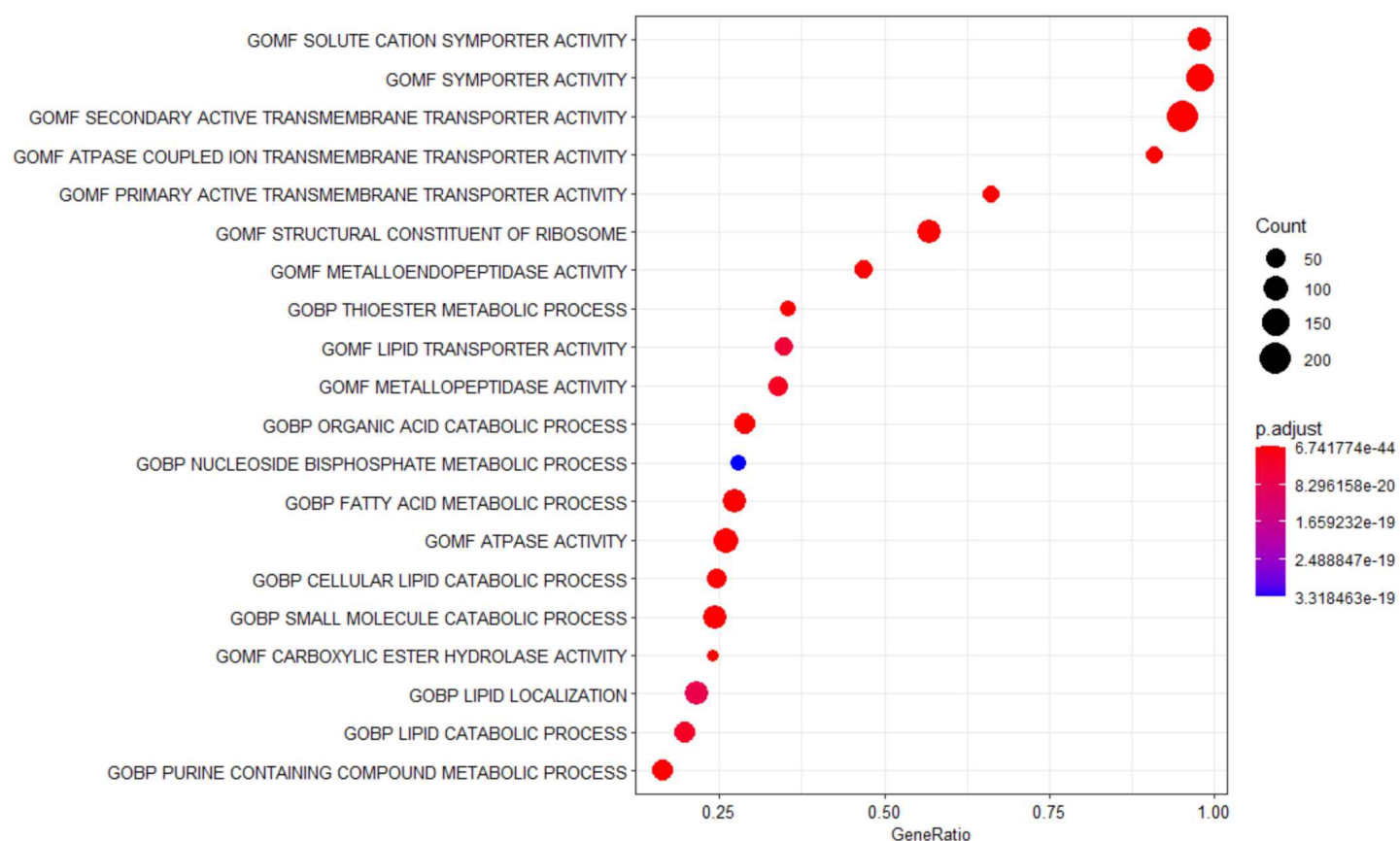
可根据函数提取出的矩阵（如本例中的c5_result）直接查看想要得知的通路信息。

13. Dotplot

```

1 # dotplot
2 dotplot(c5_result_list[[2]], showCategory = 20, font.size = 10)

```

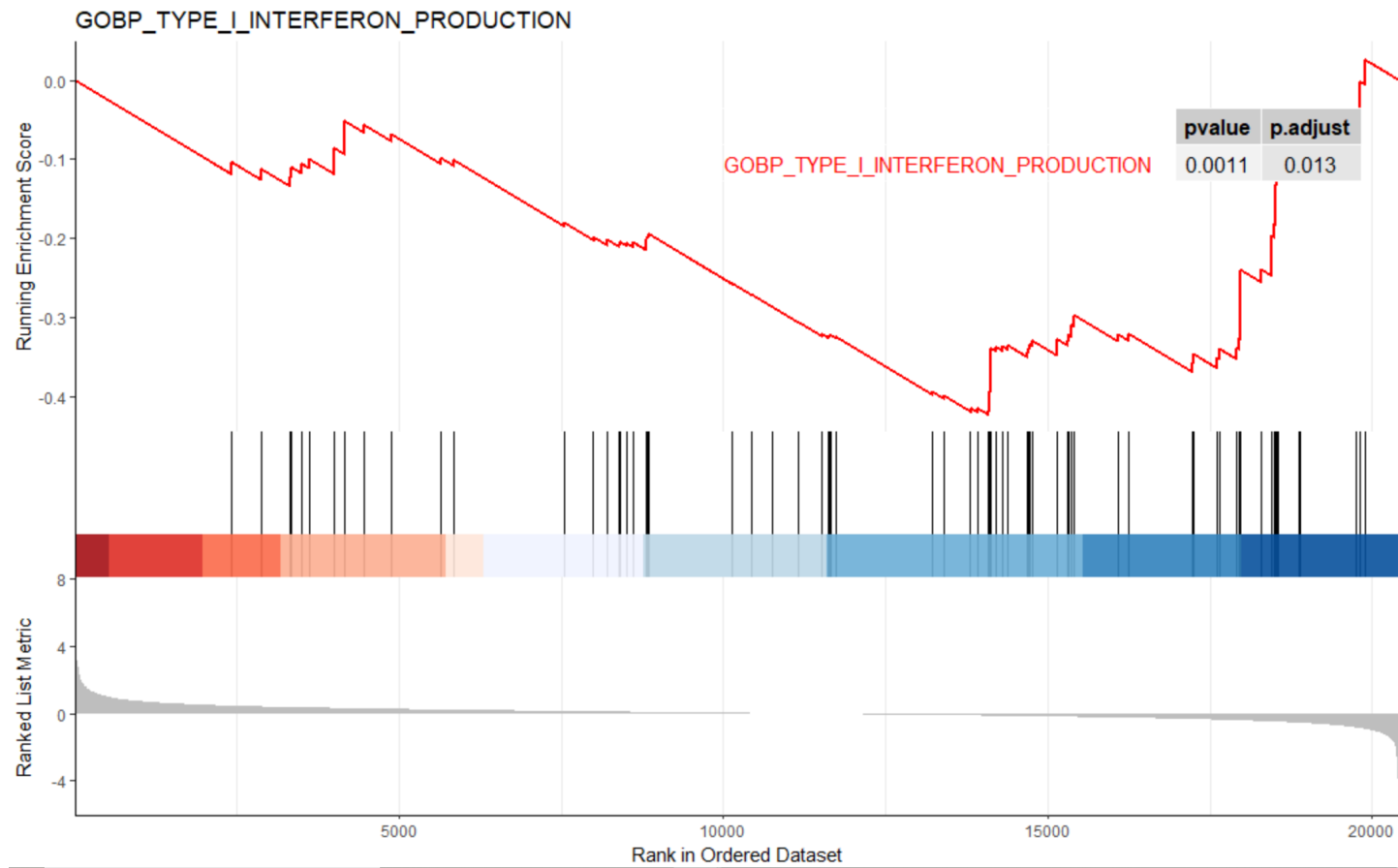


14. 富集图

```

1 # 寻找某一特定通路
2 which(c5_result$ID == 'GOBP_TYPE_I_INTERFERON_PRODUCTION') # 848
3 gseaplot2(c5_result_list[[2]],
4           geneSetID = 848,
5           color = "red",
6           rel_heights = c(1.5, 0.5, 1), # 子图高度
7           subplots = 1:3, # 显示哪些子图
8           pvalue_table = T, # 是否显示pvalue表
9           title = c5_result_list[[1]]$Description[848],
10          ES_geom = "line") # or "dot"

```

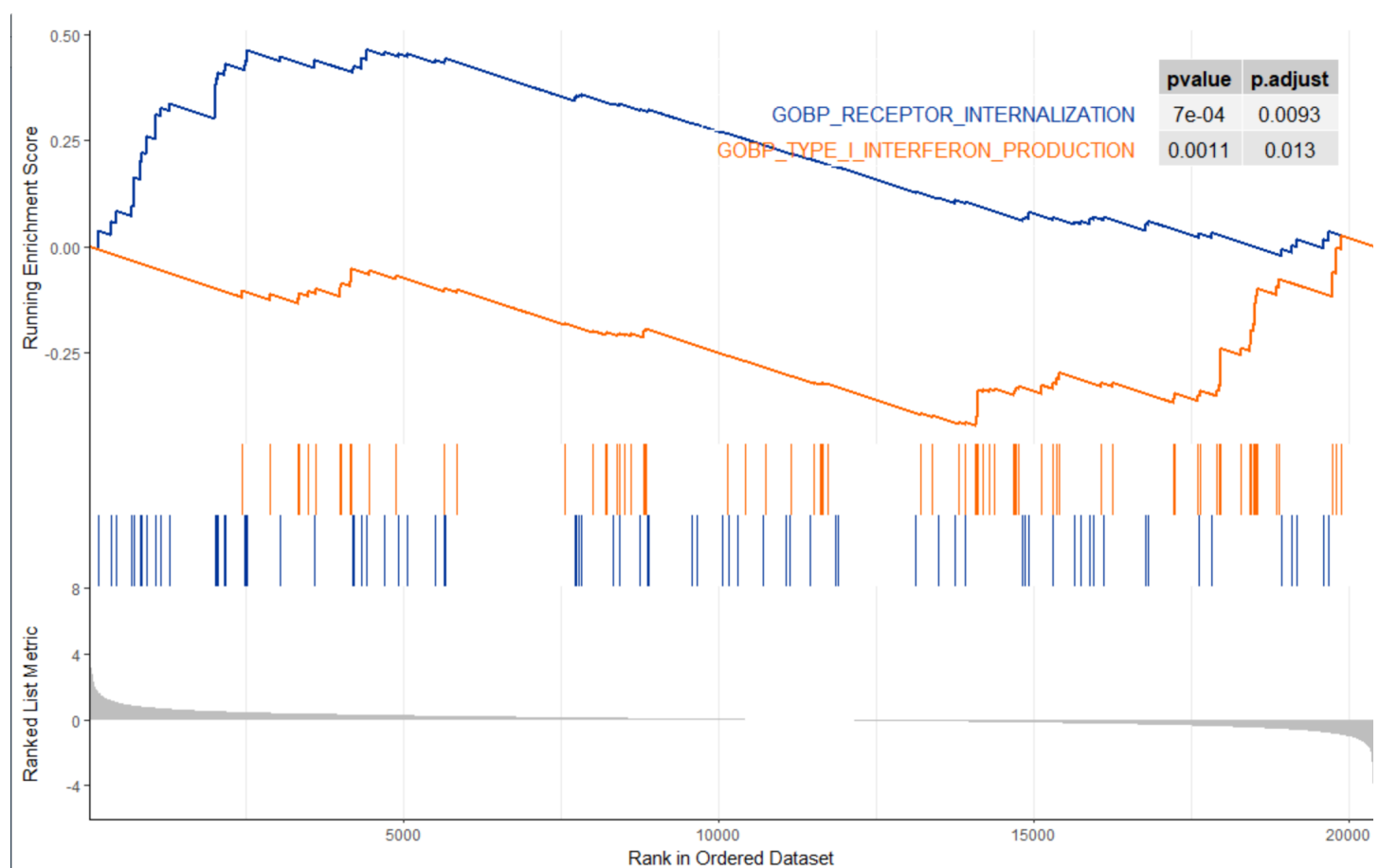


当然，也可以在一张图中同时显示多条通路

```

1 gseaplot2(c5_result_list[[2]],
2           geneSetID = c(772, 848),
3           color = c("#003399", "#FF6600"),
4           pvalue_table = TRUE,
5           ES_geom = "line")

```

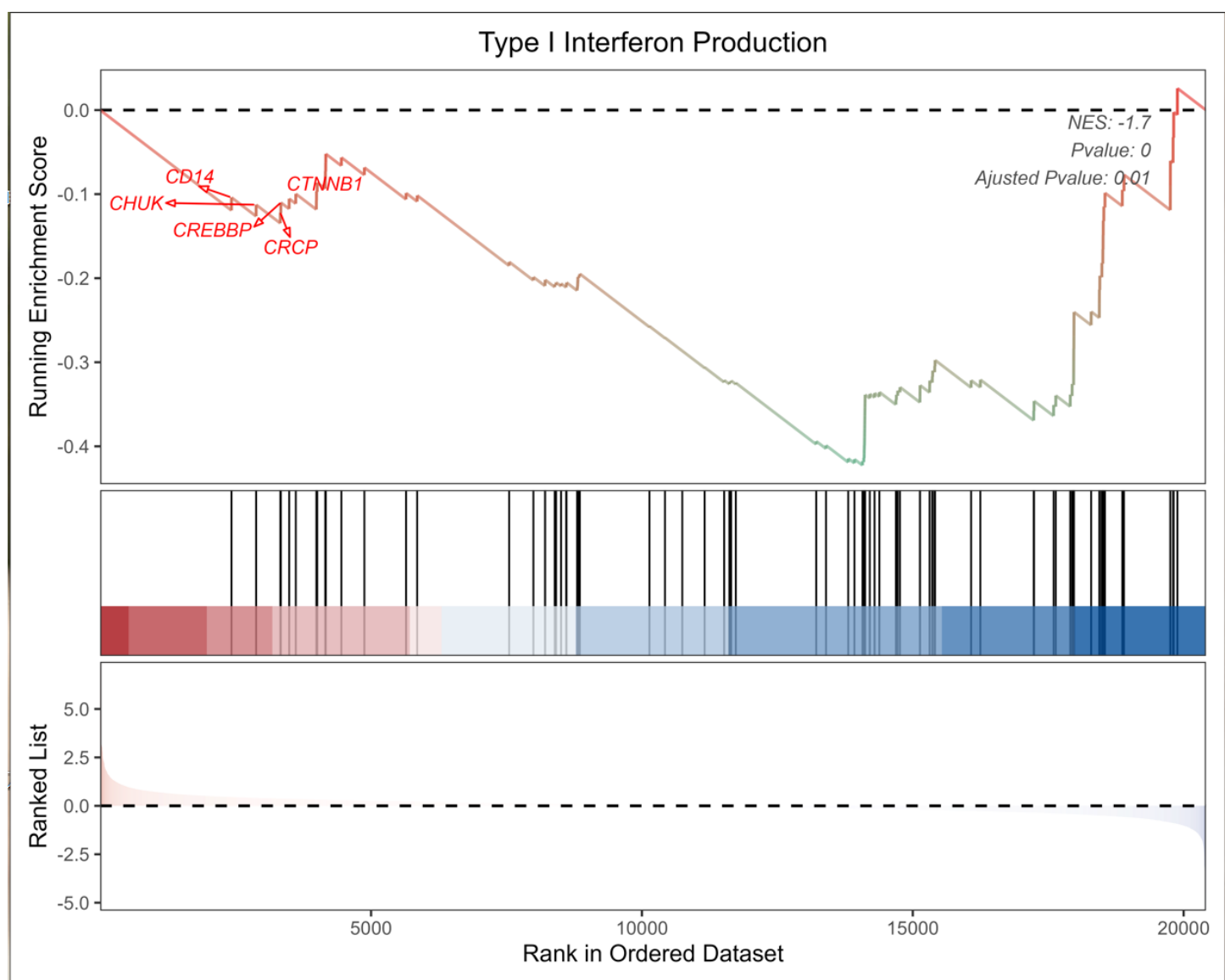


一个最近发现的新包：GseaVis，可以在图上注释基因（虽然有些乱）

```

1 # 标注前5名的基因名字
2 gseaNb(
3   object = c5_result_list[[2]],
4   geneSetID = c5_result_list[[2]]@result$ID[848],
5   addPval = T,
6   pvalX = 0.95,
7   pvalY = 0.8,
8   addPoint = F,
9   newCurveCol = c("#001871", "#b9b4ad", "#f99f1c"),
10  newHtCol = c("#001871", "white", "#f99f1c"),
11  addGene = T,
12  markTopgene = T, # 是否标注Top基因
13  topGeneN = 5, # 标注前多少个gene
14  geneCol = 'red', # 基因名标签颜色更改
15  rmSegment = T
16 )

```

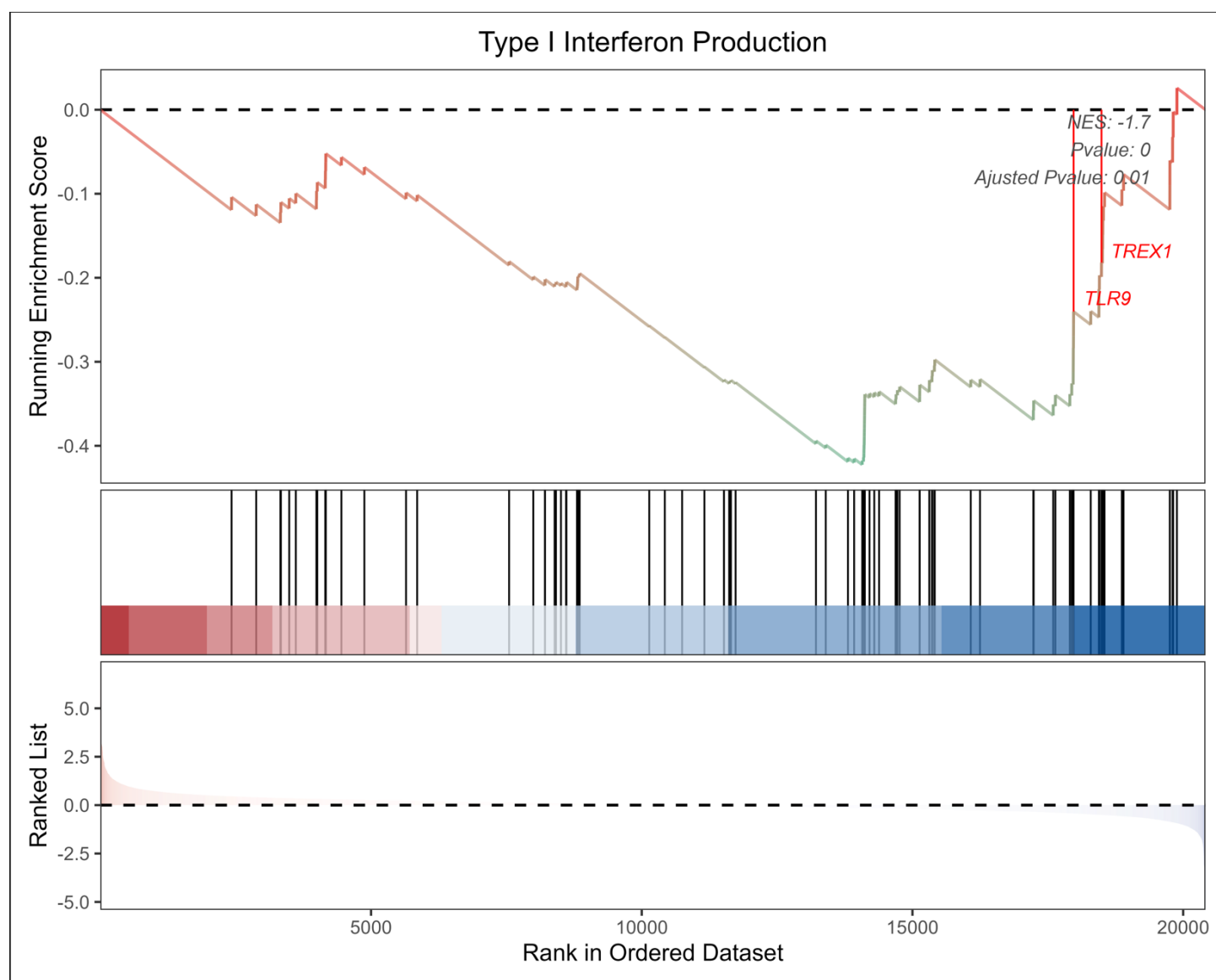


也可以按需求标注出基因及其位置。

```

1 gseaNb(
2   object = c5_result_list[[2]],
3   geneSetID = c5_result_list[[2]]@result$ID[848],
4   addPval = T,
5   pvalX = 0.95,
6   pvalY = 0.8,
7   addPoint = F,
8   newCurveCol = c("#001871", "#b9b4ad", "#f99f1c"),
9   newHtCol = c("#001871", "white", "#f99f1c"),
10  addGene = c('TREX1', 'TLR9'),
11  geneCol = 'red',
12  rmSegment = F
13 )

```



四、Kyoto Encyclopedia of Genes and Genomes (KEGG)

KEGG分析流程与GO十分相似，只不过用了不同的数据库。因此不做过多说明。

事先说明：R语言连接KEGG十分不稳定，我经查询和多次尝试仍未得到结果，网上大多有的方法是该连接随机性太强（有的时候能得到结果，有的时候则不行）。另一种方法是在进行KEGG分析前进行一下设置。

```
1 R.utils::setOption("clusterProfiler.download.method", "auto")
```

15. 分析+画图

```
1 # KEGG
2 KEGG_diff <- enrichKEGG(gene = diff_entrez$ENTREZID,
3                           organism = "mmu",
4                           pvalueCutoff = 0.05,
5                           qvalueCutoff = 0.05)
6
7 KEGG_result <- KEGG_diff@result
```

其中organism参数的设置可参考：https://www.genome.jp/kegg/catalog/org_list.html
(人hsa, 小鼠mmu.....)

```
1 # 探索/调出选定的KEGG通路
2 # 直接跳转到对应网页，红色为富集到该通路的差异基因
```



```
3 browseKEGG(KEGG_diff, "hsa04610")
```

```
1 ## Figure 1
2 barplot(
3   KEGG_diff,
4   x = "Count",
5   color = "p.adjust",
6   showCategory = 20,
7   font.size = 12,
8   title = "KEGG enrichment barplot",
9   label_format = 30
10 )
11
12 ## Figure 2
13 dotplot(
14   KEGG_diff,
15   x = "GeneRatio",
16   color = "p.adjust",
17   title = "Top 20 of Pathway Enrichment",
18   showCategory = 20,
19   label_format = 30
20 )
```

如果有任何疑问（包括但不限于不清晰的地方，解释/代码错误，更有意义的作图方式），麻烦联系卢建璋进行修改。祝大家科研顺利。