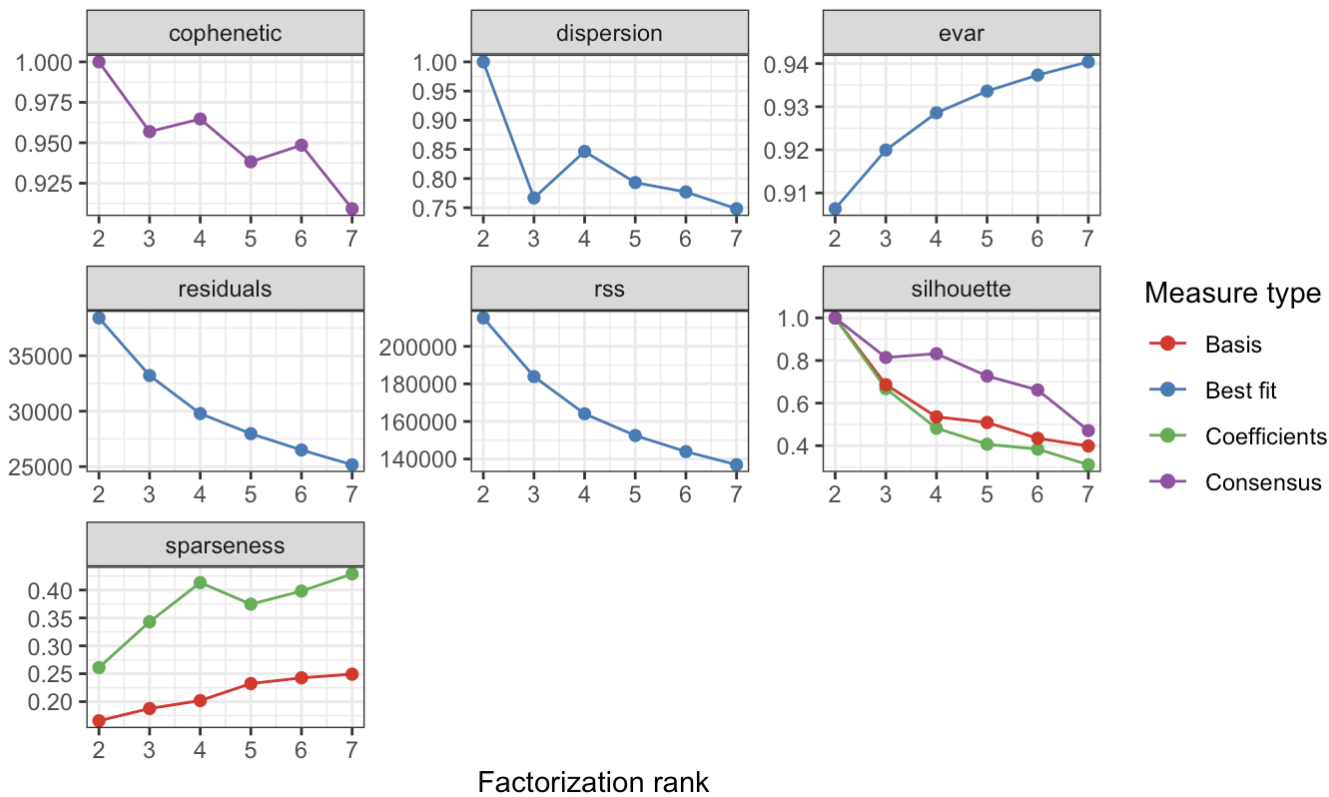
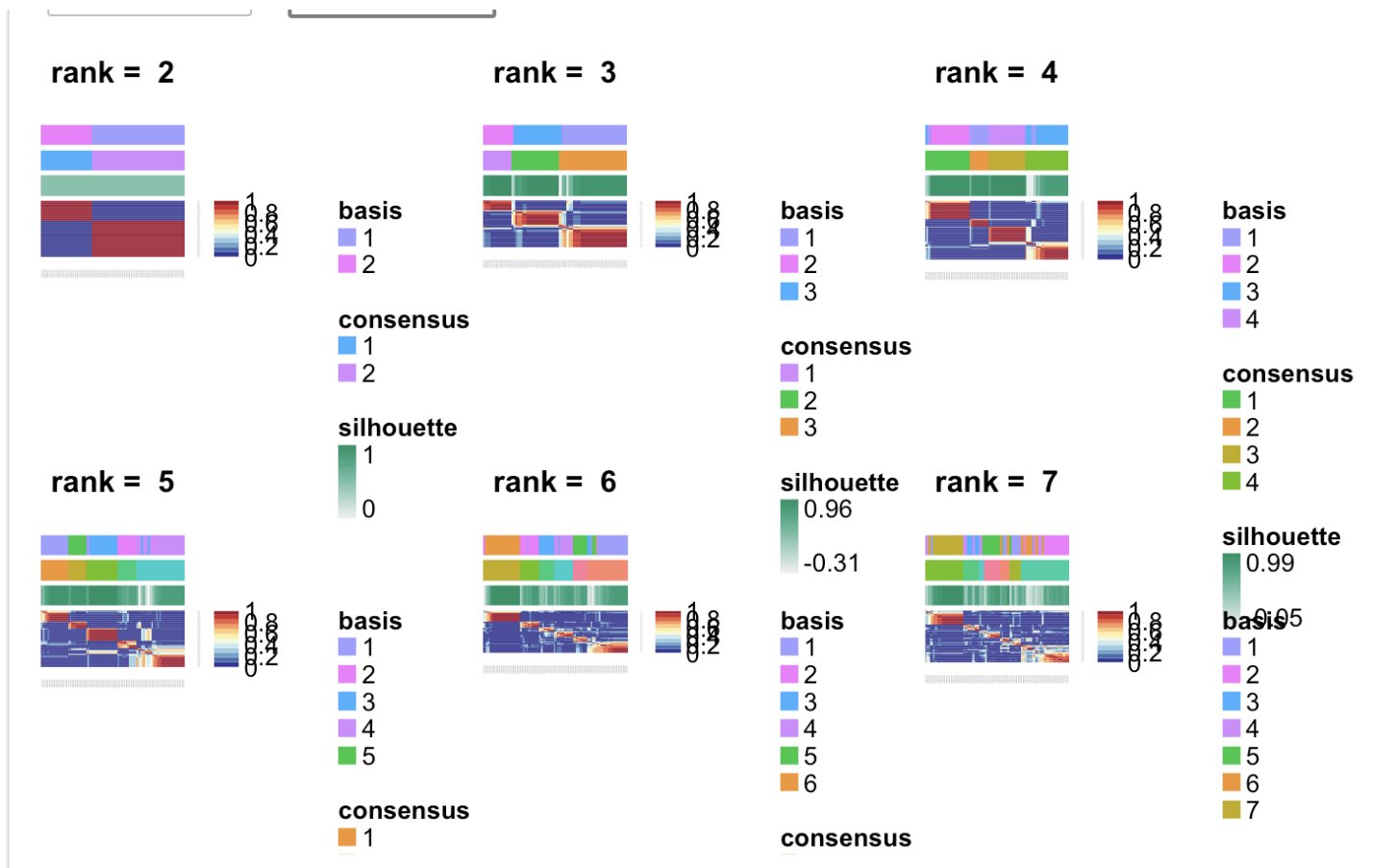


# ML\_LAQ2\_submit

Q1

NMF rank survey





k=4 gives the best performance based on cophenetic coefficient, and now we perform NMF using k=4

Justification for selecting k = 4 as the optimal number of clusters:

Based on the NMF rank survey plots, we evaluated several metrics including the cophenetic coefficient, dispersion, sparseness, and silhouette scores. While k = 2 gives the highest cophenetic and silhouette values, it represents a very coarse clustering structure that may overlook biologically meaningful subgroups. Among k = 3 to 7, k = 4 consistently shows local peaks or plateau behavior in key metrics (especially cophenetic, dispersion, and sparseness), suggesting a stable and interpretable clustering solution. Therefore, we selected k = 4 as the optimal number of clusters to balance stability, separation, and biological interpretability.

## Q2

代码块

```
1 library(limma)
2
3 cluster_assignment <- apply(h, 2, which.max)
4 cluster_assignment <- as.factor(cluster_assignment)
```

```

5
6 top_genes_list <- list()
7
8
9 for (k in 1:4) {
10   group <- ifelse(cluster_assignment == k, "cluster", "others")
11   design <- model.matrix(~ factor(group))
12
13   fit <- lmFit(d_exp, design)
14   fit <- eBayes(fit)
15
16   top <- topTable(fit, coef=2, number=20, sort.by="P")
17   top_genes_list[[k]] <- rownames(top)
18 }

```

代码块

```

1 top_genes_list[[1]]
2
3 top_genes_list[[2]]
4 top_genes_list[[3]]
5 top_genes_list[[4]]

```

output

```

1  [1] "SFN"      "GJB6"      "FGFBP1"    "KRT6A"     "DSG3"      "GJB2"      "PPP2R2C"
   "KRT6B"    "CXCL14"    "IL36G"     "FABP5"     "ANXA6"
2  [13] "CERS3"    "S100A2"    "PGLYRP3"   "SLC2A1"    "ELN"       "CLIP3"     "SPRR1B"
   "HAS3"
3
4  [1] "REPIN1"    "ALDH1A1"   "ITGA3"     "MIR936"    "PLEK2"     "ARHGEF26"
   "COL17A1"   "ATP6V0E2"  "KRT14"     "LAMC2"     "LRIG1"
5  [12] "LAMA3"     "HLF"       "EPS8"      "EPHB6"     "RGS20"     "MYLIP"
   "FAM171A1"  "CAV1"      "PBX1"
6
7  [1] "FUT3"      "PRSS27"    "CLCA4"     "SPRR2C"    "SCNN1B"    "KRT78"
   "PSCA"      "SCEL"      "DUOX2"     "PPL"
8  [11] "TMPRSS11F" "C2orf54"   "SPINK7"    "A2ML1"     "TTC9"      "CCDC64B"
   "CEACAM5"   "NCCRP1"    "FM02"      "CD24"
9
10 [1] "OSGIN1"     "CYP4F3"    "ALDH3A1"    "LOC344887"  "CYP4F11"
   "DMRT2"     "AKR1C3"    "ABCC1"
11 [9] "TSPAN7"     "TRIM16L"   "MDGA1"      "MRAP2"      "CBR1"
   "PTGR1"     "LOC100133286" "UGT1A1"
12 [17] "AK126334"   "TXNRD1"    "ADAM23"     "SLC7A11"

```

Q3

I selected upregulated genes based on these two criteria:

```
$logFC > 0 & $adj.P.Val < 0.05 .
```

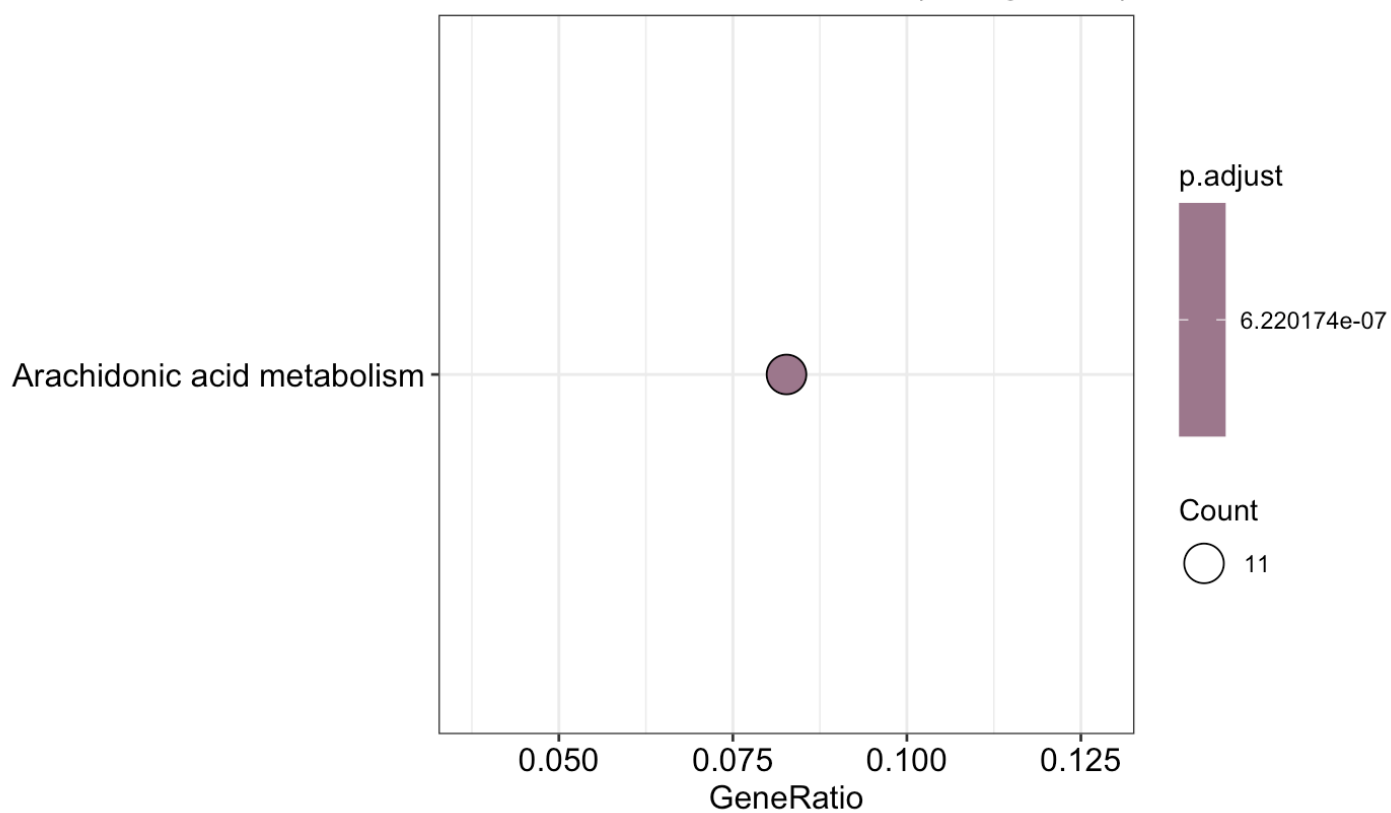
I selected downregulated genes based on these two criteria.

```
$logFC < 0 & $adj.P.Val < 0.05
```

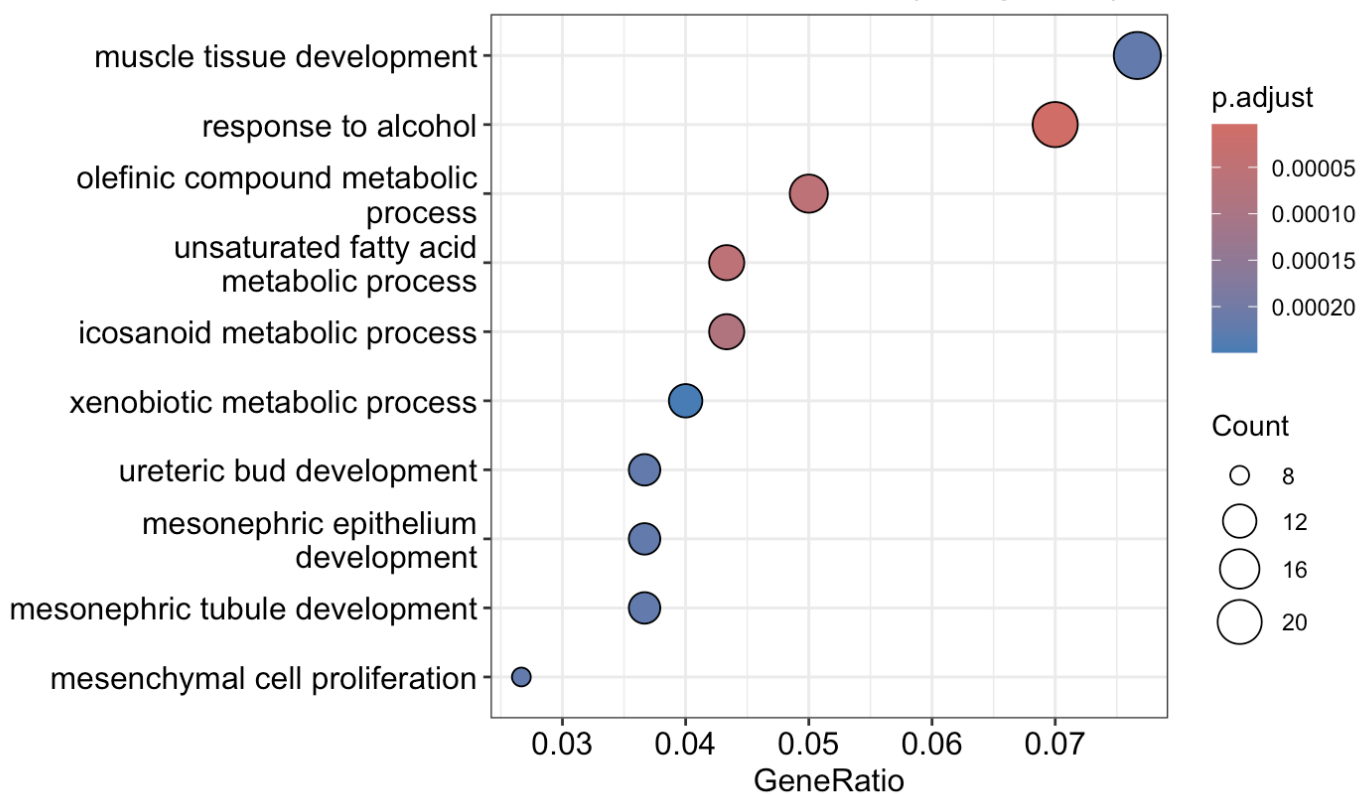
up



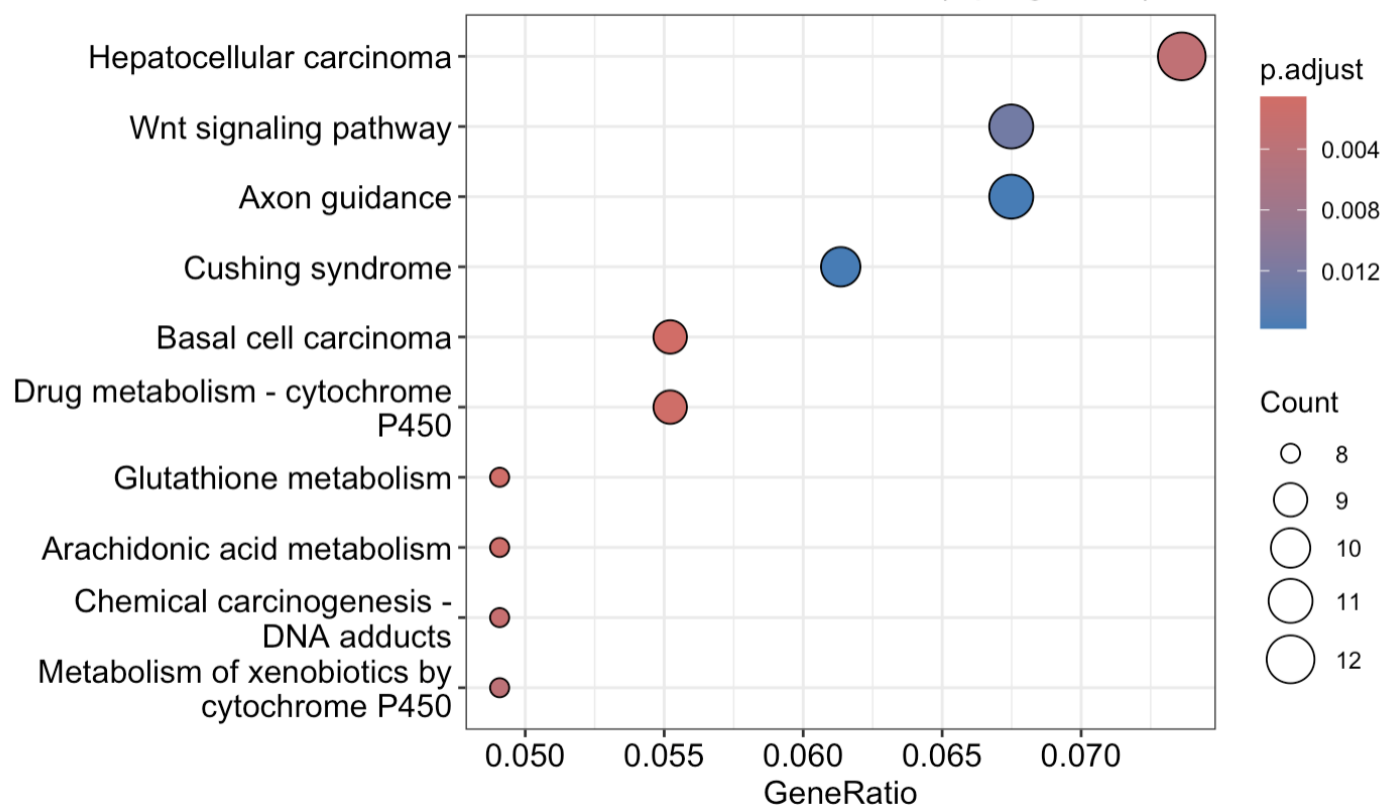
Cluster 1 - KEGG enrichment (Upregulated)



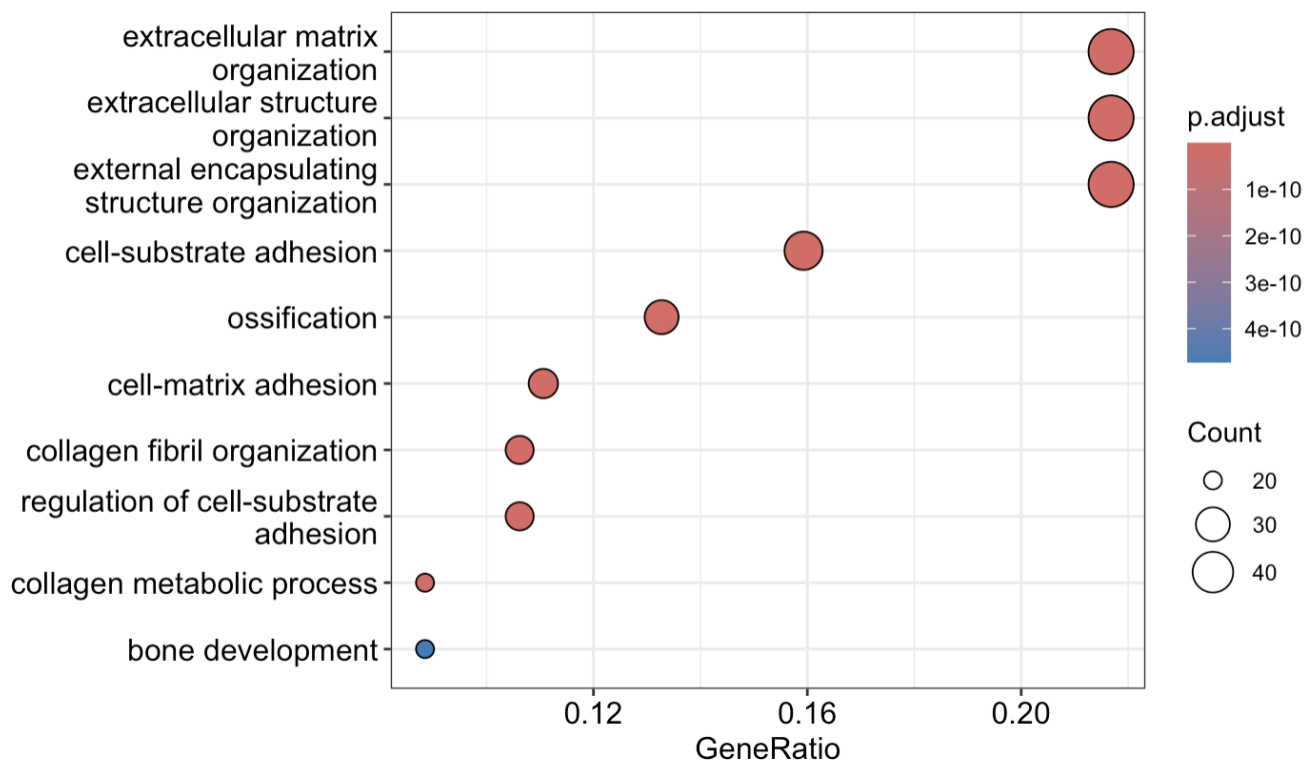
Cluster 2 - GO enrichment (Upregulated)



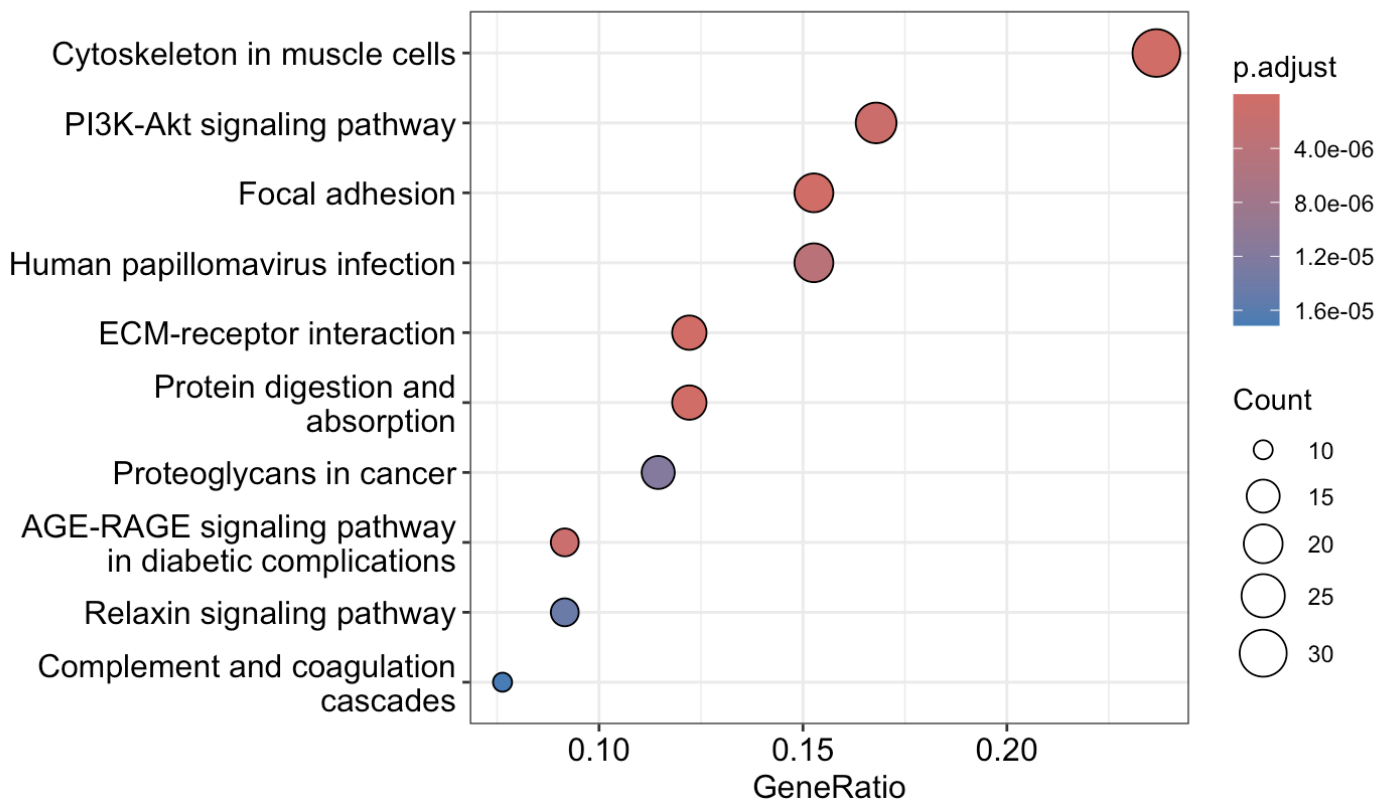
Cluster 2 - KEGG enrichment (Upregulated)



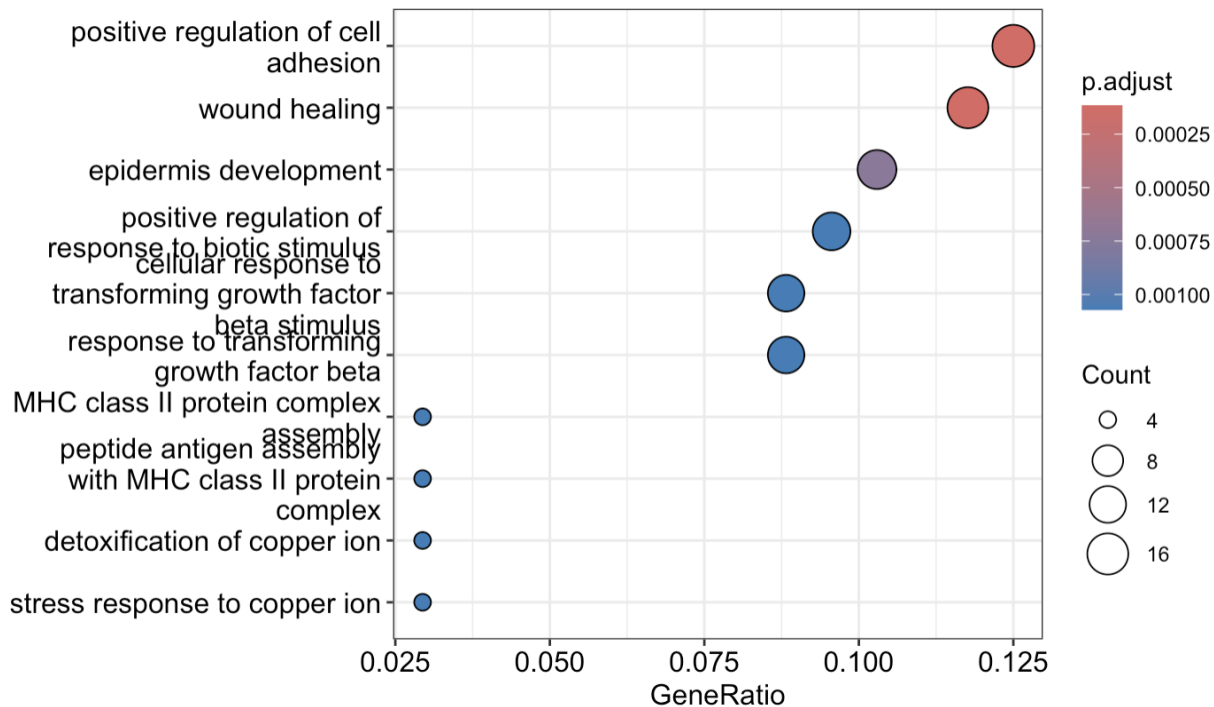
Cluster 3 - GO enrichment (Upregulated)

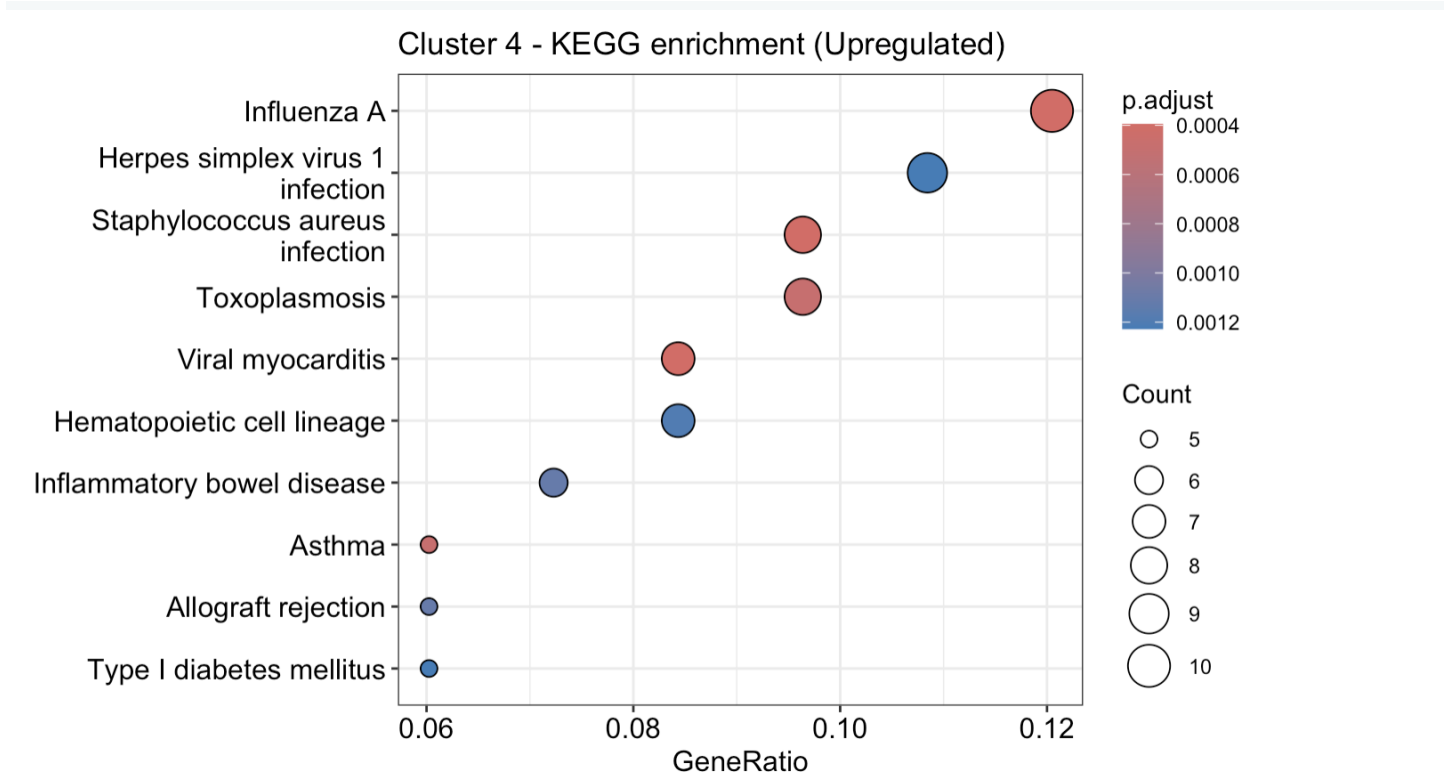


Cluster 3 - KEGG enrichment (Upregulated)



Cluster 4 - GO enrichment (Upregulated)



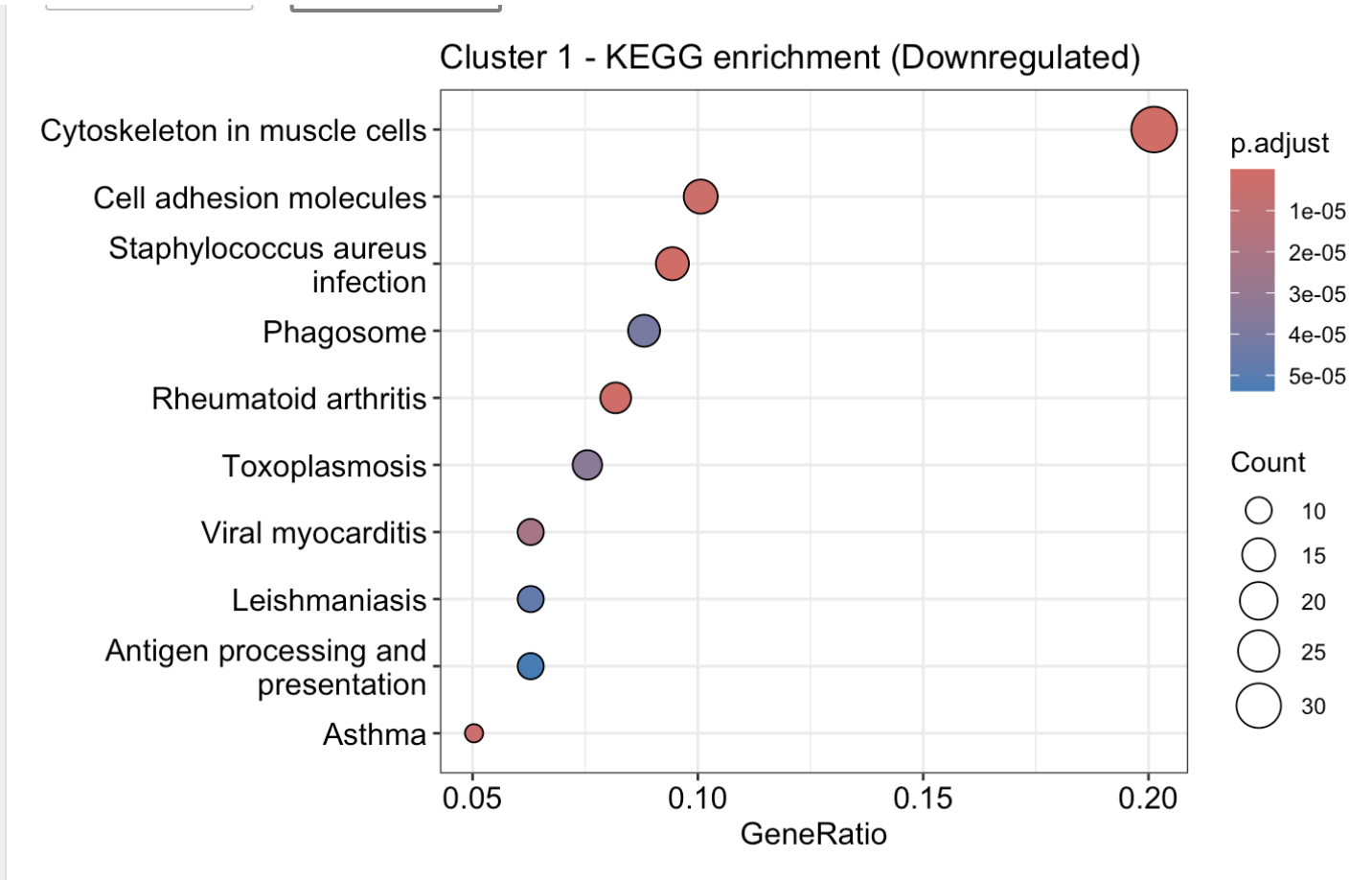
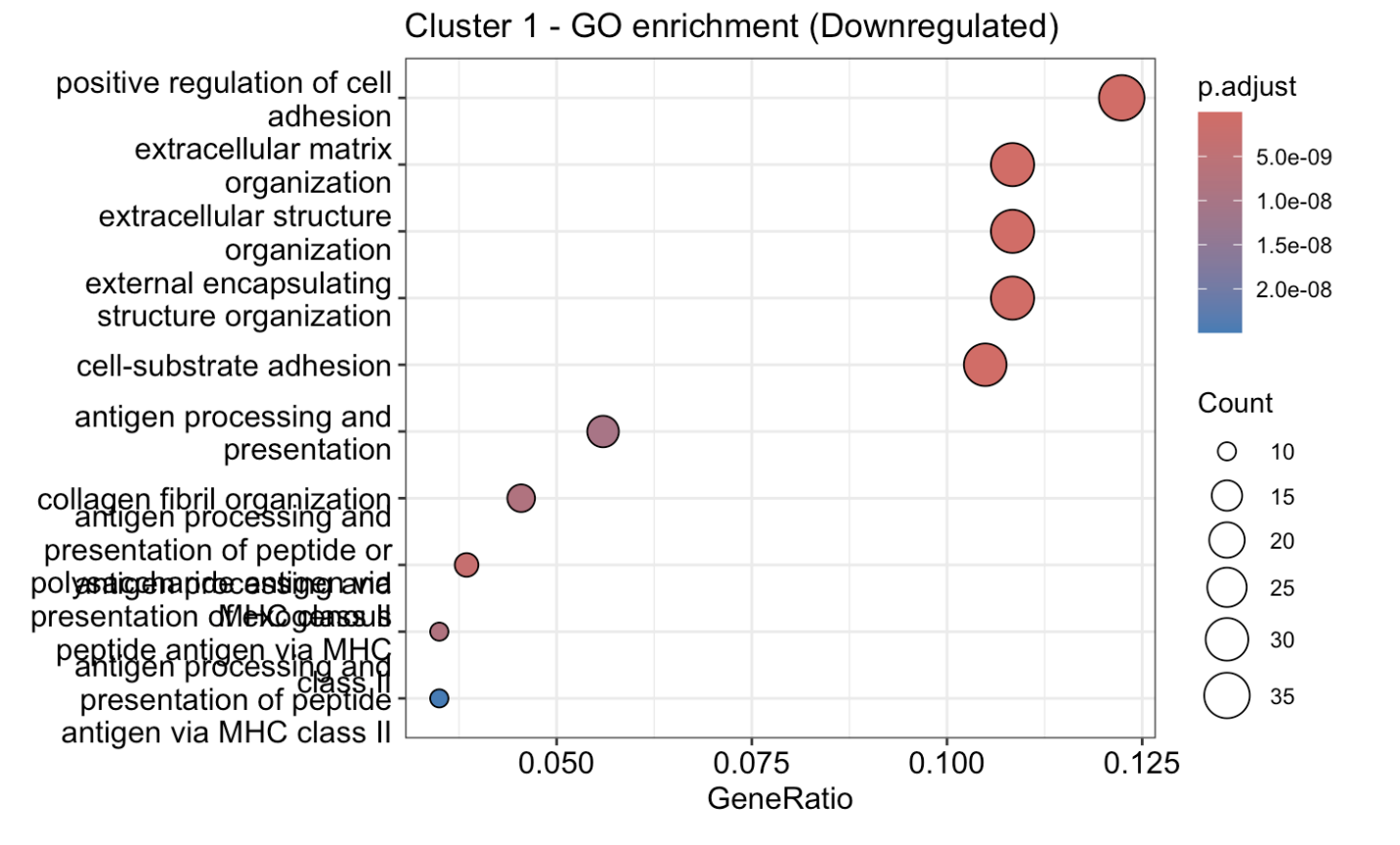


conclusion

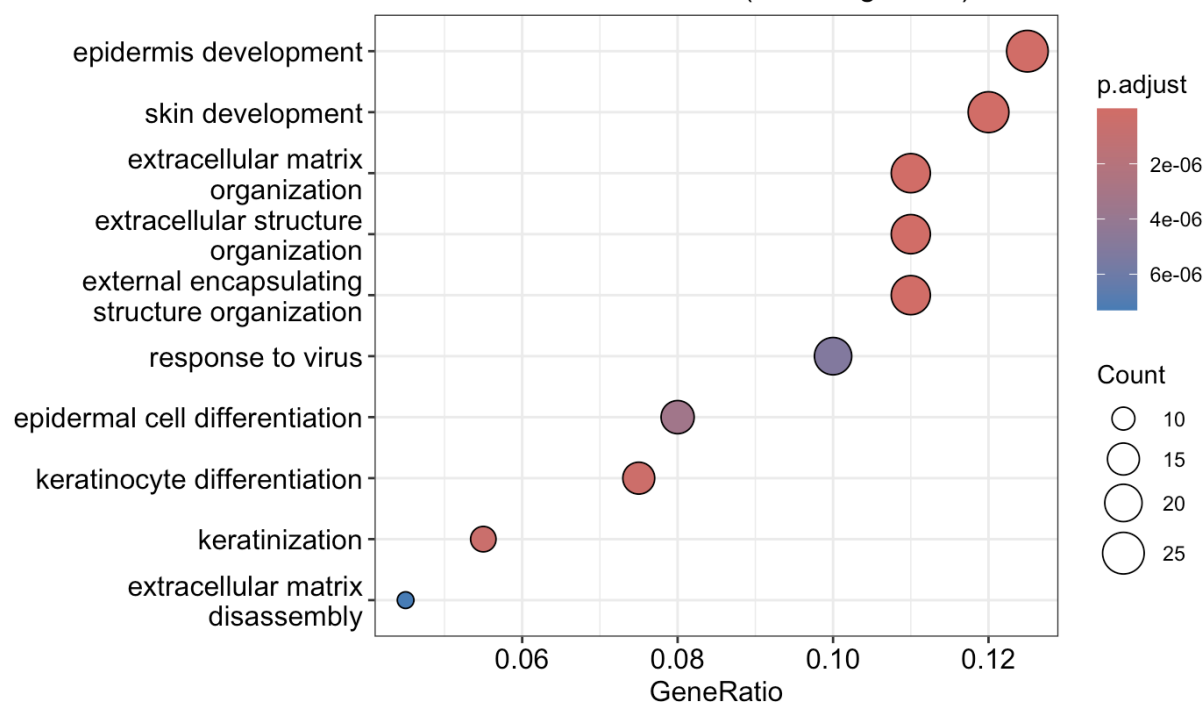
Cluster	Top 3 GO Biological Process (Upregulated)	Top 3 KEGG Pathways (Upregulated)
Cluster 1	epidermis development skin development epidermal cell differentiation	Arachidonic acid metabolism
Cluster 2	muscle tissue development response to alcohol olefinic compound metabolic process	Hepatocellular carcinoma Wnt signaling pathway Axon guidance
Cluster 3	extracellular matrix organization extracellular structure organization external encapsulating structure organization	Cytoskeleton in muscle cells PI3K-Akt signaling pathway Focal adhesion
Cluster 4	positive regulation of cell adhesion wound healing epidermis development	Influenza A Herpes simplex virus 1 infection Staphylococcus aureus infection



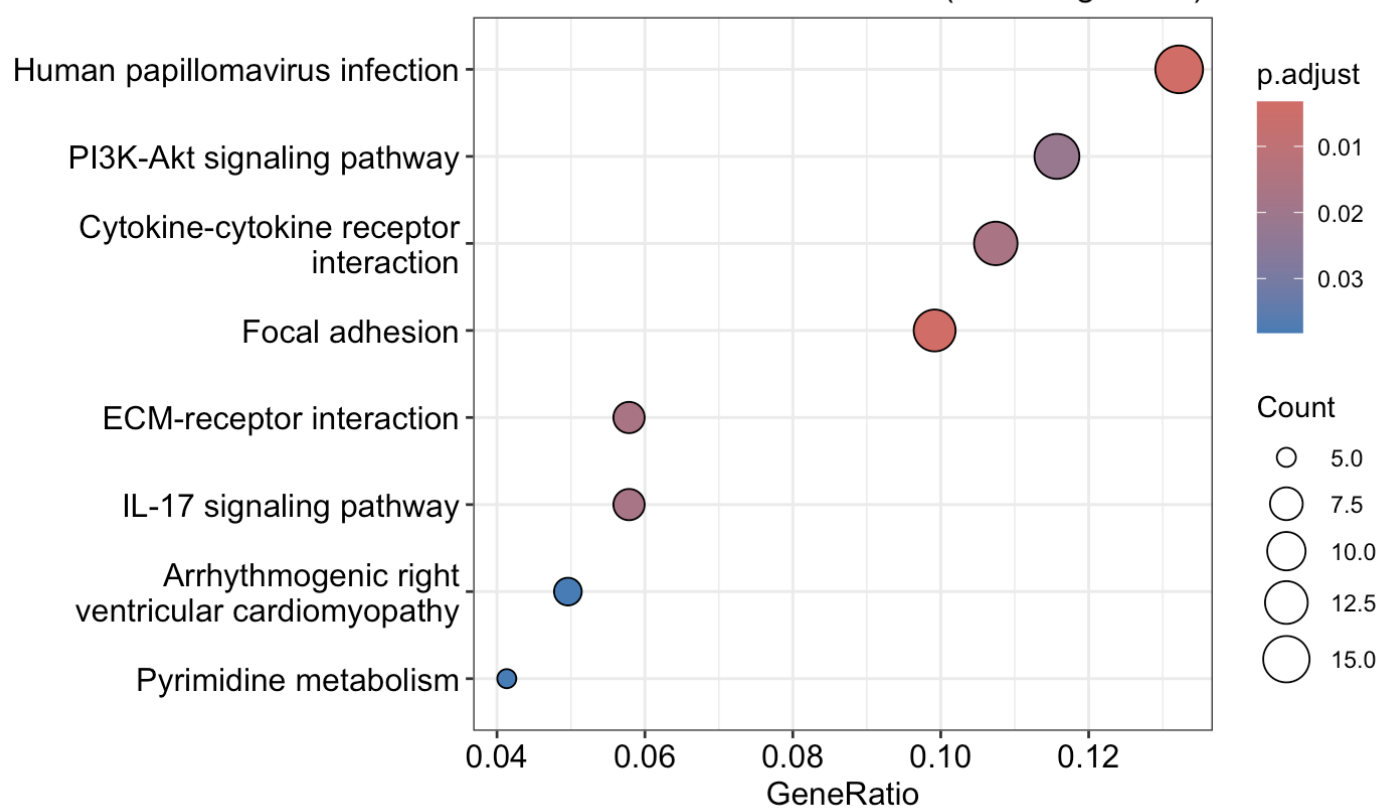
down



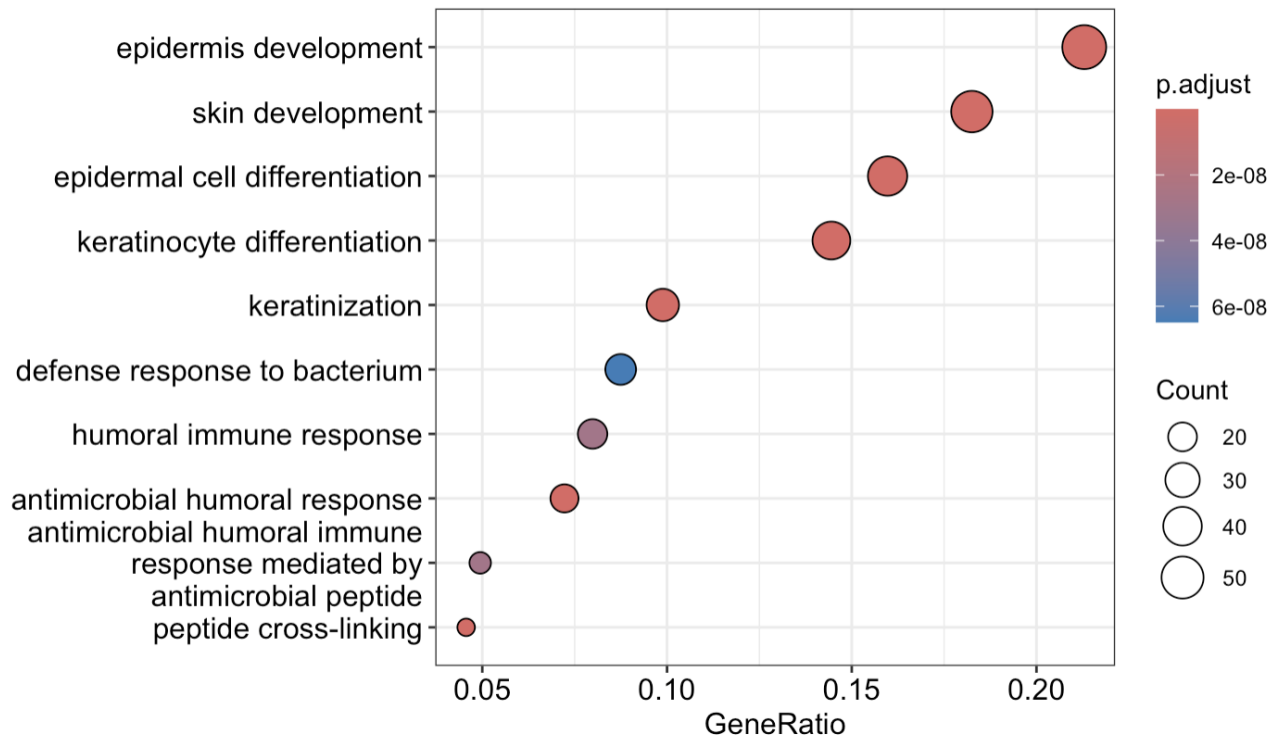
Cluster 2 - GO enrichment (Downregulated)



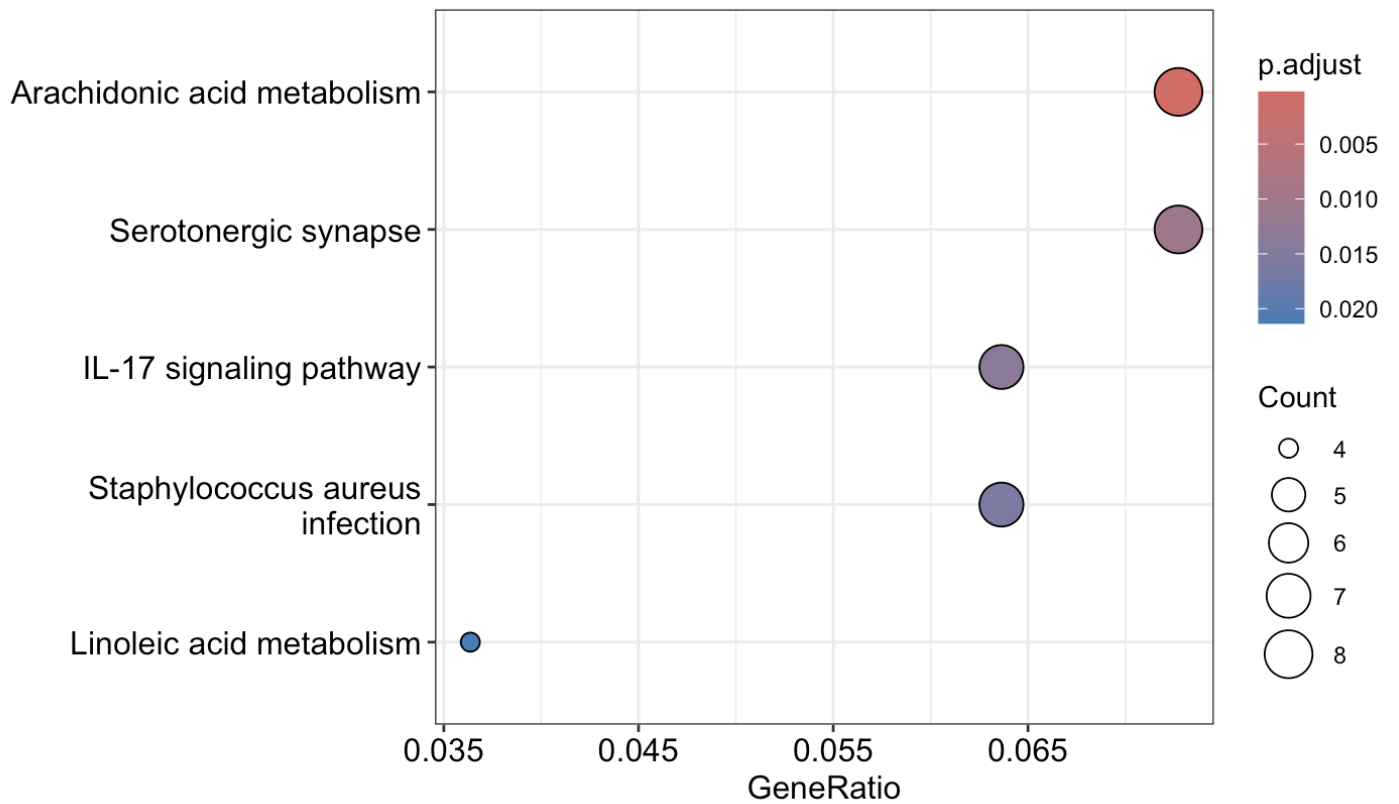
Cluster 2 - KEGG enrichment (Downregulated)



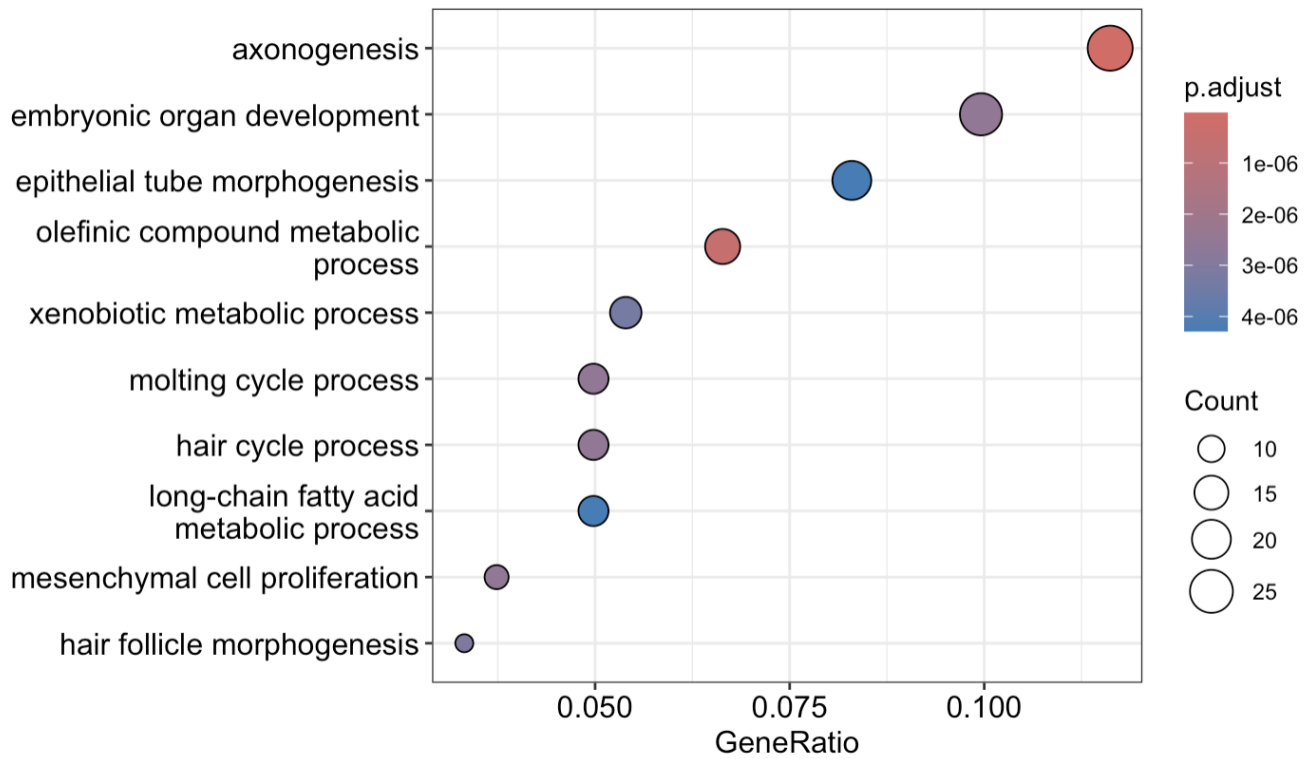
Cluster 3 - GO enrichment (Downregulated)



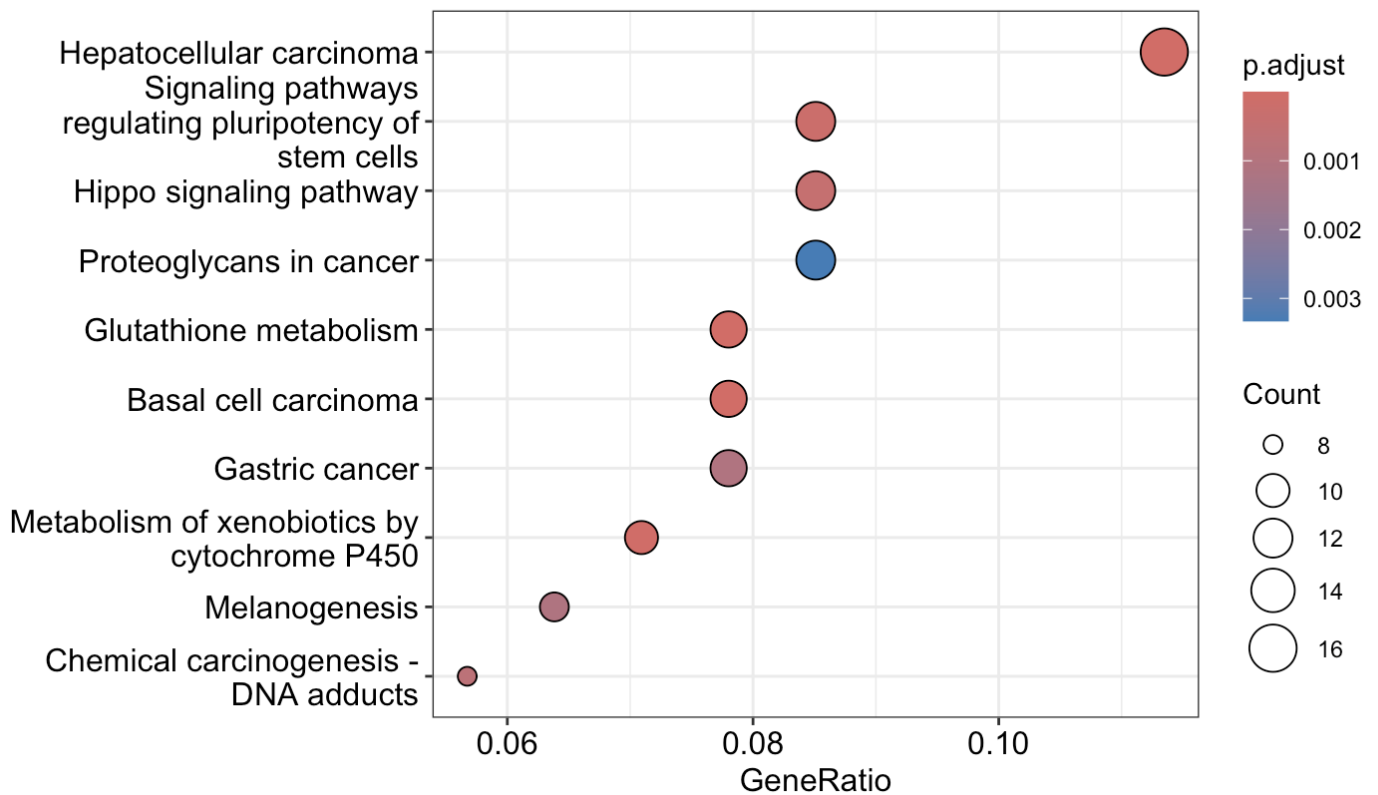
Cluster 3 - KEGG enrichment (Downregulated)



Cluster 4 - GO enrichment (Downregulated)



Cluster 4 - KEGG enrichment (Downregulated)



conclusion

Cluster	Top 3 GO Biological Process (Downregulated)	Top 3 KEGG Pathways (Downregulated)
Cluster 1	positive regulation of cell adhesion extracellular matrix organization extracellular structure organization	Cytoskeleton in muscle cells Cell adhesion molecules Staphylococcus aureus infection
Cluster 2	epidermis development skin development extracellular matrix organization	Human papillomavirus infection PI3K-Akt signaling pathway Cytokine-cytokine receptor interaction
Cluster 3	epidermis development skin development epidermal cell differentiation	Arachidonic acid metabolism Serotonergic synapse IL-17 signaling pathway
Cluster 4	axonogenesis embryonic organ development epithelial tube morphogenesis	Hepatocellular carcinoma Signaling pathways regulating pluripotency of stem cells Hippo signaling pathway

## Q4

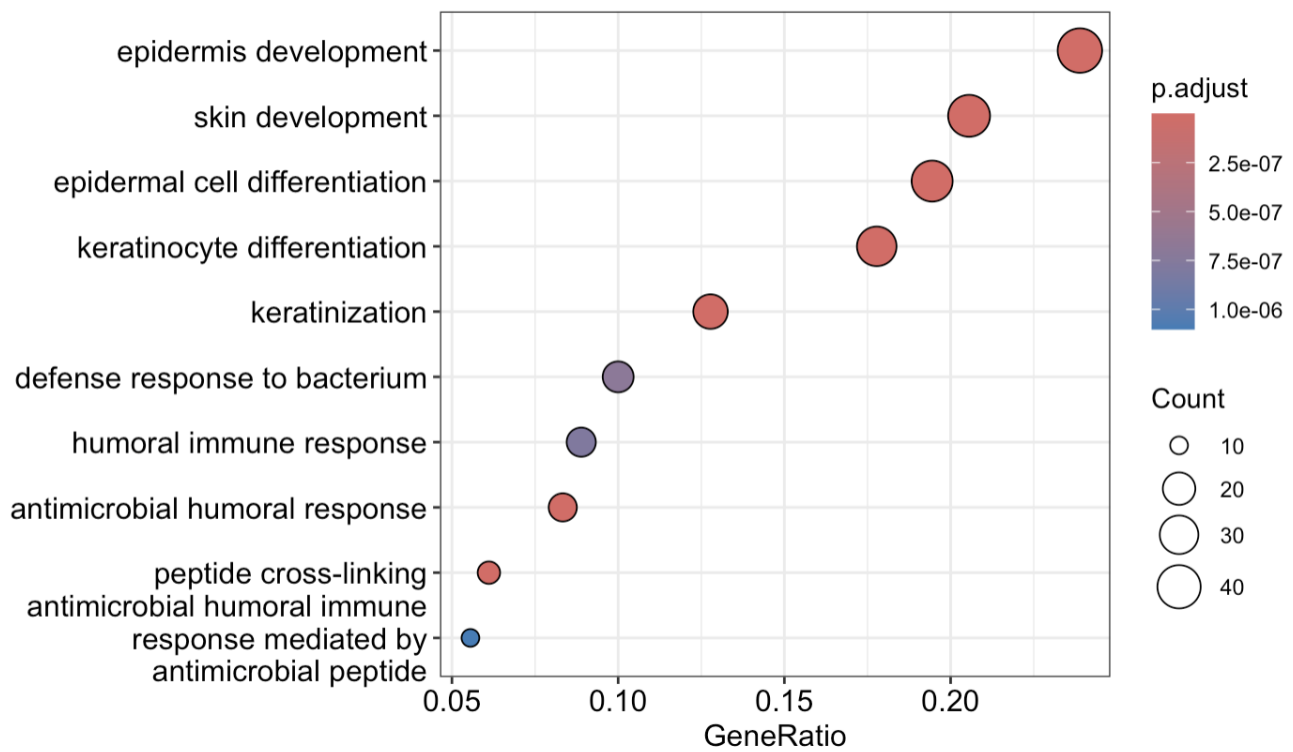
I selected upregulated and downregulated genes based on these two criteria separately

```
$logFC > 0 & $adj.P.Val < 0.05
```

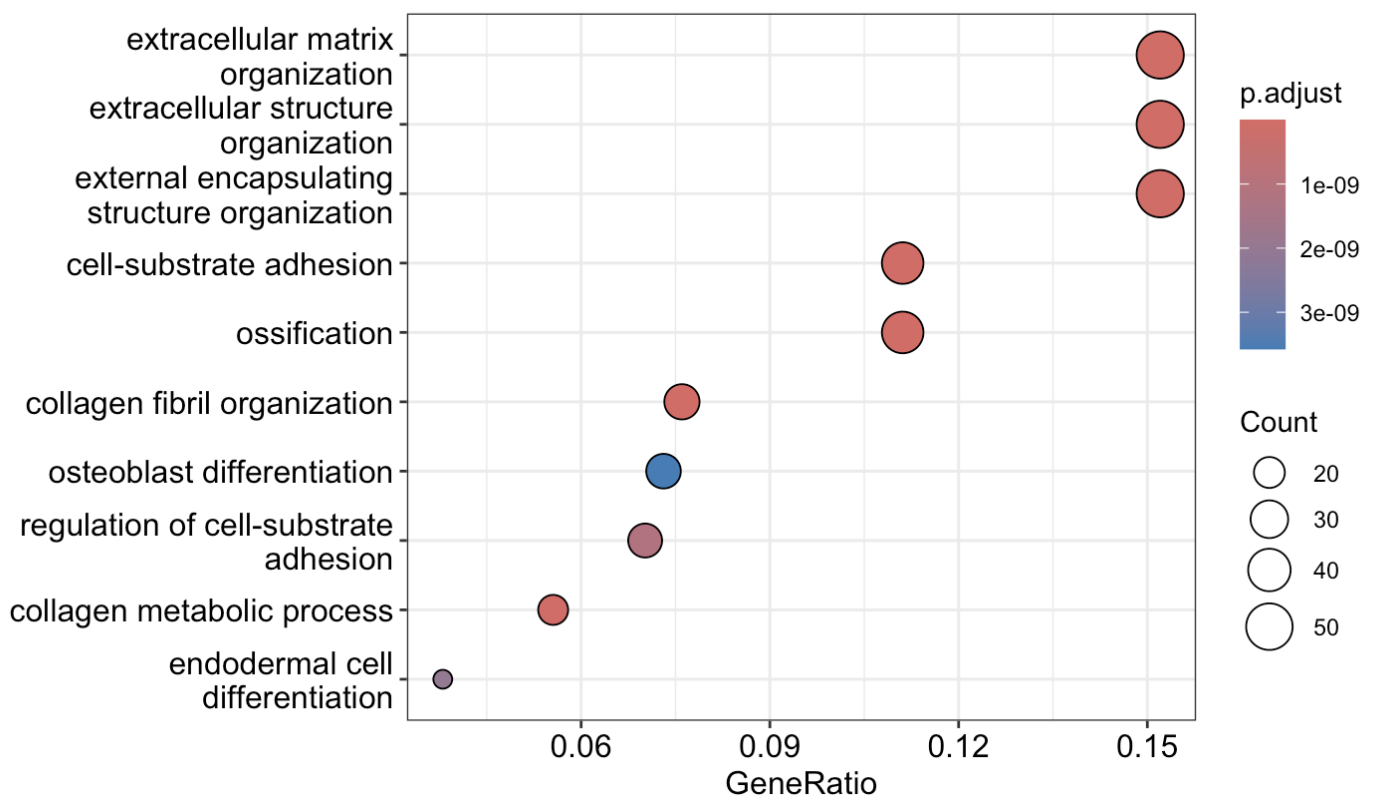
```
$logFC < 0 & $adj.P.Val < 0.05
```

## Cluster 1

Consensus Clustering - Cluster 1 Upregulated GO

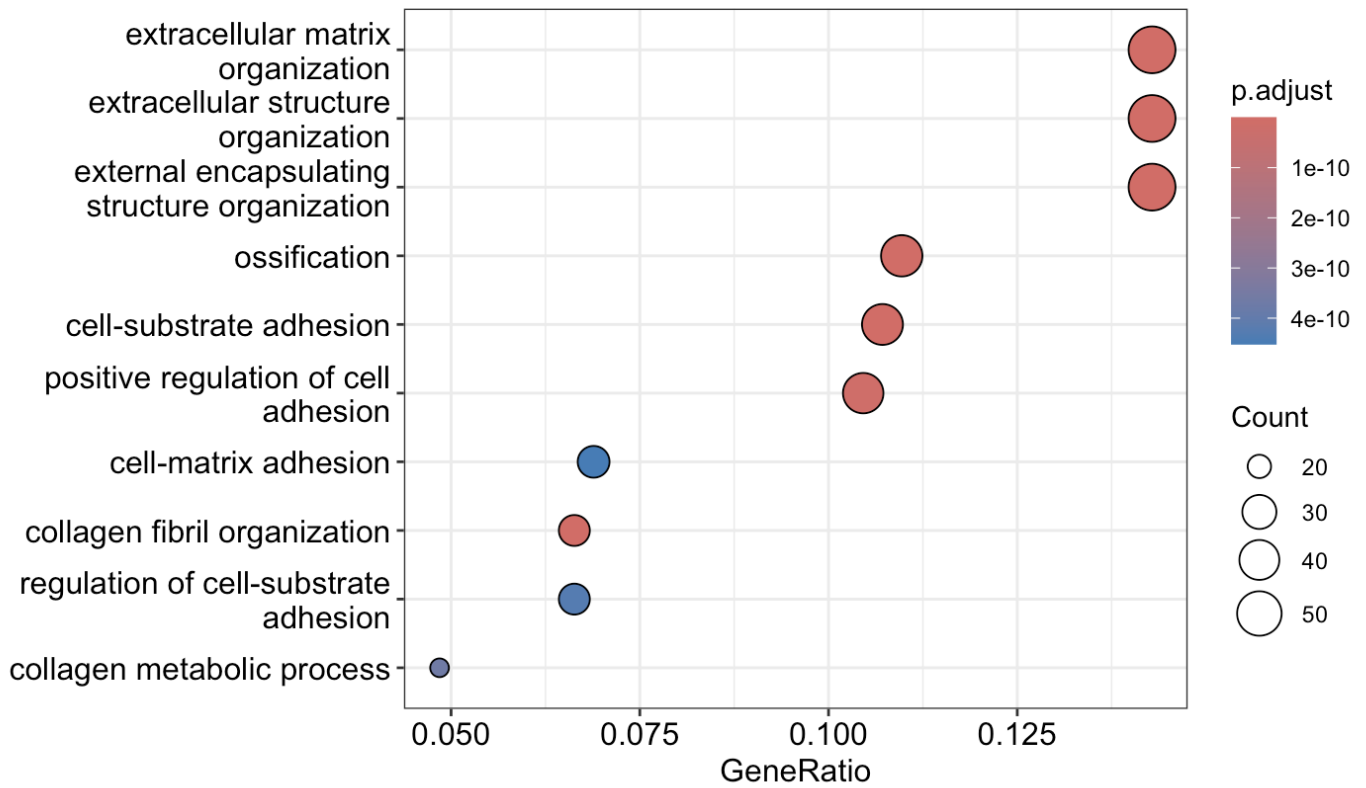


Consensus Clustering - Cluster 1 Downregulated GO

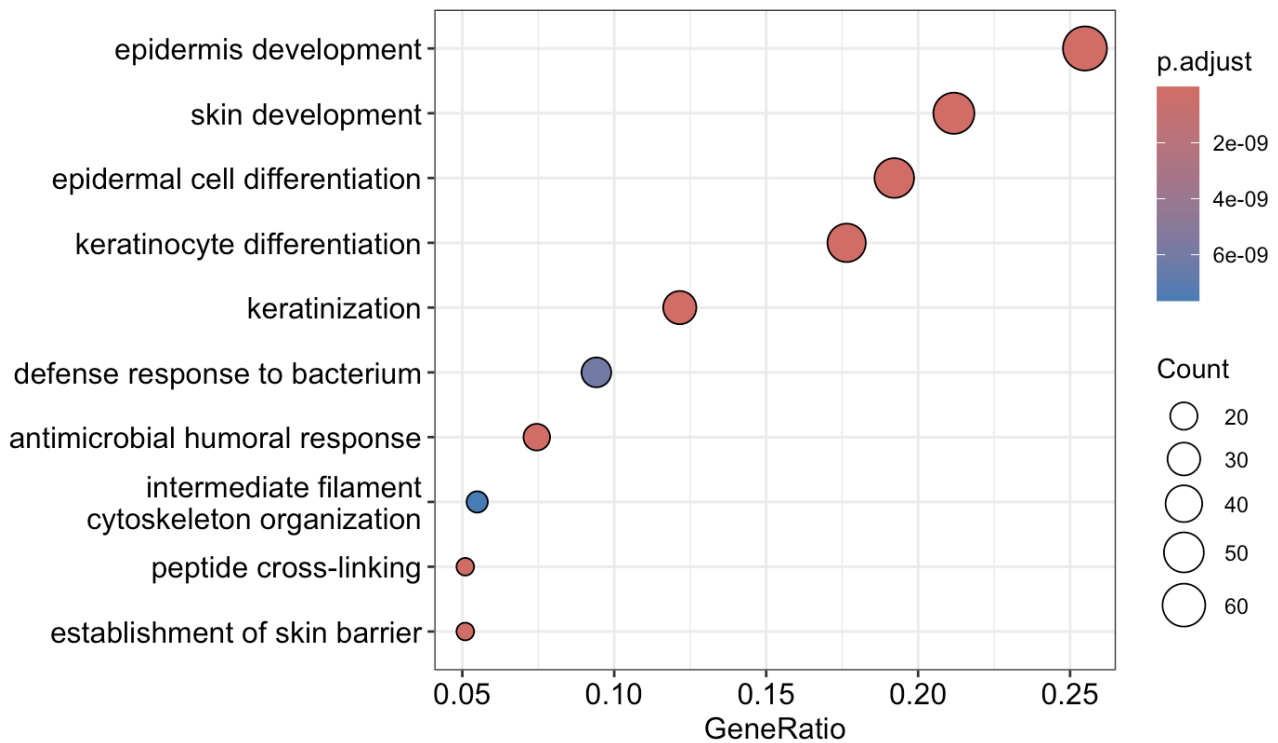


Cluster 2

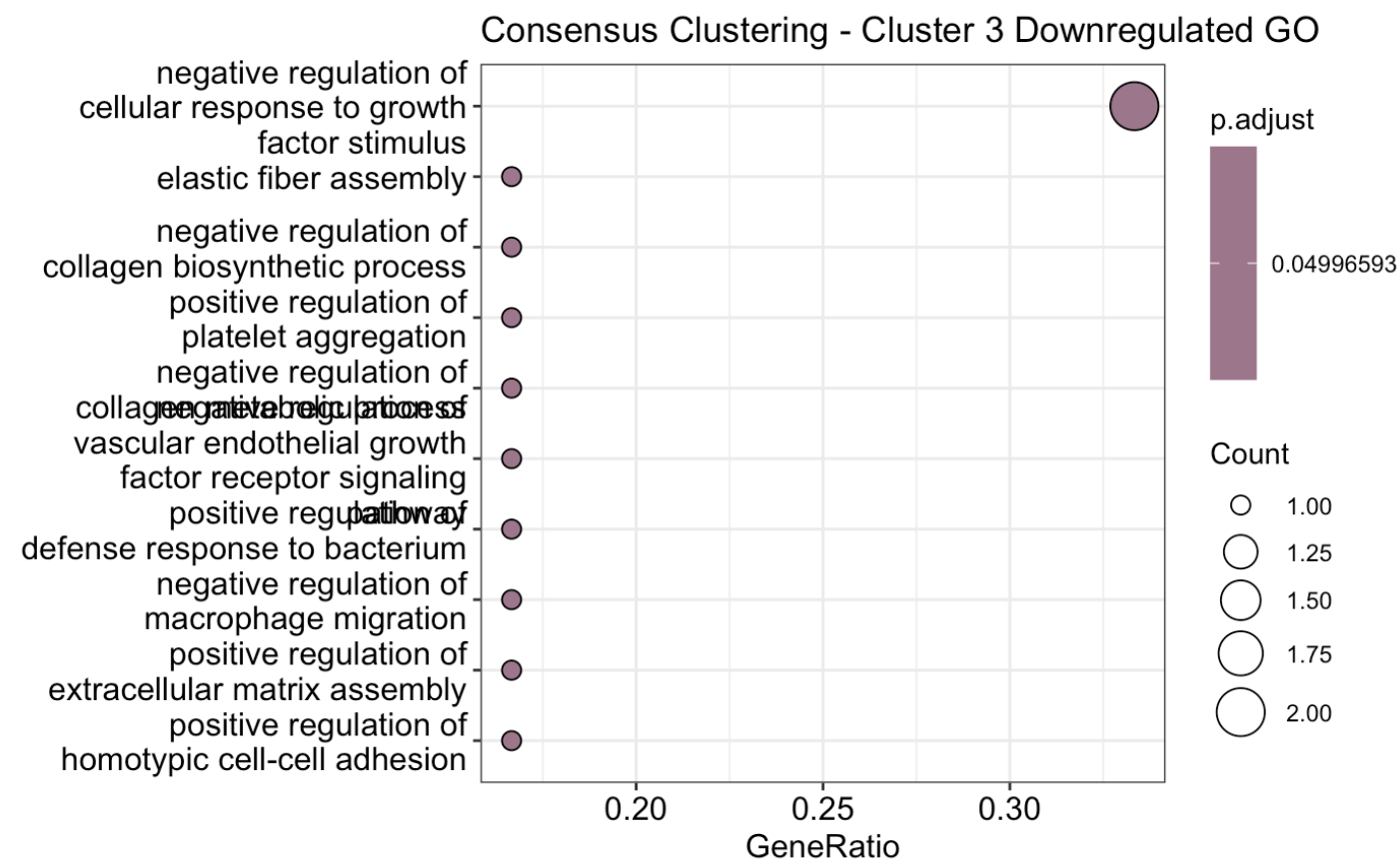
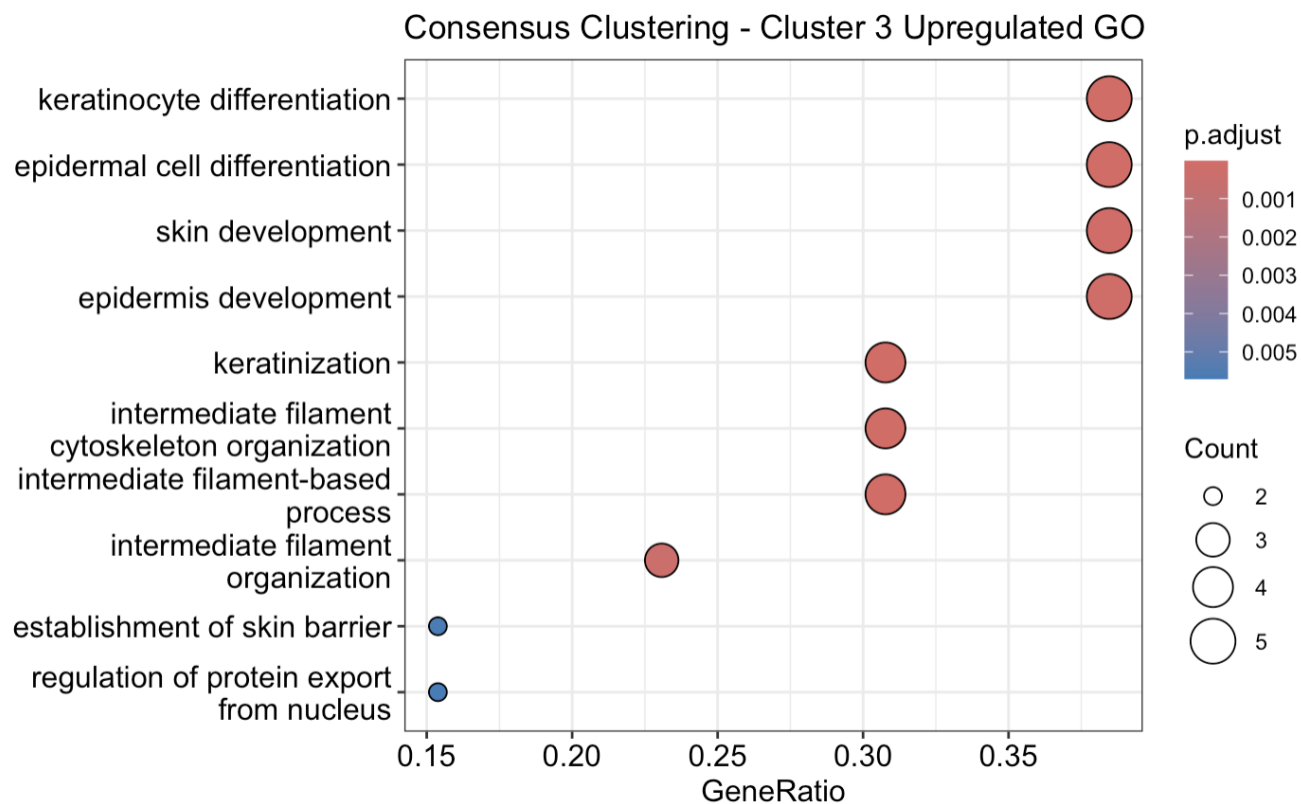
Consensus Clustering - Cluster 2 Upregulated GO



Consensus Clustering - Cluster 2 Downregulated GO



### Cluster 3

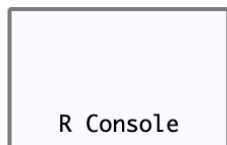


Cluster 4



```
`{r}
```

```
dotplot(go_up_ccp[[4]], title = "Consensus Clustering - Cluster 4 Upregulated GO")
dotplot(go_down_ccp[[4]], title = "Consensus Clustering - Cluster 4 Downregulated GO")
`{r}
```



Error in ans[ypos] <- rep(yes, length.out = len)[ypos] :  
replacement has length zero

[Show Traceback](#)

No plot here

## conclusion

Cluster	Top 3 Upregulated GO Terms	Top 3 Downregulated GO Terms
Cluster 1	epidermis development skin development epiderma l cell differentiation	extracellular matrix organization extracellul ar structure organization external encapsulating structure organization
Cluster 2	extracellular matrix organization extracellul ar structure organization external encapsulating structure organization	epidermis development skin development epiderma l cell differentiation
Cluster 3	keratinocyte differentiation epiderm al cell differentiation skin development	negative regulation of cellular response to growth factor stimulus elastic fiber assembly negative regulation of collagen biosynthetic process
Cluster 4	<i>No significant GO enrichment</i>	<i>No significant GO enrichment</i>

## Q5

### up

To compare the clustering results obtained from NMF and Consensus Clustering, we examined both cluster membership and functional enrichment outputs.

The cluster assignments from both methods show considerable overlap. Cross-tabulation of cluster membership revealed near one-to-one correspondence, suggesting high structural consistency across algorithms.

Functionally, GO enrichment analysis of upregulated genes showed strong agreement between corresponding clusters. For example, both methods identified "epidermis development", "skin development", and "epidermal cell differentiation" as the top enriched biological processes in Cluster 1. Similarly, both methods captured extracellular matrix-related signatures in another cluster.

However, discrepancies were observed in Cluster 4 from the Consensus Clustering, which lacked significant upregulated or downregulated GO terms, unlike NMF. This may indicate that Consensus Clustering grouped less coherent or less biologically distinct samples in that cluster.

Overall, the clustering structure and biological interpretations between the two methods were largely consistent, especially for key functional clusters.

### down

When comparing the GO enrichment results of downregulated genes across NMF and Consensus Clustering, Clusters 1 and 2 demonstrated strong consistency. Both methods identified processes related to extracellular matrix organization and epidermal development, respectively.

Cluster 3 showed discrepancies: while NMF revealed skin-related pathways, the Consensus Clustering method highlighted immune-related and regulatory processes, possibly reflecting differences in sample partitioning.

Cluster 4 from the Consensus Clustering lacked significantly enriched downregulated pathways, whereas NMF detected nervous system development-related terms. This indicates that NMF may

better capture certain biologically relevant subgroups in this case.

Overall, the GO enrichment results from both methods are largely consistent for the major clusters, with minor differences in specific subgroups.

## Q6

To generate the most robust and reliable clusters from gene expression data in an unsupervised setting, the following procedure is recommended:

1. **Feature Selection:** Start with a variance-based filtering strategy, such as MAD (median absolute deviation), to select the most informative genes (e.g., top 1500–3000 genes) and reduce noise.
2. **Apply Multiple Clustering Algorithms:** Use at least two complementary methods, such as:
  - **Non-negative Matrix Factorization (NMF):** Captures parts-based representation and is particularly suitable for gene expression data.
  - **Consensus Clustering:** Evaluates the stability of clustering results through resampling, useful for determining robust cluster membership.
3. **Cluster Number Estimation:**
  - Use internal metrics (e.g., cophenetic coefficient, silhouette width, residuals) to determine the optimal number of clusters.
  - Avoid choosing  $k = 2$  unless biologically justified, as it often oversimplifies the heterogeneity.
4. **Validation:**
  - Perform **biological pathway enrichment (GO/KEGG)** to interpret each cluster.
  - Assess whether clusters are biologically meaningful and non-overlapping in terms of functional enrichment.
5. **Cross-Comparison:**

- Compare results from multiple clustering methods.
- Use consensus or overlapping clusters as a confidence measure for stability.

## 6. Visualization and Interpretability:

- Visualize consensus matrices, silhouette plots, and enrichment bubble plots.
- Clear visualization helps reveal outliers and cluster boundary ambiguity.

## 7. Biological Input and Reproducibility:

- Consult domain knowledge for validating clusters.
- Ensure reproducibility via random seeds, transparent parameter reporting, and standard pipelines.

This integrative and comparative strategy helps identify **robust, reproducible, and biologically relevant clusters** rather than relying on a single method or metric.