

Multi-Omics analysis identifies signature genes to predict bladder cancer survival

Group 2: Pham Gia Cuong, Matanat Mammadli, Samra Hamidovic



Overview

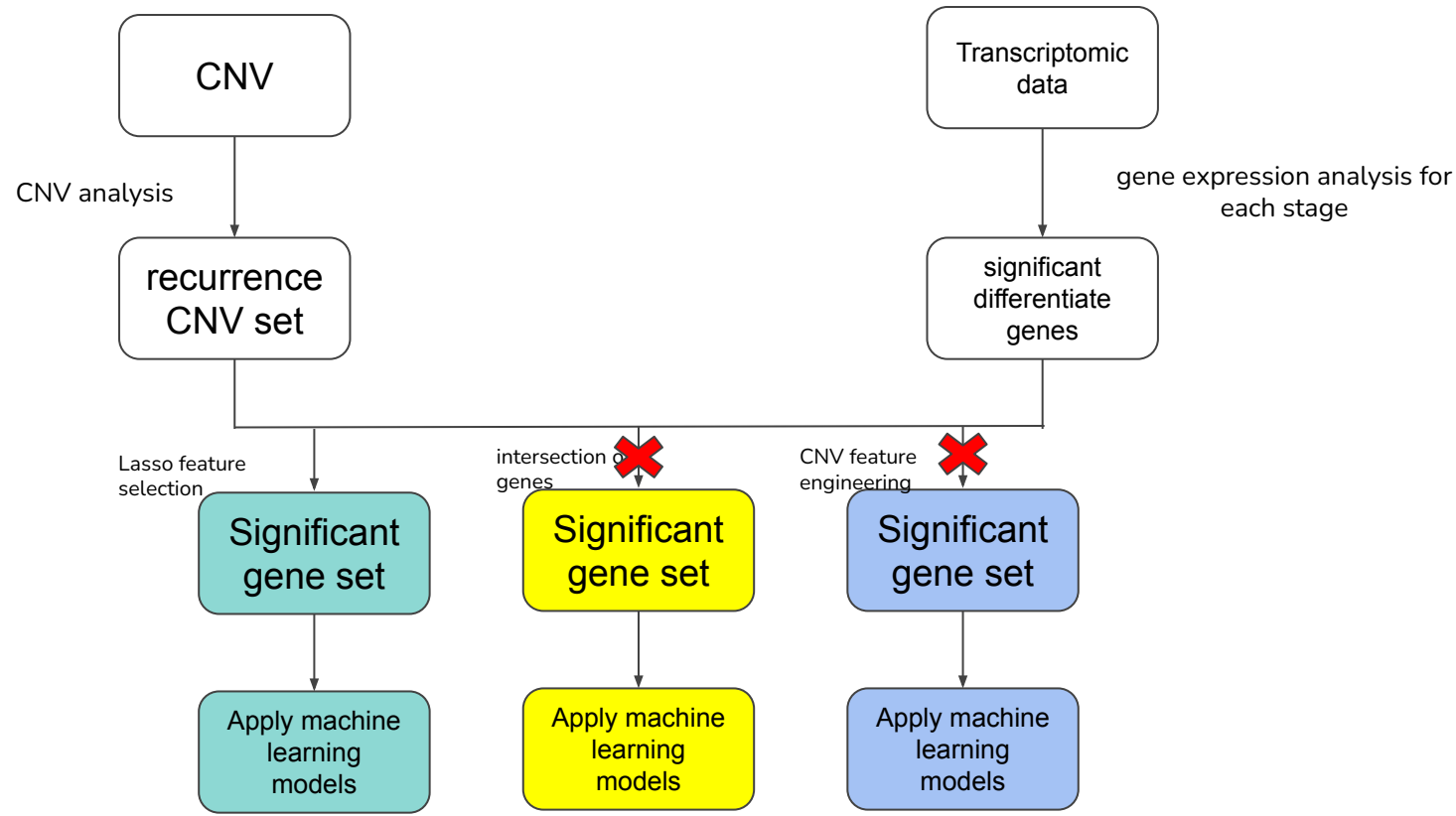
1. Remind
2. Workflow
3. CNV analysis
4. Gene expression analysis
5. Feature selection
6. Apply machine learning model
7. Survival analysis
8. GO enrichment analysis



Remind:

- Aim:
 - Using CNV data and transcriptomic data to predict survival of Bladder Cancer patients
 - Key genes contain CNV and play a big role in survival of patient
 - Functional annotation
- Data:
 - TCGA Bladder cancer data
 - CNV data: CNVs from tumor and normal cells of different samples (within-sample homogeneity)
 - Transcriptomic data: 442 samples
 - Survival data: survival label and survival time.
 - Phenotype data:
 - tumor diagnose stage: 422 samples
 - lost of follow up patient: 150 samples are lost of follow up

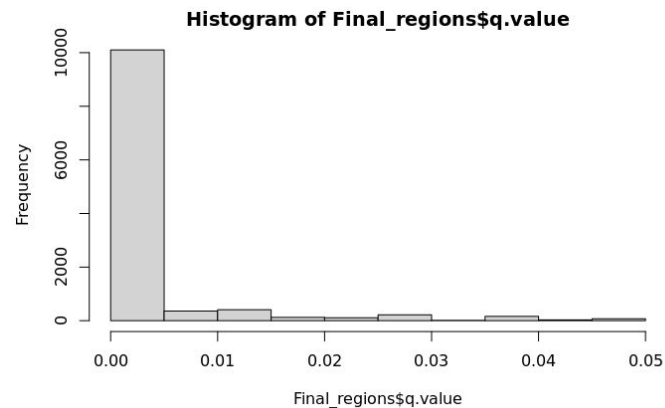
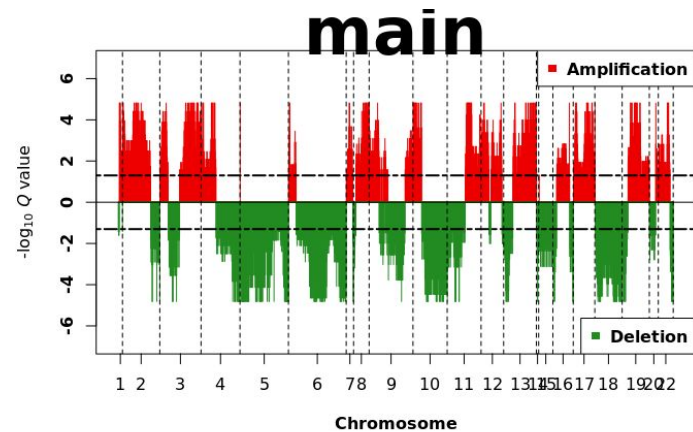
Workflow





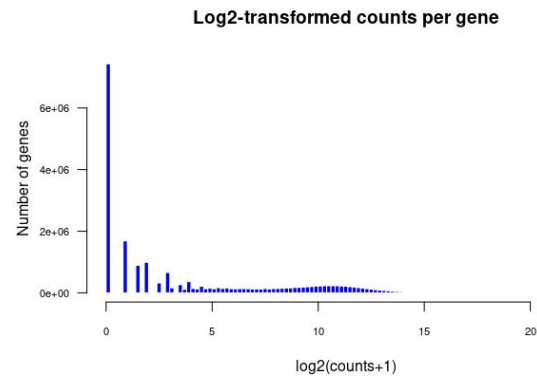
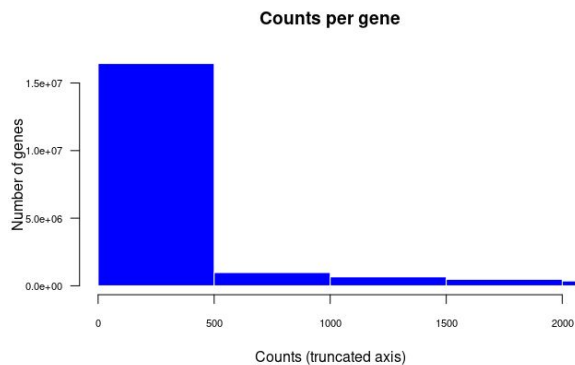
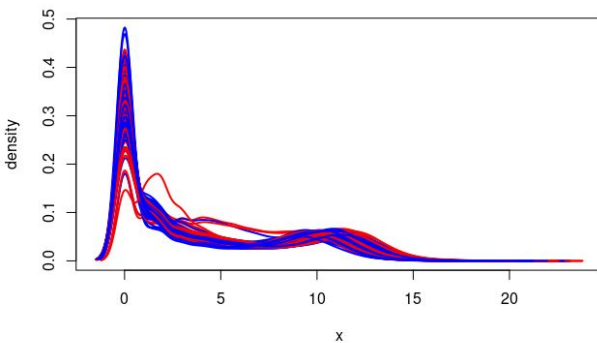
CNV analysis

- 412 samples , each sample contains tumor and normal cell types
- CNVs: chromosome, start, end, probe and segmean
- Aim: finding independent and recurrent copy number abbreviations
 - r-package: GAIA (genomic analysis of significant chromosomal aberrations)
- Result: 3448 segments ($q.value < 0.01$) -> 11591 genes
- Cross checking with gene profile from gene expression data -> 9254 genes





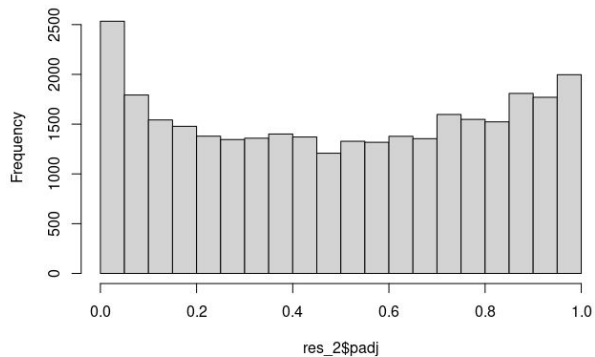
Gene expression analysis





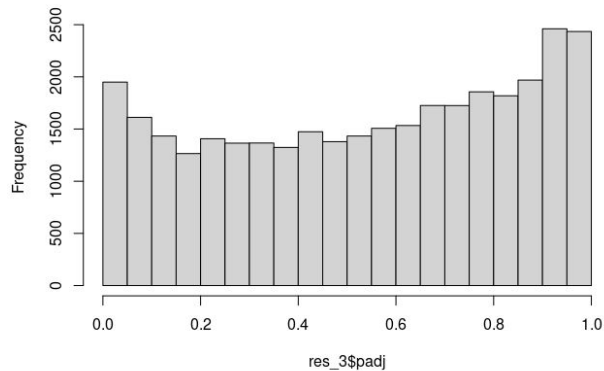
Gene expression analysis

Histogram of res_2\$padj



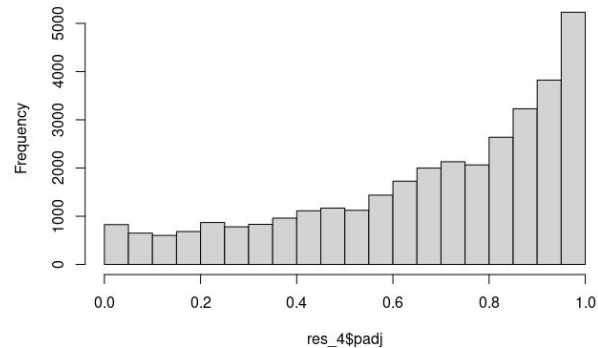
Stage 2

Histogram of res_3\$padj



Stage 3

Histogram of res_4\$padj

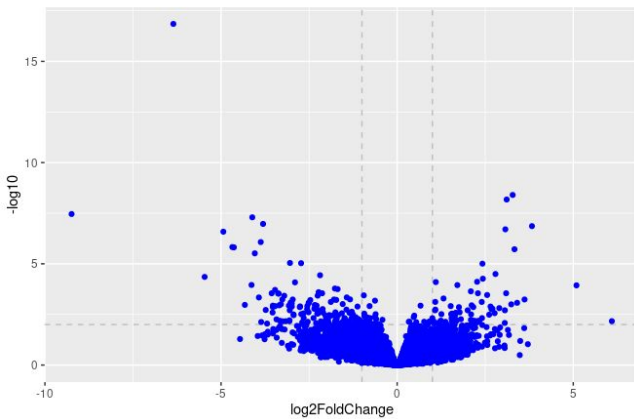


Stage 4

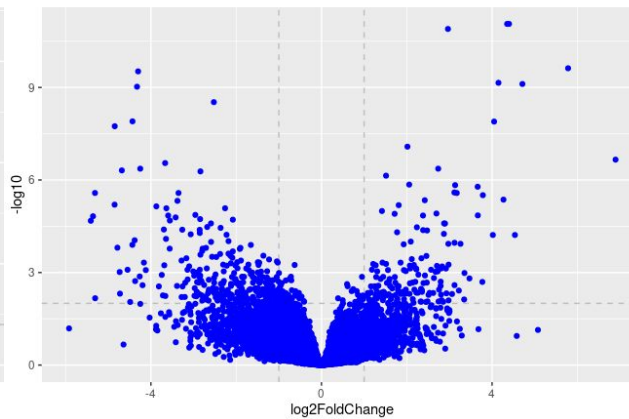


Gene expression analysis

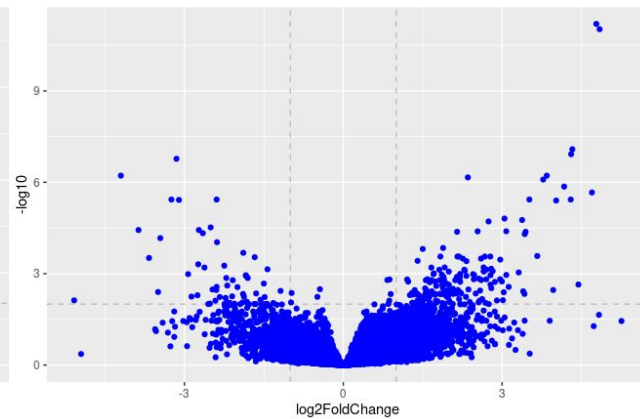
Volcano Plot of stage 2



Volcano Plot of stage 3

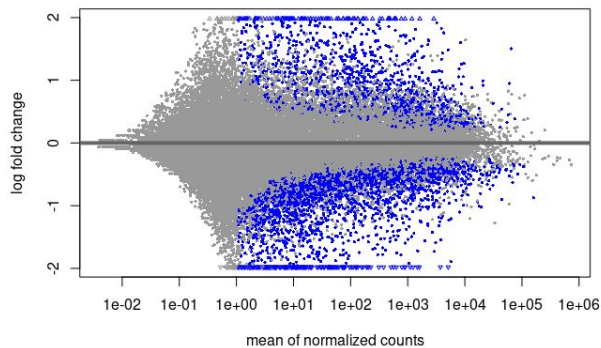


Volcano Plot of stage 4

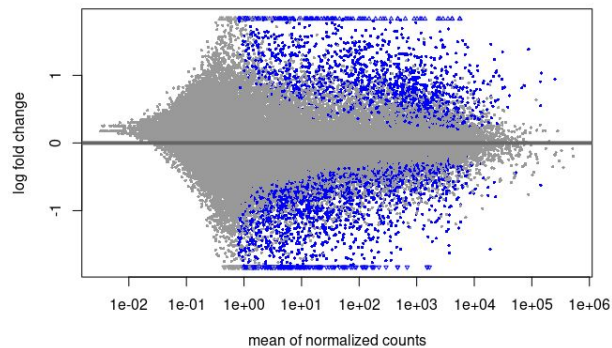




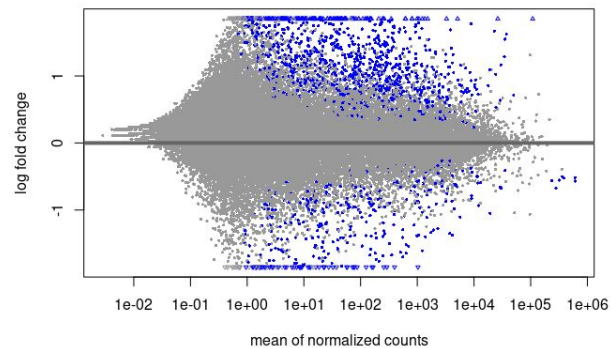
Gene expression analysis



Stage 2



Stage 3



Stage 4

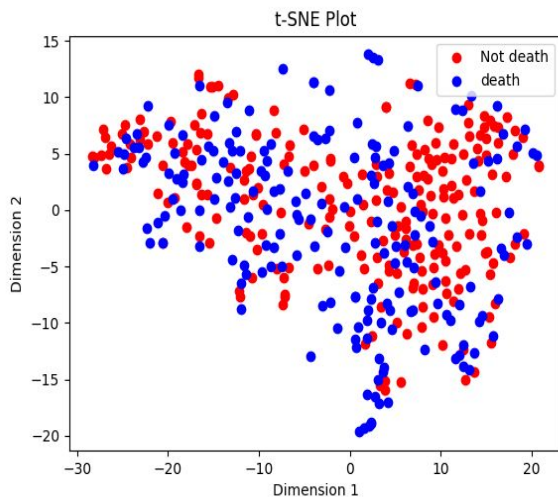


Gene expression analysis

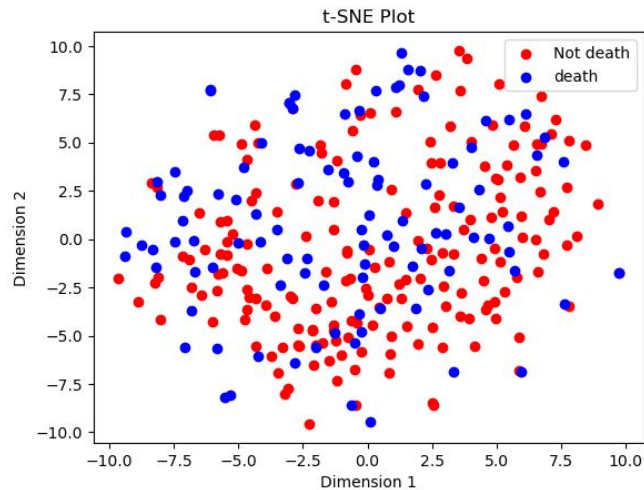
	significant genes	samples
stage 2	619	111
stage 3	1003	97
stage 4	607	84



tSNE plot

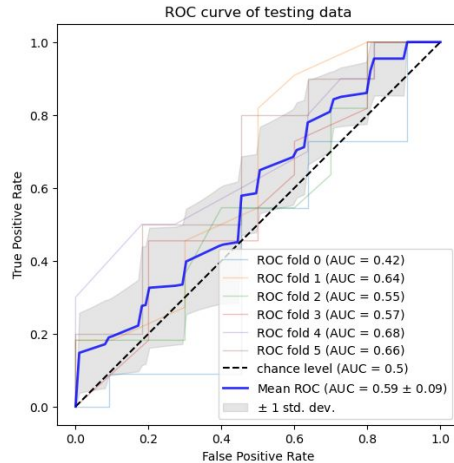


t-SNE on original raw count

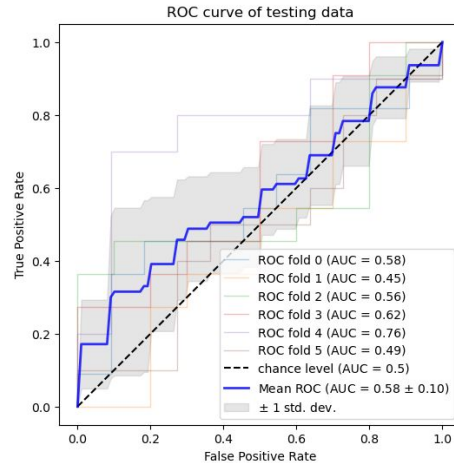


t-SNE on significant genes raw count

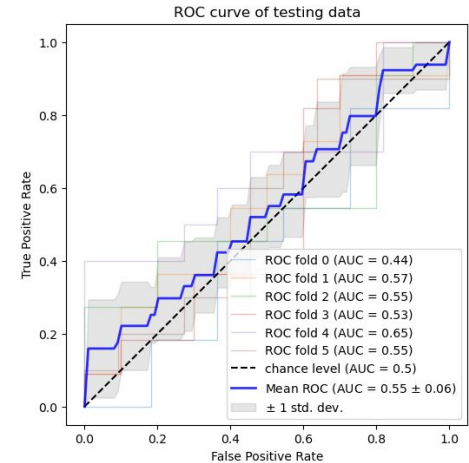
Apply machine learning models



Random forest



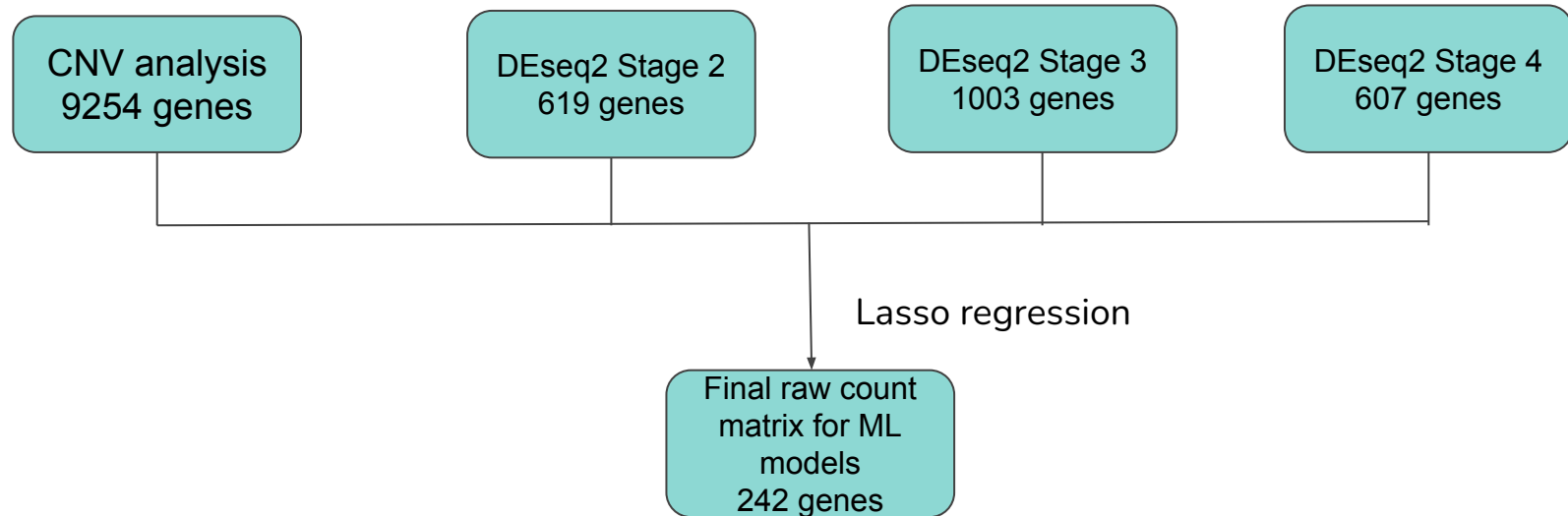
Ridge regression



XGBoost

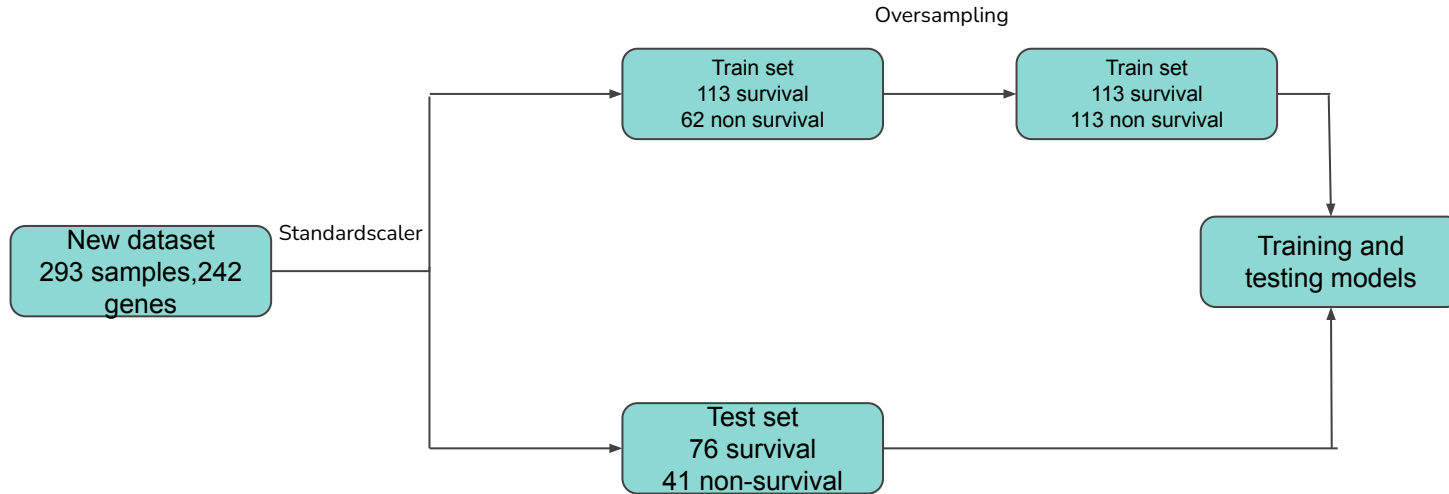
Result with the features from intersection
of CNV gene set and DESeq2

Apply machine learning model

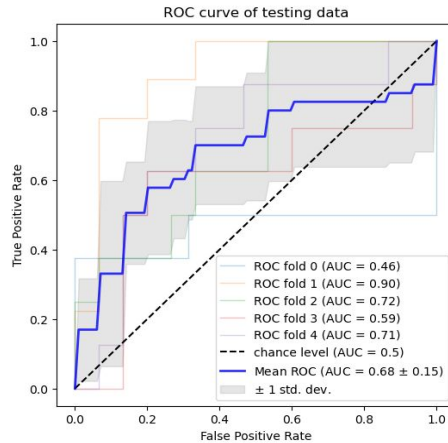




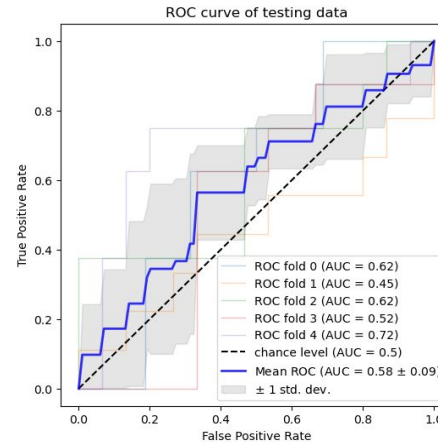
Apply machine learning model



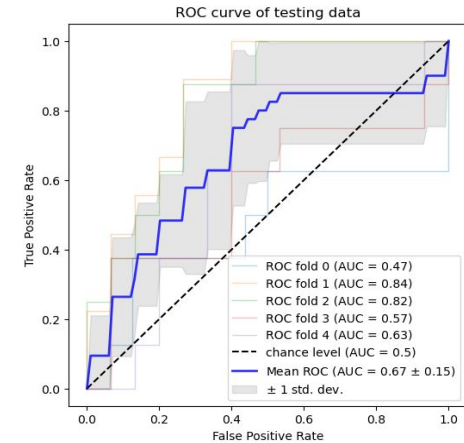
Apply machine learning model



Random forest



Ridge regression



XG Boost



Random forest

Hyperparameters:

- max features : sqrt, log2
- max depth of each tree : 10, 12, 14, 16, 20
- criterion: gini, entropy

Gridsearchcv -> max_depth=10, max_features='log2'



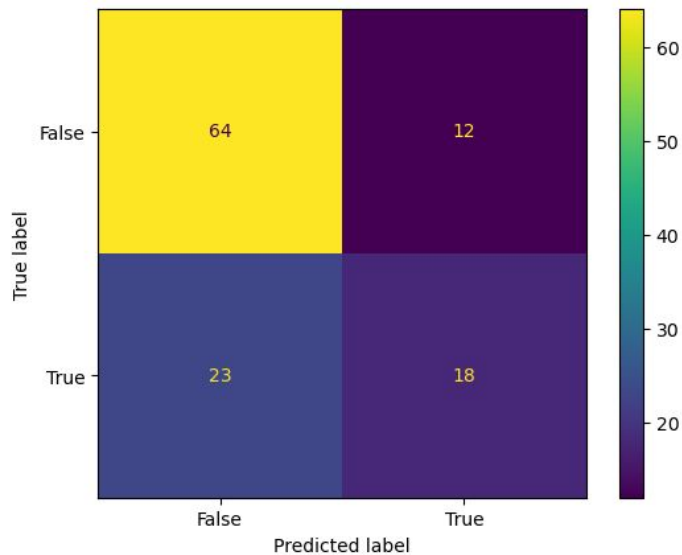
Random forest

Test set

	precision	recall	f1-score	support
0	0.74	0.84	0.79	76
1	0.60	0.44	0.51	41
accuracy			0.70	117
macro avg	0.67	0.64	0.65	117
weighted avg	0.69	0.70	0.69	117

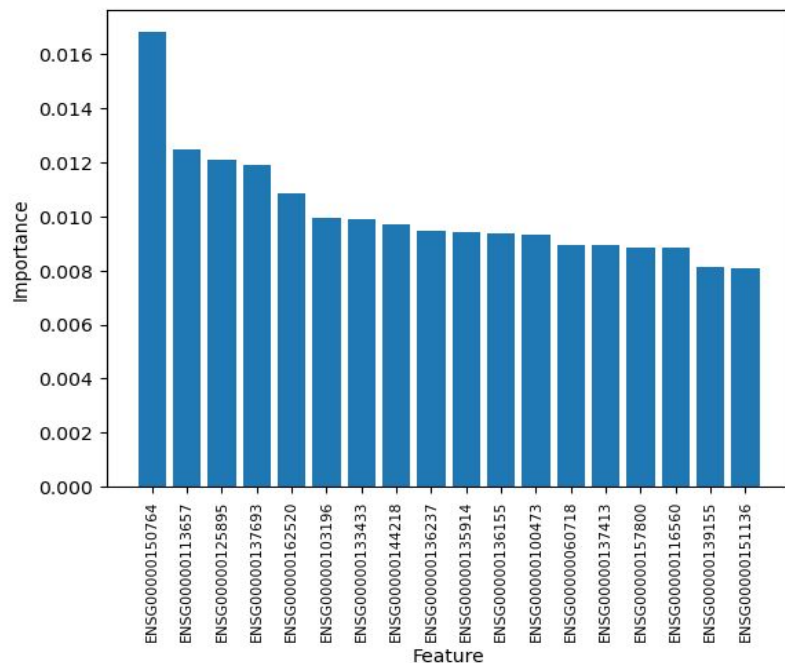
Train set

	precision	recall	f1-score	support
0	1.00	1.00	1.00	113
1	1.00	1.00	1.00	113
accuracy			1.00	226
macro avg	1.00	1.00	1.00	226
weighted avg	1.00	1.00	1.00	226



OVERFITTED :(

Features importance

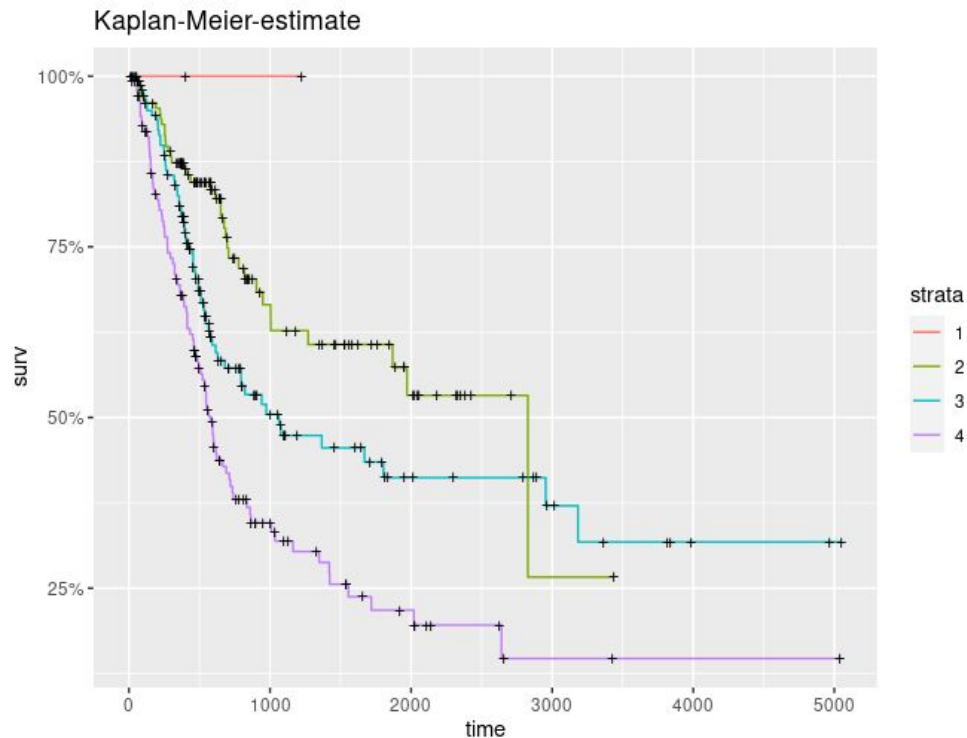


Ensemble ID	Gene name	Cancer related information
ENSG00000150764	<i>DIXDC1</i>	Participating in growing of tumor (1)
ENSG00000113657	<i>DPYSL3</i>	high DPYSL3 expression predicted a higher bladder tumor recurrence rate in patients (2)
ENSG00000125895	<i>TMEM74</i>	High expression of TMEM74 significantly shortens the surviving periods of patients in several types of cancer (3)
ENSG00000137693	<i>YAP 1</i>	Yap1 also plays an important role in the development of bladder and the deregulation of Yap1 is significantly associated with the development and metastasis of human bladder cancer (4)



Survival analysis

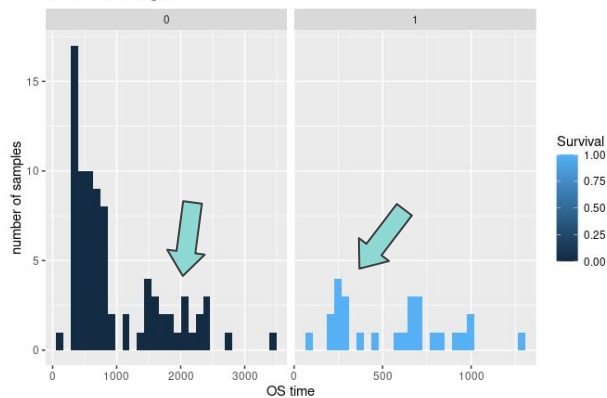
- **Reminder:**
Kaplan-Meier-curve for each stage
- Stage = **Confounder**



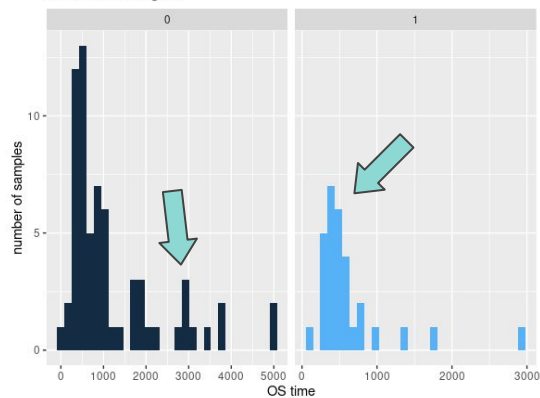


Survival analysis

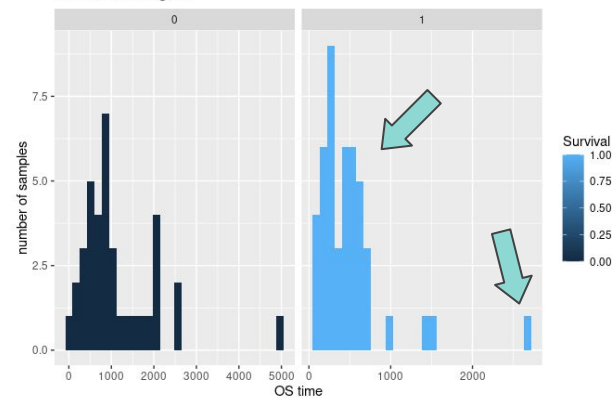
Survival of stage 2



Survival of stage 3



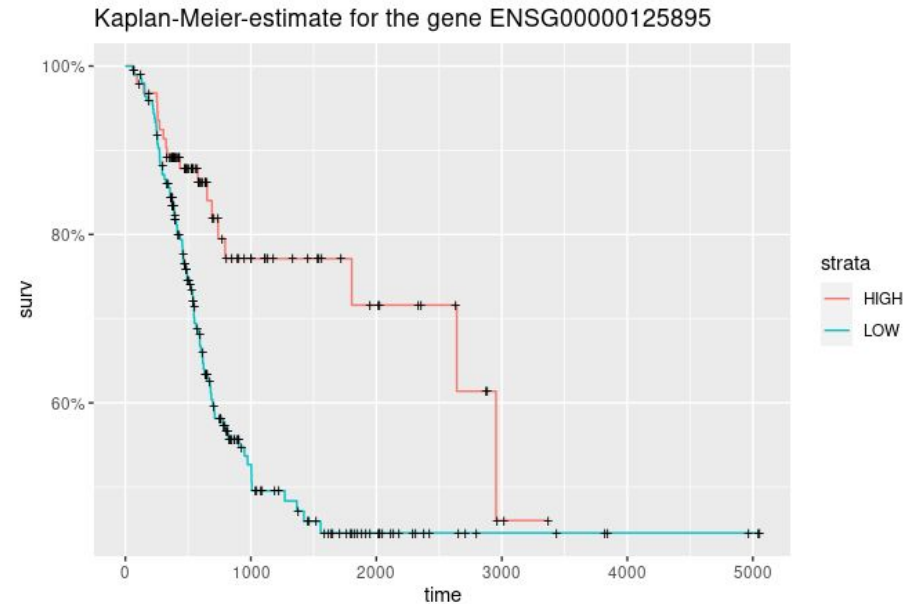
Survival of stage 4





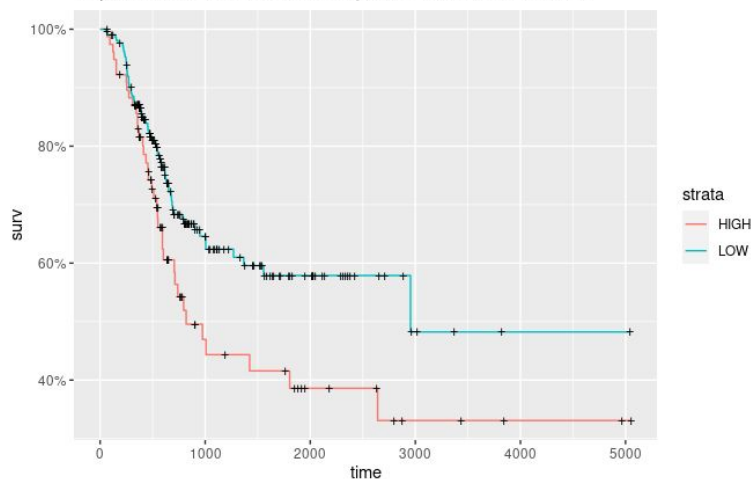
Survival analysis

- highest and lowest counts of one gene
- upregulated and downregulated
- Gene ***TMEM74***
- High expression shortens the surviving periods of patients



Survival analysis

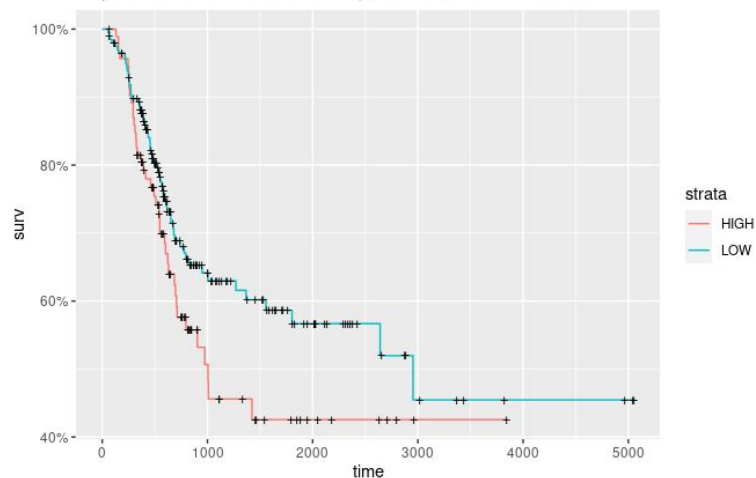
Kaplan-Meier-estimate for the gene ENSG00000113657



DPYSL3

high DPYSL3 expression predicted a higher bladder tumor recurrence rate

Kaplan-Meier-estimate for the gene ENSG00000137693

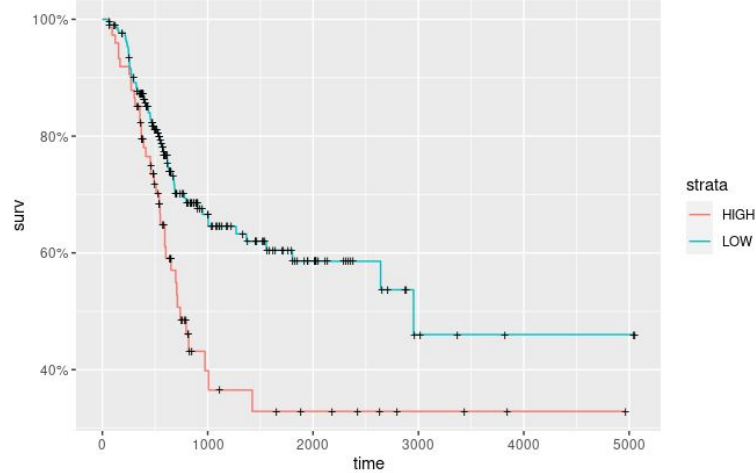


YAP1

deregulation of Yap1 is significantly associated with the development and metastasis of BLCA

Survival analysis

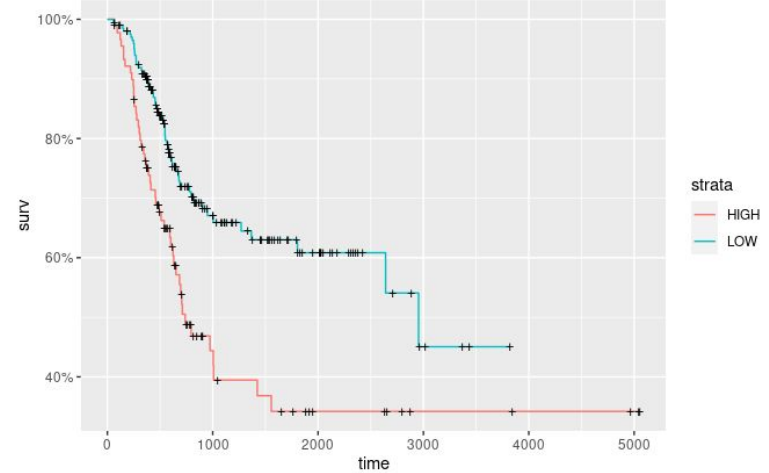
Kaplan-Meier-estimate for the gene ENSG00000150764



DIXDC1

high DPYSL3 expression predicted a higher bladder tumor recurrence rate

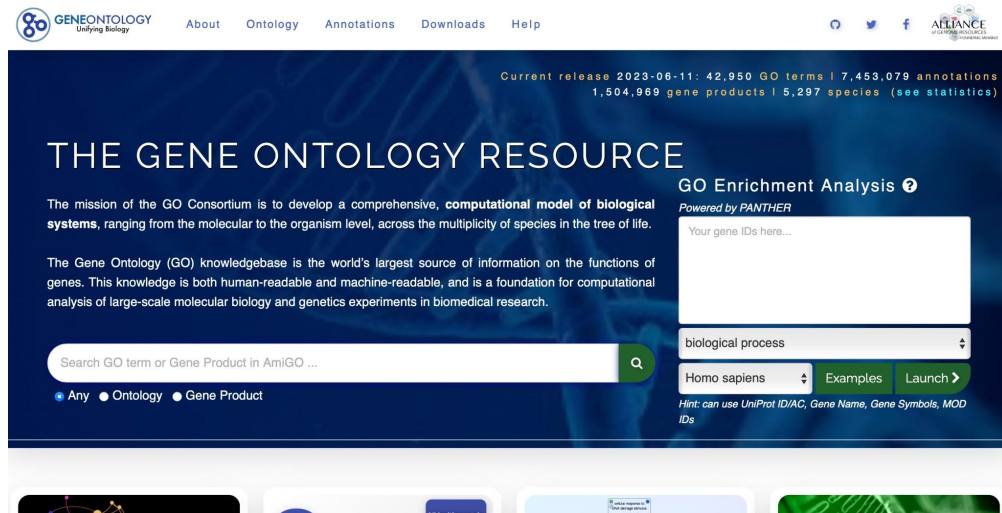
Kaplan-Meier-estimate for the gene ENSG00000162520



SYNC

GO Enrichment analysis

Gene Ontology (GO) enrichment analysis is used for interpreting high throughput molecular data and generating hypotheses about underlying biological phenomena of experiments.



The screenshot displays the Gene Ontology Resource website. At the top, the Gene Ontology logo and navigation links (About, Ontology, Annotations, Downloads, Help) are visible. A status bar indicates the current release (2023-06-11) with 42,950 GO terms, 7,453,079 annotations, 1,504,969 gene products, and 5,297 species. The main heading is "THE GENE ONTOLOGY RESOURCE". Below this, the mission statement and the description of the GO knowledgebase are provided. The "GO Enrichment Analysis" tool is highlighted, powered by PANTHER. It features a search bar for gene IDs, a dropdown menu for "biological process", and buttons for "Homo sapiens", "Examples", and "Launch". A hint at the bottom suggests using UniProt ID/AC, Gene Name, Gene Symbols, or MOD IDs.



GO Enrichment Analysis

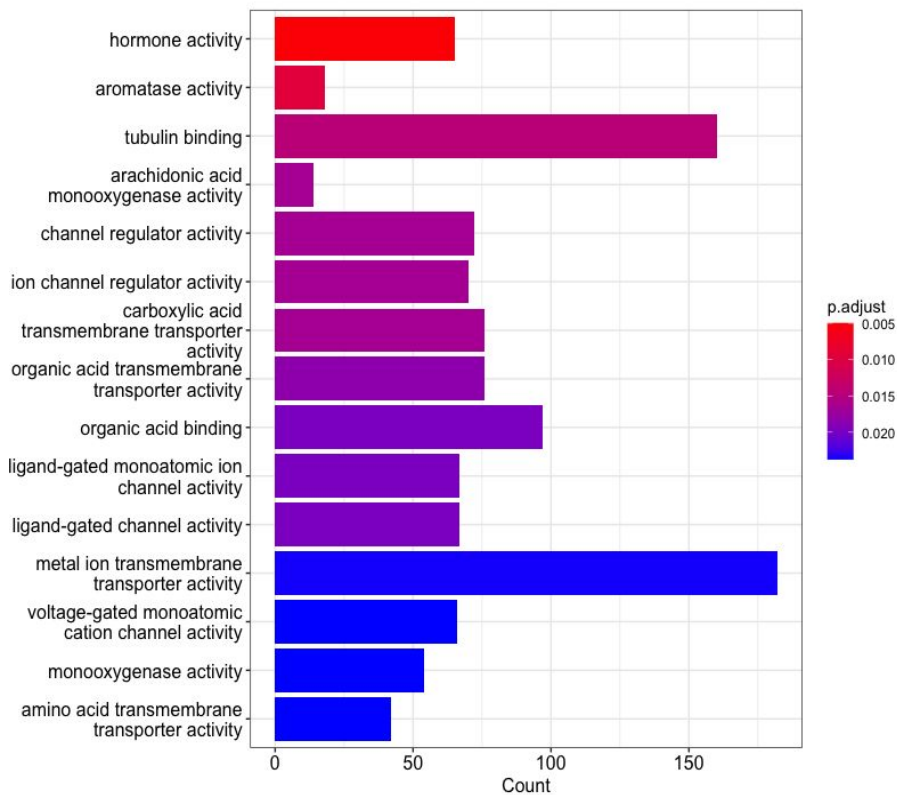
- One of the main uses of the GO (Gene Ontology) is to perform enrichment analysis on gene sets.
- For example, given a set of genes that are up-regulated under certain conditions, an enrichment analysis will find which GO terms are over-represented (or under-represented) using annotations for that gene set.

```
12 GO_results <- enrichGO(gene = genes_to_test,  
13                        OrgDb = "org.Hs.eg.db",  
14                        keyType = "ENSEMBL",  
15                        ont = "MF",  
16                        readable = TRUE  
17                        )  
18 GO_results
```



GO Enrichment analysis (Results)

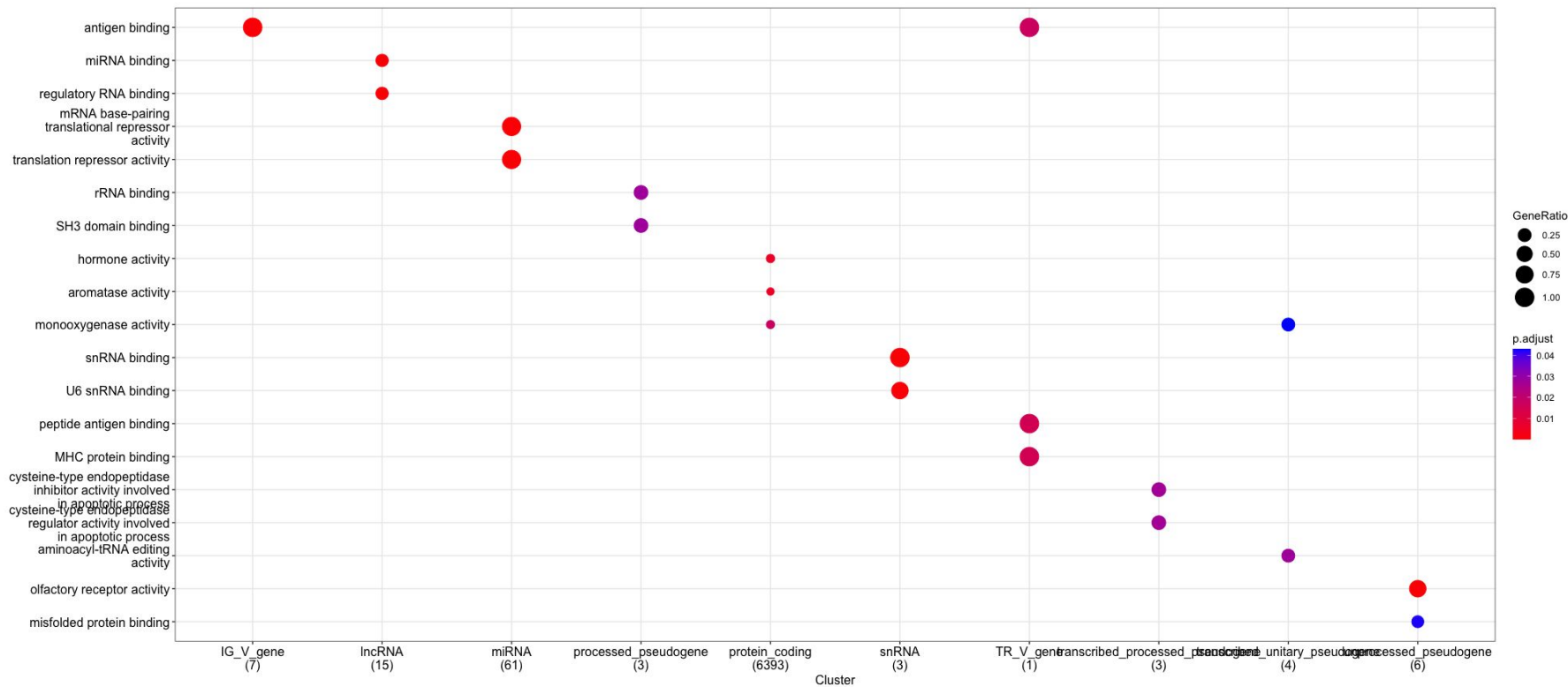
Barplot of our GO results





GO Enrichment analysis (Results)

clusterProfiler::dotplot()

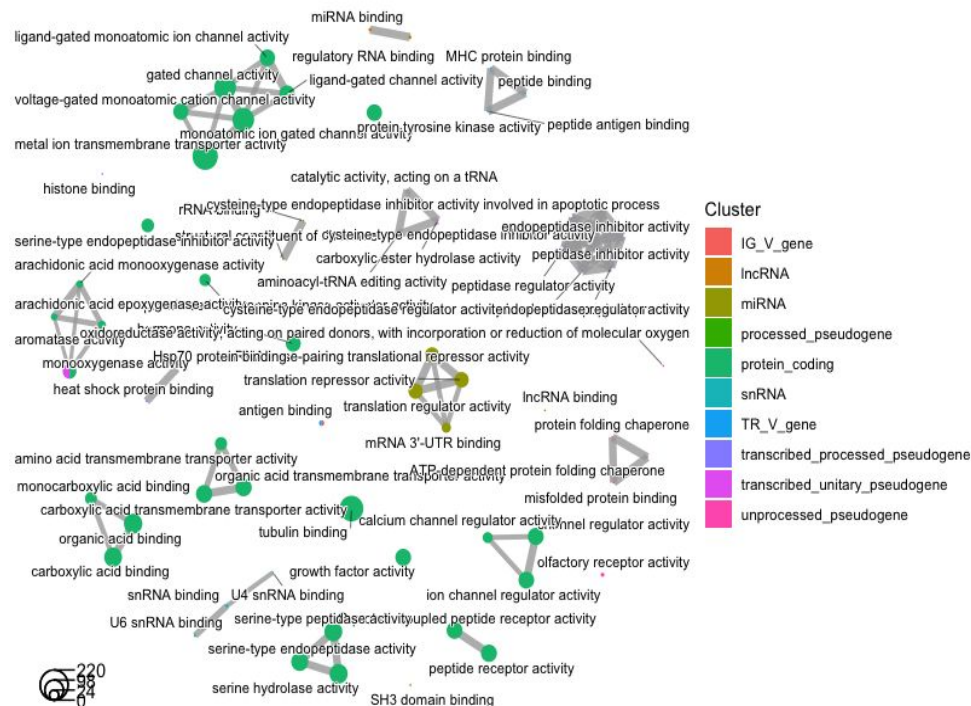


GO Enrichment analysis (Results)

`compareCluster()`

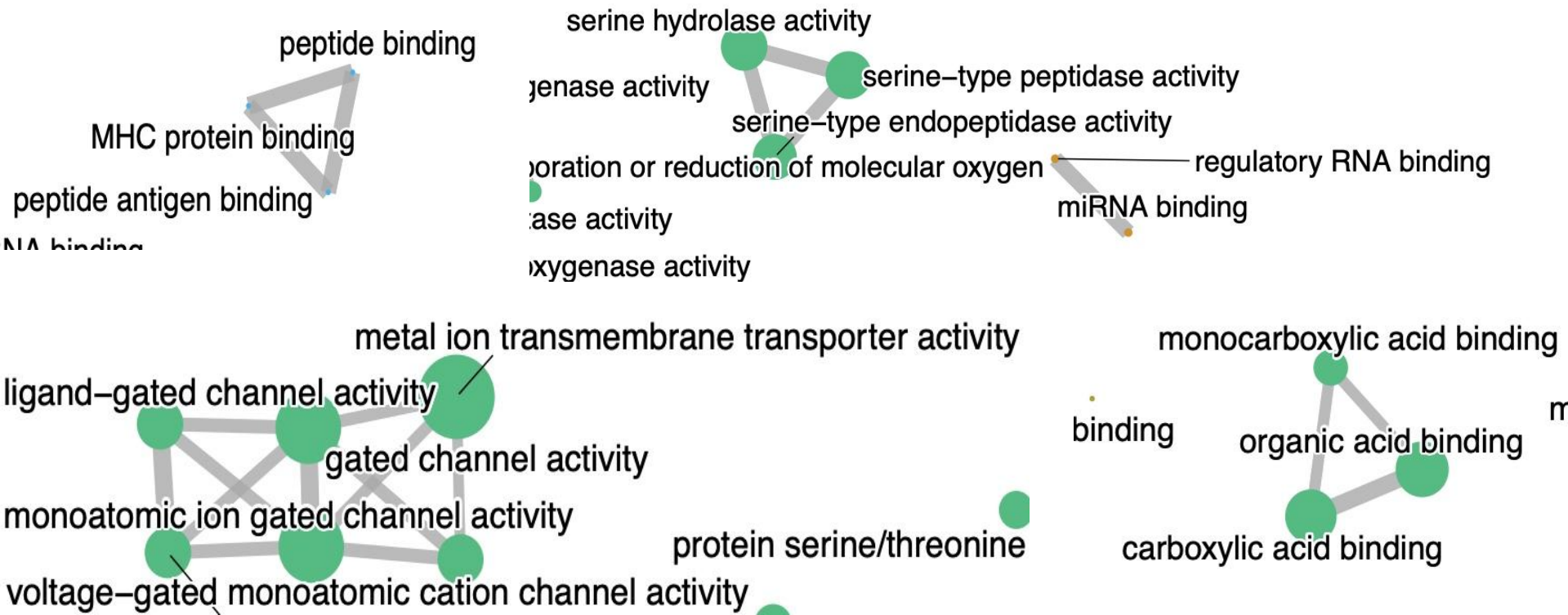
`enrichplot::emapplot()`

`cowplot::plot_grid()`





Clusters close-up





Conclusion

What difficulties have we faced:

- finding dataset
- preprocessing took longer than we thought (no LIMMA, but DESeq2)
- accuracy of our data was 70% and overfitted
- confounder: stages, lost of follow up
- Install outdated R packages

What have we learned?

- CNV analysis, DESeq2 analysis
- GO enrichment analysis

What could be improved?

- Find confounders and insights of metadata before starting
- Always checking for overfitting