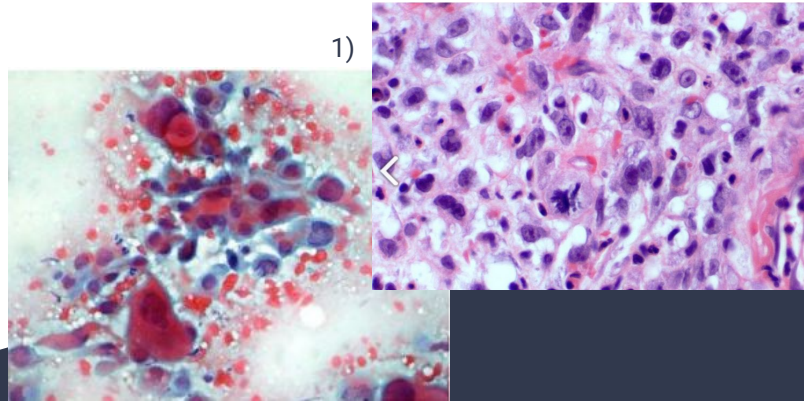


Transcriptomic and epigenomic biomarkers to differentiate between smokers with lung or laryngeal cancer

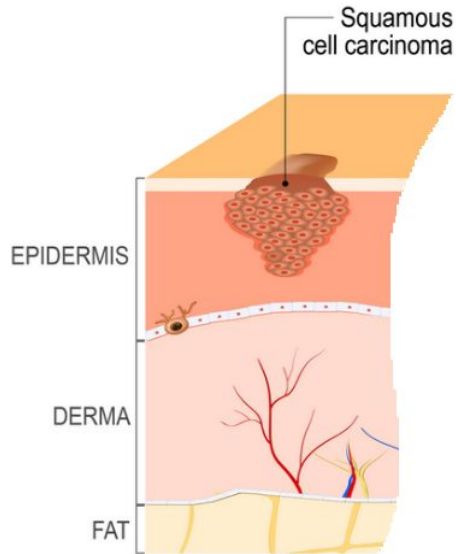
Final presentation



Group 5
Mia Anscheit
Matanat Mammadli
Friederike Wohlfarth

1) Squamous Cell Carcinoma of the lung (left) and of the larynx (right). Images from were obtained from <https://www.pathologyoutlines.com/topic/larynxcarcinomageneral.html>, <https://www.cellnetpathology.com/>, last accessed 2024/07/07.

Project introduction

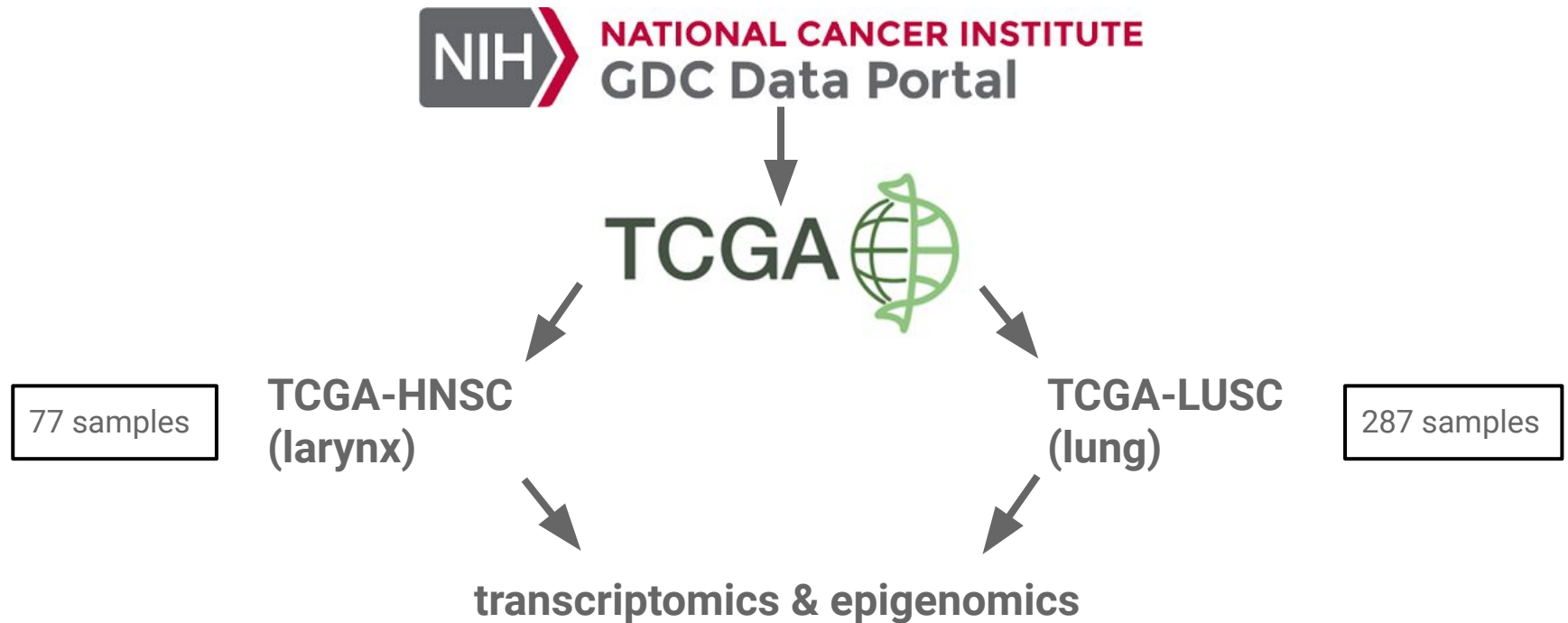


→ in respiratory organs mostly caused by smoking

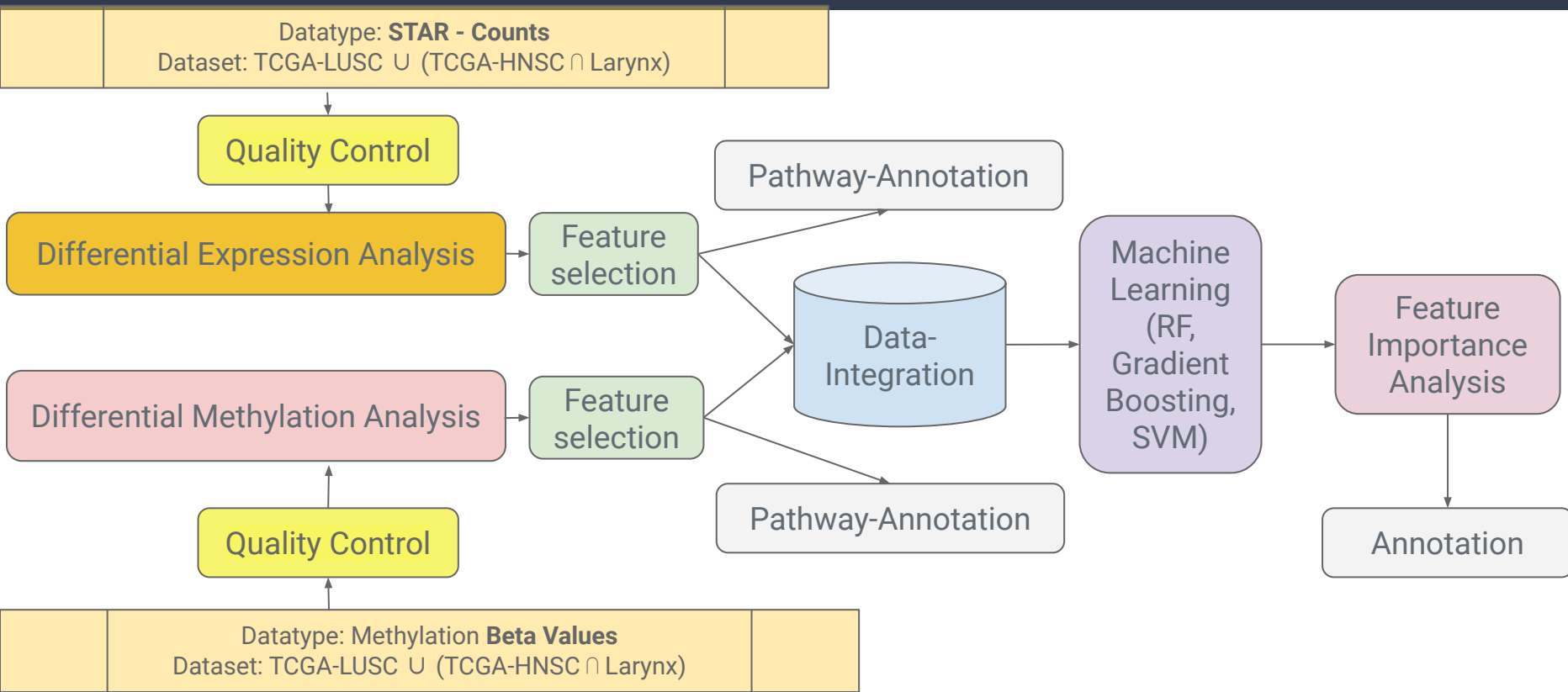


Which transcriptomic and epigenomic biomarkers are useful for the differentiation between squamous cell carcinoma in the lung and in the larynx?

Data



Overview



Quality control / Preprocessing Expression

Raw Counts in 60,660 genes



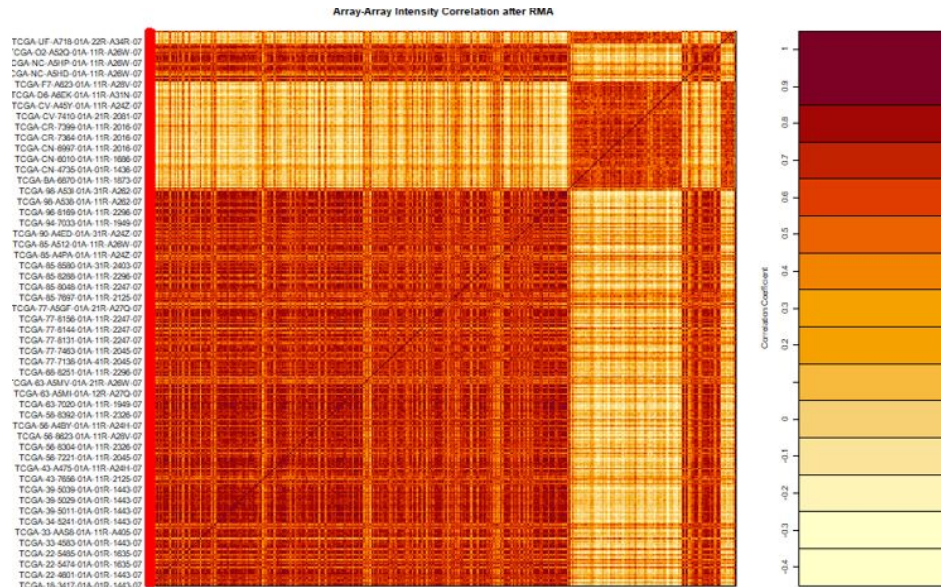
Preprocessing,
Normalization for
GC Content

60,660 genes



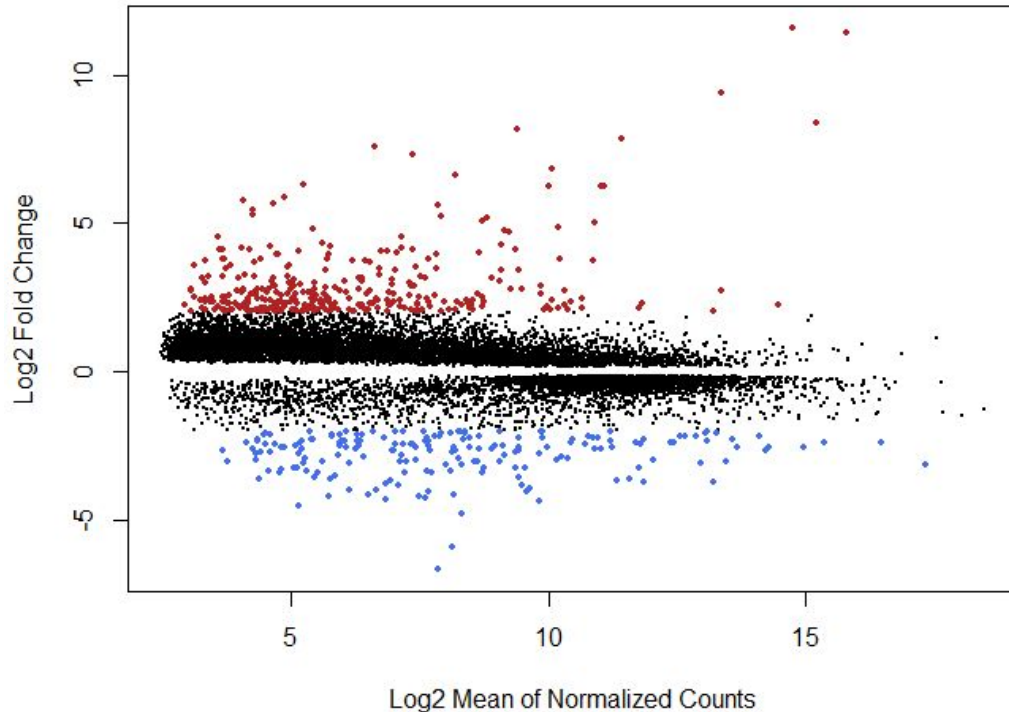
0.25 Quantile
Filtering

45,266 genes



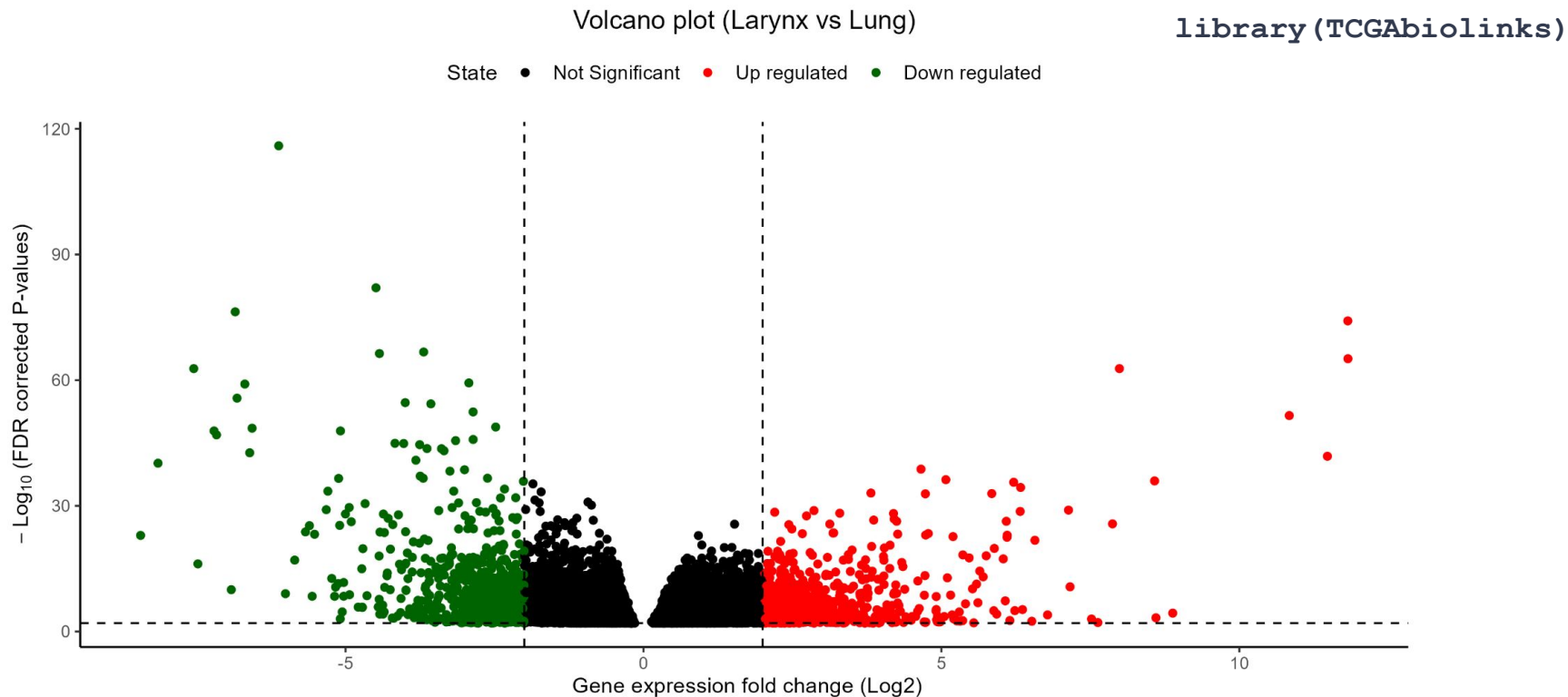
Differential gene expression analysis (DESeq2)

DESeq2 - Volcano Plot

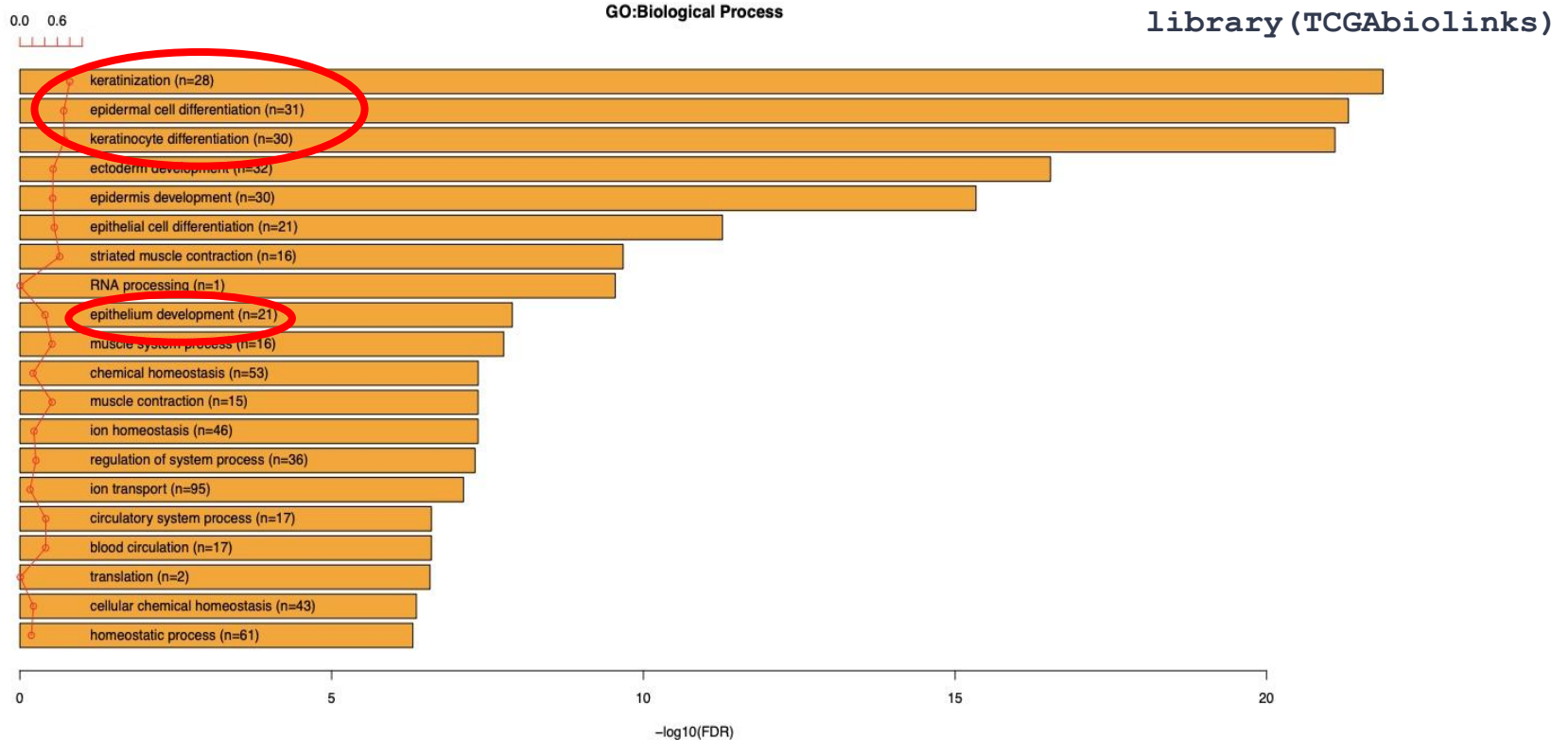


512 genes
with adjusted $p < 0.01$ &
 $\text{abs}(\log_2\text{FoldChange}) > 2$

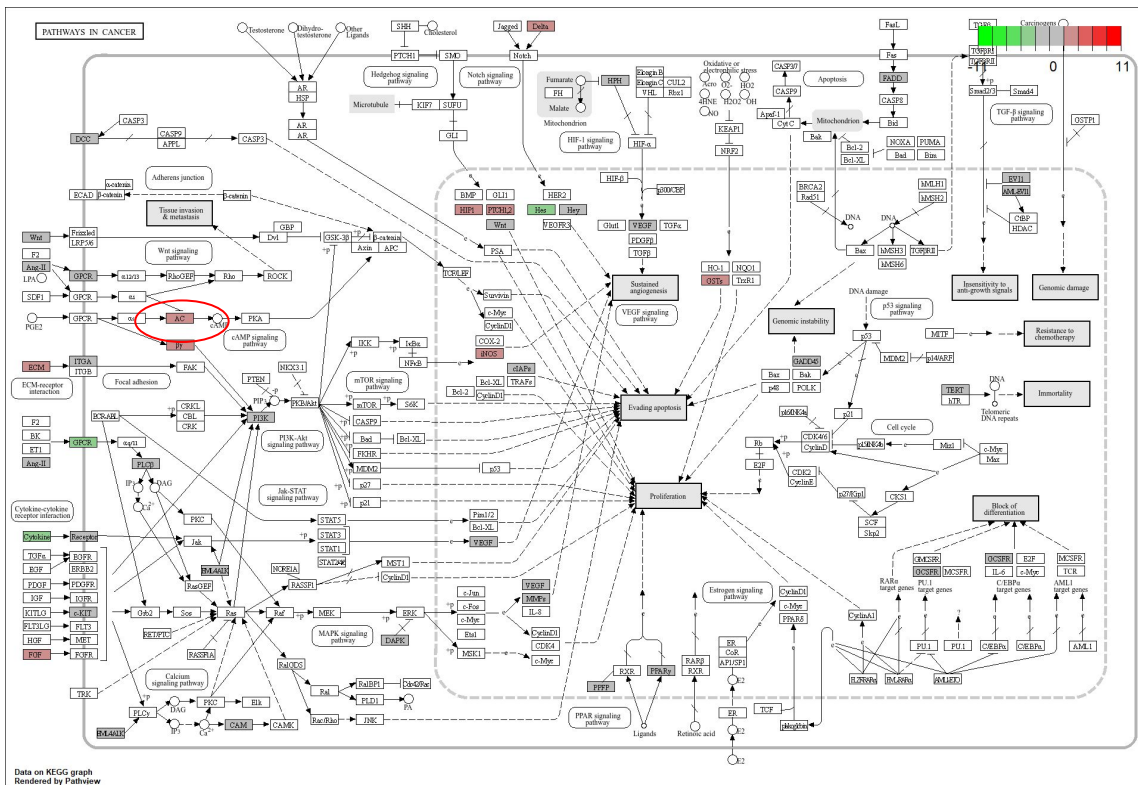
Differential gene expression analysis



Annotation

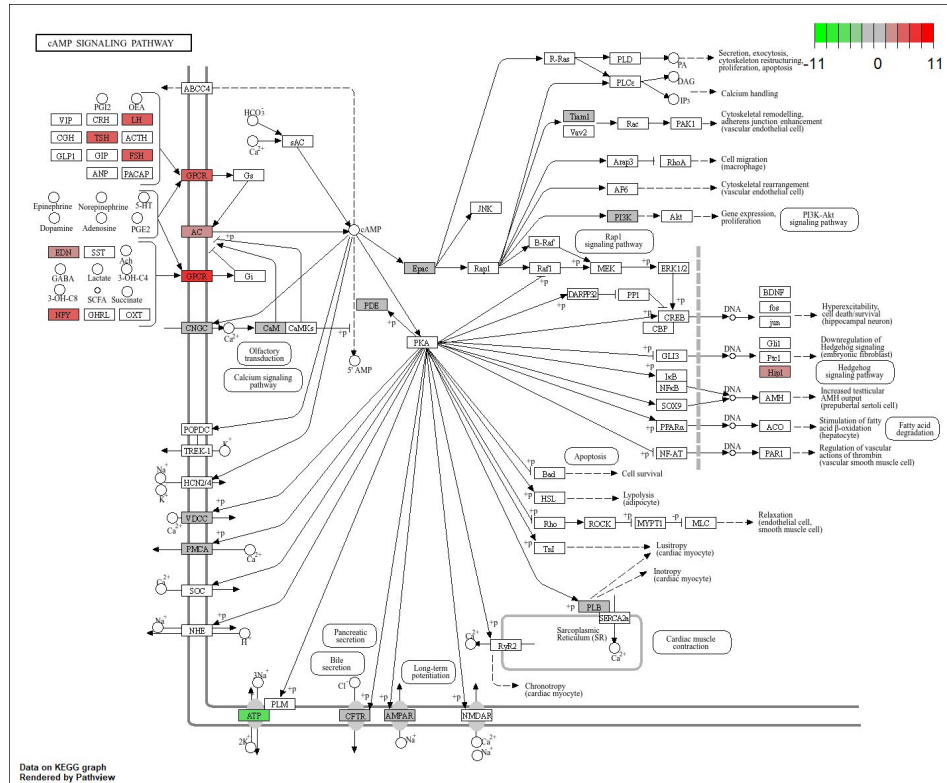


Annotation



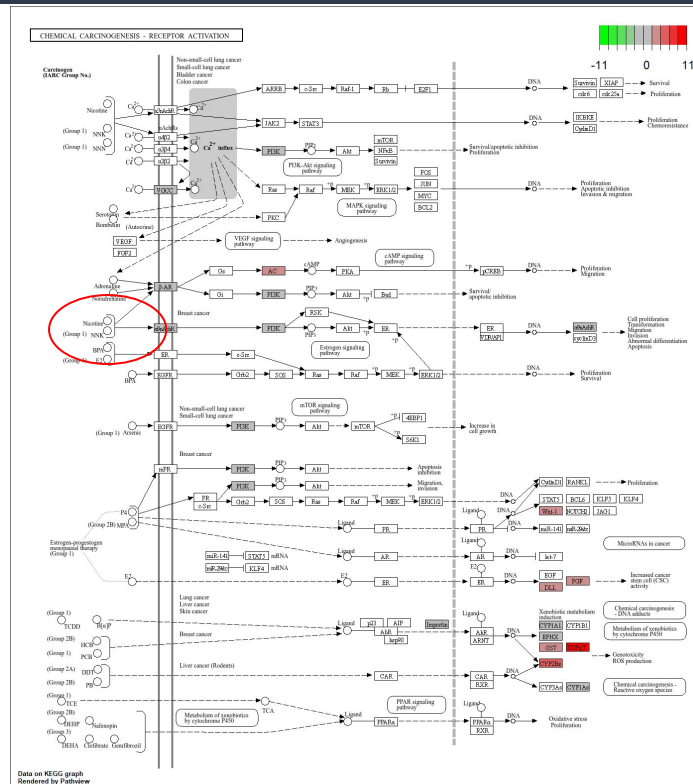
```
library(TCGAbiolinks)
```

Annotation



```
library(pathview)
```

Annotation



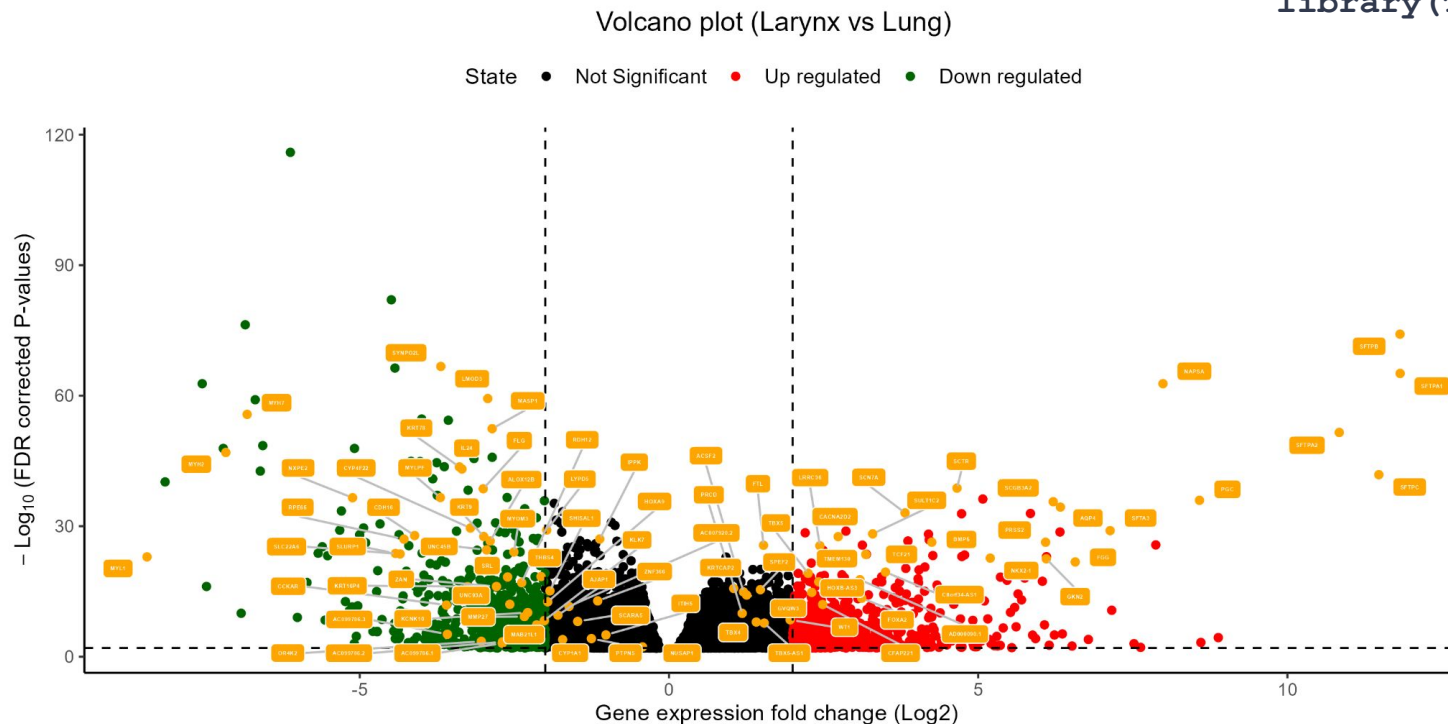
```
library(pathview)
```

DEG Analysis – Feature selection

| | | | | | library (Caret) |
|-------------|----------|---------|---------|----------|-----------------|
| AC007920.2 | FGG | LMOD3 | PRCD | SPEF2 | |
| AC099786.1 | FLG | LRRC36 | PRSS2 | SRL | |
| AC099786.2 | FOXA2 | LYPD5 | PTPN5 | SULT1C2 | |
| AC099786.3 | FTL | MAB21L1 | RDH12 | SYNPO2L | |
| ACSF2 | GKN2 | MASP1 | RPE65 | TBX4 | |
| AD000090.1 | GVQW3 | MMP27 | SCARA5 | TBX5 | |
| AJAP1 | HOXA9 | MYH2 | SCGB3A2 | TBX5-AS1 | |
| ALOX12B | HOXB-AS3 | MYH7 | SCN7A | TCF21 | |
| AQP4 | IL24 | MYL1 | SCTR | THBS4 | |
| BMP5 | IPPK | MYLPF | SFTA3 | TMEM130 | |
| C8orf34-AS1 | ITIH5 | MYOM3 | SFTPA1 | UNC45B | |
| CACNA2D2 | KCNK10 | NAPSA | SFTPA2 | UNC93A | |
| CCKAR | KLK7 | NKX2-1 | SFTPAB | WT1 | |
| CDH16 | KRT16P4 | NUSAP1 | SFTPC | ZAN | |
| CFAP221 | KRT78 | NXPE2 | SHISAL1 | ZNF366 | |
| CYP1A1 | KRT9 | OR4K2 | SLC22A6 | | |
| CYP4F22 | KRTCAP2 | PGC | SLURP1 | | |

DEA – Feature selection

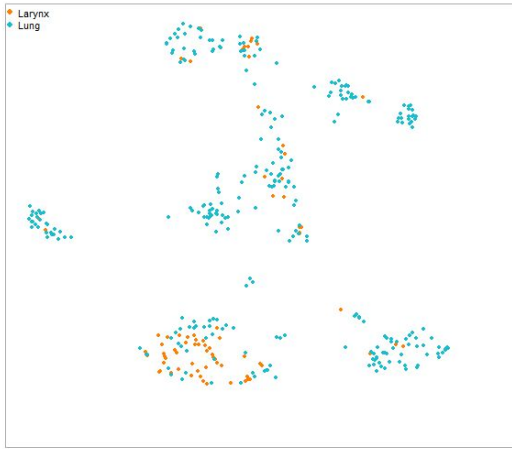
library(TCGAbiolinks)



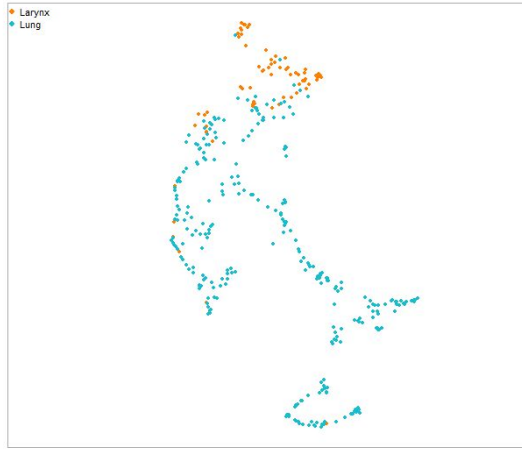
DEA – Feature selection

- Before feature selection: 60,660 genes
- After feature selection: 83 genes

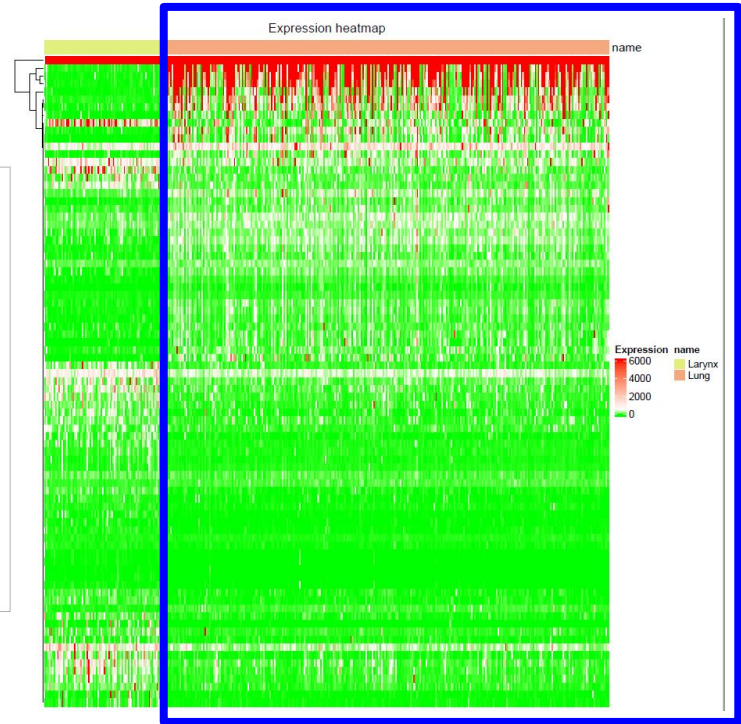
A UMAP visualization of the expression data: Larynx vs Lung before feature selection



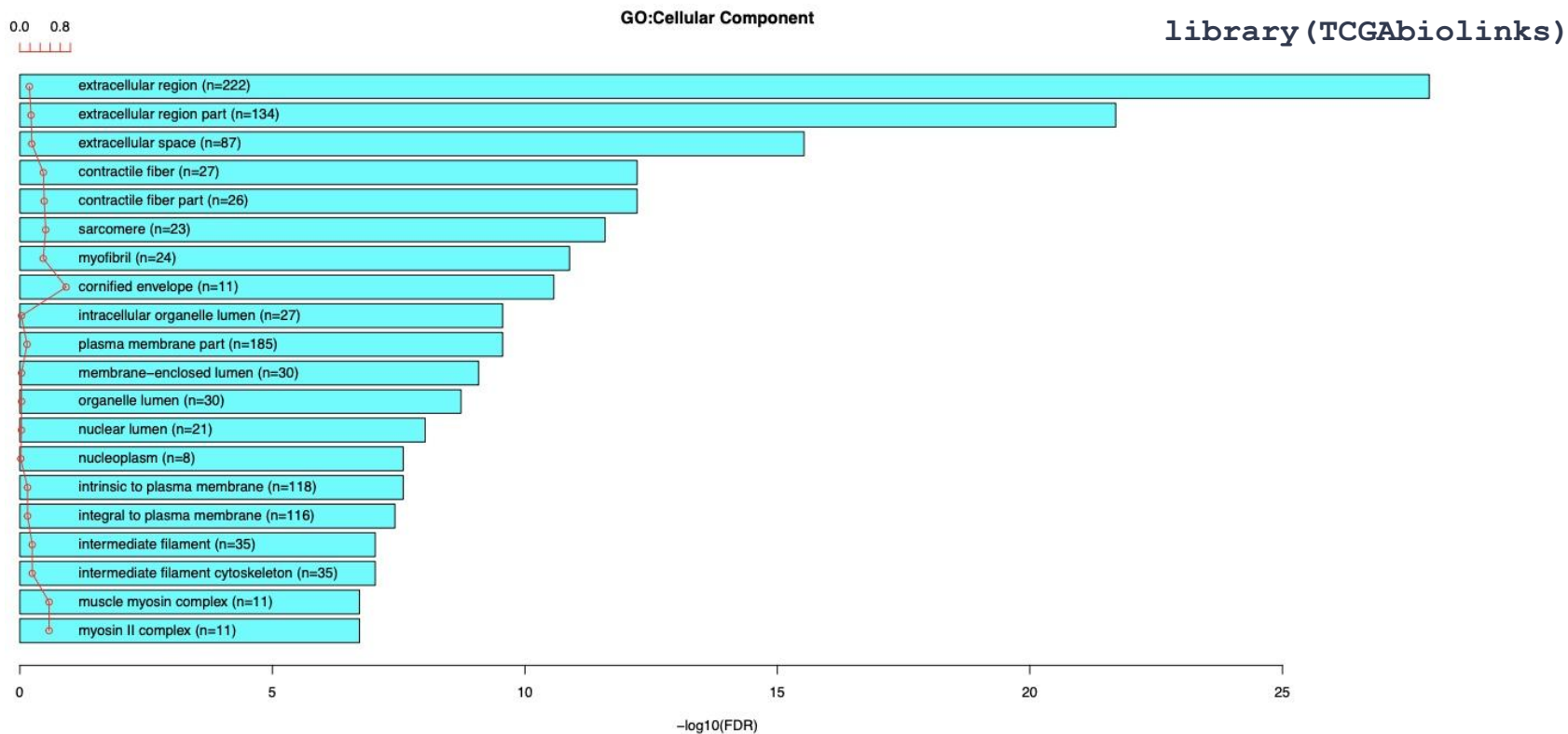
A UMAP visualization of the expression data: Larynx vs Lung after feature selection



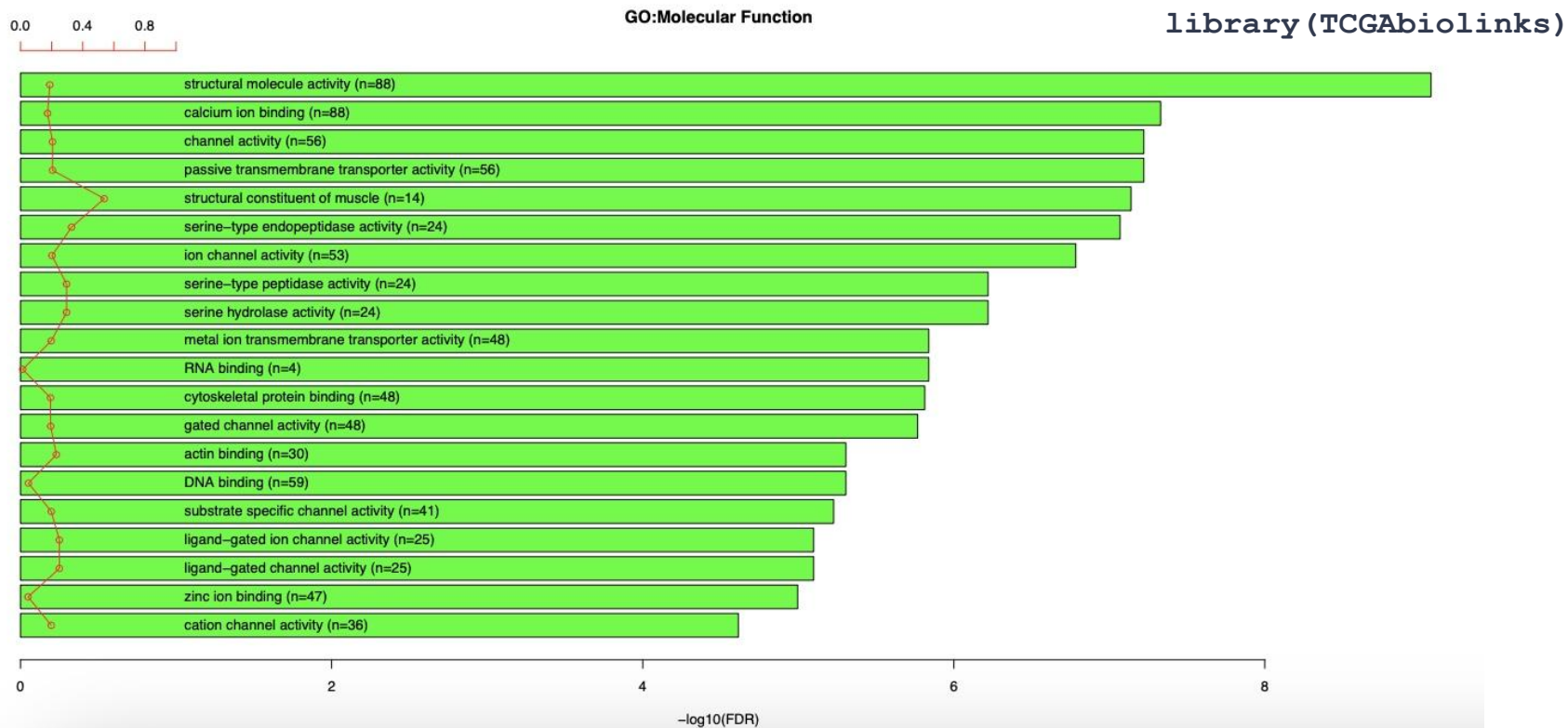
```
library(TCGAbiolinks)
library(umap)
```



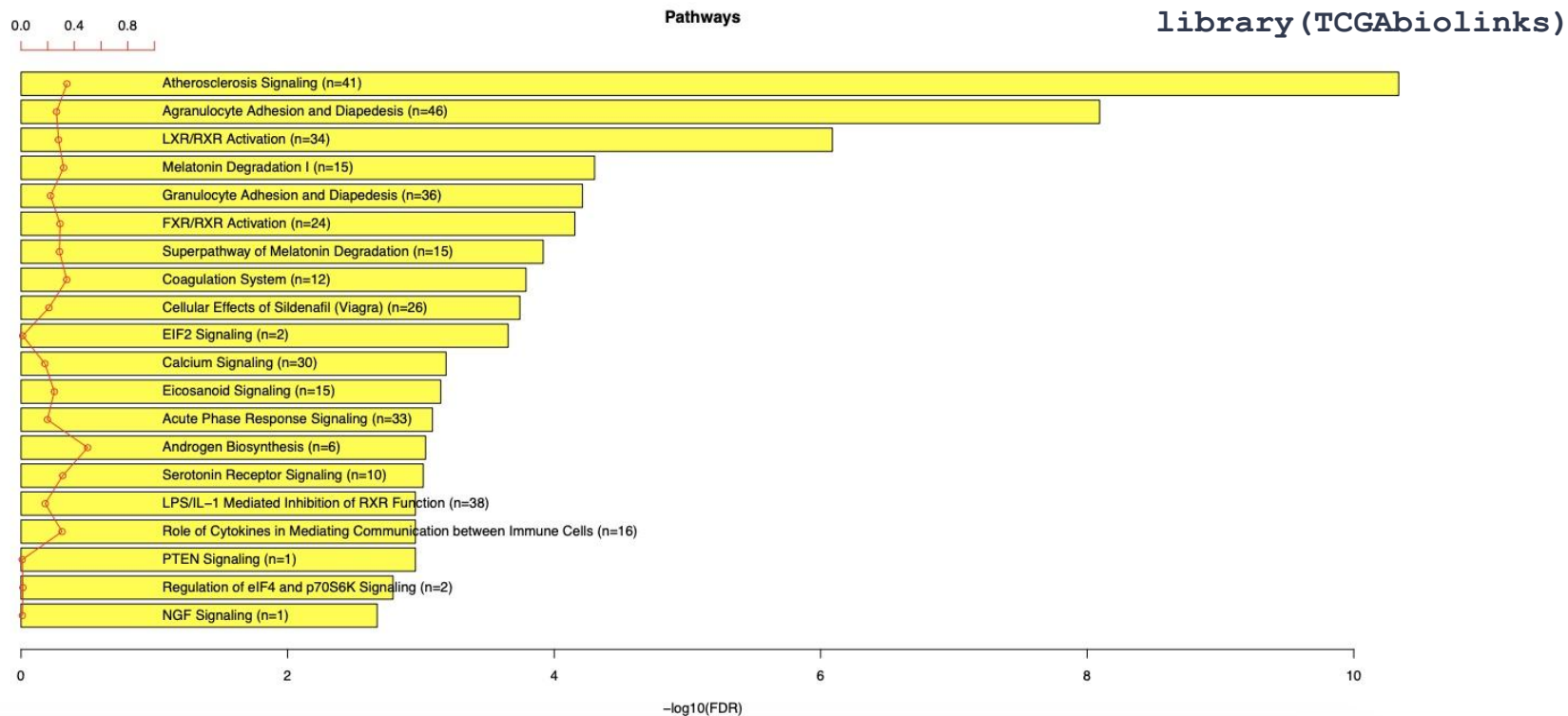
Annotation



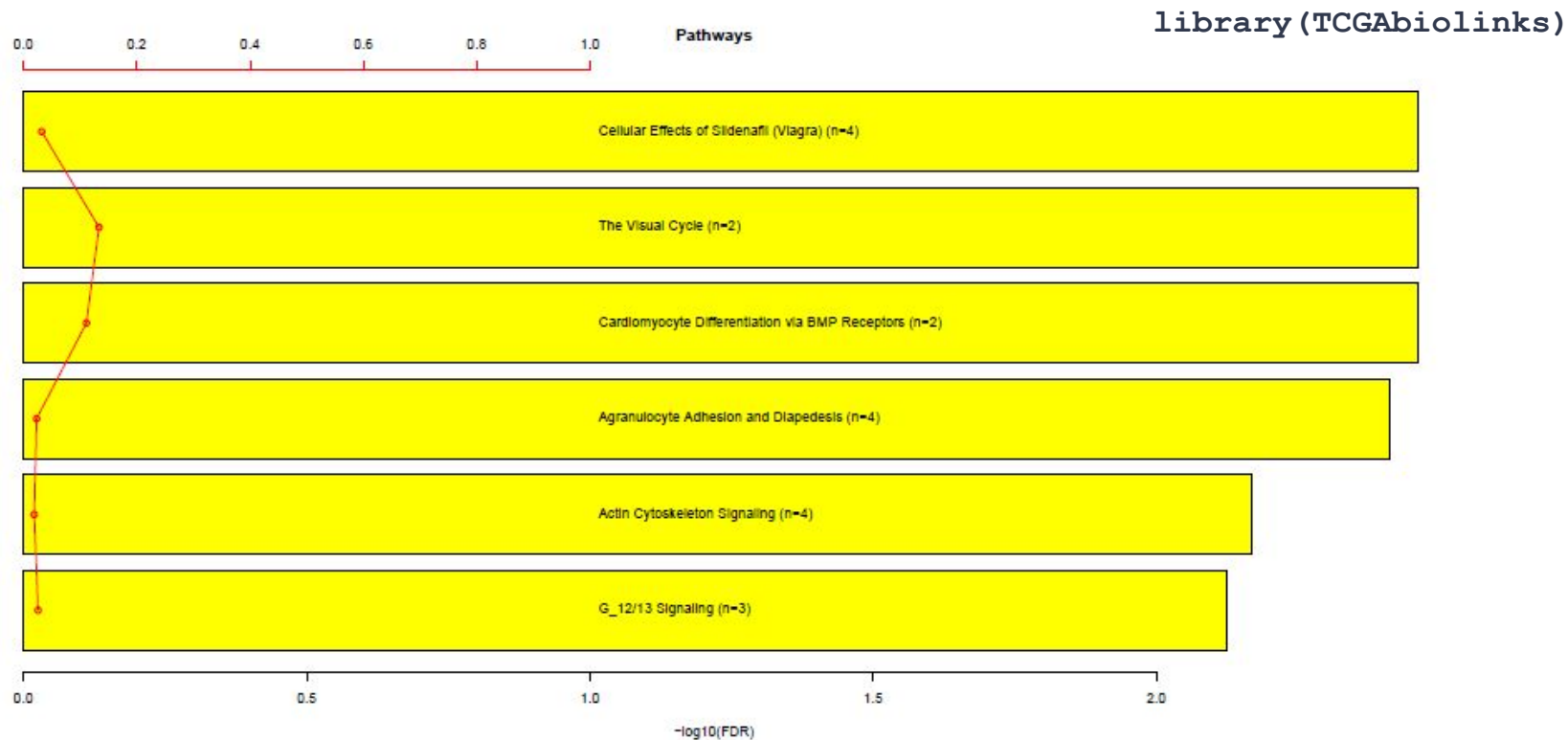
Annotation



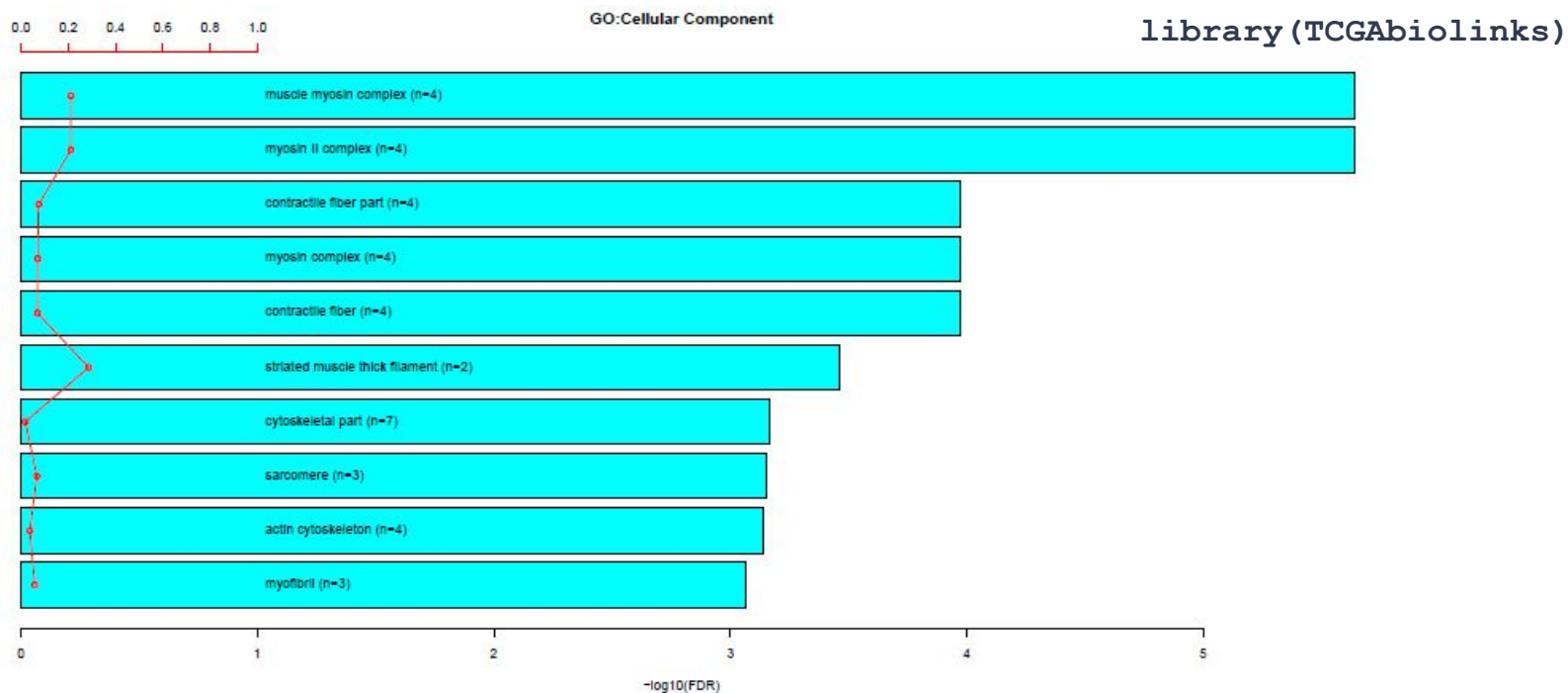
Annotation



DEG Analysis – Feature selection



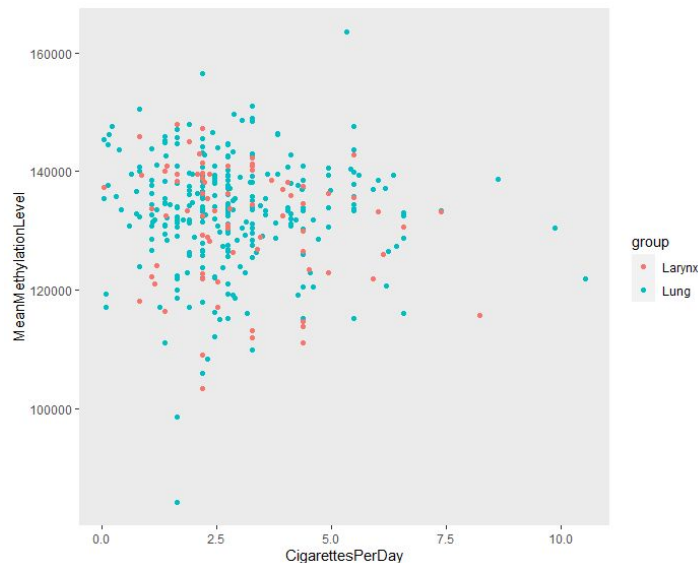
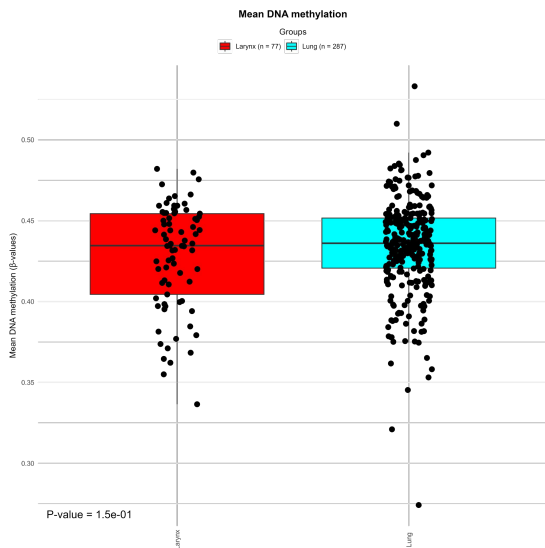
DEG Analysis – Feature selection



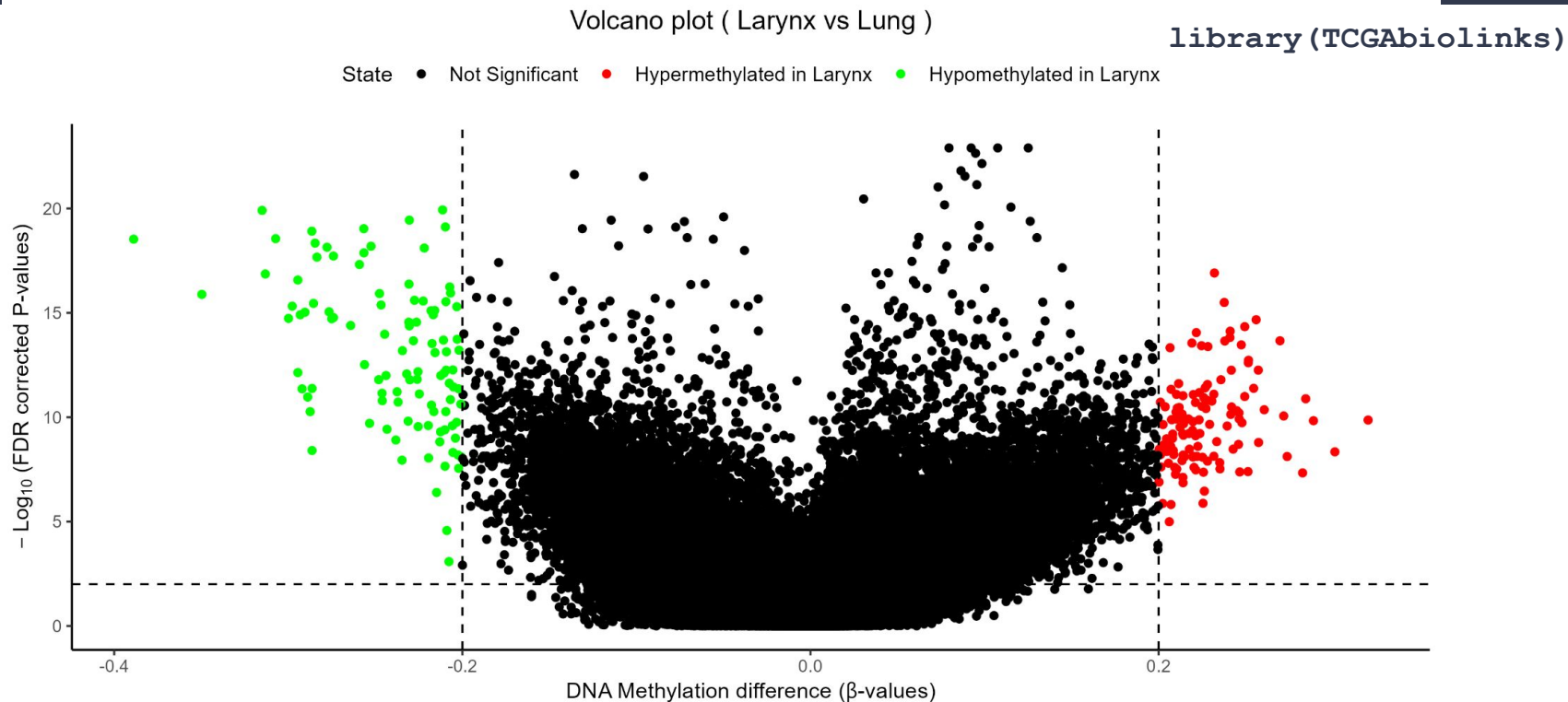
Quality control / Preprocessing Methylation

- 485,577 cpg islands in 35,557 genes, beta value > 0 \rightsquigarrow 312,864 cpg islands
- excluded X, Y and unknown chromosomes \rightsquigarrow 306,903 cpg islands

library(TCGAbiolinks)



Differential methylation analysis



DM Analysis – Feature selection

[Cells](#). 2020 Jul; 9(7): 1613.

Published online 2020 Jul 3. doi: [10.3390/cells9071613](https://doi.org/10.3390/cells9071613)

Methylation in *HOX* Clusters and Its Applications in Cancer Therapy

[Ana Paço](#),¹ [Simone Aparecida de Bessa Garcia](#),² and [Renata Freitas](#)^{2,3,*}

CLASP1
COLEC11
CTB-49A3.4
CTB-57H20.1
CTD-2555C10.3
CYR61
DHCR7
DNHD1
EPDR1
ESR2
FLOT1
FOXC1
FOXQ1

HOXA-AS2
HOXA-AS3
HOXA2
HOXA3
HOXA4
HOXA5
HOXA7
HOXB-AS3
HOXB3
HOXB4
HOXB5
HOXB6
IER3
INPP5D

MIR6080
MRO
NEURL1B
NRG2
PACS1

PCDHA1
PCDHA10
PCDHA11
PCDHA12
PCDHA13
PCDHA2
PCDHA3
PCDHA4

PCDHA5
PCDHA6
PCDHA7
PCDHA8
PCDHA9
PCDHAC1

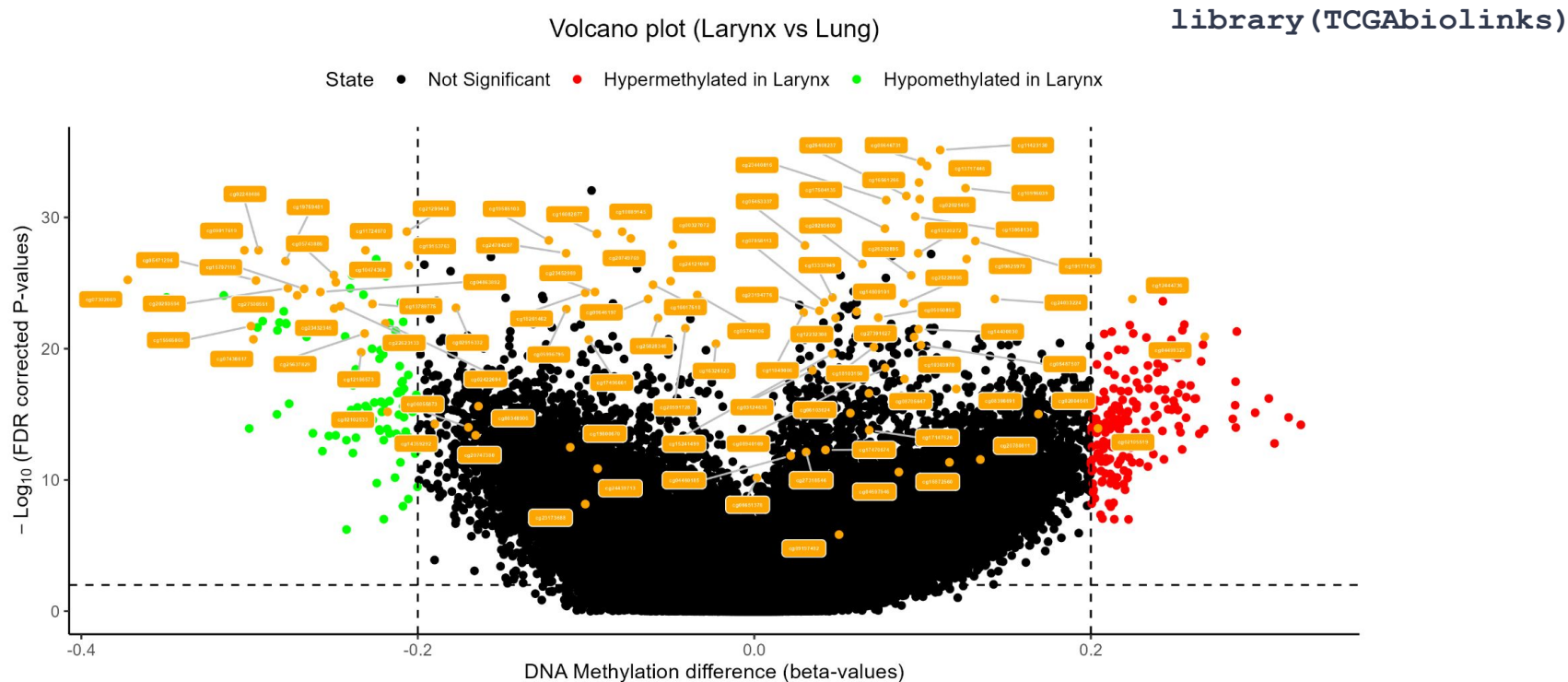
PLEKHM1P
PPP1R9A
RANBP17
RBM20
RP1-170019.22
RP1-170019.24
RP11-399K21.11
RP11-660I16.2

SCARF2
SFRP4
SKAP1
SMAD6
SMYD4
ST14
SYNE2
TARID
TIMM10B
TMC03
TMEM201
ZFPM1
ZNF467

Clustered protocadherins methylation alterations in cancer

[Ana Florencia Vega-Benedetti](#),^{#1} [Eleonora Loi](#),^{#1} [Loredana Moi](#),¹ [Sylvain Blois](#),¹ [Antonio Fadda](#),¹ [Antonella Arcella](#),³ [Manuela Badiali](#),⁴ [Felice Giangaspero](#),^{2,3} [Isabella Morra](#),⁵ [Amedeo Columbano](#),⁶ [Luigi Zorcolo](#),⁷ [Viviana Gismondi](#),⁸ [Liliana Varesco](#),⁸ [Sara Erika Bellomo](#),⁹ [Silvia Giordano](#),^{9,10} [Matt](#)

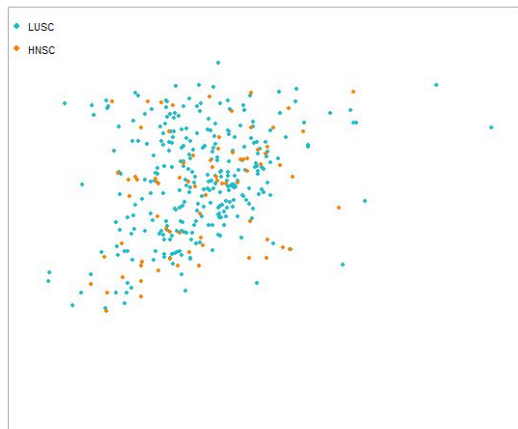
DM Analysis – Feature selection



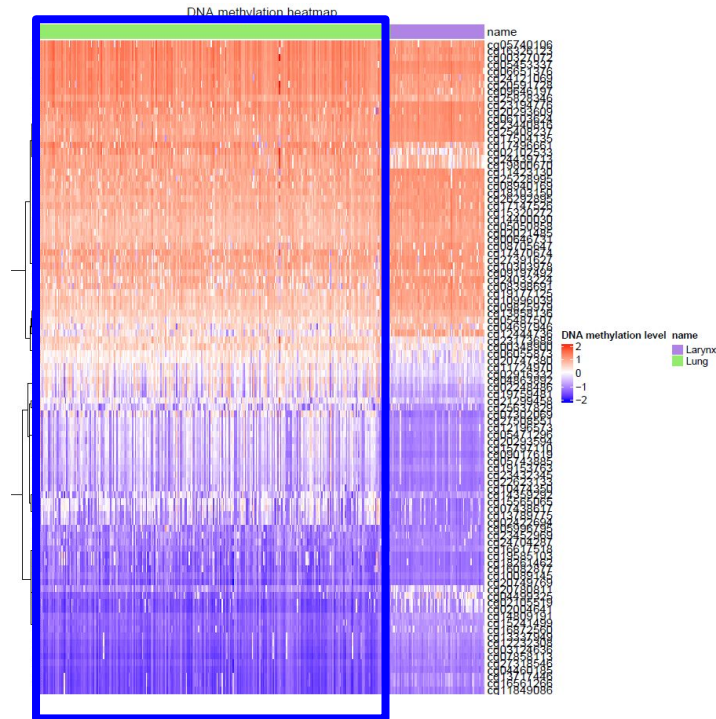
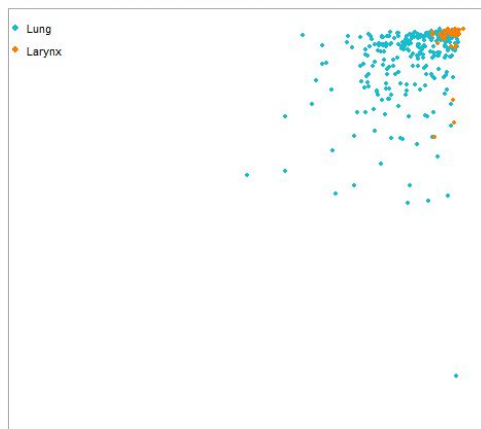
DM Analysis – Feature selection

- Before feature selection: 306,903 cpG islands
- After feature selection: 97 cpG islands

A UMAP visualization of the methylation data: Larynx vs. Lung before feature selection



A UMAP visualization of the methylation data: Larynx vs. Lung after feature selection

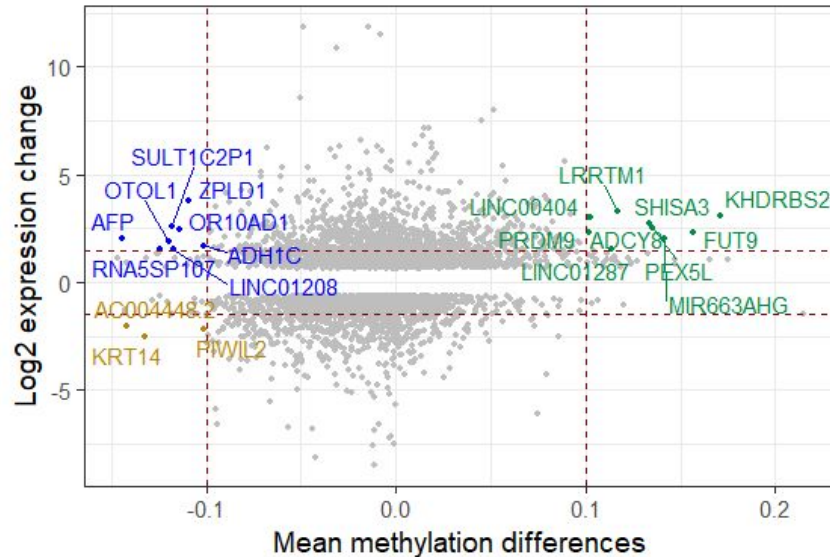


```
library(TCGAbiolinks)
library(umap)
```


Data integration Connecting Expression and Methylation

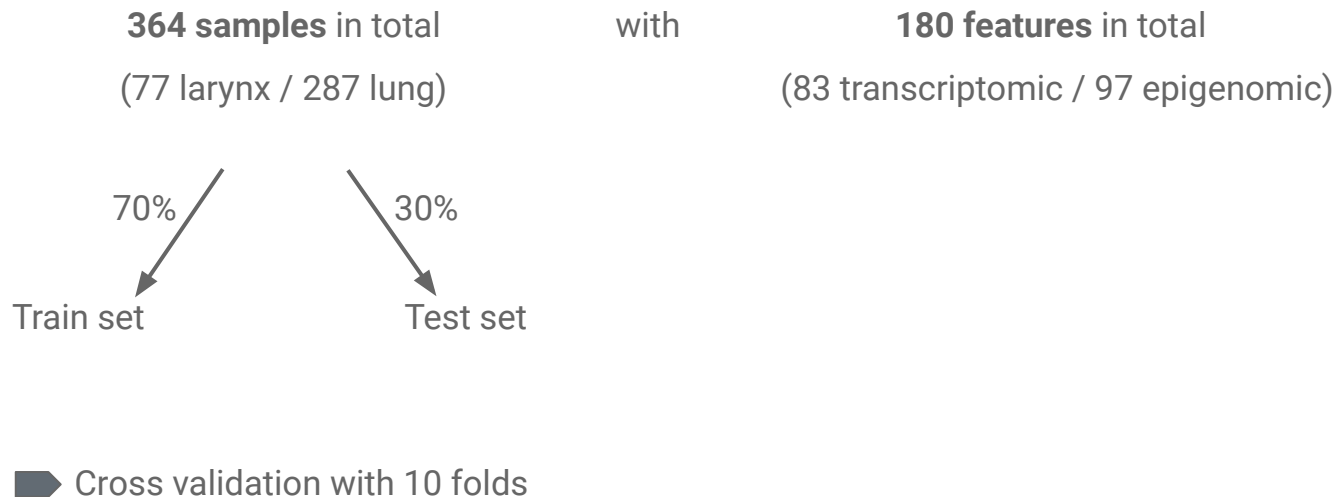
- aggregated mean beta values for methylation and counts for expression for each gene **library (TCGAbiolinks)**

Correlation between DNA methylation and Gene expression



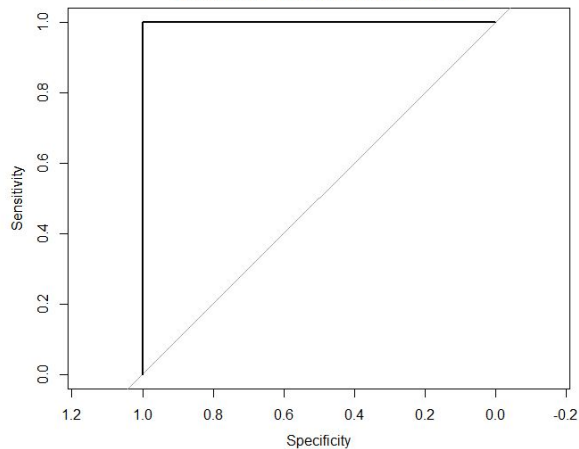
KR
ZPLD1
ADCY8
PIWIL2
KHDRBS2
SHISA3
PEX5L
FUT9
SULT1C2P1
OR10AD1
MIR663AHG
PRDM9
AFP
ADH1C
OTOL1

Machine Learning (with caret)

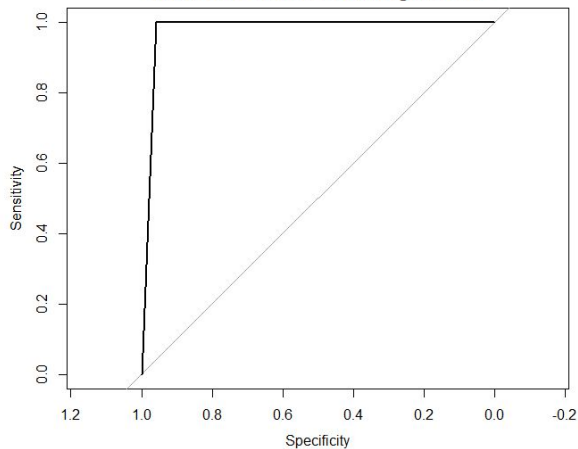


Machine Learning (with caret)

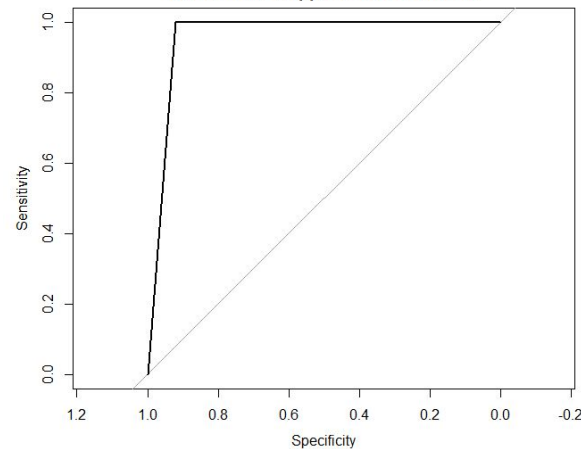
ROC Curve - Random Forest Model



ROC Curve - Gradient Boosting Machine



ROC Curve - Support Vector Machine



```
Reference
Prediction larynx lung
larynx      23    0
lung         0   86

Accuracy : 1
95% CI : (0.9667, 1)
No Information Rate : 0.789
P-value [Acc > NIR] : 6.037e-12
```

```
Reference
Prediction larynx lung
larynx      23    1
lung         0   85

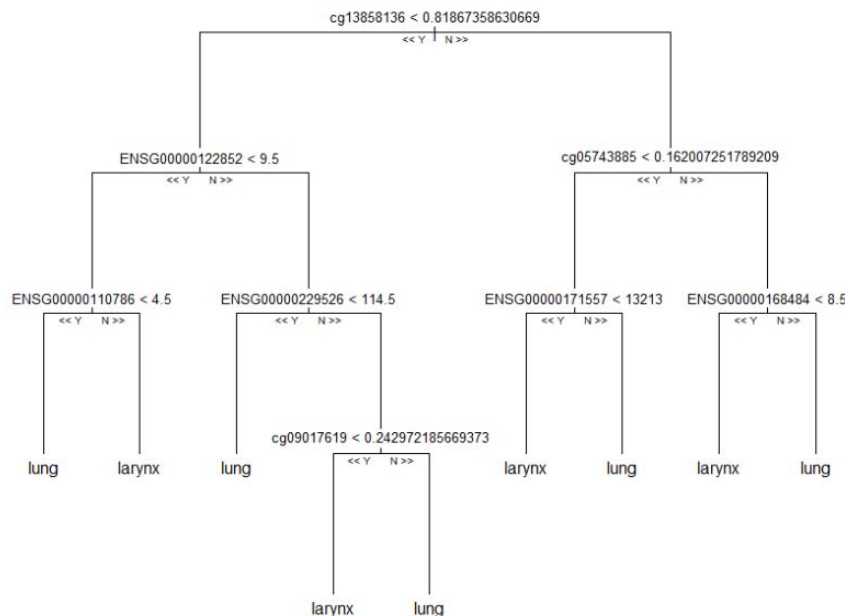
Accuracy : 0.9908
95% CI : (0.9499, 0.9998)
No Information Rate : 0.789
P-value [Acc > NIR] : 1.82e-10
```

```
Reference
Prediction larynx lung
larynx      23    2
lung         0   84

Accuracy : 0.9817
95% CI : (0.9353, 0.9978)
No Information Rate : 0.789
P-value [Acc > NIR] : 2.724e-09
```

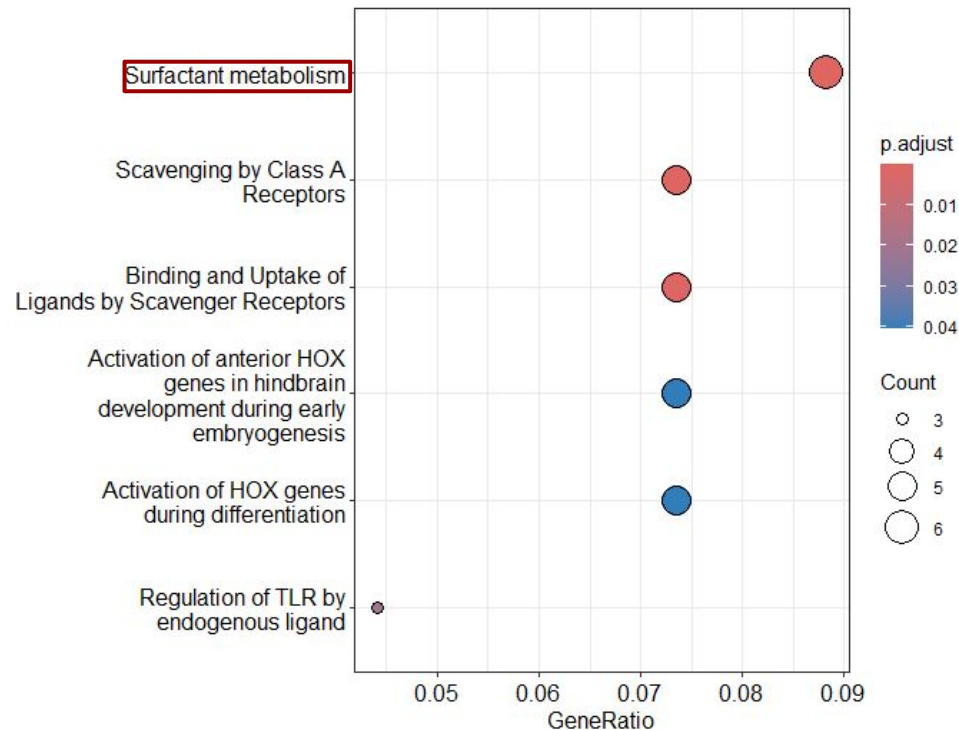
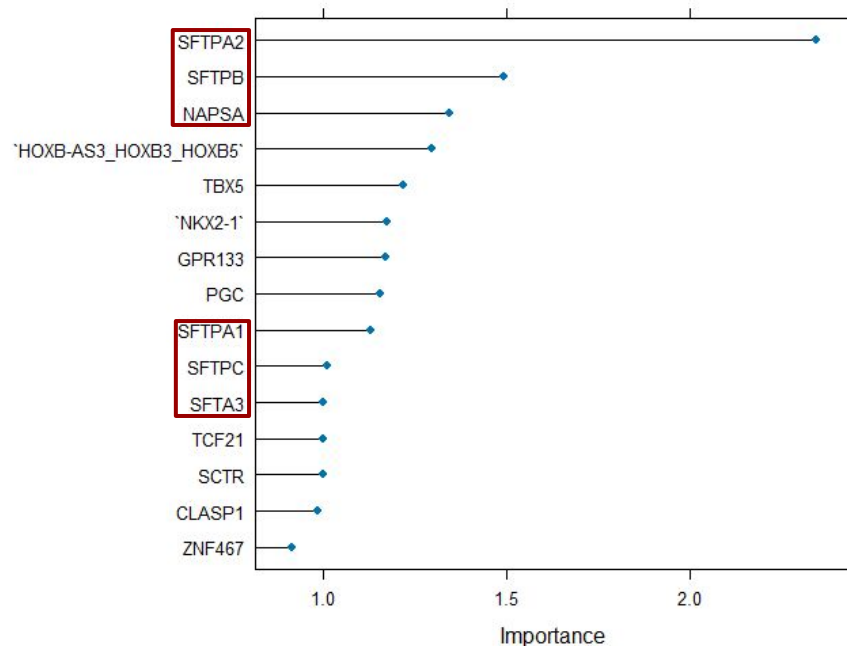
Example of the Random Forest decision trees

library(reprtree)



Feature importance (with caret)

Most important features



Summary & Conclusion

- we trained three models to predict the site of squamous cell carcinoma
- differences in expression and methylation level may be characterized by tissue of origin
- the application could be diagnostic of patients with Cancer of unknown primary (CUP)
- Here, the cancer spreads within lymph nodes in the and the primary cannot be found, the ML can be used to predict where the tumor originates from
- the model could be applied to predict the primary site of the tumor

References

- **Maximilian Leitheiser et al.** “Machine learning models predict the primary sites of head and neck squamous cell carcinoma metastases based on DNA methylation”. en. In: The Journal of Pathology 256.4 (2022). Eprint: <https://onlinelibrary.wiley.com/doi/pdf/10.1002/path.5845>, pp. 378–387. issn:1096-9896. (visited on 07/06/2024).
- **Philipp Jurmeister et al.** “Machine learning analysis of DNA methylation profiles distinguishes primary lung squamous cell carcinomas from head and neck metastases”. In: Science Translational Medicine 11.509 (Sept. 2019). Publisher: American Association for the Advancement of Science, eaaw8513. url: <https://www.science.org/doi/10.1126/scitranslmed.aaw8513> (visited on 07/06/2024).
- **Anil Vachani et al.** “A 10-Gene Classifier for Distinguishing Head and Neck Squamous Cell Carcinoma and Lung Squamous Cell Carcinoma”. In: Clinical Cancer Research 13.10 (May 2007), pp. 2905–2915. issn: 1078-0432. url: <https://doi.org/10.1158/1078-0432.CCR-06-1670> (visited on 07/06/2024).
- **Yawei Li et al.** “Machine Learning for Lung Cancer Diagnosis, Treatment, and Prognosis”. In: Genomics Proteomics Bioinformatics 20.5 (2022). issn: 850-866. doi: 10.1016/j.gpb.2022.11.003.
- **Yin Li et al.** “Transcriptomic and functional network features of lung squamous cell carcinoma through integrative analysis of GEO and TCGA data”. In: Sci Rep 8.1 (2018). doi: 10.1038/s41598-018-34160-w.
- **Alana Sorgini et al.** “Analysis of the TCGA Dataset Reveals that Subsites of Laryngeal Squamous Cell Carcinoma are Molecularly Distinct”. In: Cancers 13.1 (Dec. 31, 2020), p. 105. issn: 2072-6694. doi: 10.3390/cancers13010105.
- **Joshua D. Campbell et al.** “Genomic, Pathway Network, and Immunologic Features Distinguishing Squamous Carcinomas”. In: Cell Reports 23.1 (Apr. 3, 2018), 194–212.e6. issn: 2211-1247. Doi: 10.1016/j.celrep.2018.03.063.
- **Katherine A. Hoadley et al.** “Cell-of-Origin Patterns Dominate the Molecular Classification of 10,000 Tumors from 33 Types of Cancer”. In: Cell 173.2 (Apr. 5, 2018), 291–304.e6. issn: 1097-4172. doi: 10.1016/j.cell.2018.03.022