

Materials and Methods

For the means of this project two separate analysis are performed: a pan-cancer analysis focusing on differences between cancer types and a focused analysis investigating THCA.

Description of the underlying data

For the analysis four data sets were provided. For pan-cancer analysis a gene expression data frame with normalized and log2 transformed bulk RNA-seq expression data for 60,489 genes in 9741 patients with 33 different forms of cancer was used. The data was derived from The Cancer Genome Atlas (TCGA). Complementing the TCGA expression data is an annotation dataframe with 37 clinical annotations regarding tumor type, tumor stage, gender, age, etc. for all patients.

The third object is a list containing five lists for the focused analysis, one list for each tumor type (BRCA, KIRC, LUAD, PRAD, THCA). For our focused analysis, the only the THCA data were used. The THCA list consists of three data frames: The first two contain normalized and log2 transformed bulk RNA-seq expression data for 19,624 genes in 59 THCA patients for carcinogenic and homeostatic tissue. The third dataframe complements the data with the respective clinical annotations.

The last object contains 46 pathways associated with the hallmarks of cancer in form of a list of string vectors.

Metabolic pathway selection

To perform enrichment analysis later on, 6366 canonical pathways were selected from the Molecular Signatures Database (MSigDB)

@msigdb

with the `msigdb::msigdb()` function. As not to introduce a bias during enrichment analysis, the similarity of MSigDB pathways among themselves as well as with the hallmark pathways was computed with the Jaccard index. Pathways with a Jaccard index greater than the 1σ range were discarded.

Preprocessing of expression data

Data cleaning

All expression data were checked for missing values with the `na.omit()` function. Subsequently, low variance filtering was performed for TCGA and THCA tumor expression data. The variances of expression were computed for every gene across all samples and then, genes with variances below a threshold were discarded to reduce dimensionality.

Biotype filtering

Next, biotype filtering was performed for pan-cancer and THCA expression data to reduce dimensionality further. Only genes sharing biotypes with the hallmark pathways were kept for the the following analysis. The biotypes of the genes were retrieved using the `biomart::getBM()` function from the `biomaRt` package [1]. To allow for an appropriate comparison within all pathways, only MSigDB pathways where over 99% of their respective genes were present in the filtered expression data were selected as final pathways.

Methods for descriptive analysis

Mean-variance plot

In a mean-variance plot the variance is plotted over the mean of expression values of single genes across all patients. Thus, the variance and mean were calculated for each gene in the THCA expression data. The final plot was created with the package `ggplot2` ??.

KANN RAUS m.M.n Violin plot

To check the distribution of a data set and compare it with other data sets violin plots are used. Based on how similar the violin plots are, it can be implied that the data is normalized. Violin plots are tilted and mirrored density plots of gene expression values. The y-axis shows the gene expression value and the x-axis shows the amount of genes with a certain gene expression value.

Jaccard-Index

The Jaccard-Index is a method to describe the similarity between two quantities. To compute it, the intersection of all gene ENSEMBL-IDs from two compared pathways was divided by their union. We used this method to determine the degree in which pathways are similar to each other.

Volcano plot

A volcano plot is used to identify genes displaying significantly different expression in carcinogenic versus homeostatic tissues. First, the log2 fold change (Log2FC) is calculated for each gene across all samples in the THCA expression data in the following way:

$$\log_2FC = \text{mean}(\text{normaltissue}) - \text{mean}(\text{tumortissue})$$

Next, a two-sided t-test was performed with the `t.test()` function to determine the significance of a difference in expression. To avoid the accumulation of type one errors, a Bonferroni correction was performed. `n` is the number of genes in the cleaned data set for focused analysis:

$$\alpha = \frac{0.025}{n}$$

In the volcano plot the $-\log_{10}$ of the calculated p-values is plotted against the Log2FC. Genes with a lower p-value than the corrected significance level α are significantly differently expressed. If the Log2FC is additionally positive, the genes are significantly overexpressed in tumor tissue, if the Log2FC is negative, the genes are significantly underexpressed in tumor tissue.

Dimension reduction and pathway enrichment analysis

The GSEA was used to identify enriched pathways in THCA tumor tissue. Here, GSEA was performed with the package “fgsea” [fgsea]. First, the expression values were ranked decreasingly by log2FC for every patient. Log2FC was chosen as the ranking metric as it is easy to compute and shows a high sensitivity. **xxx Quelle: Ranking metrics in gene set enrichment analysis: do they matter?** Secondly, using the ranked log2FC vectors, the enrichment score of each pathway was calculated for each patient with the `fgseamultilevel()` function.

As no normal tissue reference data was provided for the TCGA expression data, pathway activities were computed via GSVA. The analysis was performed with the `gsva()` function from the “GSVA” package

[@gsva]. To give a general overview over the differences in expression of THCA and homeostatic thyroid tissue GSVA, the THCA expression data were also analysed by GSVA. To do so, tumor and normal expression data were combined into a singular dataframe of which enrichment scores were computed with `gsva()`. Then, the GSVA data was split again and the log2FC between the two matrices was computed and taken as pathway activity.

PCA was performed to provide an uncorrelated dataset for the subsequent UMAP. For the TCGA GSVA pathway activity data the `prcomp()` function was used. To verify the results, PCA was performed on TCGA expression data, as well. In this case `Seurat::RunPCA()` from the Seurat package was used to minimized computation times. **xxx Quelle: seurat package**

UMAP analysis was used to identify and visualize clusters in TCGA GSVA and expression data. This was achieved with the `umap()` function from the package “umap”

`@umap`

running on all PCs from TCGA GSVA and expression data.

Regression analysis

Linear Regression

A linear regression analysis is performed to predict the activity of xxx based on the activity of the other pathways.

Firstly, the correlation of the pathways for predicting is checked, only pathways with a low correlation were kept. In the next step, the variance is checked, 80% of the genes with low variance were omitted.

For the regression analysis only 20% of the pathways were used, to only use significant pathways.

The regression analysis was tested by

Neuronal Network

A neural network was used to predict the activity of REACTOME_INTERLEUKIN_36_PATHWAY based on the activity of other pathways. Therefore, the network was trained with the pathway activity of 45 xxx patients from the THCA data for focused analysis. The other 15 patients were used to validate the network, obtaining a mean squared error (MSE) value, to evaluate the precision of the network.

For identification of the best initial conditions, 25 different networks are generated, each one with 2 hidden layers and different combinations of neurons per layer. For each combination the MSE is calculated and the 3 combinations with the lowest MSE are selected for selection of the best initial conditions regarding the weights and biases. For each of the 3 networks 100 random initial conditions are tested, resulting in one network with the lowest MSE.

Packages

```
## Warning: Paket 'readxl' wurde unter R Version 4.1.3 erstellt
```

Table 1: Packages used in the analysis.

Package	Localisation	Usage	Link
biomart	pre_02, pre_03, pre_05	renaming the genenames from the hallmarkpathways-dataframe into ensembleIDs	https://bioconductor.org/packages/release/bioc/html/biomart.html
msigdb	pre_03	downloading all of the canonical pathways and the genes which they include in homo sapiens from the msigbdr data base	https://bioconductor.org/packages/release/bioc/html/msigdb.html
dplyr	pre_04, pre_05	tidying and manipulating of dataframes	https://cran.r-project.org/web/packages/dplyr/index.html
ggplot2	pre_04, pre_05, descr_03, descr_04, THCA_01, THCA_02, pan_02, pan_04	allows for the creation of plots with more detailed options	https://cran.r-project.org/web/packages/ggplot2/index.html
pheatmap	descr_01, pan_01, neu_02, neu_04	allows for the creation of heatmaps with more detailed options	https://cran.r-project.org/web/packages/pheatmap/pheatmap.html
vioplot	descr_02	creation of violinplots	https://cran.r-project.org/web/packages/vioplot/index.html
VennDiagram	descr_05	creation of VENN-diagrams	https://cran.r-project.org/web/packages/VennDiagram/VennDiagram.html
dplyr	THCA_01, pan_01		
fgsea	THCA_01, pan_01	to do a GSEA	https://bioconductor.org/packages/release/bioc/html/fgsea.html
GSVA	THCA_01, pan_03	to do a GSVA	https://bioconductor.org/packages/release/bioc/html/GSVA.html
ComplexHeatmap	THCA_01, pan_03, pan_04	allows for the creation of heatmaps with more detailed options	https://bioconductor.org/packages/release/bioc/html/ComplexHeatmap.html
metaplot	THCA_02, pan_02, pan_04	data-driven plots	https://cran.r-project.org/web/packages/metaplot/index.html
gridExtra	THCA_02, pan_02, pan_04	"implementation of "grid" graphics "	https://cran.r-project.org/web/packages/gridExtra/index.html
umap	THCA_02, pan_02, pan_04	to do a UMAP	https://cran.r-project.org/web/packages/umap/index.html
gage	pan_01	application of GSEA	https://bioconductor.org/packages/release/bioc/html/gage.html
psych	pan_02	iterative factor analysis	https://cran.r-project.org/web/packages/psych/index.html
cluster	pan_04	cluster analysis	https://cran.r-project.org/web/packages/cluster/cluster.pdf
MASS	neu_00	implementation of neural network	https://cran.r-project.org/web/packages/MASS/index.html
neuralnet	neu_03	training of neural networks	https://cran.r-project.org/web/packages/neuralnet/neuralnet.pdf
AnnotationDbi	descr_03	translating ensemble ids into genenames	https://bioconductor.org/packages/release/bioc/html/AnnotationDbi.html
org.Hs.eg.db	descr_03	translating ensemble ids into genenames	https://bioconductor.org/packages/release/bioc/html/org.Hs.eg.db.html

Table 2: Packages used in the analysis.

Package	Localisation	Usage	Link
biomart	pre_02, pre_03, pre_05	renaming the genenames from the hallmarkpathways-dataframe into ensembleIDs	https://bioconductor.org/packages/release/bioc/html/biomart.html

Package	Localisation	Usage	Link
msigdb	pre_03	downloading all of the canonical pathways and the genes which they include in homo sapiens from the msigbdr data base	https://bioconductor.org/packages/release/data/experiment/html/msigdb.html
dplyr	pre_04, pre_05	tidying and manipulating of dataframes	https://cran.r-project.org/web/packages/dplyr/index.html
ggplot2	pre_04, pre_05, descr_03, descr_04, THCA_01, THCA_02, pan_02, pan_04	allows for the creation of plots with more detailed options	https://cran.r-project.org/web/packages/ggplot2/index.html
pheatmap	descr_01, pan_01, neu_02, neu_04	allows for the creation of heatmaps with more detailed options	https://cran.r-project.org/web/packages/pheatmap/pheatmap.pdf
vioplot	descr_02	creation of violinplots	https://cran.r-project.org/web/packages/vioplot/index.html
VennDiagram	pan_05	creation of VENN-diagrams	https://cran.r-project.org/web/packages/VennDiagram/VennDiagram.pdf
dplyr	THCA_01, pan_01	NA	NA
fgsea	THCA_01, pan_01	to do a GSEA	https://bioconductor.org/packages/release/bioc/html/fgsea.html
GSVA	THCA_01, pan_03	to do a GSVA	https://bioconductor.org/packages/release/bioc/html/GSVA.html
ComplexHeatmap	THCA_01, pan_03, pan_04	allows for the creation of heatmaps with more detailed options	https://bioconductor.org/packages/release/bioc/html/ComplexHeatmap.html
metaplot	THCA_02, pan_02, pan_04	data-driven plots	https://cran.r-project.org/web/packages/metaplot/index.html
gridExtra	THCA_02, pan_02, pan_04	implementation of “grid” graphics	https://cran.r-project.org/web/packages/gridExtra/index.html
umap	THCA_02, pan_02, pan_04	to do a UMAP	https://cran.r-project.org/web/packages/umap/index.html
gage	pan_01	application of GSEA	https://bioconductor.org/packages/release/bioc/html/gage.html
psych	pan_02	iterative factor analysis	https://cran.r-project.org/web/packages/psych/index.html
cluster	pan_04	cluster analysis	https://cran.r-project.org/web/packages/cluster/cluster.pdf
MASS	neu_00	implementation of neural network	https://cran.r-project.org/web/packages/MASS/index.html

Package	Localisation	Usage	Link
neuralnet	neuro_03	training of neural networks	https://cran.r-project.org/web/packages/neuralnet/neuralnet.pdf
AnnotationDbi	bioc_03	translating ensemble ids into gennames	https://bioconductor.org/packages/release/bioc/html/AnnotationDbi.html
org.Hs.eg.db	data_03	translating ensemble ids into gennames	https://bioconductor.org/packages/release/data/annotation/html/org.Hs.eg.db.html

Table 3: Packages used in the analysis.

Package	Localisation	Usage
biomart	pre_02, pre_03, pre_05	renaming the genes
msigdb	pre_03	downloading all of t
dplyr	pre_04, pre_05	tidying and manipu
ggplot2	pre_04, pre_05, descr_03, descr_04, THCA_01, THCA_02, pan_02, pan_04	allows for the creati
pheatmap	descr_01, pan_01, neu_02, neu_04	allows for the creati
vioplot	descr_02	creation of violinplo
VennDiagram	descr_05	creation of VENN-d
dplyr	THCA_01, pan_01	NA
fgsea	THCA_01, pan_01	to do a GSEA
GSVA	THCA_01, pan_03	to do a GSVA
ComplexHeatmap	THCA_01, pan_03, pan_04	allows for the creati
metaplot	THCA_02, pan_02, pan_04	data-driven plots
gridExtra	THCA_02, pan_02, pan_04	implementation of "
umap	THCA_02, pan_02, pan_04	to do a UMAP
gage	pan_01	application of GSEA
psych	pan_02	iterative factor anal
cluster	pan_04	cluster analysis
MASS	neu_00	implementation of n
neuralnet	neu_03	training of neural ne
AnnotationDbi	descr_03	translating ensemble
org.Hs.eg.db	descr_03	translating ensemble