
CAnCER bioMarker Prediction Pipeline (CAMPP)

Thilde Bagger Terkelsen
Danish Cancer Society Research Center
thilde@cancer.com, thildebate@gmail.com

2019-04-10

Introduction

The CAnCER bioMarker Prediction Pipeline (CAMPP) is a simple bioinformatics tool intended to automatize identification of potential diagnostic and prognostic cancer biomarkers. The pipeline is versatile and may be used for analysis of a variety of quantitative biological data from high throughput platforms, including genes, proteins, small RNAs, lipids and glycans. CAMPP currently supports; differential expression/abundance analysis, LASSO/Elastic Net regression, Weighed Gene Co-expression Network Analysis, Correlation analysis and Survival analysis (Cox proportional hazard regression). CAMPP is written in R [1] and runs via a linux command-line with flags specifying arguments.

Contents

1	About CAMPP	2
2	Requirements	4
3	Download	5
4	Installation of R-packages	5
5	Running CAMPP	7
5.1	Mandatory Data Input	7
5.2	Arguments	8
6	N-glycan Serum Markers for BC Diagnostics	14
6.1	Data Normalization, Transformation and Distributional Checks	14

6.2	Differential Abundance Analysis and LASSO/Elastic- Net Regression	16
6.3	Correlation of abundances in TIF and serum	18
6.4	Weighed Gene Co-expression Network Analysis	20
6.5	Survival Analysis	22
6.6	Protein-Protein / Gene-miR Interaction Networks	25

Workflow Overview

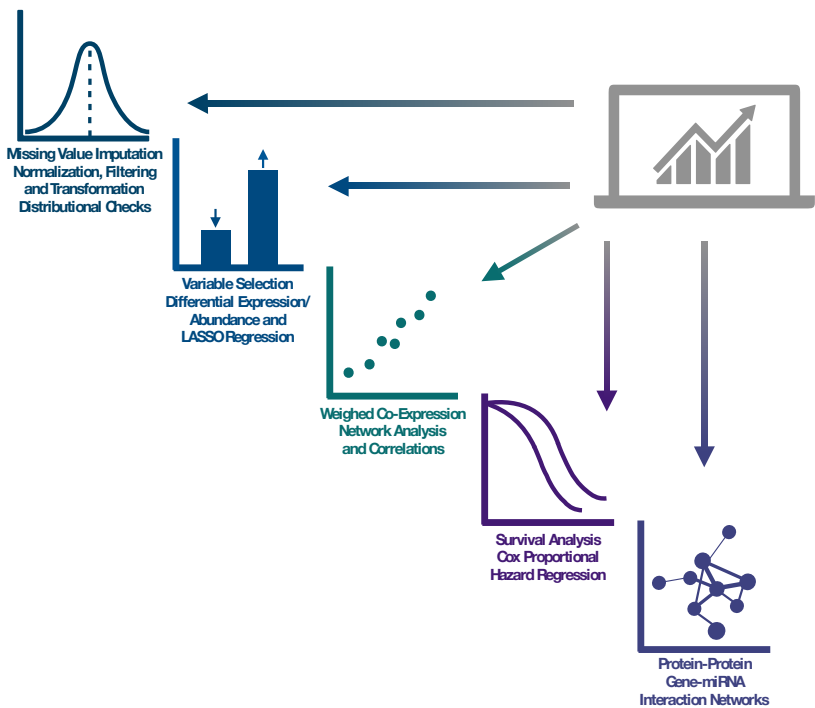


Figure 1: Cancer bioMarker Prediction Pipeline Analysis Flow.

1 About CAMPP

The CAnCer bioMarker Prediction Pipeline was developed for internal use at the Danish Cancer Society Research Center (DCRC) providing researchers with an easy way to identify biomarkers for cancer diagnostics and prognostics. In its original form the pipeline was implemented using interstitial fluid samples and tissue samples from patients with

breast cancer but it may be applied to any type of data with comparable properties.

The pipeline can perform the following types of analysis (I) Differential expression/abundance analysis (limma [26]), (II) LASSO/Elastic-Net regression (glmnet [12]), (III) Weighed Gene Co-expression Network Analysis (WGCNA [23]), (IV) Correlation analysis (Pearson/Spearman), (V) Survival analysis (Cox proportional hazard regression, survcomp [29]) and (VI) protein-protein/gene-miRNA interaction network analysis (multimiR [27] and the STRING [20]). In addition to different types of analysis the pipeline performs missing value imputation, normalization and transformation, along with data distributional checks. Differential expression/abundance analysis can be performed with as few as 3 biological replicates in each group, whereas LASSO/Elastic-Net regression, WGCNA, correlation and survival analysis is only advisable with a large(r) number of samples. CAMPP may be run with a variety of biological molecules (genes, miRNAs, proteins etc.) from various platforms (high-throughput sequencing, microarray data, liquid chromatography-mass spectrometry, etc.). For differential expression/abundance analysis the limma package [26] for R is utilized. Though limma was originally designed for microarray data, and more recently RNAseq, a number of studies have shown the versatility of this software for the analysis of other -omics data [4, 21, 24]. LIMMA has few underlying statistical assumptions and is known to be powerful for small sample sizes as a result of shrinkage of feature-specific variances [31].

N.B The user should be careful with LASSO/Elastic-Net regression, as this type of analysis needs a good number of samples, in a balanced group design, to yield reliable results. Recommended is a minimum of 30 samples for each group in the design [12]. To perform WGCNA [23] at least 15 samples must be available for analysis (section 5 in FAQ here: <https://horvath.genetics.ucla.edu/html/CoexpressionNetwork/Rpackages/WGCNA/faq.html>), see description of WGCNA on page X.

Survival (over-all, relapse-free) of patients based on the abundance/expression of a given marker is predicted using cox proportional-hazard model. Check for proportional hazard and linearity of continuous covariates are automatically performed. The model is corrected for age at diagnosis (surgery, entry into trial) by default, however, if correction for additional confounders is desired these should be specified by the user.

N.B For every parameter added to the model more events are needed for appropriate statistical power. A rule of thumb is ~ 10 events for each parameter, however, this will somewhat depend on the number of levels within a given parameter. Generally results of the survival analysis should be interpreted with caution as the pipeline is unable to account for all possible options and assumptions related to this type of analysis.

The check for proportional hazard and linearity of continuous covariates should be OK before interpreting any results!

2 Requirements

To run CAMPP, a working R version 3.5.1 (or newer) is required [1]. In addition, macbook users must have Xcode installed, as well as the Xcode command-line tools. Install Xcode as shown below.

Open a terminal on your computer and type:

Install Xcode

```
xcode-select -install
```

The Pipeline relies on a variety of R-packages, see list below. CAMP-PIInstall.R script checks whether these packages are installed and helps the user download, install and load them (see section on 4).

Table 1: R-packages

Analysis	Packages
Missing value imputation	impute (v 1.56.0) [17]
Distributional checks	fitdistrplus (v 1.0.11) [8]
Excel formatting	openxlsx (v 4.1.0) [33]
Plotting	heatmap.plus (v 1.3), squash (v 1.0.8) viridis (v 0.5.1), ggplot2 (v 3.1.0) [7, 11, 14, 34]
Data management	data.table (v 1.11.8) stackoverflow (v 0.1.2), plyr (v 1.8.4) scales (v 1.0.0) [9, 13, 35, 36]
K-means Clustering	mclust (v 5.4.3) [30]
DE/DA Analysis	sva (v 3.30.0), limma (v 3.38.2) [19, 26]
LASSO/Elastic-Net	glmnet (v 2.0.16) [12]
Co-expression analysis	WGCNA (v 1.66) [23]
Interaction Networks	igraph (v 1.2.4), biomaRt (v 2.38.0) multiMiR (v 1.4.0), devtools (v 2.0.1) [6, 10, 27, 37]
Survival analysis	survminer (v 0.4.3), survcomp (v 1.32.0) [22, 29]

To obtain information on R-packages above, e.g. developers and release information (articles) go to `R` and type `citation("package_name")`.

3 Download

The CAncer bioMarker Prediction Pipeline is easily installed by following the steps outlined below.

- (1) Download the GitHub repository as a .zip locally from the address below. (2) If you wish to work externally on a server you can upload the repository from your computer to the server using scp from the command-line or (3) use git for a direct download:

Download CAMPP

```
(1) https://github.com/ELELAB/CAncer-bioMarker-Prediction-
-Pipeline-CAMPP

(2) scp CAncer-bioMarker-Prediction-Pipeline-CAMPP-master.zip
name@login:/path-to-dir.

(3) git clone git://github.com/ELELAB/CAncer-bioMarker-
Prediction-Pipeline-CAMPP.git
```

If you are using the "git clone" option, make sure that the Git software has been installed.

- Unzip the repository. From the command-line use "unzip" to decompress.

Decompress CAMPP

```
unzip CAncer-bioMarker-Prediction-Pipeline-CAMPP-master.zip
```

The unzipped repository should contain the following: `Readme.md` `LICENSE.md` `CAMPPInstall.R`, `CAMPPFunctions.R` `CAMPP.R`.

4 Installation of R-packages

The CAncer bioMarker Prediction Pipeline consists of three R-scripts; `CAMPPInstall.R`, `CAMPPFunctions.R` and `CAMPP.R`. The `CAMPPInstall.R` only needs to be used the first time the pipeline is run - this script ensures that all required R-packages will be installed. The `CAMPPFunctions` is a script containing custom functions used in the analysis and must therefore be located in the directory from where the

pipeline is run.

1. In order to ensure that all R-packages needed have been installed, run the CAMPPInstall.R script from the command-line.

Install R-packages from Mac OS X and Ubuntu terminal

```
Rscript CAMPPInstall.R
```

Install R-packages from Windows terminal

```
"C:\Program Files\R\R-3.5.1\bin\Rscript" CAMPPInstall.R
```

In the windows example above the R-version is 3.5.1, this should of course match the R-version installed on the computer!

2. Running CAMPPInstall.R will generate a script named CAMPPmissingpackages.R. This script contains the names of any packages not installed. If no CAMPPmissingpackages.R script is generated it means that the required R-packages already exist.
3. To require any missing packages open R in the command-line and source the CAMPPmissingpackages.R script. R will suggest a library to place the packages in and will query about a "CRAN-mirror" - pick the mirror closest to your geographical location.

Install Missing R-packages from Mac OS X and Ubuntu terminal

```
R  
source("CAMPPmissingpackages.R")
```

Install Missing R-packages from Windows terminal

```
C:\Program Files\R\R-3.5.1\bin\Rscript" R  
source("CAMPPmissingpackages.R")
```

4. Although the CAMPPmissingpackages.R will aid the user in installing the required R-packages from both CRAN and Bioconductor, there may be instances where a package is not updated and compatible with the running version of R.

Warning message:

```
In install.packages("name_of_package") : package 'name_of_pack-
```

age' is not available...

In this case the package may need to be installed directly from the github repository. In order to do this check that the R-package devtools is installed and loaded. Next, use `install.github()` to install:

Install Missing R-packages

```
install.packages("devtools")
library(devtools)
install_github(Bioconductor/name_of_package)
install_github(Bioconductor-mirror/name_of_package)
```

5 Running CAMPP

After installation of all R-packages the CAnCER bioMarker Prediction Pipeline is ready to be used.

It is essential that the script `CAMPPFunctions.R` is located in the directory from which the pipeline is run.

The pipeline is run from the command-line using flags. The user may find the example files used for running the example below at <https://github.com/ELELAB/N-glycan-TIF> along with the original publication. Files are located in the *Data/DataExamples* folder.

5.1 Mandatory Data Input

1. **Data (-d):** An .xlsx (or .txt) file containing feature expression/abundance. With rows as variables, and columns as samples, e.g. columns are N-glycan, protein, (mi)RNA identifiers and rows are sample IDs. The repository <https://github.com/ELELAB/N-glycan-TIF/tree/master/Data/DataExamples> contains an example with N-glycans named: `glycandata.xlsx`. If two datasets are provided (for correlation and/or network analysis), this option should be specified as a comma separated list (without quotes or parenthesis!) of length two, first entry being data file 1 and second entry data file 2.
2. **Data Variant (-v):** The user must specify what type of data is provided in order for the pipeline to pick the appropriate normalization and/or transformation. Options include; array (microarray data), seq (high throughput sequencing data), ms (mass spectrometry data) or other (other type). If two datasets are provided (for correlation and/or network analysis), this option should be specified as a comma separated list (without quotes or parenthesis!) of length two, first entry referring to data file 1 and second entry referring to the data file 2.

- Sequencing data (**-v** is set to **seq**): Variables with low counts over all groups (tissue, treatment) are filtered out, library sizes are scaled (normalization method is weighted trimmed mean of M-values, TMM) and data are voom transformed.
- Microarray data (**-v** is set to **array**): Data are log transformed and either quantile normalized (`normalizeBetweenArrays`) or standardized using mean or median (specify option **-z**).
- Mass spectrometry data (**-v** is set to **ms**): **IF** option **-t** is specified, then data will be log transformed (log2, log or logit as specified).

It should be noted that CAMPP does NOT perform within-array-normalization (`normalizeBetweenArrays`), which is standard for two color intensity data, e.g. this must be done before hand (see limma manual for more information [26]).

3. **Metadata (-m)**: An .xlsx file (or .txt) containing metadata. This file must contain at least two columns, one with identifiers matching the column names in the data file and one with groups to contrast in analysis e.g. diagnosis (tumor or normal), tumor stage (1,2 or 3), drug treatment (A, B C) ect. If two datasets are provided (for correlation and/or network analysis), this option should be specified as a comma separated list (without quotes or parenthesis!) of length two, first entry being metadata file 1 and second entry metadata file 2. The repository <https://github.com/ELELAB/N-glycan-TIF/tree/master/Data/DataExamples> contains an example with N-glycans named: glycanmetadata.xlsx.
- **Ids and Groups (-g)**: The user must specify which columns in the metadata file corresponds to the sample ids and groups for comparison, receptively. This is done by providing a list of two strings separated by a comma (without quotes or parenthesis!) , indicating the names of the columns which should be used. If two datasets are provided, and both datasets are to be corrected for experimental batch, this option should be specified as a comma separated list (without quotes or parenthesis!) of length four. The first two entries in the list specifying names of columns (ids and groups) to use from metadataset 1 and third and fourth entry specifying names of columns (ids and groups) to use from metadataset 2.

5.2 Arguments

1. **A Simple Run**: To run the pipeline, arguments (flags) must be set to specify the input. In the simplest case we are interested in performing differential abundance/expression analysis and for this we only need to specify the mandatory inputs.

We want a custom name for our results, so we set **-n** and specify a string, in this case we will call it "SimpleRun".

Mandatory Arguments

- d: Quantitative data
- m: Metadata
- v: Data variant
- g: Ids and Groups

The box below shows an example of the simplest possible CAMPP run.

First Run Example

```
Rscript CAMPP.R -d data.xlsx -m metadata.xlsx -v seq  
-g ids,group -n SimpleRun
```

2. The output of the command above will be an .txt (tabular) file with the identifiers, test-statistics, fold changes (logFC) and p-values of any differentially expressed/abundant variables.

Table 2: Example of .txt Output

logFC	t	P.Val	adj.P.Val	B	name	dir.	comp.
2.1	11.6	7.7e-22	4.9e-20	38.9	pA	up	T-N
1.3	8.3	8.4e-14	1.1e-12	20.6	pK	up	T-N
.
.
.
-1.2	-10.2	4.1e-18	1.3e-16	30.5	pH	down	T-N

N.B. The pipeline logs while running and produces a text file, CAMPPlog.txt, with any errors or comments accumulated during analyses - this file should be viewed carefully.

3. **Help:** In addition to the two mandatory arguments above there are a range of optional arguments which may be utilized. The -h (help) option allows the user to obtain information about all available flags.

CAMPP Help

```
Rscript CAMPP.R -h
```

4. **Other Arguments:** The table below show other available arguments.

Optional Arguments

- s: Multidimensional Scaling Plot
- t: Data Transformation
- b: Batch Correction
- j: Distributional Checks
- c: Color Scheme
- r: Covariates
- f: Cut-offs for logFC and corrected p-value
- k: Kmeans Clustering
- l: LASSO/Elastic-Net Regression
- w: Weighed Gene Co-expression Analysis
- o: Correlation Analysis
- u: Survival Analysis (Cox Regression)
- p: Protein-Protein Interaction Networks
- i: miRNA-Gene Interaction Networks

- **Multidimensional Scaling plot (-s):** The flag -s may be used to generate a preliminary multidimensional scaling (MDS) plot. Multidimensional scaling (MDS) reduces high-dimensional data to two (or more) "observable" dimensions (M1 and M2) in such a way that the inter-sample distance relationship in the simplified dimensions mimic the distance relationship between the original values of samples in n-dimensional space (n = number of variables). MDS was performed with the euclidean distances as the distance metric. A MDS plot may help to determine whether there appears to be any "grouping" of data in a desired way (tumor samples together and normal samples together) or in an undesired way (experimental batches) before performing the DE/DA analysis.
- **Data Transformation(-t):** The flag -t may be set if the user desires the expression/abundance data to be transformed before analysis. A logarithmic transformation is recommended as the variance of measurements, from most platforms, depend on the expression/abundance level itself. The log-transformation reduces this dependency and additionally pushes the negative binomial distribution, displayed by count data (miRNA, mRNA), towards a normal distribution. The user may choose between log₂, log₁₀, logit or voom transformation. For RNA-seq the voom transformation is recommended [26], while a log transformation may be more appropriate for proteomics and N-glycan abundances. If another sample paired dataset of expression values are provided, this option should be specified as a comma separated list (without quotes or parenthesis!) of length two, first entry referring to data file 1 and second entry referring to the data file 2.
- **Batch (-b):** If the data comes from experimental batches and the user wants to correct for this, a column specifying which batch each sample belongs, should also be included in the metadata file.

The argument `-b` takes a string (no quotes!) referring to the name of the column in the metadata file denoting batches (e.g. A, B, C, or batch1, batch2, batch3, etc.). Batch type must be noted as a character, meaning numbers alone are not allowed. If two datasets are provided, this option should be specified as a comma separated list (without quotes or parenthesis!) of length two, first entry matching the name of a column in metadata file 1 and the second entry matching the of a column in metadataset 2.

- **Distributional Checks (-j):** The flag `-j` may be set to FALSE to remove the default check of variable distributions. If the flag is not specified (or set to TRUE), the pipeline will produce plots including histograms, quantile-quantile plots and probability plots. By default 10 random variables are picked from the dataframe for plotting (here we are assuming that most variables in a gene expression matrix or protein abundance matrix will belong to the same family of distributions.) **We heavily recommend always running the data checks and to NOT ignore the output, as this vital to whether or not results are reliable!**
- **Colors (-c):** Flag to change group color scheme. Accepted R-colors must be specified in a comma separated list (without quotes or parenthesis!) of a length matching the number of groups.
- **Covariates(-r):** The user may specify flag `-r` if covariates should be included in the differential expression/abundance analysis and/or the survival analysis. This argument takes a comma separated list (without quotes or parenthesis!). The first element in this list must be either TRUE or FALSE. If TRUE is specified then covariates will be included in both DE/DA analysis and Survival Analysis. If FALSE is specified covariates will ONLY be used for Survival Analysis. All other elements of the list after element one (TRUE/FALSE) must be strings matching one or more column names in the metadata file. **Age** is automatically added as a covariate for survival analysis and should therefore not be specified with `-r`!
- **LogFC and FDR (-f):** Cut-offs for log2 fold change and corrected p-value (fdr). Defaults are $\log_2FC > 1$ or $\log_2FC < -1$ and $fdr < 0.05$. This argument takes a comma separated list of length two (without quotes or parenthesis!). The first element specifying cut-off for logFC and the second element specifying cut-off for FDR. If two datasets are provided, the list must have length four, e.g. cut-offs for both sets.
- **Kmeans Clustering (-k):** The flag `-k` is set to specify K-means clustering. This argument takes a string specifying which column in the metadata file should be used to label the samples in the returned plot(s). If `-k` is set but left empty, no labels are added to the plot(s). The number of clusters tested will be based on number of samples, fewer samples will result in fewer kmeans

tested. A folder with MDS plots will be returned for the best n number of kmeans, based on the bayesian information criterion (BIC) [30]. If the dataset has many variables i.e. RNAseq with many genes, multiple samples of 3000 variables will be generated and tested to overcome issues with computational time and the consensus of best n kmeans will be returned.

Clustering may only be performed one dataset at a time!

- **LASSO/Elastic-Net Regression (-l):** The flag -l may be set to specify least absolute shrinkage and selection operator (-l 1.0) or Elastic-Net (0.0 < -l < 1.0) regression. LASSO/EN is performed using the R-package **glmnet** [12]. K-fold (default is 10) cross validation (cv.glmnet) is used to estimate the optimal value for the hyperparameter, lambda. LASSO/EN may be performed in two ways, (I) the dataset is split into training and testing subsets, k-fold cross validation is performed on the training dataset, followed by estimation of specificity and sensitivity (area under the curve = AUC) using the test dataset, or (II) k-fold cross validation is performed using the full dataset, no AUC is reported. CAMPP will automatically estimate whether the input dataset is large enough to split into training and test subsets and whether EN/LASSO is advisable to perform altogether. Specifying the flag -l will produce a list of variables selected by LASSO and a file with the overlap between differential expression/abundance analysis and LASSO regression. LASSO is run with n (default is 10) different random seeds and the consensus set of variables is returned.

N.B LASSO it not appropriate for heavily unbalanced group designs! When running LASSO bar-plots with cross-validation errors and AUCs are returned. The user should evaluate whether these are reasonable.

LASSO may only be performed on one dataset at a time!

- **Weighed Gene Co-expression Network Analysis (-w):** The flag -w must be set in order to perform Weighed Gene Co-expression Network Analysis with the R-package WGCNA [23]. WGCNA, despite its name, is in this case not exclusive to gene expression data but may be applied to any quantitative expression data. Minimum module size is 10 variables and modules with less than 25% dissimilarity will be merged (default values). The cutoff for top most interconnected variables (genes, proteins ect.) in an identified module is set to the 75th quantile by default. The user may specify different cutoffs for minimum module size, module dissimilarity merging and and % of interconnected variables to report with the flag -x.

N.B the softpower plot generated by a WGCNA run should be inspected before interpretation of results. If the data is very heterogeneous, or if there are too few variables or samples (minimum 15 samples for WGCNA) the scale-free topology fit index (printed to the screen while running) might fail to reach values above 0.8 for reasonable powers (see point 5 in WGCNA tuto-

rial: <https://horvath.genetics.ucla.edu/html/CoexpressionNetwork/Rpackages/WGCNA/faq.html>). If this is the case, the dataset may not be appropriate for WGCNA.

WGCNA may only be performed with one dataset at a time!

- **Correlation Analysis (-o):** The flag -o must be set to perform correlation analysis between two matched datasets. For option -o the user must provide a string specifying which subset of variables should be included in the correlation analysis; "ALL" = all variables (not advisable, unless dataset is small), "DA" / "DE" = Differentially Expressed/Abundant, "LASSO" / "EN" = LASSO / Elastic Net results or "Consensus" = Overlap between DE and LASSO/EN. Naturally two datasets must be input in order to perform correlation analysis. Both datasets are given to the argument (-d), with the names of the files separated by a comma (no quotes or parenthesis!). The two datasets do not need to have the same dimensions, but there must be at least a partial overlap in both variables and samples (column names). Column names of these datasets should match the IDs in the `ids` column in the two metadata files.
- **Survival Analysis (-u):** The flag -u must be set in order to perform survival analysis using cox proportional hazard model. Options for survival analysis are; "ALL" (not advisable!), "DA" / "DE" = Differentially Expressed/Abundant, "LASSO" / "EN" = LASSO / Elastic Net results or "Consensus" = Overlap between DE and LASSO/EN, referring to the set of variables used for cox-regression. For survival analysis the metadata file must contain at least four columns in addition to the sample IDs named; 'age' (age in years at diagnosis, surgery or entry into trial), 'outcome.time' (time until end of follow-up, censoring or death in weeks, months or years), 'outcome' (numeric 0 = censoring, 1=dead) and 'survival' (numeric 0 = no survival info, 1=survival info available). If the user wishes to correct for potential confounders (e.g. tumor grade, hormone levels, drug-treatment ect.) these should also be included in the `metadata.xlsx`. CAMPP checks two underlying assumptions of the cox model before performing survival analysis (I) a linear relationship of continuous covariates with log hazards and (II) proportional hazards of categorical and continuous covariates, e.i. constant relative hazard. If the requirement of linearity is not fulfilled, cubic splines will be added to the covariate(s) in question.
- **P-P Interactions (-p):** Flag -p may be specified to perform protein-protein interaction network analysis. Input for this argument is a comma separated list of length two (without quotes or parenthesis!), where the first element specifies the type of gene IDs in the data file, accepted IDs are: uniprotswissprot, ensembl_peptide_id, hgnc_symbol, ensembl_gene_id or ensembl_transcript_id. The second element specifies the protein-protein

interaction database to use, currently the only supported database is STRING [20], accepted inputs are: stringdatabase.

- **Gene-miRNA Interactions (-i):** Flag -i may be specified to perform gene-miRNA interaction network analysis. Input for this argument is a comma separated list of length two (without quotes or parenthesis!), where the first element specifies the type of miRNA IDs in the data file, accepted IDs are: mature_mirna_ids or mature_mirna_accession. The second element specifies the Gene-miRNA interaction database to use (validated, predicted or both), accepted inputs are: targetscan, mirtarbase or tarscanbase.

N.B: If both -p and -i are set, CAMPP will integrate P-P and Gene-miR networks. Naturally two data files must be provided to use both -p and -i, it is assumed that data file 1 contains genes and data file 2 contains miRNAs!

6 N-glycan Serum Markers for BC Diagnostics

This section contains an example of how CAMPP is run using different flags in the command-line and what output the user can expect. The files used for running this example may be found in the repository <https://github.com/ELELAB/N-glycan-TIF/tree/master/Data/DataExamples>.

The case below uses N-glycans abundances measured using high resolution quantitative Ultra-Performance Liquid Chromatography (UPLC) [28] from interstitial samples [15] and matched serum. Tumor interstitial fluid (TIF), normal interstitial fluid (NIF) and serum samples were collected from ~ 90 women diagnosed with breast cancer (BC). A total of 165 N-glycan groups were identified [32].

Briefly, the involvement of N-glycosylation in development and progression of BC has been documented by both in vitro and in vivo studies [3, 5, 16]. Several circulating N-glycan patterns with altered glycan structures, possibly originating from a primary tumor or from other organs, in response to a neoplastic process, have recently been described in a number of studies by using high-throughput N-glycan profiling [2, 18, 25, 28].

6.1 Data Normalization, Transformation and Distributional Checks

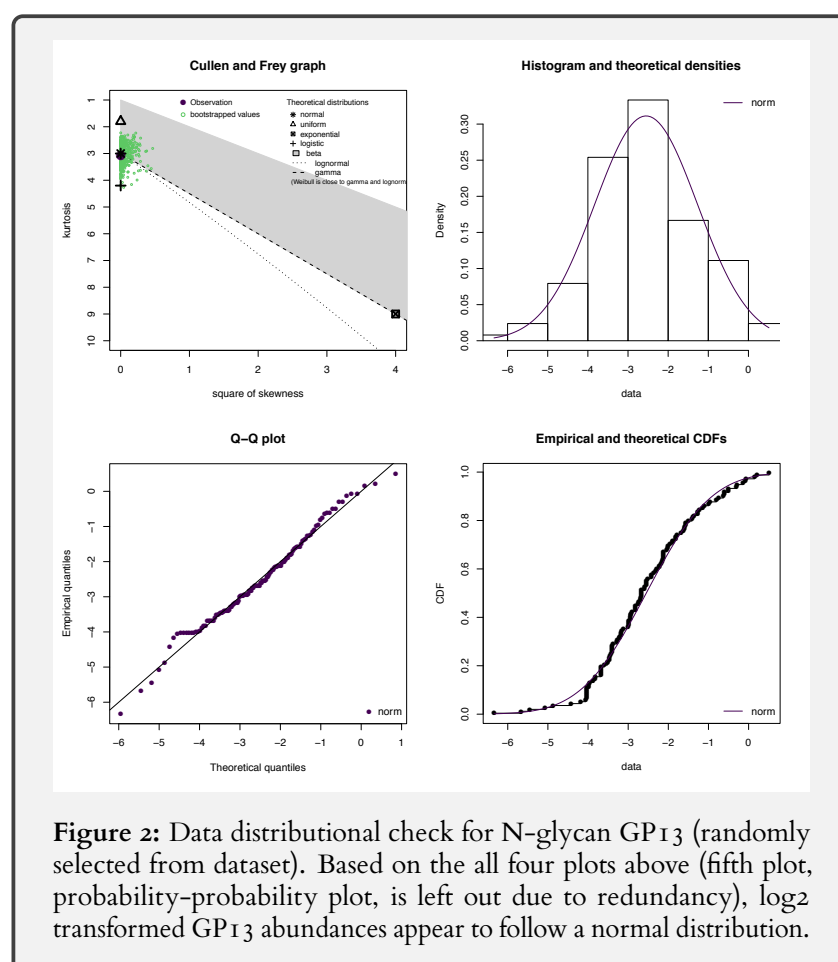
If an input dataset contains missing values CAMPP will automatically impute these, unless missing per column > 70%. Depending on which data input is given CAMPP will perform data normalization and transformation (**argument -v, data variant, see section on mandatory input**). The quantitative N-glycan data (used in this example) are relative (fractional) meaning that the sum of all values within one sample yields 100. Abundances of N-glycans were quantified using

liquid chromatography tandem mass spectrometry (LC-MS/MS). The dataset had already been standardized by the MS-software, therefore the option `-v` was set to `ms` and the option `-t` was set to `log2`. As the N-glycan data were quantified over three LC-MS/MS runs, the argument `-b` was set, e.g. data were corrected for experimental batch.

Before performing any analysis it is advisable to evaluate the distribution of the normalized data. CAMPP automatically generates distributional plots for `n` (default is 6) randomly selecting input variables for the user to evaluate - to skip this step the argument `-j` may be set to `FALSE`).

Output:

- I. A 5-panel plot for each of the `n` (default is 6) randomly selected variables. The plot contains, quantile -and probability plots, a histogram of expression/abundance values with fitted distributions and a skewness-kurtosis plot (Cullen and Frey graph, see `fitdistrplus` manual [8]) - See **Figure 2**.



6.2 Differential Abundance Analysis and LASSO/Elastic-Net Regression

Differential Abundance Analysis and LASSO regression was run with correction for experimental batch (p1, p2, p3), in this particular case the column name specifying batch (in metadata file) was "pool". The data were log₂ transformed (flag -t) and a preliminary MDS plot was generated (flag -s.) The cut-off for a significant hit was an FDR ≤ 0.05 (default) - no logFC cut-off was set to define N-glycan differential abundance (flag -f). The -l was set to 1.0, specifying LASSO (for Elastic-Net 0.0 < -l < 1.0). Files used for this example may be found at <https://github.com/ELELAB/N-glycan-TIF/tree/master/Data/DataExamples>.

The box below shows an example of DAA and LASSO regression with CAMPP.

Differential Abundance Analysis.

```
Rscript CAMPP.R -d glycandata.xlsx -m glycanmetadata.xlsx  
-v ms -g sids,cn -t log2 -b pool -l 1.0 -s TRUE  
-f 0,0.05 -n FirstRunOutput
```

Output:

The run above generates three outputs:

1. A multidimensional scaling plot (FirstRunOutput_MDSplot.pdf) using the abundance/expression of variables (in this case N-glycans). The components M1 and M2 in the plot below are those which best retained the distance relationship between samples in two dimensions. - See **Figure 3**.
2. A tabular .txt file (FirstRunOutput_DE.txt) with the stats for significant differentially abundant/expressed variables (in this case N-glycans).
3. A heatmap (FirstRunOutput_heatmap.pdf) visualizing the power of the identified DA/DE variables (in this case N-glycans) for partitioning of the samples in groups/condition/treatment (in this case NIF and TIF samples) - See **Figure 4**.
4. A tabular .txt file (FirstRunOutput_LASSO.txt) with the LASSO selected variables for group classification.
5. A tabular .txt file (FirstRunOutput_DEA_LASSO_Consensus.txt) containing the overlap of DAA results and LASSO results.
6. Bar-plots (FirstRunOutput_CrossValidationPlot.pdf, FirstRunOutput_AUCTestDataClassification.pdf) with cross-validation errors and AUCs. - **Figures not shown.**

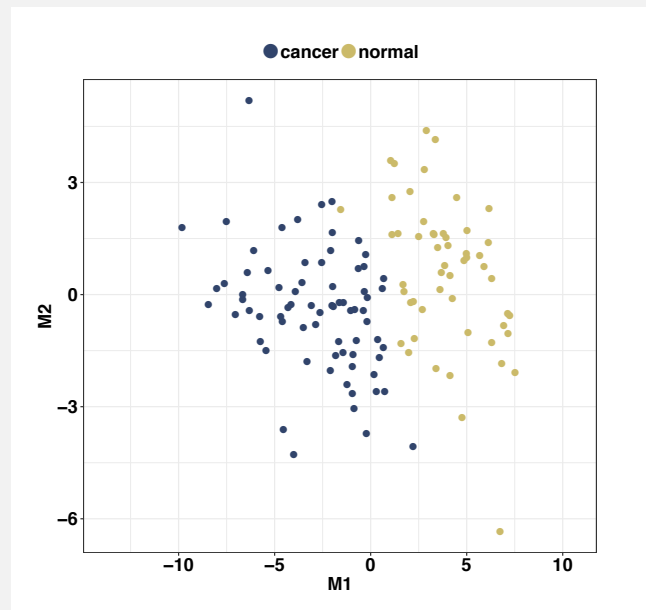


Figure 3: Multidimensional Scaling Plot showing the separation of tumor interstitial and normal interstitial fluid samples based on N-glycan abundances. The data was batch corrected before plotting.

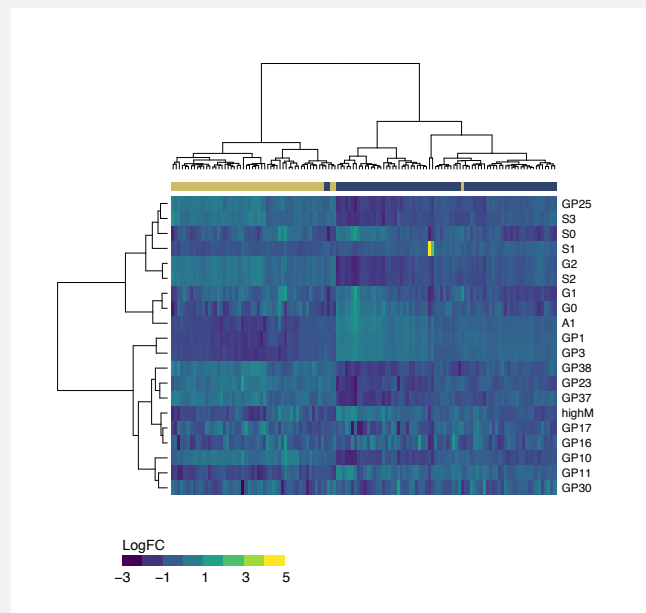


Figure 4: Heatmap showing the partitioning of TIF and NIF samples based on 20 N-glycan peaks/groups identified as differentially abundant.

The multidimensional scaling plot in **Figure 3** indicates that N-glycan abundance patterns may confer a separation of tumor interstitial fluid and normal interstitial fluid samples from patients with breast cancer – at least when all N-glycans variables are retained. The MDS plot may be used as an indicator of whether the differential expression analysis is likely to yield any results, e.g. no clustering of samples from different groups/conditions/treatments implies that the biological data does NOT reflect the sample classification.

6.3 Correlation of abundances in TIF and serum

In order to determine whether N-glycans with differential abundances in tumour –and normal interstitial fluids displayed corresponding patterns in serum, correlation was performed using TIF samples and matched serum.

For correlation analysis with CAMPP, the user must provide:

1. Two .xlsx (or .txt) files with expression/abundance data (rows as variables and columns as samples). See section 5.2.4 for specifications. Examples may be found in the repository <https://github.com/ELELAB/N-glycan-TIF/tree/master/Data/DataExamples>.
2. Two .xlsx (or .txt) files with metadata for each dataset. Examples may be found in the repository <https://github.com/ELELAB/N-glycan-TIF/tree/master/Data/DataExamples>.
3. A string specifying which subset of variables to use for correlation analysis, options are; "ALL", "DE" / "DA", "LASSO"/ "EN" or "Consensus" – see section X.

Files used in the example below, can be found in the repository <https://github.com/ELELAB/N-glycan-TIF/tree/master/Data/DataExamples>.

Command-line box 6.2.3 shows an example of a CAMPP correlation analysis run with TIF and serum N-glycan abundances. Both data and serumdata were produced by tandem mass spectrometry and -v was therefore set to ms,ms. Also, both sets were log2 transformed (argument -t) before analysis. Only differentially abundant N-glycans were included in the correlation analysis As serum data were not produced in batches argument -e was unneeded.

Correlation Analysis.

```
Rscript CAMPP.R -d glycandata.xlsx, glycanSdata.xlsx
-m glycanmetadata.xlsx, -m glycanSmetadata.xls
-b pool -v ms,ms -g sids,cn,sids,cn -t log2,log2 -o DA
-n SecondRunOutput
```

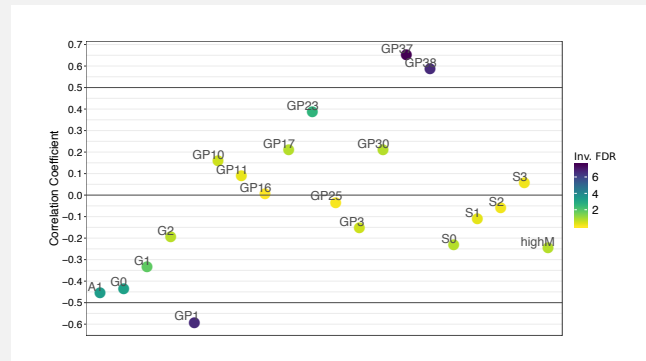


Figure 5: Correlation plot depicting correlation coefficients associated with abundances of N-glycans in TIF and matched serum. Correlation was performed with the 20 N-glycans identified as DA in normal and tumor fluids. Dots are colored in accordance with inverse (scaled) FDRs, e.g. darker shade indicates smaller adjust p-value and vice versus.

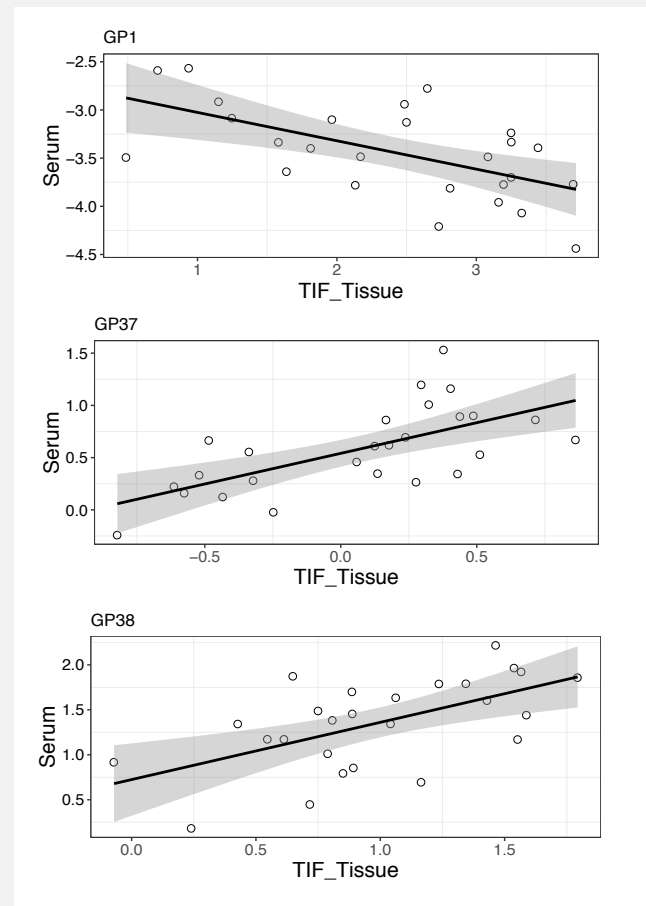


Figure 6: Scatter plots of the 3 N-glycans (GP1, GP37 and GP38) with significant correlations between abundances in TIF and serum. Shaded area indicates confidence of regression line.

Output:

The run above generates three outputs:

1. A tabular .txt file (SecondRunOut_corr_serum.txt) with the stats for variables (in this case N-glycans) displaying significant correlation of abundance/expression in tumor sample (interstitial fluid) and matched serum.
2. A correlation plot (SecondRunOutput_corrplot.pdf) visualizing the correlation coefficients and adjusted p-values of all differentially expressed/abundant variables - in this case N-glycans found to partition TIF and NIF. See **Figure 5**.
3. Individual scatter plots (SecondRunOutput_individual_corrplots.pdf) for the variables displaying a significant correlation between expression/abundance in TIF/tissue and serum (FDR <= 0.05 is default). See **Figure 6**.

Based on the .txt file with coefficients and adjusted p-values (not shown), three N-glycans, GP1, GP37 and GP38 were found to display significant correlations of abundances in TIF and serum (see Figure 5). It is worth noting that levels of GP1 in TIF and serum are negatively correlated indicating an accumulation of this N-glycan within the tumor interstitium, perhaps as a result of primary tumor response. GP37 and GP38 have positive correlation coefficients and may be considered as potential diagnostic serum markers.

Figure 6 shows the scatter plots (TIF abundance plotted against serum abundance) of GP1, GP37 and GP28, with confidence shading of regression lines.

6.4 Weighed Gene Co-expression Network Analysis

In order to explore the abundance relationship between N-glycans in an unsupervised way, Weighed Gene Co-expression Network Analysis was performed with WGCNA [23]. The cutoff for top most interconnected variables (genes, proteins ect.) in an identified module is set to 25% by default. The user may specify a different cutoff with the flag -x. Cutoff for module similarity merging is 0.25 and cutoff for minimum module size is 10 (defaults from the WGCNA tutorial here: <https://horvath.genetics.ucla.edu/html/CoexpressionNetwork/Rpackages/WGCNA/faq.html>). These default parameters cannot be changed with a flag but may easily be modified in the source code of the CAMPP.R script if desired.

Weighed Gene Co-expression Network Analysis.

```
Rscript CAMPP.R -d glycandata.xlsx -m glycanmetadata.xlsx  
-v ms -g sids,cn -t log2 -b pool -w TRUE -n ThirdRunOutput
```

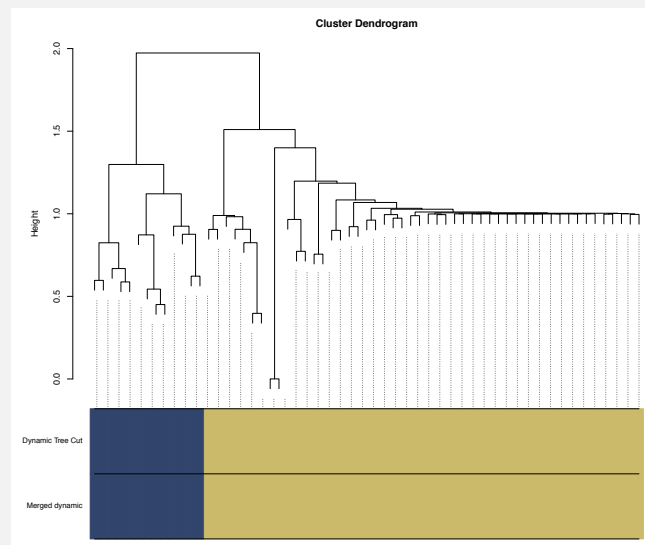


Figure 7: WGCNA module dendrogram showing hclust of N-glycans and which modules each N-glycan was assigned to. The plot shows the modules, in this case two, before and after merging on module similarity. As similarity was not enough to merge, the two original modules were retained in this example.

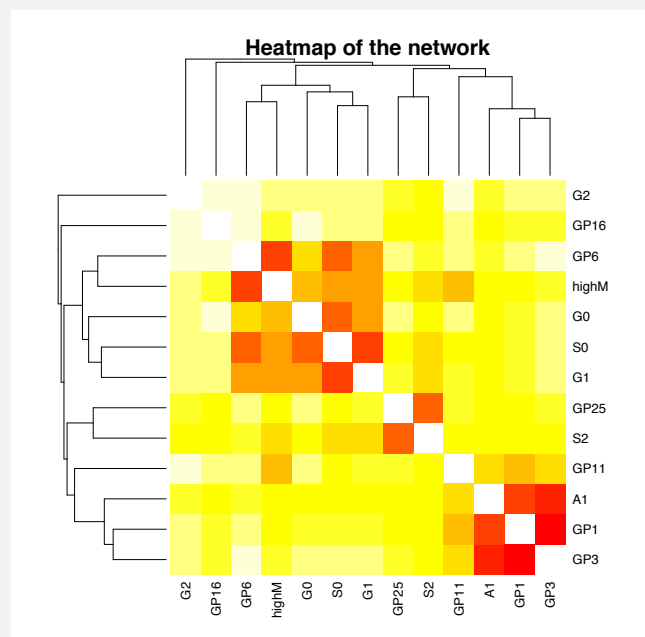


Figure 8: Heatmap of the blue module (heatmap of turquoise module not shown).

Output:

The run above generates three outputs:

1. A tabular .txt file (`ThirdRunOut_WGCNAres.txt`) with the stats for variables (in this case N-glycans) which were the among the top n % (default 25%) most interconnected variables from each identified module, e.g. those with the highest `kWithin` values.
2. WGCNA dendrogram, e.g. clustering of variables (N-glycans), with modules colored before and after merging (`ThirdRunOut_WGCNA_ModuleTree.pdf`). See **Figure 7**.
3. A heatmap for each module network, colored according to variable (N-glycan) co-expression (`ThirdRunOut_moduleHM.pdf`). See **Figure 8**.

The heatmap in **Figure 8** below shows that within the blue module N-glycan features `A1`, `GP1` and `GP3` were the most highly co-expressed, followed by N-glycan features `So`, `G1`, `GP6` and `highM`. These results are consistent with both the results of the TIF and serum correlation analysis (**Figure 5**), as well as the variables selected from LASSO regression and differential expression analysis.

6.5 Survival Analysis

Cox proportional hazard model was used to assess whether any differentially abundant N-glycans had potential as prognostic biomarkers for the survival of breast cancer patients.

CAMPP only provides cox proportional hazard models for overall survival and simple relapse-free survival, e.g. one time entry per patient/sample.

CAMPP checks two underlying assumptions of the cox model before performing survival analysis (I) a linear relationship of continuous covariates with log hazards and (II) proportional hazards of categorical and continuous covariates, e.i. constant relative hazard. If the requirement of linearity is not fulfilled, cubic splines will be added to the covariate(s) in question, and analysis is continued. If the proportional hazard assumption is violated for a categorical variable the user may apply stratification and re-run the pipeline. However, the pipeline does not currently handle continuous variables with non-proportional hazards, e.i. that the user should ignore any cox models where this is the case.

Mandatory Columns for Survival Analysis (metadata.xlsx):

In order to perform cox proportional hazard regression with CAMPP, the user must ensure that the metadata file contains at least four columns in addition to the sample IDs named; 'age' (age in years at diagnosis, surgery or entry into trial), 'outcome.time' (time until end of follow-up, censoring or death in weeks, months or years), 'outcome' (numeric 0 = censoring, 1=dead) and 'survival' (numeric 0 = no survival info,

i=survival info available). If the user wishes to correct for potential confounders (e.g. tumor grade, hormone levels, drug-treatment ect.) these should also be included in the `metadata.xlsx`. The flag `-u` is set to "DA" specifying that the variables which were differentially expressed should be used for survival analysis, alternatives are; "LASSO", "EN" or "Consensus".

Files used in the example below, can be found in repository <https://github.com/ELELAB/N-glycan-TIF/tree/master/Data/DataExamples>. Command-line box 6.3.2 shows an example of a CAMPP survival analysis run. TIF N-glycan abundances were corrected for experimental batch and log2 transformed before regression analysis.

Survival Analysis.

```
Rscript CAMPP.R -d glycandata.xlsx -m glycanmetadata.xlsx  
-v ms -g sids,cn -t log2 -b pool -u DA -n FourthRunOutput
```

In the example above the cox model only contains patient age at diagnosis (default), however, the user may specify other confounders to include. covariates are included with the flag `-r`. If multiple confounders are added these should be separated by commas and their names should match those of the desired columns within `metadata.xlsx`. **N.B** the first argument specified in the `-r` list should always be either TRUE or FALSE. TRUE means that covariates should be included both in the design matrix for differential expression analysis and survival analysis, whereas FALSE specifies that covariates should only be used for cox regression. In the run below, survival analysis was performed with tumor infiltrating lymphocyte status (name: TILS, type: categorical) and tumor size (name: TSize, type: continuous) as covariates – See the `metadata.xlsx` sheet in the `DataExamples` folder.

N.B It is important to note that each time a covariate is added to the model we increase the degrees of freedom. It is therefore essential to have a dataset of an appropriate size for the number of covariates added. The good rule of thumb is 10 events per parameter/variable.

Survival Analysis (confounders).

```
Rscript CAMPP.R -d glycandata.xlsx -m glycanmetadata.xlsx  
-v ms -g sids,cn -t log2 -b pool -u DA -r FALSE,TILS,TP  
-n FourthRunOutput
```

If a CAMPP run produces the warning message below, it means that

some of the specified covariates violate the proportional hazard assumption:

WARNING: The following variables and/or covariates failed the test of proportional hazard: TILS.

If the covariates that failed are categorical you may use strata by re-running the pipeline adding flag -y followed by the names of the categorical covariates to stratify (if multiple, separate by comma). N.B, this pipeline does not handle continuous variables that violate the proportional hazard assumption, if any of these failed PH test, the hazard ratios of these should NOT be evaluated.

If the s are categorical we can re-run the pipeline with stratification using the flag -y:

Survival Analysis (confounders and stratification).

```
Rscript CAMPP.R -d glycandata.xlsx -m glycanmetadata.xlsx -v ms
-g sids,cn -t log2 -b pool -u DA -r FALSE,TILS,TP -y TILS
-n FourthRunOutput
```

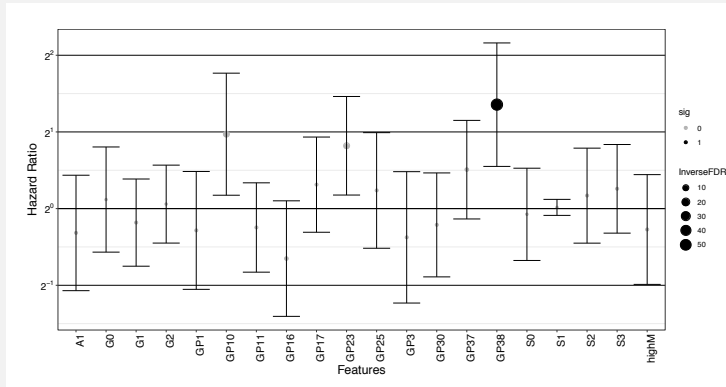


Figure 9: Summary plot of hazard ratios (and confidence intervals) for cox regression with each of the 20 N-glycans found to be differentially abundant between TIF and NIF. A hazard ratio (HR) < 1 indicates that a high level of a given N-glycan is associated with a positive outcome (longer overall survival after diagnosis), while a HR > 1 suggests that a high level of a given N-glycan predicts poorer overall survival. The dot size denotes inverse FDR, e.g. a larger dot is associated with a smaller FDR. Significant N-glycans (FDR <= 0.05) have blue dots, while non-significant N-glycans have orange dots.

Output:

The runs above generate two outputs:

1. A tabular .txt file (FourthRunOut_survival.txt) with the stats for variables (in this case N-glycans) displaying significant associations between abundance/expression in tumor sample (interstitial fluid) and patient outcome.
2. A plot of the hazard ratios for variables (in this case N-glycans) with confidence intervals and significance (FourthRunOut_survivalplot.pdf). See **Figure 8**.

Figure 9 shows that one N-glycan, GP38, was found to be significantly associated with survival, e.i. a high level of this N-glycan was predictive for poor overall survival. GP38 was one of the three N-glycans displaying a correlation between abundances in TIF and serum, suggestion that this glycan structure may have potential as a prognostic serum biomarker.

6.6 Protein-Protein / Gene-miR Interaction Networks

CAMPP may be used to perform protein-protein and/or gene-miRNA interaction network analysis. Interactions are returned for differentially expressed genes/miRNAs. In this case LASSO/EN may not be specified instead, as this type of analysis does not yield the statistics necessary for the analysis.

P-P interactions are extracted from the STRING database [20], with a lower score cut-off > 25th quantile. MiRNAs are retrieved using the multiMiR R-package [27], the user may specify whether to use predicted miRNA targets (TargetScan, lower score cut-off > 25th quantile), validated gene-miRNA pairs (miRTarBase) or both of these. In the following section we given an example of how to run this analysis with miRNA and gene expression data. Here we cannot use the N-glycan dataset as N-glycans are neither genes nor miRNAs. Instead we will show a fictive example. Note that gene -and miRNA identifiers must be in the approved list of inputs - See specifics under the Arguments section above.

The box below shows an example of protein-protein interaction network analysis with CAMPP. Here, the data file contains gene symbols and so the first element in the list given to argument `-p` is `hgnc_symbol`. The second element is the p-p interaction database, in this case the STRING database.

Protein-Protein Network Interaction Analysis.

```
Rscript CAMPP.R -d Genedata.xlsx -m Genemetadata.xlsx
-v seq -g gID,gCN -b gBatch
-p hgnc_symbol,stringdatabase -n FifthRunOut
```

The box below shows an example of gene-miRNA interaction network

analysis with CAMPP. Here, the data file contains miRNA IDs and so the first element in the list given to argument *-i* is *mature_mirna_ids*. We decide that we want predicted gene-miRNA interaction pairs, so we set the second element in list *-i* to *targetscan*.

Protein-Protein Network Interaction Analysis.

```
Rscript CAMPP.R -d miRNAdata.xlsx -m miRNAmetadata.xlsx  
-v seq -g miRID,miRCN -b miRBatch  
-i mature_mirna_ids,targetscan -n FifthRunOut
```

The box below shows an example of CAMPP run where both gene expression data and paired miRNA expression data are available. Here, we can integrate protein-protein and gene-miRNA interactions into one network.

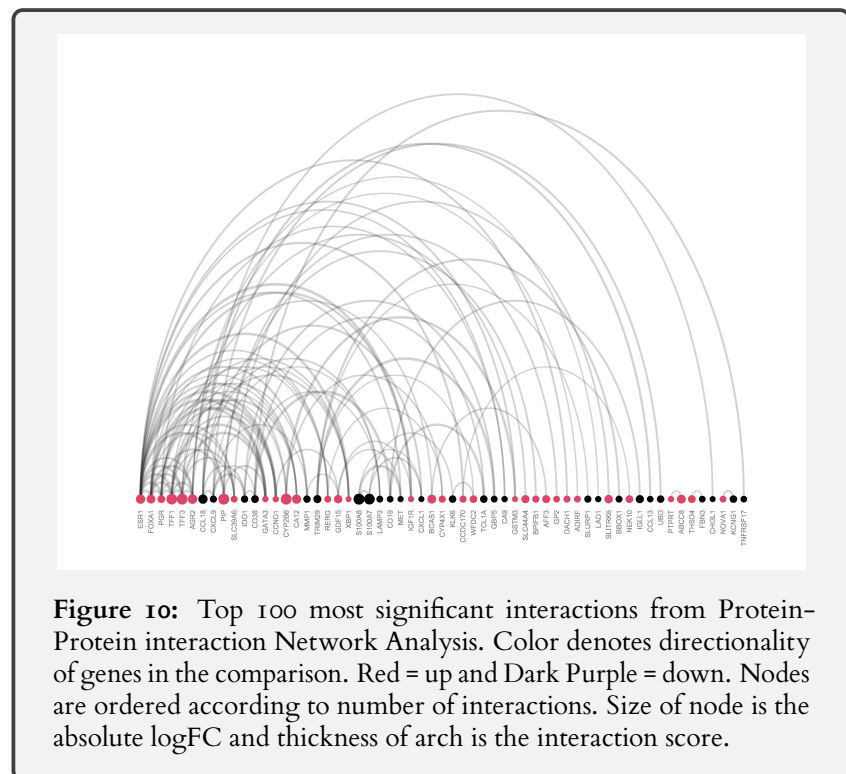
Protein-Protein Network Interaction Analysis.

```
Rscript CAMPP.R -d Genedata.xlsx,miRNAdata.xlsx  
-m Genemetadata.xlsx,miRNAmetadata.xlsx  
-v seq,seq -g gID,gCN,miRID,miRCN -b gBatch,miRBatch  
-p hgnc_symbol,stringdatabase -i mature_mirna_ids,targetscan  
-n FifthRunOut
```

Output:

The runs above will return the following outputs:

1. A tabular .txt file (*FifthRunOut_AllInteractions.txt*) which contains all interaction pairs with accompanying logFCs, FDRs and interaction scores. This file may be used to visualize networks in Cytoscape or with another similar tool. If more than two groups were contrasted in the DE/DA analysis, then a file for each pairwise contrast will be returned.
2. A plot of the top 100 most significant interactions, based on absolute logFCs and inter-connectivity. (*FifthRunOut_TopInteractions.pdf*). If more than two groups were contrasted in the DE/DA analysis, then a plot for each pairwise contrast will be returned. See **Figure 10**.



References

- [1] R Core Team. *R: A language and environment for statistical computing*. R Foundation for Statistical Computing, Vienna, Austria. URL <http://www.R-project.org/>. (2014).
- [2] Abd Hamid, U.M et al. *A strategy to reveal potential glycan markers from serum glycoproteins associated with breast cancer progression*. *Glycobiology* 18, (2008): 1105–1118.
- [3] Abbott, K.L. et al. *Targeted glycoproteomic identification of biomarkers for human breast carcinoma*. *Journal of Proteome Research* 7, (2008):1470–1480.
- [4] Castello, A., et al. *Insights into RNA biology from an atlas of mammalian mRNA-binding proteins*. *Cell* 149.6 (2012): 1393–1406.
- [5] Christiansen M.N et al. *Cell surface protein glycosylation in cancer*. *Proteomics* 14, (2014):525–546.
- [6] Csardi G. and Nepusz T. *The igraph software package for complex network research* InterJournal, Complex Systems 1695. 2006. <http://igraph.org>.
- [7] Day, A. *heatmap.plus: Heatmap with more sensible behavior*. R package (2012) version 1.3. <https://CRAN.R-project.org/package=heatmap.plus>

- [8] Delignette-Muller, M. L., and Dutang, C. *fitdistrplus: An R package for fitting distributions*. Journal of Statistical Software 64.4 (2015): 1–34.
- [9] Dowle, M. and Srinivasa, A. *data.table: Extension of 'data.frame'*. R package (2019) version 1.12.0. <https://CRAN.R-project.org/package=data.table>
- [10] Durinck, S. et al. *Mapping identifiers for the integration of genomic datasets with the R/Bioconductor package biomaRt*. Nature Protocols 4, 1184–1191 (2009).
- [11] Eklund, A. C. *squash: Color-Based Plots for Multivariate Visualization*. R package (2017) version 1.0.8. <https://CRAN.R-project.org/package=squash>
- [12] Friedman, J., Hastie, T., and Tibshirani, R. *Regularization paths for generalized linear models via coordinate descent*. Journal of statistical software 33.1 (2010): 1.
- [13] Fultz, N. and the StackOverflow.com community. *stackoverflow: Stack Overflow's Greatest Hits*. R package (2018) version 0.3.0. <https://CRAN.R-project.org/package=stackoverflow>
- [14] Garnier, S. *viridis: Default Color Maps from 'matplotlib'*. R package (2018) version 0.5.1. <https://CRAN.R-project.org/package=viridis>
- [15] Gromov, P. et al. *Tumor interstitial fluid—a treasure trove of cancer biomarkers*. Biochimica et Biophysica Acta (BBA)—Proteins and Proteomics 1834.11 (2013): 2259–2270.
- [16] Guo, H.B. et al. *Specific posttranslational modification regulates early events in mammary carcinoma formation*. Proceedings of the National Academy of Sciences 107, (2010): 21116–21121.
- [17] Hastie, T. et al. *impute: impute: Imputation for microarray data*. R package (2018) version 1.56.0.
- [18] Haakensen, V.D. et al. *Serum N-glycan analysis in breast cancer patients – Relation to tumour biology and clinical outcome*. Journal of Molecular Oncology, (2015):1–14.
- [19] Jeffrey T. et al. *sva: Surrogate Variable Analysis*. R package (2019) version 3.30.1.
- [20] Jensen, L. J., et al. *STRING 8 - a global view on proteins and their functional interactions in 630 organisms*. Nucleic acids research 37 (suppl-1), D412–D416, (2008).
- [21] Kammers, K., et al. *Detecting significant changes in protein abundance*. EuPA open proteomics 7 (2015): 11–19.
- [22] Kassambara, A. and Kosinski, M. *survminer: Drawing Survival Curves using 'ggplot2'*. R package (2018) version 0.4.3. <https://CRAN.R-project.org/package=survminer>

- [23] Langfelder, P. and Horvath, S. *WGCNA: an R package for weighted correlation network analysis*. BMC Bioinformatics (2008), 9:559.
- [24] Pineda, A.L., et al. *On Predicting lung cancer subtypes using ‘omic’ data from tumor and tumor-adjacent histologically-normal tissue*. BMC cancer 16.1 (2016): 184.
- [25] Potapenko, I.O. et al. *Glycan-related gene expression signatures in breast cancer subtypes; relation to survival*. Journal of Molecular Oncology 9, (2013): 861–876.
- [26] Ritchie, M.E., et al. *limma powers differential expression analyses for RNA-sequencing and microarray studies*. Nucleic Acids Research 43(7), e47.
- [27] Ru, Y. et al. *The multiMiR R package and database: integration of microRNA–target interactions along with their disease and drug associations*. Nucleic Acids Res 42 (17) (2014). e133. doi: 10.1093/nar/gku631
- [28] Saldova, R., et al. *Association of N-glycosylation with breast carcinoma and systemic features using high-resolution quantitative UPLC..* Journal of proteome research 13.5 (2014): 2314–2327.
- [29] Schroeder MS., et al. *survcomp: an R/Bioconductor package for performance assessment and comparison of survival models*. Bioinformatics 27(22): 3206–3208. (2011)
- [30] Scrucca et al. *mclust 5: clustering, classification and density estimation using Gaussian finite mixture models*. The R journal 8.1 (2016): 289.
- [31] Soneson, C., and Delorenzi, M. et al. *A comparison of methods for differential expression analysis of RNA-seq data*. BMC bioinformatics 14.1 (2013): 91.
- [32] Terkelsen, T. and Haakansen, V.D. et al. *N-glycan signatures identified in tumor interstitial fluid and serum of breast cancer patients: association with tumor biology and clinical outcome*. Molecular oncology 12.6 (2018): 972–990.
- [33] Walker, A. *openxlsx: Read, Write and Edit XLSX Files*. R package version 4.1.0. <https://CRAN.R-project.org/package=openxlsx>
- [34] H. Wickham. *ggplot2: Elegant Graphics for Data Analysis*. Springer-Verlag New York, (2016).
- [35] Wickham, H. *The Split-Apply-Combine Strategy for Data Analysis*. Journal of Statistical Software, 40(1), 1–29, (2011). URL: <http://www.jstatsoft.org/v40/i01/>.
- [36] Wickham, H. *scales: Scale Functions for Visualization*. R package (2018) version 1.0.0. <https://CRAN.R-project.org/package=scales>.
- [37] Wickham, H. et al. *devtools: Tools to Make Developing R Packages Easier*. R package version 2.0.1. (2018). <https://CRAN.R-project.org/package=devtools>