# fgczgseaora: unifying methods on gene (protein) set enrichment

European Bioconductor Meeting 2019 - Brussels

Lucas Kook and Witold Wolski (wew@fgcz.ethz.ch)

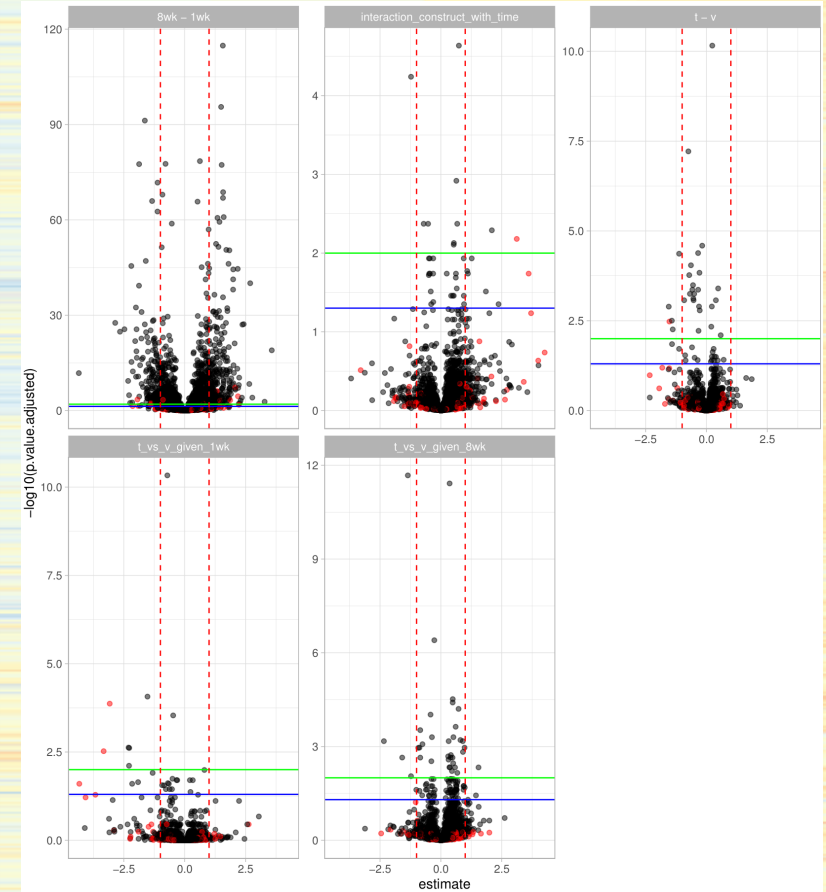[Proteome Informatics - Functional Genomics Center Zurich]

09 December, 2019

# Overview

- Pathway analysis for proteomics quantification experiments
- fgczgseaora
- Outlook

# Protein quantification experiments

- determine protein foldchanges
  for various contrasts (comparisons of treatments)
- up to thousands of proteins
- only *abundant* proteins quantifed (detection bias)

# Pathway analysis

- Over-Represenation Analysis (ORA)
- Gene Set Enrichment Analysis (GSEA)

Pathway analysis uses a priori gene sets that have been grouped together by their involvement in the same biological pathway, or by proximal location on a chromosome. Examples of gene set database are Gene Ontology (GO), KEGG, Reactome and many more.

# Over-Representation Analysis (ORA)

- Dychotomize list of proteins
  (e.g. using a *threshold* into overexpressed - Yes/No).
- Test if a geneset is *over-represented*
  in on of the sublists
  (e.g. Fischers Exact Test).
- how to choose the threshold?

```
## Pathway GO:0003091

##                    Differentially expressed
## GO Term               Yes          No
##    Contained           12           3
##    Not Contained        7          24

## p-value: 0.00034
```
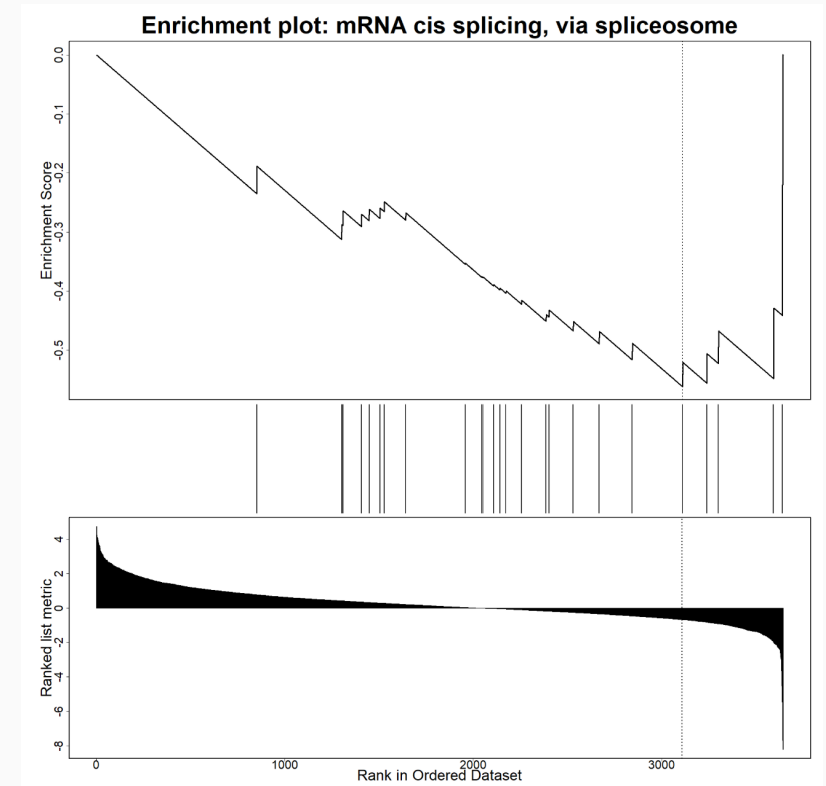
# Gene Set Enrichment Analysis (GSEA)

- Ranked list (no threshold required)
- locate genes of genesets in ranked list
- compute enrichment score



Enrichment plot: mRNA cis splicing, via spliceosome

Gene Sets can be highly correlated, because they contain the same proteins. Multiplicity adjustment assumes indpendence (FDR).

# fgczgseaora

- Easily generate reports to be delivered to biologists.
- For ORA We can only use tools which allow to specify detection background.
- Map identifiers - support for *sp* identifiers
- Ideally run packages locally
- Provide a similar **R** and command line interface to run ORA GSEA.

# Many R packages are available

## R packages for pathway analysis

| Package | Repo | Maintenance | offline | ID.Mapping | ORA | GSEA |
|---|---|---|---|---|---|---|
| WebGestaltR | CRAN | + | - | + | + | + |
| FGNet | Bioc | + | (-) | (-) | - | + |
| HTSanalyzeR | Bioc | - | (-) | - | + | + |
| sigora | CRAN | + | + | (-) | + | - |
| SetRank | CRAN | - | (-) | - | - | + |
| STRINGdb | Bioc | + | - | (-) | + | + |
| enrichR | CRAN | + | - | + | (+) | + |
| TopGO | Bioc | ... | | | | |

- We did integrate:
  - `WebgestaltR` (online only)
  - `sigORA` (offline)

WebgesaltR - Various gene set databases, id mapping, allows for downloading html results. sigORA - uses gene pair signatures. Searches background and pathways for protein pairs unique to a given pathway. By this it decreases the correlation among gene sets.

# Common R interface

```r
runWebGestaltGSEA(
    data = dd,
    fpath = "",
    ID_col = "UniprotID",
    score_col =  "estimate",
    organism =  "hsapiens",
    target = "geneontology_Biological_Process",
    nperm = 500,
    outdir = file.path(odir, "WebGestaltGSEA")
)
```

```r
runWebGestaltORA(
    data = dd,
    fpath = "",
    ID_col = "UniprotID",
    score_col =  "estimate",
    organism =  "hsapiens",
    threshold = 1,
    greater = TRUE,
    target = "geneontology_Biological_Process",
    nperm = 500,
    outdir = file.path(odir, "WebGestaltORA")
)
runSIGORA(
    data = dd,
    score_col = "estimate",
    threshold = 1,
    greater = TRUE,
    target = "GO",
    outdir = file.path(odir, "sigORA")
)
```

# Command line interface

```
Rscript lfq_multigroup_gsea.R ./foldchange_estimates.xlsx -o hsapiens
Rscript lfq_multigroup_ora.R ./foldchange_estimates.xlsx -t uniprotswissprot
```

The enrichment methods in this package (ORA, GSEA sigORA) come with a `docopt` based command line tool to facilitate analysing batches of files.

# Command line interface

```
"WebGestaltR GSEA for multigroup reports

Usage:
  lfq_multigroup_gsea.R <grp2file> [--organism <organism>] [--outdir <outdir>] [--

Options:
  -o --organism <organism> organism [default: hsapiens]
  -r --outdir <outdir> output directory [default: results_gsea]
  -t --idtype <idtype> type of id used for mapping [default: uniprotswissprot]
  -i --ID_col <ID_col> Column containing the UniprotIDs [default: UniprotID]
  -n --nperm <nperm> number of permutations to calculate enrichment scores [defau
  -e --score_col <score_col> column containing fold changes [default: pseudo_estir
  -c --contrast <contrast> column containing fold changes [default: contrast]

Arguments:
  grp2file  input file
" → doc
library(docopt)
```

# HTML outputs - Multiple Contrasts and Targets

- creates folder structure with HTML files
  visualizing the ORA and GSEA results:
  - For all contrasts
    e.g. t - v, 8wk - 1wk etc.
  - and all selected target
    e.g. GO Bioprocess, GO Molecular Function
- These files are linked from an `index.html`
- can easily be stored and delivered as part of analysis.

# Outlook

## Outlook

- Standardize R-API interface
- Standardize return values and reports.
- add one or two more packages (`edgeR`, `topGO`, ?)

## THANK YOU!

## Acknowledgments:

Paolo Nanni, Christian Panse, Ralph Schlapbach, Tobias Kockmann