

# Estimating Gene-Specific Phenotypes with **gespeR**

Fabian Schmich

April 25, 2023

## Contents

### 1 The **gespeR** Model

This package provides algorithms for deconvoluting off-target confounded phenotypes from RNA interference screens. The package uses (predicted) siRNA-to-gene target relations in a regularised linear regression model, in order to infer individual gene-specific phenotype (GSP) contributions. The observed siRNA-specific phenotypes (SSPs) for reagent  $i = 1, \dots, n$  as the weighted linear sum of GSPs of all targeted genes  $j = 1, \dots, p$

$$Y_i = x_{i1}\beta_1 + \dots + x_{ip}\beta_p + \epsilon_i, \quad (1)$$

where  $x_{ij}$  represents the strength of knockdown of reagent  $i$  on gene  $j$ ,  $\beta_j$  corresponds to the GSP of gene  $j$  and  $\epsilon_i$  is the error term for SSP  $i$ . The linear regression model is fit using elastic net regularization:

$$\hat{\beta} = \underset{\beta}{\operatorname{argmin}} \left\{ \sum_{i=1}^n \left( y_i - \sum_{j=1}^p x_{ij}\beta_j \right)^2 + \lambda \sum_{j=1}^p (\alpha\beta_j^2 + (1-\alpha)|\beta_j|) \right\}. \quad (2)$$

Here  $\lambda$  determines the amount of regularization and  $\alpha$  is the mixing parameter between the ridge and lasso penalty with  $0 \leq \alpha \leq 1$ . The elastic net penalty selects variables like the lasso and shrinks together the coefficients of correlated predictors like ridge. This allows for a sparse solution of nonzero GSPs, while retaining simultaneous selection of genes with similar RNAi reagent binding patterns in their respective 3' UTRs. For more information and for citing the **gespeR** package please refer to:

Schmich F (2023). *gespeR: Gene-Specific Phenotype Estimator*. <https://doi.org/10.18129/B9.bioc.gespeRdoi:10.18129/B9.bioc.gespeR>, R package version 1.32.0, <https://bioconductor.org/packages/gespeR>.

### 2 Working Example

In this example, we first load simulated phenotypic readout and siRNA-to-gene target relations. The toy data consists of four screens (A, B, C, D) of 1,000 siRNAs and a limited gene universe of 1,500 genes. Detailed description of how the data was simulated can be accessed using `?simData`. First, we load the package:

```
library(gespeR)

## Warning: replacing previous import 'utils::findMatches' by 'S4Vectors::findMatches' when
loading 'AnnotationDbi'
```

Now the phenotypes and target relations can be initialised using the `Phenotypes` and `TargetRelations` commands. First, we load the four phenotypes vectors:

```

phenos <- lapply(LETTERS[1:4], function(x) {
  sprintf("Phenotypes_screen_%s.txt", x)
})
phenos <- lapply(phenos, function(x) {
  Phenotypes(system.file("extdata", x, package="gespeR"),
    type = "SSP",
    col.id = 1,
    col.score = 2)
})
show(phenos[[1]])

## 1000 SSP Phenotypes

## Warning: 'tbl_df()' was deprecated in dplyr 1.0.0.
## i Please use 'tibble::as_tibble()' instead.
## i The deprecated feature was likely used in the gespeR package.
## Please report the issue to the authors.
## This warning is displayed once every 8 hours.
## Call 'lifecycle::last_lifecycle_warnings()' to see where this warning was
## generated.

## # A tibble: 1,000 x 2
##   ID          Scores
##   <chr>         <dbl>
## 1 siRNAID_0001 -0.930
## 2 siRNAID_0002 -1.13
## 3 siRNAID_0003 -1.05
## 4 siRNAID_0004  0.808
## 5 siRNAID_0005 -1.42
## 6 siRNAID_0006  1.64
## 7 siRNAID_0007 -0.157
## 8 siRNAID_0008  0.748
## 9 siRNAID_0009 -0.959
## 10 siRNAID_0010 -0.0440
## # i 990 more rows

```

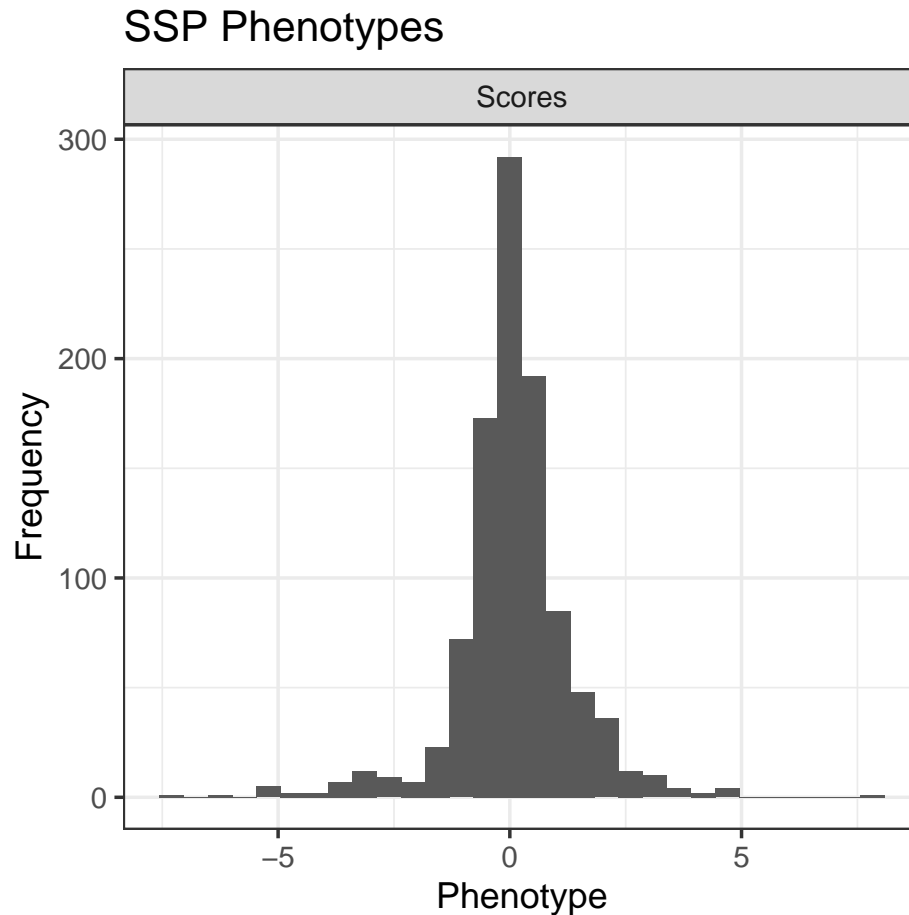
A visual representation of the phenotypes can be obtained with the `plot` method:

```

plot(phenos[[1]])

## 'stat_bin()' using 'bins = 30'. Pick better value with 'binwidth'.

```



Now, we load the target relations for all four screens using the constructor of the `TargetRelations` class:

```
tr <- lapply(LETTERS[1:4], function(x) {
  sprintf("TR_screen_%s.rds", x)
})
tr <- lapply(tr, function(x) {
  TargetRelations(system.file("extdata", x, package="gesperR"))
})
show(tr[[2]])
```

```
## 1000 x 1500 siRNA-to-gene relations.
## 10 x 5 sparse Matrix of class "dgCMatrix"
##               colnames
## rownames      geneID_0001 geneID_0002 geneID_0003 geneID_0004 geneID_0005
## siRNAID_0001      .      .      .      .      .
## siRNAID_0002      .      .      .      .      .
## siRNAID_0003      .  0.4385757      .      .      .
## siRNAID_0004      .      .      .      .      .
## siRNAID_0005      .      .      .      .      .
## siRNAID_0006      .      .      .      .      .
## siRNAID_0007      .      .      .      .      .
## siRNAID_0008      .      .      .      .      .
## siRNAID_0009      .      .      .      .      .
## siRNAID_0010      .      .      .      .      .
## ...
```

For large data sets, e.g. genome-wide screens, target relations objects can become very big and the user may not want to keep all values in the RAM. For this purpose, we can use the `unloadValues` method. In this example, we write the values to a temp-file, i.e. not the original source file, which may be required, when we do not want to overwrite existing data, after, for instance, subsetting the target relations object.

```
# Size of object with loaded values
format(object.size(tr[[1]]), units = "Kb")

## [1] "717.7 Kb"

tempfile <- paste(tempfile(pattern = "file", tmpdir = tmpdir()), ".rds", sep="")
tr[[1]] <- unloadValues(tr[[1]], writeValues = TRUE, path = tempfile)

# Size of object after unloading
format(object.size(tr[[1]]), units = "Kb")

## [1] "178.4 Kb"

# Reload values
tr[[1]] <- loadValues(tr[[1]])
```

In order to obtain deconvoluted gene-specific phenotypes (GSPs), we fit four models on the four separate data sets using cross validation by setting `mode = "cv"`. We set the elastic net mixing parameter to 0.5 and use only one core in this example:

```
res.cv <- lapply(1:length(phenos), function(i) {
  gesperR(phenotypes = phenos[[i]],
    target.relations = tr[[i]],
    mode = "cv",
    alpha = 0.5,
    ncores = 1)
})
```

The `ssp` and `gsp` methods can be used to obtain SSP and GSP scores from a `gesper` object:

```
ssp(res.cv[[1]])

## 1000 SSP Phenotypes
##
## # A tibble: 1,000 x 2
##   ID          Scores
##   <chr>         <dbl>
## 1 siRNAID_0001 -0.930
## 2 siRNAID_0002 -1.13
## 3 siRNAID_0003 -1.05
## 4 siRNAID_0004  0.808
## 5 siRNAID_0005 -1.42
## 6 siRNAID_0006  1.64
## 7 siRNAID_0007 -0.157
## 8 siRNAID_0008  0.748
## 9 siRNAID_0009 -0.959
## 10 siRNAID_0010 -0.0440
## # i 990 more rows

gsp(res.cv[[1]])
```

```
## 1500 GSP Phenotypes
##
## # A tibble: 1,500 x 2
##   ID          Scores
##   <chr>        <dbl>
## 1 geneID_0001      NA
## 2 geneID_0002      NA
## 3 geneID_0003      NA
## 4 geneID_0004      NA
## 5 geneID_0005      NA
## 6 geneID_0006      NA
## 7 geneID_0007      NA
## 8 geneID_0008      NA
## 9 geneID_0009      NA
## 10 geneID_0010     NA
## # i 1,490 more rows

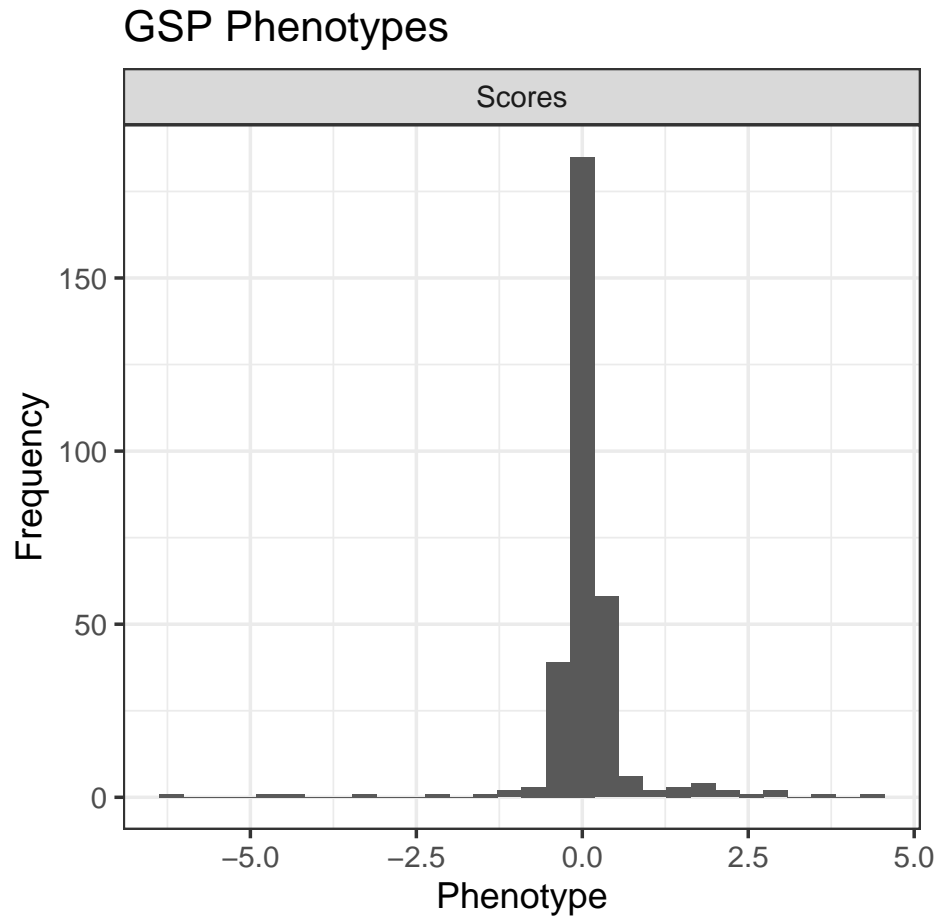
head(scores(res.cv[[1]]))

## # A tibble: 6 x 2
##   ID          Scores
##   <chr>        <dbl>
## 1 geneID_0001      NA
## 2 geneID_0002      NA
## 3 geneID_0003      NA
## 4 geneID_0004      NA
## 5 geneID_0005      NA
## 6 geneID_0006      NA
```

The fitted models can also be visualised using the `plot` method:

```
plot(res.cv[[1]])

## 'stat_bin()' using 'bins = 30'. Pick better value with 'binwidth'.
## Warning: Removed 1185 rows containing non-finite values ('stat_bin()').
```

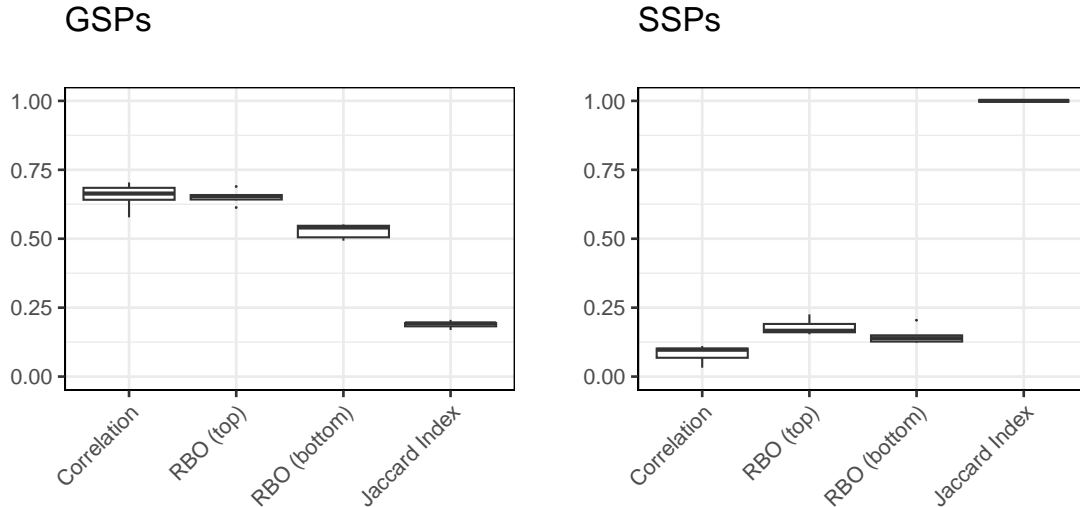


The `concordance` method can be used to compute the concordance between ranked lists of phenotypes. Here we compute concordance between all pairs of GSPs, as well as between all pairs of SSPs, from all four data sets:

```
conc.gsp <- concordance(lapply(res.cv, gsp))
conc.ssp <- concordance(lapply(res.cv, ssp))
```

We can visualise the `concordance` objects using the `plot` method:

```
plot(conc.gsp) + ggtitle("GSPs\n")
plot(conc.ssp) + ggtitle("SSPs\n")
```



### 3 sessionInfo()

- R version 4.3.0 RC (2023-04-13 r84269 ucrt), x86\_64-w64-mingw32
- Locale: LC\_COLLATE=C, LC\_CTYPE=English\_United States.utf8, LC\_MONETARY=English\_United States.utf8, LC\_NUMERIC=C, LC\_TIME=English\_United States.utf8
- Time zone: America/New\_York
- TZcode source: internal
- Running under: Windows Server 2022 x64 (build 20348)
- Matrix products: default
- Base packages: base, datasets, grDevices, graphics, methods, stats, utils
- Other packages: gesper 1.32.0, ggplot2 3.4.2, knitr 1.42
- Loaded via a namespace (and not attached): AnnotationDbi 1.62.0, Biobase 2.60.0, BiocFileCache 2.8.0, BiocGenerics 0.46.0, BiocManager 1.30.20, Biostrings 2.68.0, Category 2.66.0, DBI 1.1.3, GSEABase 1.62.0, GenomeInfoDb 1.36.0, GenomeInfoDbData 1.2.10, IRanges 2.34.0, KEGGREST 1.40.0, Matrix 1.5-4, MatrixGenerics 1.12.0, R6 2.5.1, RBGL 1.76.0, RColorBrewer 1.1-3, RCurl 1.98-1.12, RSQLite 2.3.1, Rcpp 1.0.10, S4Vectors 0.38.0, XML 3.99-0.14, XVector 0.40.0, affy 1.78.0, affyio 1.70.0, annotate 1.78.0, biomaRt 2.56.0, bit 4.0.5, bit64 4.0.5, bitops 1.0-7, blob 1.2.4, cachem 1.0.7, cellHTS2 2.64.0, cli 3.6.1, codetools 0.2-19, colorspace 2.1-0, compiler 4.3.0, crayon 1.5.2, curl 5.0.0, dbplyr 2.3.2, digest 0.6.31, doParallel 1.0.17, dplyr 1.1.2, evaluate 0.20, fansi 1.0.4, farver 2.1.1, fastmap 1.1.1, filelock 1.0.2, foreach 1.5.2, genefilter 1.82.0, generics 0.1.3, glmnet 4.1-7, glue 1.6.2, graph 1.78.0, grid 4.3.0, gtable 0.3.3, highr 0.10, hms 1.1.3, httr 1.4.5, hwriter 1.3.2.1, iterators 1.0.14, labeling 0.4.2, lattice 0.21-8, lifecycle 1.0.3, limma 3.56.0, locfit 1.5-9.7, magrittr 2.0.3, matrixStats 0.63.0, memoise 2.0.1, munsell 0.5.0, parallel 4.3.0, pillar 1.9.0, pkgconfig 2.0.3, plyr 1.8.8, png 0.1-8, preprocessCore 1.62.0, prettyunits 1.1.1, progress 1.2.2, rappdirs 0.3.3, reshape2 1.4.4, rlang 1.1.0, scales 1.2.1, shape 1.4.6, splines 4.3.0, splots 1.66.0, stats4 4.3.0, stringi 1.7.12, stringr 1.5.0, survival 3.5-5, tibble 3.2.1, tidyselect 1.2.0, tools 4.3.0, utf8 1.2.3, vctrs 0.6.2, vsn 3.68.0, withr 2.5.0, xfun 0.39, xml2 1.3.3, xtable 1.8-4, zlibbioc 1.46.0