

GRridge: Adaptive group-regularized ridge regression by use of co-data

Mark A. van de Wiel and Putri W. Novianti

April 30, 2018

Department of Epidemiology & Biostatistics
VU University Medical Center
Amsterdam, The Netherlands

`mark.vdwiel@vumc.nl`

Contents

1	Overview	2
1.1	Elements of the GRridge package	2
1.2	Getting started	3
2	Example 1: DNA-Methylation data	4
3	Example 2: mRNA sequencing data	7
3.1	A partition containing overlapping groups	8
3.2	Merge groups in a partition	8
3.3	A procedure to select and to order multiple partitions	9

1 Overview

Predicting binary or continuous response from high-dimensional data is a well-addressed problem nowadays. Many existing methods have been adapted to cope with high-dimensional data, in particular by means of regularization. Adaptive group regularized ridge regression was introduced to improve the predictive performance of logistic ridge regression by incorporating external and/or internal auxiliary information on the features: the co-data. More formally, co-data can be described as any nominal or quantitative feature-specific information, obtained independently of the response variable. Four types of co-data are distinguished:

1. Response-independent summaries in the primary data (e.g. standard deviation).
2. Feature-specific summaries from an independent study (e.g. p-values).
3. Feature groupings from public data bases (e.g. pathways).

GRridge package implements adaptive group-regularized (survival, linear, logistic) ridge regression by use of co-data. The package includes convenience functions to convert such co-data to the correct input format. In addition, it includes functions for evaluating the predictive performance.

GRridge package was based on these following publications:

Mark van de Wiel, Tonje Lien, Wina Verlaat, Wessel van Wieringen, Saskia Wilting. (2016). Better prediction by use of co-data: adaptive group-regularized ridge regression. *Statistics in Medicine*, 35(3), 368-81. [[Wiel et al., 2016](#)]

Putri W. Novianti, Barbara C. Snoek, Saskia Wilting, Mark van de Wiel. (2017). Better diagnostic signatures from RNAseq data through use of auxiliary co-data. *Bioinformatics*. 33, 1572-1574. [[Novianti et al., 2017](#)]

1.1 Elements of the GRridge package

Key elements of the **GRridge** package are:

1. An automatic function to create a partition of features, namely **CreatePartition** and **matchGeneSets** function for non-overlapping and overlapping groups, respectively. The package also provides a function to regroup a considerably number of overlapping groups by **mergeGroups** function.
2. Comparison of the performance of **GRridge** model with, ordinary ridge regression, lasso and non-penalized regression, using the **grridge** function:
 - ridge regression: the function automatically estimated this ridge regression model.

- lasso: `comparelasso=TRUE`.
 - non-penalized regression: `compareunpenal=TRUE`.
3. The `grridgeCV` function provides predicted classes and predicted probabilities that are estimated by cross-validation.
 4. Post-hoc feature selection of X-relevant features. Two features selection approaches are available, namely forward selection ("`selectionForward=TRUE`", "`maxsel=X`" in the `grridge` function) and feature selection via L1-penalization method ("`selectionEN=TRUE`", in the `grridge` function). The forward selection may perform reasonably well when the group-penalties are strong, in particular for multiple partitions. The second feature selection approach, on the other hand, is likely superior when group-penalties are less strong, because it inherits the superior selection properties of the elastic-net in such a setting. As it will perform similar as forward selection selection approach, we strongly suggest to use the group-regularized elastic-net for post-hoc feature selection ("`selectionForward=FALSE`", "`selectionEN=TRUE`").
 5. Evaluation of the performance of classification models is visualized by receiver operating characteristics (ROC) curve ("`roc`" function) and is quantified by area under the curve ("`auc`" function).

1.2 Getting started

The GRridge package depends on these following R packages: `penalized` [Goeman, 2010], `glmnet` [?], `survival` (Therneau, 2015) and `Iso` (Turner, 2015).

2 Example 1: DNA-Methylation data

A cervical cancer study measures DNA methylation level on normal healthy controls (control, n=20) and high-grade precursor lesions (precursor, n=17) tissue biopsies [Farkas et al., 2013]. A popular platform for measuring methylation is the Infinium HumanMethylation450 BeadChip (Illumina, San Diego, CA, USA), which contains 450,000 probes per individual, where each probe renders a so-called beta-value. The preprocessing process rendered 40,000 methylation probes [Wiel et al., 2016].

Load the GRridge library and its dependencies

```
> library(GRridge)
```

Load the primary data set

```
> data(dataFarkas)
```

It contains these following objects

- datcenFarkas: methylation data for cervix samples (arcsine-transformed beta values).
- respFarkas: binary response (Normal and Precursor).
- CpGannFarkas: annotation of probes according to location (CpG-Island, North-Shelf, South-Shelf, North-Shore, South-Shore, Distant).

We first create a partition based on location of the probes (CpGannFarkas). For nominal input (factor), `CreatePartition(vec)` creates a partition of features (probes) with groups according to the levels of `vec`.

```
> firstPartition <- CreatePartition(CpGannFarkas)
```

```
[1] "Summary of group sizes:"  
Distant  Island  N_Shelf  N_Shore  S_Shelf  S_Shore  
 14047   12858    2006    5262    1765    4062
```

A practical issue when applying penalized regression is the need or 'no need' for standardization of the features. A potential of GRridge method is that it can let the data decide how the variance of features should impact the penalties. More discussion about this issue can be found in [Wiel et al., 2016]. A partition based on standard deviation (sd) of each feature is then created. For numeric input it creates a partition according to ranking, here into uniformly-sized groups based on sds. The argument `decreasing=FALSE` implies here that groups of probes with smaller sds may potentially be penalized less when using the monotone argument below in the `grridge` function (which implicitly also happens when standardizing the data).

```

> sdsF <- apply(datcenFarkas,1,sd)
> secondPartition <- CreatePartition(sdsF,decreasing=FALSE,
+                                 uniform=TRUE,grsize=5000)

[1] "Group size 5000"
[1] "Sorting vec in increasing order, assuming small values are MORE relevant"
[1] "Summary of group sizes:"
  group1 group2 group3 group4 group5 group6 group7 group8 group9 group10
    5000   5000   5000   5000   5000   5000   5000   5000   5000   5000

```

Concatenate two partitions into one list

```

> partitionsFarkas <- list(cpg=firstPartition,sds=secondPartition)

```

A list of monotone functions from the corresponding partition.

```

> monotoneFarkas <- c(FALSE,TRUE)

```

monotoneFarkas indicates that monotone increasing group-penalties are desired for the 2nd partition (sd-based), and not for the first one.

`grridge()` function applies group-regularized ridge to data `datcenFarkas`, response `respFarkas` and probe grouping `partitionFarkas`. It recognizes automatically whether survival, linear or logistic (here) regression should be performed. Here, it saves the prediction objects from ordinary and group-regularized ridge. Includes non-penalized intercept by default.

```

> grFarkas <- grridge(datcenFarkas, respFarkas, partitionsFarkas,
+                   optl=5.680087, monotone= monotoneFarkas)

```

The group-penalties from the GRridge model (`grFarkas`) is shown as follow

```

> grFarkas$lambdaults

```

`grridgeCV()` function performs 10-fold cross-validation to assess predictive performances of the predictors saved in the `grFarkas` object. Invokes `grridge()` using the same arguments as used by the above call to `grridge()` to create `grFarkas`. The result is a matrix with 3 columns containing the true response, and the predictions by ordinary and group-regularized logistic ridge.

```

> grFarkasCV <- grridgeCV(grFarkas, datcenFarkas,
+                        respFarkas, outerfold=10)

```

The performance of probabilistic classifiers is visualized by ROC curves and is measured by AUC.

```

> cutoffs <- rev(seq(0,1,by=0.01))
> rocridgeF <- roc(probs=grFarkasCV[,2],

```

```
+           true=grFarkasCV[,1], cutoffs=cutoffs)
> rocrridgeF <- roc(probs=grFarkasCV[,3],
+           true=grFarkasCV[,1], cutoffs=cutoffs)
> plot(rocridgeF[1,], rocridgeF[2,], type="l", lty=1, ann=FALSE, col="grey")
> points(rocrridgeF[1,], rocrridgeF[2,], type="l", lty=1, col="black")
> legend(0.6, 0.3, legend=c("ridge", "GRridge"),
+       lty=c(1,1), lwd=c(1,1), col=c("grey", "black"))
>
```

3 Example 2: mRNA sequencing data

Blood platelets extracted from patients with breast cancer (breast, n=40) and colorectal cancer (CRC, n=41) were used to profile their RNA markers for the purpose of early cancer detection [Best et al., 2015]. The raw sequencing data set is publicly available in GEO database (GEO: GSE68086). The raw data was preprocessed, rendering 18,410 transcripts (or features).

Load the GRridge library and its dependencies

```
> library(GRridge)
```

Load the primary data set

```
> data(dataWurdinger)
```

The object contains

- `datWurdingerBC` : A matrix containing preprocessed mRNA sequencing data (quasi-gaussian scale, normalized). Columns are samples (81 samples with Breast Cancer and Colorectal Cancer) and rows are features (18410 features).
- `respWurdinger` : A factor containing responses for samples with Breast cancer (n=40) and colorectal cancer (n=41)
- `annotationWurdinger` : A list containing `ensembleID`, `geneSymbol`, `entrezID` and chromosome location.
- `coDataWurdinger` : A list containing co-data (i) sets from immunologic signature pathway and (ii) platelets expressed genes

In this second example, we focus on the application of GRridge method by using multiple external co-data, namely (1) immunologic signature pathway (2) transcription factor based pathway (3) platelet expressed genes and (4) genomic annotation based on chromosomal location. The first two co-data are based on the gene set enrichment analysis (GSEA) from the Molecular Signatures Database (MSigDB). We created a list of platelet-expressed genes is based on the joint lists of two studies, i.e. [Gnatenko et al., 2009] and [Bugert et al., 2003]. The last partition was based on chromosomal location taken from biomaRt databases [Durinck et al., 2009].

Here, we focus on the binary classification case between breast cancer (breast) and colorectal cancer (CRC).

First, the preprocessed primary data set were transformed and standardized

```
> # Transform the data set to the square root scale
> dataSqrtWurdinger <- sqrt(datWurdinger_BC)
> #
> #Standardize the transformed data
```

```

> datStdWurdinger <- t(apply(dataSqrtWurdinger,1,
+                           function(x){(x-mean(x))/sd(x)}))
> #
> # A list of gene names in the primary RNAseq data
> genesWurdinger <- as.character(annotationWurdinger$geneSymbol)

```

3.1 A partition containing overlapping groups

We first show an example of GRridge classification model by using overlapping groups, i.e. pathway-based grouping. Transcription factor based pathway was extracted from the MSigDB (Section C3: motif gene sets; subsection: transcription factor targets; file's name: "c3.tft.v5.0.symbols.gmt"). The gene sets are based on TRANSFAC version 7.5 database (<http://www.gene-regulation.com/>).

A partition based on the GSEA object (TFsym) is then created. Some features may belong to more than one group. The argument `minlen=25` implies the minimum number of features in a gene set. If `remain=TRUE`, gene sets with less than 25 members are grouped to the "remainder" group. "genesWurdinger" is an object containing gene names from the mRNA sequencing data set. See `help(matchGeneSets)` for more detail information. The TFsym can be downloaded from: <https://github.com/markvdwiel/GRridgeCodata/tree/master/Transcription-factor-binding-site-pathway>

```

> gseTF <- matchGeneSets(genesWurdinger,TFsym,minlen=25,remain=TRUE)

```

The output value of the `matchGeneSets` function can be used directly as an input in the `grridge` function (`partitions=gseTF`). There is no need to create a partition via the `CreatePartition` function. A similar approach can be done for other pathways based partition.

3.2 Merge groups in a partition

Pathway-based partition often contains a considerable number of gene sets (or groups). There are 615 and 4871 groups in the transcription factor and immunological based pathway, respectively (per July 4, 2016). Overfitting may be an issue in the GRridge predictive modeling on such a large number of groups. As a solution, a data driven re-grouping based on hierarchical clustering analysis can be applied. The GRridge package provides a function to merge groups in a partition, i.e. `mergeGroups`. In this example the initial gene sets will be re-grouped into 5 groups (`maxGroups=5`).

```

> gseTF_newGroups <- mergeGroups(highdimdata= datStdWurdinger,
+                               initGroups=gseTF, maxGroups=5)

```

To extract indices of new groups,

```

> gseTF2 <- gseTF_newGroups$newGroups

```


The `gseTF2` object is a list of the components of which contain the indices of the features belonging to each group. This object is in the same format as the output from the `CreatePartition` function. Hence, the result can be used directly as an input in the `grridge` function.

To observe members of the new groups,

```
> newClustMembers <- gseTF_newGroups$newGroupMembers
```

3.3 A procedure to select and to order multiple partitions

Although there is no harm including multiple co-data sets, we have shown in the miRNAseq data that the more co-data used does not guarantee the better predictive performance of a GRridge model. We introduce a procedure to optimize the use of co-data in a GRridge model. A co-data set will be included in the predictive model, if it gives a significant improvement to the model. The contribution of a co-data is evaluated by cross-validation likelihood (cvl) value. The partitions selection procedure is similar with the forward selection in a classical regression model. The `grridge` function gives an option to optimize the use of co-data. This selection reassures each and every set gives additive positive affect to the predictive model.

We apply the partitions selection procedure on the three available co-data, described as follows

1. co-data 1: a partition based on immunologic signature pathway This following object was obtained by the same approach as transcription factor based pathway mentioned in the previous sections. We merged the initial gene sets (groups) into five new groups following the procedure mentioned in section 3.2.

```
> immunPathway <- coDataWurdinger$immunologicPathway  
> parImmun <- immunPathway$newClust
```

2. co-data 2: a partition based on a list of platelets expressed genes.

```
> plateletsExprGenes <- coDataWurdinger$plateletgenes  
> # Group genes in the primary data based on the list  
> # The genes are grouped into  
> # either "NormalGenes" or "Non-overlapGenes"  
> is <- intersect(plateletsExprGenes, genesWurdinger)  
> im <- match(is, genesWurdinger)  
> plateletsGenes <- replicate(length(genesWurdinger), "Non-overlapGenes")  
> plateletsGenes[im] <- "NormalGenes"  
> plateletsGenes <- as.factor(plateletsGenes)  
> parPlateletGenes <- CreatePartition(plateletsGenes)
```

3. co-data 3: a partition based on chromosomal location. A list of chromosomal location based on biomaRt data bases.

```

> ChromosomeWur0 <- as.vector(annotationWurdinger$chromosome)
> ChromosomeWur <- ChromosomeWur0
> idC <- which(ChromosomeWur0=="MT" | ChromosomeWur0=="notBiomart" |
+             ChromosomeWur0=="Un")
> ChromosomeWur[idC] <- "notMapped"
> table(ChromosomeWur)
> parChromosome <- CreatePartition(as.factor(ChromosomeWur))

```

Concatenate all partitions into one list.

```

> partitionsWurdinger <- list(immunPathway=parImmun,
+                             plateletsGenes=parPlateletGenes,
+                             chromosome=parChromosome)

```

A list of monotone functions from the corresponding partitions,

```

> monotoneWurdinger <- c(FALSE, FALSE, FALSE)

```

Start selecting and ordering partitions.

```

> optPartitions <- PartitionsSelection(datStdWurdinger, respWurdinger,
+                                     partitions=partitionsWurdinger,
+                                     monotoneFunctions=monotoneWurdinger,
+                                     optl=160.527)

```

The output of the `PartitionsSelection` function is a numeric vector containing the order of the selected partition(s). We may plug-in this object to the `ord` argument in the `grridge` function. To reduce the computational time, we may use the optimum `lambda2` to the GRridge predictive modeling, resulted from the `optPartitions`.

As comparison, lasso model is also built (`comparelasso=TRUE`) and post-hoc feature selection via L1 penalization method is performed (`selectionEN=TRUE`) by predetermined the number of selected markers (`maxsel=10`).

```

> # To reduce the computational time, we may use the optimum lambda2
> # (global lambda penalty) in the GRridge predictive modeling,
> # optl=optPartitions$optl
> # GRridge model by incorporating the selected partitions
> partitionsWurdinger_update = partitionsWurdinger[optPartitions$ordPar]
> monotoneWurdinger_update = monotoneWurdinger[optPartitions$ordPar]
> grWurdinger <- grridge(datStdWurdinger, respWurdinger,
+                       partitions=partitionsWurdinger_update,
+                       monotone= monotoneWurdinger_update,
+                       innfold = 3, comparelasso=TRUE,
+                       optl=optPartitions$optl, selectionEN=TRUE,
+                       maxsel=10)

```

`grWurdinger$resEN$whichEN` contains indexes (and the features' name) of the selected features based on the group- weighted elastic net and `grWurdinger$resEN$betaEN` has the information of the beta value of the corresponding selected feature.

Asses the performance of the GRridge model by performing 10-fold CV

```
> grWurdingerCV <- grridgeCV(grWurdinger, datStdWurdinger,  
+                             respWurdinger, outerfold=10)
```

The performance of probabilistic classifiers is visualized by ROC curves and is measured by AUCs

```
> cutoffs <- rev(seq(0,1,by=0.01))  
> rocridge <- roc(probs= grWurdingerCV[,2],  
+                 true= grWurdingerCV[,1],cutoffs)  
> rocGRridge <- roc(probs= grWurdingerCV[,3],  
+                  true= grWurdingerCV[,1],cutoffs)  
> rocLasso <- roc(probs= grWurdingerCV[,4],  
+                 true= grWurdingerCV[,1],cutoffs)  
> rocGRridgeEN <- roc(probs= grWurdingerCV[,5],  
+                     true= grWurdingerCV[,1],cutoffs)  
> plot(rocridge[1,],rocridge[2,],type="l",lty=2,ann=FALSE,col="grey")  
> points(rocGRridge[1,],rocGRridge[2,],type="l",lty=1,col="black")  
> points(rocLasso[1,],rocLasso[2,],type="l",lty=1,col="blue")  
> points(rocGRridgeEN[1,],rocGRridgeEN[2,],type="l",lty=1,col="green")  
> legend(0.6,0.35, legend=c("ridge","GRridge", "lasso","GRridge+varsel"),  
+        lty=c(1,1), lwd=c(1,1),col=c("grey","black","blue","green"))
```

References

- Best, M. G., Sol, N., Kooi, I., Tannous, J., Westerman, B. A., Rustenburg, F., Schellen, P., Verschueren, H., Post, E., Koster, J., et al. (2015). Rna-seq of tumor-educated platelets enables blood-based pan-cancer, multiclass, and molecular pathway cancer diagnostics. *Cancer Cell*, 28(5):666–676.
- Bugert, P., Dugrillon, A., Günaydin, A., Eichler, H., and Klüter, H. (2003). Messenger rna profiling of human platelets by microarray hybridization. *Thrombosis and haemostasis*, 90(4):738–748.
- Durinck, S., Spellman, P. T., Birney, E., and Huber, W. (2009). Mapping identifiers for the integration of genomic datasets with the r/bioconductor package biomart. *Nature protocols*, 4(8):1184–1191.
- Farkas, S. A., Milutin-Gašperov, N., Grce, M., and Nilsson, T. K. (2013). Genome-wide dna methylation assay reveals novel candidate biomarker genes in cervical cancer. *Epigenetics*, 8(11):1213–1225.
- Gnatenko, D. V., Dunn, J. J., Schwedes, J., and Bahou, W. F. (2009). Transcript profiling of human platelets using microarray and serial analysis of gene expression (sage). *DNA and RNA Profiling in Human Blood: Methods and Protocols*, pages 245–272.
- Goeman, J. J. (2010). L1 penalized estimation in the cox proportional hazards model. *Biometrical journal*, 52(1):70–84.
- Novianti, P., Snoek, B., Wilting, S., and van de Wiel, M. (2017). Better diagnostic signatures from rnaseq data through use of auxiliary co-data. *Bioinformatics*.
- Wiel, M. A., Lien, T. G., Verlaat, W., Wieringen, W. N., and Wilting, S. M. (2016). Better prediction by use of co-data: adaptive group-regularized ridge regression. *Statistics in Medicine*, 35(3):368–381.