

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

Mark A. van de Wiel^{1,2}, Tonje G. Lien³, Wina Verlaet⁴,
Wessel N. van Wieringen^{1,2}, Saskia M. Wilting⁴

1. Department of Epidemiology and Biostatistics, VU University Medical Center, PO Box 7057, 1007 MB Amsterdam, The Netherlands
2. Department of Mathematics, VU University, Amsterdam, The Netherlands
3. Department of Mathematics, University of Oslo, Oslo, Norway
4. Department of Pathology, VU University Medical Center, Amsterdam, The Netherlands

Keywords: Classification, logistic ridge regression, empirical Bayes, Random Forest, feature selection, methylation

Supplementary Material is available from: www.few.vu.nl/~mavdwiel/grridge.html

Abstract

For many high-dimensional studies, additional information on the variables, like (genomic) annotation or external p-values, is available. In the context of binary and continuous prediction, we develop a method for adaptive group-regularized (logistic) ridge regression, which makes structural use of such ‘co-data’. Here, ‘groups’ refer to a partition of the variables according to the co-data. We derive empirical Bayes estimates of group-specific penalties, which possess several nice properties: i) they are analytical; ii) they adapt to the informativeness of the co-data for the data at hand; iii) only one global penalty parameter requires tuning by cross-validation. In addition, the method allows use of *multiple* types of co-data at little extra computational effort.

We show that the group-specific penalties may lead to a larger distinction between ‘near-zero’ and relatively large regression parameters, which facilitates post-hoc variable selection. The method, termed **GRidge**, is implemented in an easy-to-use R-package. It is demonstrated on two cancer genomics studies, which both concern the discrimination of precancerous cervical lesions from normal cervix tissues using methylation microarray data. For both examples, **GRidge** clearly improves the predictive performances of ordinary logistic ridge regression and the group lasso. In addition, we show that for the second study the relatively good predictive performance is maintained when selecting only 42 variables.

1 Introduction

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 and new ones, e.g. based on feature extraction, have been devised (Hastie et al., 2008). These methods have in common that the input is a response vector

of length n and a numerical $n \times p$ design matrix, where n is the number of independent samples and $p > n$ is the number of variables. Then, the predictor is usually learned solely from this input, possibly in combination with or followed by variable selection.

Co-data comprises of all information on the measured variables other than their numerical values for the given study. A few examples in the context of genomics are: a) Data or summaries like p -values from an external study with a related objective on the same set of variables (or highly overlapping); b) Database information that summarizes the (a priori) importance of genes for a class of diseases, e.g. the Cancer Gene census (Futreal et al., 2004); c) Genomic annotation, e.g. the chromosome on which a gene is located. Co-data of type a), also referred to as ‘historical data’, has been demonstrated to potentially benefit the analysis of a given clinical trial, in particular when sample size n is small (Neuenschwander et al., 2010). For such low-dimensional data, assigning weight(s) to the co-data, e.g. by choice of the prior in a Bayesian setting, is a difficult issue, because it usually implies a subjective setting. In a high-dimensional setting like ours, however, we show that one can use empirical Bayes principles to let the data decide how informative the co-data should be.

The empirical Bayes approach sets our approach apart from other ones that use co-data to improve prediction or variable selection, like the group-lasso (Meier et al., 2008), a general multi-penalty approach (Tai and Pan, 2007) or a weighted lasso approach (Bergersen et al., 2011). In addition, unlike those methods, our approach is able to handle co-data of many different types: the external information on the variables can be binary, nominal, ordinal or continuous plus it can manage *multiple* sources of co-data iteratively.

We focus mostly on logistic ridge regression to present our approach, but also demonstrate the generality of the approach by an extension to random forest classification. We start out by recapitulating logistic ridge regression and the first two moments of the parameter estimates. These are then used to derive an empirical Bayes estimate for group-specific penalties. Next, we present a more stable iterative alternative, and also address iteration on multiple partitions of the variables. If the co-data is available as a continuous summary like a vector of p -values, we argue that one may use rank-based small groups of variables in combination with enforced monotony for the group-specific penalties.

A consequence of the use of group-specific penalties is that it can facilitate *a posteriori* variable selection. We show that effective group-regularization may result in a relatively heavy-tailed empirical distribution of the regression parameter estimates. This, as we illustrate by an example, may allow selection of a fairly sparse model with hardly any loss of predictive accuracy.

The approach is demonstrated on two cancer genomics examples. Both examples concern discriminating precancerous cervical lesions from normal cervix tissues using methylation microarray data. For the first data set, we first demonstrate that our method can automatically account for different standard deviations across variables. Next, we show that the use of annotation of the methylation probes (which are the regression variables) for group-regularization improves the prediction for 86% of the samples. The second example concerns a diagnostic setting using methylation profiles from self-collected cervico-vaginal lavages (self samples). The resulting samples are likely to be impure, which presents a challenge for discriminating the two classes. Here, we show that use of the p -values from the first study, which concerns more pure samples, as a basis for group-regularization in the second study, increases the area-under-the-ROC-curve from 67% to 74%. In addition, applying variable selection on the basis of the parameter estimates of the group-regularized approach rendered an equally accurate model with only 42 variables. Together with simulated data sets, the two examples were also used to compare our method with existing ones.

We conclude with remarks on i) conceptual differences between our approach and related methods; ii) possible extensions of our method; and iii) the corresponding R-package `GRridge` and its computational efficiency.

1.1 Logistic ridge regression

It is well known that classical ridge regression corresponds to Bayesian ridge regression: the maximum a posteriori estimate for regression parameters $\boldsymbol{\beta} = (\beta_1, \dots, \beta_p)$ corresponds to the classical estimate $\hat{\boldsymbol{\beta}}$ when using a central Gaussian prior for β_k with a variance $\tau^2 \propto 1/\lambda$, where λ is the penalty parameter in the classical ridge setting. We explore this fact to develop an empirical Bayes estimate of group-specific penalties. We explain the procedure for logistic ridge regression; the changes needed for linear ridge regression are detailed in the Supplementary Material. The results of ordinary logistic ridge regression (hence ignoring the groups) at a given value of global penalty parameter λ (e.g. obtained by cross-validation) are used as a starting point.

We first recapitulate some results for logistic ridge regression. For independent responses $Y_i \in \{0, 1\}$, $i = 1, \dots, n$, we have

$$Y_i \sim \text{Bernoulli}[\text{expit}(X_i\boldsymbol{\beta})],$$

where $X = (X_1^T, \dots, X_n^T)^T$ is the $n \times p$ design matrix and $\text{expit}(X_i\boldsymbol{\beta}) = \exp(X_i\boldsymbol{\beta}) / (1 + \exp(X_i\boldsymbol{\beta}))$. The estimate $\hat{\boldsymbol{\beta}}$ maximizes the penalized log-likelihood:

$$\sum_{i=1}^n [Y_i \log(p_i) + (1 - Y_i) \log(1 - p_i)] - \lambda \sum_{k=1}^p \beta_k^2, \quad (1)$$

where $p_i = \text{expit}(X_i\boldsymbol{\beta})$. Typically, the Newton-Raphson (NR) algorithm is used to maximize (1). To this end, given current estimate $\tilde{\boldsymbol{\beta}}$, define $X_W = W^T X$, $W = (\text{diag}(\tilde{p}_i(1 - \tilde{p}_i)))^{1/2}$ and $\tilde{p}_i = \text{expit}(X_i\tilde{\boldsymbol{\beta}})$. Moreover, let $\mathbf{z} = (z_i)_{i=1}^n$ and $z_i = \text{logit}(\tilde{p}_i) + (Y_i - \tilde{p}_i) / (\tilde{p}_i(1 - \tilde{p}_i))$. Then, the NR update (Cule et al., 2011) is:

$$\hat{\boldsymbol{\beta}} = (X_W^T X_W + 2\lambda I)^{-1} X_W^T \mathbf{z}. \quad (2)$$

We assume (2) has been applied until convergence of the NR algorithm. Note that penalization causes bias, so, with $Y = (Y_1, \dots, Y_n)$: $E_Y(\hat{\beta}_k) \neq \beta_k$. Both $E_Y(\hat{\beta}_k)$ and $V_Y(\hat{\beta}_k)$ can be approximated, as shown below. We will use these moments to derive an empirical Bayes estimate of the group-specific penalties.

The first-order approximation μ_k of $E_Y(\hat{\beta}_k)$ is (le Cessie and van Houwelingen, 1992; Cule et al., 2011):

$$\begin{aligned} \mu_k &= [(I - 2\lambda(X_W^T X_W + 2\lambda I)^{-1})\boldsymbol{\beta}]_k = [(X_W^T X_W + 2\lambda I)^{-1}(X_W^T X_W + 2\lambda I - 2\lambda I)\boldsymbol{\beta}]_k \\ &= [(X_W^T X_W + 2\lambda I)^{-1} X_W^T X_W \boldsymbol{\beta}]_k =: \sum_{\ell=1}^p c_{k\ell} \beta_\ell \end{aligned} \quad (3)$$

where $[M]_k$ denotes the k th row (component) of any matrix (vector) M . In addition, we have (le Cessie and van Houwelingen, 1992; Cule et al., 2011) for $\Sigma = \text{Cov}(\hat{\boldsymbol{\beta}})$:

$$\hat{\Sigma} \approx (X_W^T X_W + 2\lambda I)^{-1} X_W^T X_W (X_W^T X_W + 2\lambda I)^{-1}. \quad (4)$$

Calculation of both μ_k and $\hat{\Sigma}$ requires the inverse of the large $p \times p$ matrix $M_\lambda = (X_W^T X_W + 2\lambda I)^{-1}$. However, singular value decomposition (SVD) of $X_W^T = U D V^T$ reduces the calculation of M_λ to inversion of an $n \times n$ matrix and matrix multiplication.

1.2 Empirical Bayes estimation of group penalties

Assume we have a partition of the variables into G groups, $(\mathcal{G}_1, \dots, \mathcal{G}_G)$, of sizes (K_1, \dots, K_G) . Then, replace the penalty term in (1) by a generalized ridge penalty term (Hoerl and Kennard,

1970):

$$\sum_{i=1}^n [Y_i \log(p_i) + (1 - Y_i) \log(1 - p_i)] - \sum_{g=1}^G \lambda_g \sum_{k \in \mathcal{G}_g} \beta_k^2, \quad (5)$$

where $\lambda_g = \lambda'_g \lambda$ with global penalty λ known and penalty multipliers λ'_g to be estimated. Let us assume an independent Gaussian (and hence ridge-type) prior:

$$\beta_k \sim N(0, \tau_{g(k)}^2), \quad (6)$$

where $g(k)$ denotes the group that variable k belongs to. Then, $E_{Y,\beta}(\hat{\beta}_k^2) = V_{Y,\beta}(\hat{\beta}_k)$, because $E_{Y,\beta}(\hat{\beta}_k) = 0$, where expected value E is computed with respect to the random variables in the subscript. Using (3) and (5) with $v_k = \hat{\Sigma}_{kk}$, plus the facts that $E_{\beta_i, \beta_j}(\beta_i \beta_j) = 0$ and $E_{\beta_i}(\beta_i^2) = \tau_h^2$ for $i \in \mathcal{G}_h$, we have:

$$\begin{aligned} E_{Y,\beta}(\hat{\beta}_k^2) &= E_{\beta} \left[V_Y(\hat{\beta}_k) + (E_Y(\hat{\beta}_k))^2 \right] \approx E_{\beta} \left[v_k + (E_Y(\hat{\beta}_k))^2 \right] \\ &= v_k + E_{\beta}[\mu_k^2] = v_k + E_{\beta} \left[\left(\sum_{\ell=1}^p c_{k\ell} \beta_{\ell} \right)^2 \right] = v_k + \sum_{h=1}^G \sum_{\ell \in \mathcal{G}_h} c_{k\ell}^2 \tau_h^2. \end{aligned} \quad (7)$$

This second moment renders the g th estimation equation by i) substituting $E_{Y,\beta}(\hat{\beta}_k^2)$ by its estimate, $\hat{\beta}_k^2$; ii) dividing both sides of (8) by v_k and subtracting 1; and iii) aggregating over all $k \in \mathcal{G}_g$:

$$\sum_{k \in \mathcal{G}_g} (\hat{\beta}_k^2 / v_k - 1) = \sum_{k \in \mathcal{G}_g} v_k^{-1} \left[\sum_{h=1}^G \sum_{\ell \in \mathcal{G}_h} c_{k\ell}^2 \tau_h^2 \right] = \sum_{k \in \mathcal{G}_g} \left[\sum_{h=1}^G \sum_{\ell \in \mathcal{G}_h} d_{k\ell}^2 \tau_h^2 \right] = \sum_{h=1}^G \alpha_{gh} \tau_h^2, \quad (8)$$

where $d_{k\ell} = c_{k\ell} / \sqrt{v_k}$ and $\alpha_{gh} = \sum_{k \in \mathcal{G}_g} \sum_{l \in \mathcal{G}_h} d_{kl}^2$. Let $B_g = \sum_{k \in \mathcal{G}_g} (\hat{\beta}_k^2 / v_k - 1)$. Then, the empirical Bayes estimate for $\tau_1^2, \dots, \tau_G^2$ is obtained by solving the system (linear in τ_h^2):

$$\begin{cases} B_1 = \sum_{h=1}^G \alpha_{1h} \tau_h^2 \\ B_2 = \sum_{h=1}^G \alpha_{2h} \tau_h^2 \\ \vdots \\ B_G = \sum_{h=1}^G \alpha_{Gh} \tau_h^2. \end{cases} \quad (9)$$

Naive computation of the coefficients α_{gh} requires calculation of the possibly very large $p \times p$ matrix $D = (d_{k\ell})_{k,\ell=1}^p$. We experienced that this may consume considerable computing time and memory. Fortunately, a much more efficient calculation is possible. To see this, first note from $d_{k\ell} = c_{k\ell} / \sqrt{v_k}$ and (3) that D is a product of a $p \times n$ matrix, $L = \text{diag}(1/\sqrt{v_k})(X_W^T X_W + 2\lambda I)^{-1} X_W^T$, and an $n \times p$ matrix $R = X_W$, where L can be efficiently computed by SVD of X_W^T . Matrix decomposition of L and R according to the groups implies that $\alpha_{gh} = \sum_{k,\ell} (d_{k\ell}^{gh})^2$, where $d_{k\ell}^{gh}$ are the elements of $D^{gh} = L_g R_h$ with $L_g = (L_k)_{k \in \mathcal{G}_g}$ and $R_h = (R_{\ell})_{\ell \in \mathcal{G}_h}$. D^{gh} may still be a prohibitively large matrix, and hence we wish to avoid computing it. The following proposition provides an efficient work-around for this when $p \gg n$, because it enables computation of α_{gh} by element-wise multiplication of $L_g^T L_g$ and $R_h R_h^T$, where both matrix products are of dimensions $n \times n$ only.

Proposition *Let L and R be $p_1 \times n$ and $n \times p_2$ matrices and $D = LR$. Let $A \circ B$ be the Hadamard (element-wise) product of any equally-sized matrices A and B . Denote the sum of elements of A by $[A]_{\Sigma} = \sum_{k,\ell} a_{k\ell}$. Then,*

$$\alpha = \sum_{k=1}^{p_1} \sum_{\ell=1}^{p_2} (d_{k\ell})^2 = [D \circ D]_{\Sigma} = [(L^T L) \circ (R R^T)]_{\Sigma}. \quad (10)$$

Proof: For any quadruple of matrices A, B, C and E of arbitrary dimensions $q \times r, q \times r, r \times s, r \times s$, respectively, we have

$$\begin{aligned} [(A^T B) \circ (C E^T)]_{\Sigma} &= \sum_{k,\ell} \left(\sum_i (A^T)_{ki} b_{i\ell} \sum_j c_{kj} (E^T)_{j\ell} \right) = \sum_{k,\ell} \left(\sum_i a_{ik} b_{i\ell} \sum_j c_{kj} e_{\ell j} \right) \\ &= \sum_{i,j} \left(\sum_k a_{ik} c_{kj} \sum_{\ell} b_{i\ell} e_{\ell j} \right) = [(AC) \circ (BE)]_{\Sigma}. \end{aligned}$$

Substituting $A = L, C = R, B = L$ and $E = R$ completes the proof. \square

System (10) generally results in satisfactory solutions when p is not extremely large (see also the Simulation section). For very large p , however, we experienced that it may lead to extreme (and even negative) values of the estimates. Such instabilities may be caused by strong multi-collinearity between variables (likely present in high-dimensional settings), also across groups, which affects the coefficients α_{gh} in (9). This may hamper disentangling the contributions of the various groups to each of the left-sides in (10). Therefore, we provide an iterative alternative below, but we first discuss how to obtain the group penalties from $\hat{\tau}_1^2, \dots, \hat{\tau}_G^2$, the solutions of (10).

To re-estimate β the resulting group-specific variances, $\hat{\tau}_g^2$ need to be converted to group-specific penalty multipliers λ'_g , using $\lambda'_g \propto 1/\hat{\tau}_g^2$. The missing constant is obtained by calibrating the mean of the inverses of λ'_g to 1. This amounts to solving for calibration constant C , with $K_g = |G_g|$:

$$\lambda'_g = C/\hat{\tau}_g^2 \quad \text{and} \quad \frac{1}{p} \sum_{g=1}^G K_g / \lambda'_g = 1. \quad (11)$$

This calibration is useful to avoid (often time-consuming) re-cross-validation of λ . It calibrates the mean of the inverse penalty multipliers towards the mean of those inverse multipliers in the original, initial ridge regression (with multipliers all equal to 1, implying a mean of 1). In fact, we observed for the examples below that after calibration re-cross-validation hardly changes the estimate of λ and the predictive performance. Finally, the group-specific penalty equals $\lambda_g = \lambda'_g \lambda$.

1.3 Estimation for generalized logistic ridge regression

After estimating the group-specific penalties we re-estimate β , which requires maximizing (6). This is achieved by applying ordinary logistic ridge regression, i.e. iteratively applying (2), with penalty parameter λ to a new weighted design matrix $X_W^{(2)} = X_W \Lambda^{-1/2}$, where Λ is a diagonal matrix with $\Lambda_{kk} = \lambda'_{g(k)}$. To see this, write the group-specific penalty term corresponding to variable k in group $g(k)$ as

$$\lambda_{g(k)} \beta_k^2 = \lambda [(\lambda'_{g(k)})^{1/2} \beta_k]^2 =: \lambda (\beta'_k)^2.$$

Then, write the contribution of column k in $X, [X]_k$, to the penalized likelihood (6) through $p_i = \text{expit}(X_i \beta)$ as $[X]_k (\lambda'_{g(k)})^{-1/2} \beta'_k$, which determines $X^{(2)} = X \Lambda^{-1/2}$, and hence also $X_W^{(2)} = W^T X^{(2)} = X_W \Lambda^{-1/2}$. Finally, for the new estimate of β_k , we have:

$$\hat{\beta}_k^{(2)} = (\lambda'_{g(k)})^{-1/2} \hat{\beta}'_k. \quad (12)$$

Here, the upper index in $\hat{\beta}_k^{(2)}$ refers to possible iterations, which will be introduced in the next section. The variance should be scaled analogously: $v_k^{(2)} = (\lambda'_{g(k)})^{-1} v'_k$, with $v'_k = V(\hat{\beta}'_k)$, available from (5).

1.4 An iterative alternative

Here, we provide an iterative alternative to (10). The system (10) does not make use of the fact that the initial estimates, $\hat{\beta}$, were implicitly (via the correspondence between the λ and τ^2) already obtained under a Gaussian prior with common variance τ^2 . In particular for high-dimensional data, this implicit prior has a large impact on $\hat{\beta}$. The proposed iterative alternative first estimates this common prior variance τ^2 . For that, we simply collect all variables in one group, which renders only one equation in (10) with solution:

$$\hat{\tau}^2 = \frac{\sum_{k=1}^p (\hat{\beta}_k^2 / v_k - 1)}{\sum_{k,\ell=1}^p v_k^{-1} c_{k\ell}^2}. \quad (13)$$

Then, we set out to estimate τ_g by first assuming $\tau_h^2 = \hat{\tau}^2$ for all $h \neq g$, which is reasonable given the (implicit) common prior that was used to obtain the estimates. Now, splitting the right-side of the g th equation of (10) into the contributions of group g and all other groups and substituting $\tau_h^2 = \hat{\tau}^2$ renders the estimate:

$$\hat{\tau}_g^2 = \frac{\sum_{k \in \mathcal{G}_g} (\hat{\beta}_k^2 / v_k - 1) - \sum_{k \in \mathcal{G}_g} v_k^{-1} \sum_{h \neq g} \sum_{\ell \in \mathcal{G}_h} c_{k\ell}^2 \hat{\tau}^2}{\sum_{k,\ell \in \mathcal{G}_g} v_k^{-1} c_{k\ell}^2}. \quad (14)$$

In words, (15) can be considered as an estimate of τ_g^2 that quantifies how much the observed sum of squared group g parameters (scaled by their variances) deviates from the expected contributions to this summand of all variables ℓ not in group \mathcal{G}_g . The above solution is particularly attractive when iterating the estimation, because then the updated $\hat{\tau}_g^2$ estimates adapt to the most recent generalized ridge estimates $\hat{\beta}'_k$ (13). As discussed above, these are also obtained under an implicit common prior (common λ), which allows us to iteratively use (14) and (15). From $\hat{\beta}'_k$, the iterative re-scaling in (13) then computes $\hat{\beta}_k$, which is on the original scale of the covariates X . We experienced that this alternative solution is always very competitive with explicitly solving (10), and sometimes superior, in particular when p is (very) large. The Supplementary Material provides a simulation-based comparison between the two.

Such iteration requires a stopping criterion. We simply monitor the cross-validated likelihood (CVL) and stop iterating when this decreases. The cross-validation is fast, because it only requires evaluation of the CVL for *given* global penalty λ . Moreover, we use the efficient implementation by Meijer and Goeman (2013). The resulting estimates are denoted by $\hat{\beta}_k^{(L)}$, where L is the number of iterations before the CVL decreases.

1.5 Iterating on a new partition

More than one partition of the variables may be available, as illustrated in the second example. Suppose we have two partition with G_1 and G_2 groups, respectively. Then, the above method may simply be applied by cross-tabling the two partitions, rendering $G_1 G_2$ groups. However, this may render a very large number of groups and some of these groups may contain only few variables, which may deteriorate the empirical Bayes estimates. Alternatively, one may simply iterate the group-specific regularization for the second partition after the first partition was considered. A disadvantage of that approach is that the results may (somewhat) depend on the ordering of the

partitions. For the iterative re-penalization solution (15), we therefore opt to embed iteration on partitions into the re-penalization iteration. Hence, partitions are considered in alternating order. The CVL-based stopping criterion formulated above is applied to both partitions with respect to the previous fit; if CVL does not improve, that particular partition does not take part in the outer re-penalization iteration anymore. If CVL does not improve for both partitions, the outer iteration is stopped as well. The group-regularization algorithm including this double iteration is depicted in Supplementary Figure 1.

Note that the new penalty multipliers will adapt to both the data and the current penalties. This is important when the partitions are not independent. Let $\hat{\beta}_k^{(\ell,j)}$ be the estimate of β_k for re-penalization iteration ℓ and partition $j = 1, 2$. The new estimate $\hat{\beta}_k^{(\ell,2)}$ is computed by applying (13) to $\hat{\beta}_k^{(\ell,1)}$, using grouping variable $g_2(k)$. Similarly, $\hat{\beta}_k^{(\ell+1,1)}$ is computed by applying (13) to $\hat{\beta}_k^{(\ell,2)}$, using grouping variable $g(k)$. The final penalty multiplier for variable k equals $\lambda'_{g(k)}\lambda''_{g_2(k)}$, where the latter term is the penalty multiplier based on the second partition. These notions trivially extend to more than two partitions. The final group-regularized estimates of β_k are denoted by $\hat{\beta}_k^{\text{GR}}$. The iterative group-regularization is illustrated in the second example.

1.6 Ranking-based small groups

Often, the co-data consist of external data on the same variables (e.g. genes) for an analogous, but somewhat different setting. Our second data example illustrates such a case. Then, the two data sets can not simply be pooled. However, summaries like p -values or regression coefficients based on the external data may be used to define a partition of the variables into small groups which is then used as input for the group-regularized ridge on the primary data set. We enforce monotony on the penalties of those groups to avoid over-fitting, as detailed below.

First, rank the variables according to the summary, e.g. p -values. Then, create groups of size s , where group g contains the variables with ranks $s(g-1)+1, s(g-1)+2, \dots, sg$. Apply (15) to obtain initial estimates $(\hat{\tau}_g^{\text{init}})^2$ for these small groups. Due to the size of the groups these estimates may be instable and not in line with the ranking based on the external data. Therefore, we force the estimates to be monotone by applying weighted isotonic regression (Robertson et al., 1982) of $(\hat{\tau}_g^{\text{init}})^2$ on the index (and hence group rank) g , rendering regression function $\hat{f}(\cdot)$. The weights account for possibly different group sizes. Then, the new estimates are set to $\hat{\tau}_g^2 = \hat{f}(g)$, which are substituted into (12) to obtain group-specific penalty multipliers λ'_g . Enforcing monotony highly stabilizes the estimates and interpretation of the results. In fact, even $s = 1$ might be used, but, because the stabilizing effect of the isotonic regression is potentially less strong for the extreme ranks, this could lead to over-fitting. The latter is mitigated by using small, non-singular groups. The stabilizing effect is illustrated for the second data example in Supplementary Figure 2.

The software also allows for non-uniformly-sized groups, where one specifies a minimum group size, say $s = 10$, for variables corresponding to the most extreme values of the summary, and a maximum number of groups; the group size then gradually increases for variables with less extreme values of the summary. This enables the use of fewer groups (and hence faster computations) while still maintaining a good ‘resolution’ for the extreme values of the summary.

2 Generalizing the concept I: post-hoc variable selection

A nice side effect of group-specific regularization is that it may simplify post-hoc variable selection, because the empirical distribution of estimated coefficients is typically more heavy-tailed than the one from ordinary ridge regression. Hence, there is a clearer separation between $\hat{\beta}_k$ ’s close to zero from those further away from zero. This is illustrated in Supplementary Figure 5 for the

second data example. Also, it is known that ordinary ridge regression tends to spread mass of the parameter estimates over correlated variables. Group-specific regularization can prioritize such variables, in particular when the groups are small and the range of group-specific penalties is large. A posteriori selection could be based on an information criterion or a mixture model for the $\hat{\beta}_k$'s. However, since we are in a prediction setting, we suggest to select directly on the basis of predictive performance by using CVL. For the purpose of prediction, variable selection is mainly desirable for potentially developing a measurement device (e.g. based on qPCR) with much fewer variables than the original one. Hence, we allow the user to set a maximum of variables to be selected, e.g. $p_{\max} = 100$.

A simple proposal for CVL-based selection is: sort the variables with respect to $|\hat{\beta}_k^{\text{GR}}|$; select $s, 0 \leq s \leq s_{\max}$ top-ranking variables; re-fit the model using only those variables, but with the same fixed λ and λ_g 's as for the full model; compute CVL_s on this model; find $\text{CVL}_{\max} = \max_s \text{CVL}_s$; select $s_{\text{sel}} = \min\{s : \text{CVL}_s \geq \text{CVL}_{\max} - q_{\text{marg}}|\text{CVL}_{\max}|\}$, with e.g. relative margin $q_{\text{marg}} = 1\% = 0.01$. The margin favors more sparse models: the minimization finds the model with the fewest variables such that its CVL is within a, say, 1% margin of the best. Supplementary Figure 1 depicts the entire group-regularization algorithm including variable selection, whereas Supplementary Figure 6 shows the CVL profile as a function of s for the second data example.

3 Generalizing the concept II: random forest

The concept of adaptive group-regularization (or, analogous, group-weighting) can be generalized to other classifiers, also to some of very different nature than logistic ridge regression. The Supplementary Material describes the extension to the random forest classifier in detail; below we provide a summary.

A standard random forest classifier uses only $m = \mathcal{O}(\sqrt{p})$ variables (nodes) per node split. Typically, these variables are sampled uniformly from the entire set. Now, the idea is to weigh groups by increasing or decreasing the sampling probability according to the overall importance of variables in a group. Once a set of top-ranking variables across a forest is defined by a formal selection procedure (Doksum et al., 2008) or by simply using the top $k\%$ (for, say, $k = 5$), the observed number of top-ranking variables per group is modeled by a multinomial distribution per tree. Then, the variability of the multinomial proportions across trees is modeled by a Dirichlet distribution. The parameters of this prior are then estimated by use of empirical Bayes. The result is then used for weighted sampling of variables in the trees in a new random forest. The process of random forest classification, variable ranking, selection, estimation and weighted sampling is repeated, until the out-of-bag error does not or hardly decrease anymore.

4 Simulation results: summary

We performed simulations to compare the performances of the systems-based solution (10) and the iterative solution (15). In addition, `GRridge` is compared with i) ordinary logistic ridge and ii) the group lasso (Meier et al., 2008). We also compared with the adaptive logistic ridge, which is the ridge version of the adaptive logistic lasso (Zou, 2006), simply amounting to using variable-specific penalty multipliers that are inverse proportional with the initial squared ridge parameter estimates, $\hat{\beta}_k^2$. However, since we found that the predictive performance of the adaptive logistic ridge was generally inferior to that of the ordinary logistic ridge, we do not present those results in detail.

We studied a number of scenarios where we varied the number of groups G , the size of the groups p_g , the correlation strength in X , the differential signal between the two classes of samples

across groups, and the sparsity (i.e. proportion of groups without predictive signal). Performance was evaluated by computing AUC and mean Brier residuals on a large test data set ($n_{\text{test}} = 1000$), which was generated under the same settings as the training set ($n = 100$). These are reported in extensive tables, supplied in the Supplementary Material; here, we summarize the results.

First, we observe that the systems-based solution and the iterative solution are very competitive for $p = 2000, 5000$ ($p = G * p_g$), while the latter is superior for large $p, p = 12500$. In particular, the iterative solution is indeed more stable across repeated simulations for $p = 12500$. The non-iterative, systems-based solution relies strongly on the parameter estimates of the initial logistic ridge regression, the bias of which may be strong when p is very large. The iterative solution, however, typically finds less extreme group penalties in the first iteration, then re-estimates the regression parameters, allowing those to adapt to the new penalties.

Generally, both **GRridge** versions performed at least as good as ordinary logistic ridge. As expected, we clearly observe that the gap between the performances increases with more skewed effects across groups and with increased sparsity. In addition, group lasso outperforms ordinary logistic ridge in group-sparse settings, while the reverse often holds for the non-sparse settings. **GRridge** generally outperforms the predictive accuracy of the group lasso, in some cases with fairly large margins, e.g. with AUCs that are 0.10 to 0.15 larger on the absolute scale. The group lasso becomes more competitive for high group-sparsity, in particular for p large. Yet, it seems that **GRridge** adapts well to sparsity and maintains its relatively good performance. Note that the weaker predictive performance of the group lasso may, for some applications, be counterbalanced by its group-selection property.

5 Examples: diagnostic classification using methylation data

DNA consists of the four nucleotides A, C, G and T. Methylation refers to the addition of a methyl-group to a C preceding a G (CpG), which can influence expression of the encoded gene. As such, methylation has a so-called epigenetic effect on the functionality of the cell, and consequently on the entire organism. It is believed to be an important molecular process in the development of cancer (Laird, 2003). In addition, DNA is a well-characterized and relatively stable molecule, compared to mRNA (gene expression) and many proteins. Therefore, the use of DNA methylation for diagnostic purposes is currently heavily investigated. A popular platform for measuring methylation is the IlluminaTM 450K bead chip. This platform measures 450,000 probes per individual, where each probe corresponds uniquely to a CpG location on the genome. Each probe renders a so-called beta-value, which is the estimated proportion of methylated DNA molecules for that particular genomic location in a given tissue. Like for any microarray study, the data is preprocessed using several steps; see the Supplementary Material.

We have data sets from two similar studies on cervical cancer at our disposal. The carcinogenesis of cervical cancer is relatively well-characterized. The transformation process of normal epithelium to invasive cancer takes many years, and includes distinct stages of precursor lesions (CIN; cervical intraepithelial neoplasia). Whereas low-grade precursor lesions (CIN1/2) are known to regress back to normal, high-grade precursor lesions (CIN2/3) have a relatively high risk for progression to cancer and are usually surgically removed. Therefore, accurate detection of high-grade CIN is very important. The two studies both measure methylation for normal cervical tissue and CIN3 tissue for several independent individuals, but differ in one important aspect. The first study measures methylation on CIN3 tissue biopsies, whereas the second study considered self-collected cervico-vaginal lavages of women with underlying CIN3 lesions (Gok et al., 2010). The relatively good quality of the samples in the first study may render important information about relevant methylation markers. The quality and purity of the tissues in the second study is probably inferior. This study, however, better resembles a more realistic diagnostic setting, in particular because

many countries have implemented screening programs for cervical cancer. Our first example uses the data of the first study only, but complements this with other source(s) of information to create partition(s) of the probes into groups, which are used in the group-regularized ridge regression. The second example shows how, in addition, the results of the first study can be used in our algorithm to improve diagnostic classification for the second study. We present the results from the iterative version of our method.

5.1 To standardize or not? - An automatic solution

A practical issue when applying penalized regression is the need or ‘no need’ for standardization of the covariates. There is no consensus on this issue (Zwiener et al., 2014), because on one hand standardization has the beneficial effect of rendering a common penalty more appropriate, while on the other hand it may remove some of the (differential) signal and may lead to instabilities for variables for which the sample variances are small. Standardization is equivalent to introducing a penalty multiplier that is proportional to the variance in the unstandardized setting (Zwiener et al., 2014). A potential of our method is that it can let the data decide how the variances of the variables should impact the penalties. Below, we explore this potential for the first study.

The first study (Farkas et al., 2013) contains methylation profiles of 20 and 17 unrelated normal cervical and CIN3 tissues, respectively. To enable inclusion of these data in our complementary R-package `GRridge`, thereby allowing reproduction of the results, the computations for this and the next example were performed on a random selection of 40,000 probes. We verified, however, that all results are very similar on the entire data set, which is not surprising given the smooth nature of ridge regression and the correlations between variables.

The probes are grouped in 8 groups of 5,000 each, in increasing order of the sample variances. Note that we verified whether a different grouping (e.g. 16 groups of 2,500 or 40 groups of 1,000) would affect the results. This is not the case. In line with the argumentation above (larger penalties for probes with large variances) we imposed monotony on the 8 penalty multipliers. We observed that this constraint is not very essential here, because the estimated penalty multipliers are also nearly monotonously increasing when the constraint was not imposed. `GRridge` estimated the following penalty multipliers: $(2.75 * 10^{-2}, 6.59 * 10^{-2}, 6.92 * 10^{-2}, 8.03 * 10^{-2}, 1.21 * 10^{-1}, 3.00 * 10^{-1}, 7.36 * 10^8, 2.75 * 10^9)$. So, it effectively completely removes the impact of the probes with large variances (last two groups), allowing smaller penalties for the remaining 6 groups. Interestingly, `GRridge` with variance-based groups outperforms *both* ordinary standardized and unstandardized logistic ridge, in terms of ROC-curves (see Supplementary Figure 3), AUC (0.91, 0.86, 0.76, respectively) and mean Brier residuals (defined as $1/n \sum_{i=1}^n (Y_i - p_i)^2$; 0.14, 0.16, 0.21, respectively). Hence, `GRridge` provides a well-performing, automatic solution for the standardization issue here.

5.2 Improved classification by use of probe annotation

Our hypothesis here is that the use of a priori known annotation-based partitions of the probes may improve the classification results. This second partition, next to the variance-based one above, is based on the probe’s location in or nearby a so-called CpG-island. A CpG-island is a genomic region which is relatively rich in CG base pairs, and methylation is known to be more prevalent there than elsewhere. We used the following groups (in order of decreasing distance to CpG-islands): “CpG-island (CpG)”, “North Shore (NSe)”, “South Shore (SSe)”, “North Shelf (NSf)”, “South Shelf (SSf)”, and “Distant (D)”. If probes in CpG (or any other group of probes) are on average more important for the classification, the group-regularized ridge automatically detects this and applies a smaller penalty to all probes in this group. This may improve classification when the a priori partition was indeed informative. Note that the partition used is based on a well-accepted criterion to characterize genomic locations in methylation studies.

The group-regularized ridge used 6 iterations for re-penalizing the 6 annotation-based groups and 7 iterations for the 8 variance-based ones, which increased the CVL by 40% from -20.18 to -12.03. The final penalty multipliers for the annotation-based groups (\propto inverse weights) are: $\lambda'_{\text{CpG}} = 0.015$, $\lambda'_{\text{NSE}} = 278$, $\lambda'_{\text{SSE}} = 0.12$, $\lambda'_{\text{NSf}} = 2,986$, $\lambda'_{\text{SSF}} = 2,987$ and $\lambda'_D = 685$. The group-specific penalties clearly affect the regression parameter estimates $\hat{\beta}_k^{\text{GR}}$, because larger values of λ'_g result in smaller values of $|\hat{\beta}_k^{\text{GR}}|$. Hence, these values imply that **GRridge** effectively only uses the CpG and SSE probes for the predictions. The results confirm the importance of probes on CpG islands. The variance-based penalty multipliers are $(1.93 * 10^{-1}, 2.41 * 10^{-1}, 2.41 * 10^{-1}, 2.41 * 10^{-1}, 2.92 * 10^{-1}, 6.74 * 10^{-1}, 3.29 * 10^4, 1.11 * 10^5)$. These are largely in line with the results above, although somewhat compressed, because they adapted to the annotation-based multipliers. Please note that one should be careful with interpreting the exact values of the group penalties. As indicated above, these may depend on the presence of another partition due to overlap between groups. In addition, in the Supplementary Material we show that for the annotation-based groups above the penalties vary somewhat with respect to sizes of the groups. The order of the group-penalties seems to be fairly stable, however, so we recommend to interpret the group-penalties in terms of their ranking.

To assess whether the group-regularized ridge improves classification with respect to ordinary ridge, we computed ROC curves obtained by 10-fold cross-validation. Here, we compare with ridge regression on standardized covariates, because the latter was superior to the unstandardized version, as demonstrated above. In addition, we compare with the group lasso (Meier et al., 2008), as implemented in the R package `grpreg`, using the same annotation-based groups as for **GRridge**. The last competitor is the adaptive ridge, as discussed above. Also these two methods are applied to the standardized covariates (which were verified to be superior to their unstandardized counterparts). Group lasso selected only the SSE group, but, surprisingly, not the CpG group.

The resulting ROC curves depict the False Positive Rate (FPR) versus the True Positive Rate (FPR) for a dynamic cut-off for the predicted probability on CIN3. Figure 1(a) shows the results. We clearly observe superior performance of **GRridge**, with AUC = 0.92 (and 0.86, 0.84, 0.79 for ordinary ridge, adaptive ridge and group lasso, respectively). With respect to ordinary ridge, predictions improved for 33 out of 37 observations, as displayed in Figure 1(b). Note that adding the annotation-based groups to the variance-based groups improved the AUC only slightly, from 0.91 to 0.92, probably because AUC is a rank-based criterion. In fact, the predictions did improve for 32 out of 37 observations, leading to a relatively larger improvement (decrease) for the mean Brier residual, from 0.143 to 0.116.

5.3 Improved diagnostic classification by use of external data

The second study contains methylation profiles of self-collected cervico-vaginal lavages (or self-samples) corresponding to 15 women with an unaffected (normal) cervix and 29 women with CIN3 lesions, all unrelated. Here, it is important to note that the samples of the affected cervixes may be contaminated with normal cells and cells of other origins (mostly vaginal cells and lymphocytes), due to imprecise sampling. Hence, the differential signal may be diluted. We aim to use the results of the first study for the group penalties in the second study.

In principle, we could use the results of the group-regularized ridge regression fitted on the first study, as presented in the previous section. However, the effect of the (possible) contamination may vary considerably across probes. For example, the differential signal of probes with hypo-methylation (affected < normal) in the first study is diluted more than that of hyper-methylated probes. This can be illustrated in a simple deterministic setting. In case of hypo-methylation, consider a true ratio affected/normal = $0.4/0.8 = 1/2$. Assume a contamination of 50%, then the measured ratio will be $(0.4/2 + 0.8/2)/0.8 = 3/4$, hence the ratio is 50% too large. Using the same

numbers for hyper-methylation renders a measured ratio that is only 33% too small. In addition, it is well-known that ridge regression distributes differential signal over parameters corresponding to correlated probes. Hence, the magnitude of a particular coefficient also depends on other probes. Since the dilution in Study 2 affects probes differently, the applicability of Study 1 ridge regression results for analyzing Study 2 may be limited.

Therefore, we propose to use group penalties λ_g that are simply based on t -test p -values as obtained by applying `limma` (Smyth, 2004) on Study 1. These p -values are then used to define a ranking-based partition with 100 groups of probes of minimal size $s = 10$ (size gradually increasing with the p -value) as described above. To stabilize the estimates of λ_g weights $\hat{\tau}_g^2 \propto 1/\lambda_g$ for Study 2 are forced to be monotonously decreasing with increasing Study 1 p -values as described above. The function `pava` of the R library `Iso` is used for this purpose, which is illustrated in Supplementary Figure 2. In this setting, it is reasonable to precede our method by a mild *prior* filtering: only include those probes with $\text{FDR} \leq 0.5$ and a mean absolute difference larger than 0.1 (on log-scale) in Study 1. Then, our method applies group-specific regularization to the 9491 probes surviving these thresholds.

Given the earlier argument about a stronger dilution effect on hypo- than on hyper-methylated probes (as detected in Study 1), we also considered a second sign-based partition that distinguishes those two groups of probes. Finally, we added the variance-based and annotation-based partitions introduced in the first example. This illustrates the ability of our method to operate on multiple partitions. For this example, the adaptive group-regularized ridge used 3 re-penalization iterations. The CVL increased from -28.91 to -27.54, hence a 5% improvement. The sign-based and variance-based partitions had no effect on the results (hence rendering group-specific penalties equal to 1) on top of the p -value-based partition, illustrating the adaptive nature of the algorithm. The partition based on external p -values produced 100 group-specific penalties ranging from $1.3 \cdot 10^{-3}$ to 13.3 for $g = 1, \dots, 100$, so indeed a large range (see also Supplementary Figure 4), illustrating the relevance of this partition. The annotation-based partition rendered $\lambda_{\text{CpG}}^V = 0.17$, and much larger penalty multipliers for the other 5 classes. So, also for this data set the probes on the CpG-islands correspond to smaller penalties, which is biologically plausible.

We compared `GRridge` with: i) ordinary ridge; ii) adaptive ridge; iii) group lasso with the same 100 p -value based groups (all three also applied to the filtered probe set); and iv) ordinary ridge on the entire, non-filtered probe set. For this data set, the group lasso did not select any group. Hence, no variable was selected either, rendering inferior prediction results. Probably, the weak differential signal per variable in this challenging data set caused the absence of selections for the group lasso. The ROC curves for the other three methods were obtained by applying leave-one-out cross-validation (LOOCV). Figure 2(a) shows the results: `GRridge` has a markedly higher AUC (0.74) than the ones corresponding to i) 0.67, ii) 0.67 and iv) 0.63.

We also checked whether the order in which the four partitions are used within each re-penalization iteration matters for the results. The final CVLs for all 24 possible orderings show very little variation: the range is $[-24.58, -24.54]$. Hence, we conclude that the sensitivity of the performance with respect to the ordering is very small for this data set.

5.4 Variable selection

Supplementary Figure 5 shows that the most extreme coefficients of the group-regularized ridge regression are relatively much larger than those of ordinary logistic ridge regression. In fact, for the former, the 1% most extreme coefficients account for 61% of the total sum of absolute values of the coefficients, whereas for the latter this drops to only 3%. We applied the proposed *a posteriori* variable selection to $\hat{\beta}_k^{\text{GR}}$ which rendered a model with 42 selected variables, termed `GRridge+sel`. Figure 2(b) depicts the ROC-curves and AUCs for `GRridge+sel`, `GRridge` and `lasso`, as obtained by LOOCV. First, note that the much more parsimonious `GRridge+sel` model predicts nearly as

well as the full `GRridge` model in this case (AUC = 0.72 vs AUC = 0.74). Second, to illustrate the beneficial effect of group-specific regularization in this variable selection context, we also compare `GRridge+sel` with the `lasso` (Goeman, 2010, R package `penalized`) on the same filtered data set. The `lasso` renders a somewhat more parsimonious model with 17 variables, but performs much worse in terms of prediction: Figure 2(b) depicts the ROC-curves and AUCs. Of course, the `lasso` could possibly be improved by adapting group-regularized principles as well (see Discussion).

6 Discussion

Our method is weakly adaptive in the sense that the penalties adapt in a group-specific sense only. This is an important conceptual difference with strongly adaptive methods such as adaptive lasso (Zou, 2006) and enriched random forests (Amaratunga et al., 2008), which aim to learn variable-specific penalties from the same data as the data used for classification. Such methods strongly rely on sparsity. While this may be a fairly natural assumption for some applications, we believe it to be less realistic for complex genomic traits like cancer. In fact, we observed that for both applications the adaptive lasso did not outperform the ordinary lasso, and hence performed worse than the adaptive group-regularized ridge regression.

The adaptive group-regularized ridge shares the philosophy of accounting for group structure with the group lasso (Meier et al., 2008). The latter, however, *selects* entire groups using a lasso penalty on the group-wise sum of coefficients and then spreads the coefficients within a group using a ridge penalty within a group. The group lasso is particularly attractive for selecting relatively small, interpretable groups of variables, e.g. gene pathways. However, it is less useful and suitable when the groups tend to be large (and not necessarily homogenous) as in the first example, or when the groups have no clear biological interpretation, as for the ranking-based small groups in the second example. In addition, for both simulated and real data sets we show that the predictive performance of `GRridge` is often superior to that of the group lasso. Group-specific regularization was also discussed by Tai and Pan (2007) in the context of nearest shrunken centroids and partial least squares classifiers. Their results support our claim that such regularization can improve classification performance. Their approach, however, requires cross-validation on *all* group-penalties or, when this is too computationally demanding, *a priori* fixing of weights (inverse penalties). Also, unlike `GRridge`, their method does not make use of multiple partitions of the variables, which are often available in practice.

As discussed, group-regularization helps to better discriminate small and large coefficients, and the model after variable selection may be fairly parsimonious. Yet, extension of our method to sparse methods like lasso may be desirable in some cases. These methods usually render only few non-zero coefficients, which may lead to unstable group penalties. This may be mitigated by re-sampling or by using a power transformation of ridge-based penalties, as suggested by Bergersen et al. (2011) in another setting. Alternatively, one may consider a Bayesian set-up with a selection prior, for example a Laplace prior (Park and Casella, 2008) or a horseshoe prior (Carvalho et al., 2009). The hyper-parameters of such priors would be estimated *per group of variables*, e.g. by empirical Bayes. Then, the entire posterior of each β_k , rather than just the point estimate, impacts the penalty (represented by the group-wise prior) of the group that variable k belongs to.

It is possible to shrink β towards the corresponding estimates of the external study rather than to zero, i.e. targeting shrinkage (Gruber, 1998). However, unless the two experiments are expected to be very similar in terms of design, quality, effect size distribution, and the exact meaning of the two corresponding β_k 's, this may do more harm than good. For example, our illustration on the joint use of the two methylation studies does clearly not satisfy these conditions: due to the dilution, the β_k 's in Study 2 are bound to be weaker than those in Study 1, and likely in a non-uniform way. Yet, in very well-controlled settings targeted shrinkage may be a useful extension.

We end with some practical remarks. The adaptive group-regularized logistic and linear ridge procedures are implemented in the R-package `GRridge`, available via www.few.vu.nl/~mavdwi1/grridge.html. It depends on the package `penalized` (Goeman, 2010), which is used for model fitting and cross-validation. `GRridge` provides all functionality described in this paper, including: both versions (the non-iterative, systems-based one and the iterative one), adaptive regularization on multiple partitions, variable selection, estimation of predictive accuracy by cross-validation and convenience functions to create partitions of the variables using co-data. In addition, it allows for including non-penalized variables, e.g. clinical information. It also includes both data sets discussed here. The iterative version is the default, but based on the simulations we believe one can safely use the faster, non-iterative version for $p \leq 1000$. The iterative algorithm, however, is also fairly fast. For the first example ($p = 40,000$, 7 iterations on two partitions), constructing the group-regularized ridge classifiers took 3m01s and 3m27s, for tuning the global penalty λ by LOOCV and group-regularization, respectively. Hence, 6m28s in total on a 3GHz laptop with 3.5Mb RAM. The second example ($p = 9,491$, 3 iterations on four partitions) took 31s, 23s and 14s for λ -cross-validation, group-regularization and selection, respectively, so 1m08s in total. The code used to produce the results of `GRridge` in this paper is included in the Supplementary Material.

7 Acknowledgements

We are particularly grateful to Tristan Mary-Huard for his critical comments and suggestions on earlier versions of this manuscript. In addition, we thank Sanja Farkas for providing the raw data of her study (Farkas et al., 2013) and Carel Peeters for discussing several aspects of ridge regression. This study was partly supported by the OraMod project, which received funding from the European Community under the Seventh Framework Programme, grant no. 611425. DNA methylation data of the self-samples were obtained as part of a project supported by the European Research Council (ERC advanced 2012-AdG, proposal 322986; Mass-care), by which also Wina Verlaet and Saskia M. Wilting were supported.

References

- Amaratunga, D. et al. (2008). Enriched random forests. *Bioinformatics*, **24**, 2010–2014.
- Bergersen, L.C. et al. (2011). Weighted lasso with data integration. *Stat. Appl. Genet. Mol. Biol.*, **10**.
- Carvalho, C. et al. (2009). Handling sparsity via the horseshoe. *J. Mach. Learn. Res., W&CP*, pages 73–80.
- Cule, E. et al. (2011). Significance testing in ridge regression for genetic data. *BMC Bioinf.*, **12**, 372.
- Doksum, K. et al. (2008). Nonparametric variable selection: the EARTH algorithm. *J. Amer. Statist. Assoc.*, **103**, 1609–1620.
- Farkas, S. et al. (2013). Genome-wide DNA methylation assay reveals novel candidate biomarker genes in cervical cancer. *Epigenetics*, **8**, 1213–1225.
- Futreal, P. et al. (2004). A census of human cancer genes. *Nat. Rev. Cancer*, **4**, 177–183.
- Goeman, J. (2010). L1 penalized estimation in the Cox proportional hazards model. *Biom. J.*, **52**, 70–84.

- Gok, M. et al. (2010). HPV testing on self collected cervicovaginal lavage specimens as screening method for women who do not attend cervical screening: cohort study. *BMJ*, **340**, c1040.
- Gruber, M. (1998). *Improving Efficiency by Shrinkage: The James–Stein and Ridge Regression Estimators*. Statistics: A Series of Textbooks and Monographs.
- Hastie, T. et al. (2008). *The elements of statistical learning, 2nd ed.* Springer, New York.
- Hoerl, A.E. and Kennard, R.W. (1970). Ridge regression: Biased estimation for nonorthogonal problems. *Technometrics*, **12**, 55–67.
- Laird, P. (2003). The power and the promise of DNA methylation markers. *Nat. Rev. Cancer*, **3**, 253–266.
- le Cessie, S. and van Houwelingen, J. (1992). Ridge estimators in logistic regression. *Applied Statistics*, **41**, 191–201.
- Meier, L. et al. (2008). The group Lasso for logistic regression. *J. R. Stat. Soc. Ser. B Stat. Methodol.*, **70**, 53–71.
- Meijer, R. and Goeman, J. (2013). Efficient approximate k-fold and leave-one-out cross-validation for ridge regression. *Biom. J.*, **55**, 141–155.
- Neuenschwander, B. et al. (2010). Summarizing historical information on controls in clinical trials. *Clin. Trials*, **7**, 5–18.
- Park, T. and Casella, G. (2008). The Bayesian lasso. *J. Amer. Statist. Assoc.*, **103**, 681–686.
- Robertson, T. et al. (1982). *Order Restricted Statistical Inference*. Wiley, New York.
- Smyth, G.K. (2004). Linear models and empirical Bayes methods for assessing differential expression in microarray experiments. *Stat. Appl. Genet. Mol. Biol.*, **3**, Art. 3.
- Tai, F. and Pan, W. (2007). Incorporating prior knowledge of predictors into penalized classifiers with multiple penalty terms. *Bioinformatics*, **23**, 1775–1782.
- Zou, H. (2006). The adaptive lasso and its oracle properties. *J. Amer. Statist. Assoc.*, **101**, 1418–1429.
- Zwiener, I. et al. (2014). Transforming RNA-Seq data to improve the performance of prognostic gene signatures. *PLoS ONE*, **9**, e85150.

Figures

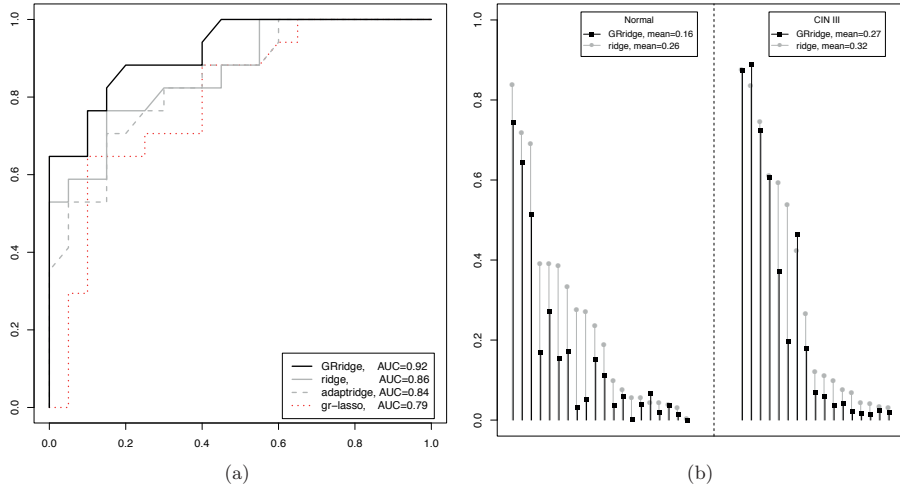


Figure 1: (a): ROC curves for first example, Group-regularized ridge (**GRridge**), ordinary **ridge**, group lasso (**gr-lasso**) and adaptive ridge (**adaptridge**). X-axis: False Positive Rate, y-axis: True Positive Rate. (b): Absolute residuals $|Y_i - p_i|$ for both classes for **GRridge** and **ridge**, in decreasing order of the **ridge** residuals

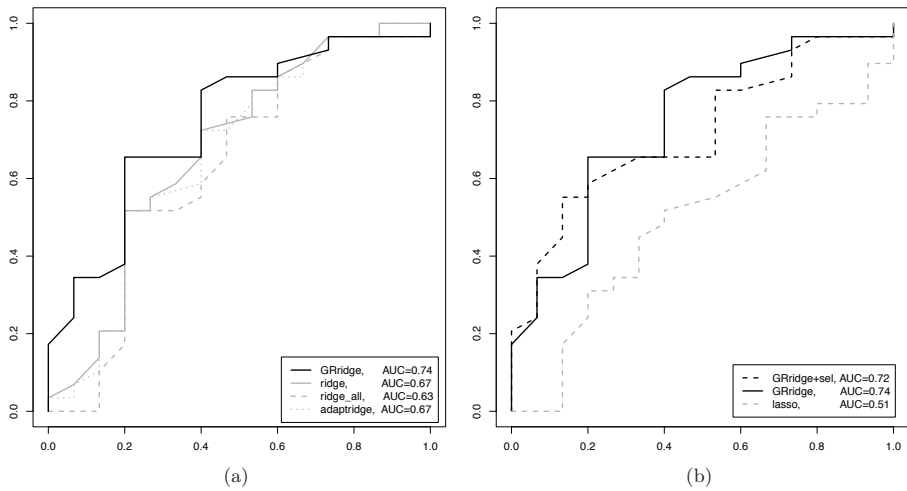


Figure 2: ROC curves for second example. Sub-figure (a): group-regularized ridge (GRridge), ridge and ridge on all variables (ridge_all); (b): Group-regularized ridge plus variable selection (GRridge+sel), GRridge, and lasso. X-axis: False Positive Rate, y-axis: True Positive Rate.

SUPPLEMENTARY MATERIAL FOR: Better prediction by use of co-data: Adaptive group-regularized ridge regression

Mark A. van de Wiel^{1,2}, Tonje G. Lien³, Wina Verlaat⁴,
Wessel N. van Wieringen^{1,2}, Saskia M. Wilting⁴

1. Department of Epidemiology and Biostatistics, VU University Medical Center, PO Box 7057, 1007 MB Amsterdam, The Netherlands
2. Department of Mathematics, VU University, Amsterdam, The Netherlands
3. Department of Mathematics, University of Oslo, Oslo, Norway
4. Department of Pathology, VU University Medical Center, Amsterdam, The Netherlands

1 Group-regularized *linear* ridge regression

Adapting the algorithm from logistic to linear ridge regression is fairly straightforward. For $\hat{\beta}$, μ_k and $c_{k\ell}$ one simply needs to substitute the weighted design matrix X_W by the ordinary $n \times p$ design matrix X and 2λ by λ in the corresponding formulas. Likewise, these substitutions render $\hat{\Sigma}$ and its diagonal elements v_k when multiplying the result with the estimated error variance, $\hat{\sigma}^2$, provided in Cule et al. (2011). Note that the expressions for the moments are exact for the linear case (Hoerl and Kennard, 1970; Gruber, 1998).

2 Effect of group sizes on group penalties

We studied the effect of the group sizes for the first example, using the annotation-based partition, complementary to the sd-based one. Table 1 displays the configurations used, with increasing unbalance in groups; the 5th configuration corresponds to the full set of 40,000 variables used in the Main Document. Table 2 shows the resulting group-penalty multipliers on a \log_{10} -scale. It should be noted that ridge regression becomes insensitive to the exact value of the group-penalty if it is very large; hence, one may expect quite some variation in the upper range of the penalties, which, however, has little impact on the actual parameter values and predictions. Nevertheless, we observe considerable variation of group-penalty multipliers across the configurations. However, the ranking of the groups, as displayed in Table 3, is fairly stable. Hence, it may be best to interpret the penalties in terms of their ranking.

3 Simulation results

We performed simulations to compare the performances of the systems-based solution and the iterative solution (see Main Document, equations (9) and (14), respectively). In addition, `GRidge` is compared with i) ordinary logistic ridge and ii) the group lasso (Meier et al., 2008), for which we used the R-packages `penalized` and `grpreg`, respectively. The latter was used with default parameter settings; we checked, however, whether use of other (reasonable) settings, in particular

Configuration	Distant	CpG	NSf	NSe	SSf	SSe
1	1500	1500	1500	1500	1500	1500
2	3000	3000	1500	3000	1500	1500
3	6000	6000	1500	3000	1500	3000
4	12000	9000	1500	3000	1500	3000
5	14047	12858	2006	5262	1765	4062

Table 1: Group-size configurations

Configuration	Distant	CpG	NSf	NSe	SSf	SSe
1	2.64	-1.09	1.33	0.25	0.50	-0.25
2	3.95	-1.71	3.96	-0.05	2.11	-0.95
3	3.64	-1.78	4.34	2.01	4.07	-0.83
4	3.64	-2.09	5.03	3.76	5.13	-1.30
5	2.81	-1.84	3.48	2.45	3.48	-0.94

Table 2: \log_{10} group-penalty multipliers for five configurations of the group sizes.

Configuration	Distant	CpG	NSf	NSe	SSf	SSe
1	6	1	5	3	4	2
2	5	1	6	3	4	2
3	4	1	6	3	5	2
4	3	1	5	4	6	2
5	4	1	6	3	6	2

Table 3: Ranks of \log_{10} group-penalty multipliers for five group size configurations.

for the **alpha** parameter (which tunes the amount of within-group ridge-type regularization), had a large effect on the performance of the group lasso in our simulations; this was not the case. We study a number of scenarios where we vary the number of groups G , the size of the groups p_g , the correlation strength in X , the differential signal between the two classes of samples across groups, and the sparsity (i.e. proportion of groups without predictive signal). Performance was evaluated by computing AUC and mean Brier residuals on a large test data set ($n_{\text{test}} = 1000$), which was generated under the same settings as the training set ($n = 100$).

Let us first describe the various simulation settings. The number of groups of variables, G was either set to 10, or to 25. The setting $G = 10$ was combined with either $p_g = 200, p_g = 500$, whereas $G = 25$ was combined with $p_g = 500$ only. This renders 3 combinations with $p = p_g * G = 2000, 5000, 12500$ number of variables. We used a multivariate Gaussian block correlation structure for X , with blocks of size 10, meaning that each variable is correlated to 9 other variables. The correlation strength ρ was set to 0.1 or 0.5. To set the differential signal between two classes of samples we largely followed Waldron et al. (2011) in the sense that we used a *constant* β_k value per group, where β_k is the parameter in the linear predictor of the logistic regression corresponding the variable $k = 1, \dots, p$. Hence, to avoid any bias, we do not generate the parameters from the Gaussian prior assumed by our method. We control the overall differential signal by fixing the mean of β_k , $\bar{\beta}$. For $\rho = 0.1$, we used $\bar{\beta} = 0.1$, whereas $\bar{\beta} = 0.01$ was used in combination with $\rho = 0.5$. These settings enabled the ordinary logistic ridge to pick up some signal, leading to

non-trivial prediction results (being either as bad as random or as good as perfect).

The degree to which groups of variables differ in terms of differential signal is controlled by simulation parameter f . Labeling the groups $1, \dots, G$, group g has a mean differential signal, $\bar{\beta}_g$, that is f times stronger than that of group $g+1$, $g = 1, \dots, G-1$. Hence, $\bar{\beta}_g$ follows a power law. Note that the order of the groups was only used for simulation purposes, not when running GRridge. We set $f = 1.3, 1.6, 2$. Finally, we study the effect of ‘sparsifying’ the group signals (in line with the spirit of the group lasso), by setting all parameters of a proportion q of the weakest groups to zero. Here, we use $q = 0, 0.7, 0.9$ for $G = 10$ (hence 10, 3, 1 non-zero groups, respectively), and $q = 0.88, 0.96$ for $G = 25$ (hence 3, 1 non-zero groups, respectively). So, for the latter case we only run the ‘group-sparse’ setting, mainly for comparing with the group lasso when the number of groups is fairly large. Finally, the binary response Y_i was obtained by sampling from a Bernoulli(p_i) distribution, where $p_i = \text{expit}(X_i\beta)$, for $i = 1, \dots, n$. All in all, this rendered $2*2*3*3=36$ and $2*3*2 = 12$ different simulation scenarios for $G = 10$ and $G = 25$, respectively, so 48 in total. Each of the simulations was repeated 5 times for which we report the median and worst performance in terms of AUC and mean Brier residuals (see Tables 4 to 6 and 7 to 9, respectively, at the end of the manuscript).

From Tables 4 to 9 we draw the following conclusions:

- Both GRridge methods perform at least as good as ordinary logistic ridge for all simulation settings. As expected, the improvement in performance is particularly evident for large values of f and/or q (both implying large differences between groups)
- The non-iterative, systems-based solution (GRridge System) and the iterative solution (GRridge Iterative) are very competitive for moderately large p : $p = 2000, p = 5000$. However, the latter is superior for the large p setting ($p = 12500$; Tables 6 and 9). The better stability of the latter is reflected by its larger minimum AUC-values [reported between brackets] and smaller maximum mean Brier residuals across the 5 simulation repeats.
- As expected, Group lasso outperforms ordinary logistic ridge in group-sparse settings, so when q is large; the gap in performance tends to be larger when f is large as well. For the non-sparse $q = 0$ setting ordinary logistic ridge can be superior, in particular when f is relatively small and correlation ρ is large.
- GRridge Iterative performs at least as good as the group lasso for all simulation settings. The latter is competitive (and superior to ‘GRridge System’) for very sparse settings such as $q = 0.96$ (only 1 non-zero group out of 25; see Tables 6 and 9). The gap in performance between GRridge Iterative and group lasso is larger for large correlation ($\rho = 0.5$).

4 Group-weighted random forest

The concept of adaptive group-regularization (or, analogous, group-weighting) can be generalized to other classifiers, also to some of very different nature than logistic ridge regression. Below, we discuss one of these: the random forest classifier.

For high-dimensional applications, a standard random forest uses only $m = \mathcal{O}(\sqrt{p})$ variables per node split (Hastie et al., 2008). Typically, these variables are sampled uniformly from the entire set. Now, the idea is to weigh groups by increasing or decreasing the sampling probability according to overall importance of variables in a group. From an (initial) random forest a variable ranking can be obtained by permutation variable importance (Breiman, 2001). In short, per variable (node) i

in a tree one simply computes the difference between the original out-of-bag (OOB) error and the one obtained when the values of the variable were randomly permuted. These differences are then averaged over the trees, and serve as a metric for variable ranking. Next, for use of the method below, one may either simply set a fixed proportion q (e.g. 5%) of highest ranking variables, or impose a more advanced variable selection method based on the ranking (Doksum et al., 2008). This defines the set of top-ranking variables.

We initiate the random forest by uniform sampling of the variables. An important property of the trees in a random forest is that these are de-correlated due to bootstrapping and the limited number of variables per tree (Hastie et al., 2008). Then, we model the number of top-ranking variables $Z_t = (Z_{1t}, \dots, Z_{Gt})$ per group $g = 1, \dots, G$ in tree $t = 1, \dots, T$, given fractions $\mathbf{f}_t = (f_{1t}, \dots, f_{Gt})$, as:

$$\begin{aligned} Z_t | \mathbf{f}_t &\sim \text{MN}(m_t^q, \mathbf{f}_t) \\ \mathbf{f}_t &\sim^{\text{iid}} \mathcal{D}(\alpha), \end{aligned} \tag{1}$$

where $\mathcal{D}(\alpha)$ denotes the Dirichlet distribution with parameters $\alpha = (\alpha_1, \dots, \alpha_G)$. The second level of (1) models the variability between trees. Then, α is estimated by conventional Empirical Bayes, hence maximizing the log marginal likelihood, which is known due to the conjugacy:

$$\hat{\alpha} = \operatorname{argmax}_{\alpha} \log P(\mathbf{Z} | \alpha) = \operatorname{argmax}_{\alpha} \sum_{t=1}^T \left[\frac{\Gamma(A)}{\Gamma(m_t^q + A)} + \sum_{g=1}^G (\Gamma(z_{gt} + \alpha_g) - \Gamma(\alpha_g)) \right], \tag{2}$$

where $A = \sum_g \alpha_g$ and z_{gt} is the observed number of top-ranking group g variables in tree t . Note that even when m_t^q is small for some trees, the estimate of α is likely to be stable given the large number of trees, T . Let α_g^{init} be the proportion of group g variables in the entire set of variables. Solving (2) will render $\hat{\alpha}_g > \alpha_g^{\text{init}}$ when group g is over-represented in the set of top-ranking variables.

Then, iterate the random forest construction by *weighted* sampling of the $m = c\sqrt{p}$ variables, independently per node split. For each tree, the proportions of group $g = 1, \dots, G$ variables are generated from the $\mathcal{D}(\hat{\alpha})$ prior, which, multiplied by m , renders m_g , the total number of variables to be generated from group g . Then, within a group, those m_g variables are sampled uniformly. Finally, iteration is monitored by the out-of-bag classification error. As for the ordinary random forest classifier, constant c may be set to 1 or tuned by minimizing the OOB error.

5 R-code

Below we display the R-code used to produce the results of `GRridge` in this paper. We include some comments; more extensive documentation on the package and the functions are available via the R-interface: `??GRridge` or, e.g., `?grridge`.

5.1 First example: Improved classification by use of probe annotation

```
#library() loads the GRridge library and its dependencies, the libraries penalized and Iso
library(GRridge)
```

```

#data() loads the data and objects necessary for for 1st example: Farkas data, using
#annotation: distance to CpG. Contains objects: datcenFarkas: methylation data for cervix
#samples (arcsine-transformed beta values), respFarkas: binary response (Normal and
#Precursor) and CpGannFarkas: annotation of probes according to location
#(CpG-Island, North-Shelf, South-Shelf, North-Shore, South-Shore, Distant)

data(dataFarkas)

#For nominal input (factor), CreatePartition(vec) creates a partition of variables (probes)
#with groups according to the levels of vec. 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).

firstPartition <- CreatePartition(CpGannFarkas)
sdsF <- apply(datcenFarkas,1,sd)
secondPartition <- CreatePartition(sdsF,decreasing=FALSE,uniform=T,grsize=5000)

#Concatenate two partitions into one list

partitionFarkas <- list(cpg=firstPartition,sds=secondPartition)

#grridge() applies group-regularized ridge to data datcenFarkas, response respFarkas and
#probe grouping partitionFarkas. It recognizes automatically whether 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.
#Argument monotone =c(FALSE,TRUE) indicates that monotone increasing group-penalties
#are desired for the 2nd partition (sd-based), and not for the first one.

grFarkas <- grridge(datcenFarkas,respFarkas,partitionFarkas,monotone=c(FALSE,TRUE),
savepredobj="all")

#Predictions from all models for NEW samples. Here, illustrated on 2 samples
#of the training data.

fakenew <- datcenFarkas[,3:4]
predict.grridge(grFarkas,fakenew)

#grridge.cv() 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 <- grridge.cv(grFarkas,datcenFarkas,respFarkas,outerfold=10)
save(grFarkas,grFarkascv,file="Farkasresults.Rdata")
grFarkascv

```

```

#roc() computes the ROC curve and auc() the area-under-the-roc-curve for the probabilistic
#classifiers

cutoffs <- rev(seq(0,1,by=0.01))
rocridgeF <- roc(probs=grFarkascv[,2],true=grFarkascv[,1],cutoffs=cutoffs)
auc(rocridgeF)
rocgrridgeF <- roc(probs=grFarkascv[,3],true=grFarkascv[,1],cutoffs=cutoffs)
auc(rocgrridgeF)

plot(rocridgeF[1,],rocridgeF[2,],type="l",lty=1,ann=F,col="black")
points(rocgrridgeF[1,],rocgrridgeF[2,],type="l",lty=1,col="grey")

```

5.2 Second example: Improved diagnostic classification by use of external data

```

#data() loads objects: datcenVerlaat: data; respVerlaat: binary response; pvalFarkas:
#p-values from Farkas study; diffmeanFarkas: effect size, Precursors - Normals [Controls]
#CpGann: CpG annotation as factor

data(dataVerlaat)

#Here, for numerical input, CreatePartition() creates a partition of variables (probes) with
#100 groups, sorted in increasing order of external p-values. Here, minimal group size equals
#10, group sizes gradually increase. Use decreasing = TRUE, if a large (summary) value from
#the external data indicates more relevance (eg test statistic).

firstPartition <- CreatePartition(pvalFarkas,decreasing=FALSE,mingr=10,ngroup=100)

#Create partition of variables based on sign of the effect size in external data.
#Positive means hyper-methylated in cases.

whpos <- which(diffmeanFarkas>0)
secondPartition <- list(Pos=whpos,Neg=(1:length(diffmeanFarkas))[-whpos])

#For nominal input (factor), CreatePartition(vec) creates a partition of variables (probes)
#with groups according to the levels of vec.

thirdPartition <- CreatePartition(CpGann)

#Create sd-based partition, based on uniformly-sized groups

sds <- apply(datcenVerlaat,1,sd)
fourthPartition <- CreatePartition(sds,decreasing=FALSE, uniform=T,grsize=1000)

#Concatenate four partitions into one list

partitionsVerlaat <- list(pFarkas=firstPartition,posneg=secondPartition,
cpg=thirdPartition,sds=fourthPartition)

```

```

#grridge() applies group-regularized logistic ridge; unpenal=~0 implies no intercept.
#In this study, cases are over-sampled which may bias the intercept. Argument monotone =
#c(TRUE,FALSE,FALSE,TRUE) indicates that monotone increasing group-penalties are desired
#for the pvalFarkas-based and sd-based partitions (and not for the other two). Selection =
#TRUE indicates that CV-likelihood-based post-hoc variable selection is applied.

grVerlaat <- grridge(datcenVerlaat,respVerlaat,unpenal=~0,partitionsVerlaat,
monotone = c(TRUE,FALSE,FALSE,TRUE),selection=TRUE,savepredobj="all")

#Displays number of selected markers.

length(grVerlaat$whichsel)

#grridge.cv() performs cross-validation (here leave-one-out by default) for assessing
#predictive performance.

grVerlaatcv <- grridge.cv(grVerlaat, datcenVerlaat,respVerlaat)
save(grVerlaat,grVerlaatcv,file="Verlaatresults.Rdata")

#Display true response [,1]; ridge prediction [,2]; grridge prediction [,3];
#grridge prediction after variable selection [,4].

grVerlaatcv

#Compute ROC curves and AUCs for the probabilistic classifiers

cutoffs <- rev(seq(0,1,by=0.01))
rocridgeV <- roc(probs=grVerlaatcv[,2],true=grVerlaatcv[,1],cutoffs)
auc(rocridgeV)
rocgrridgeV <- roc(probs=grVerlaatcv[,3],true=grVerlaatcv[,1],cutoffs)
auc(rocgrridgeV)
rocgrridgeSelV <- roc(probs=grVerlaatcv[,4],true=grVerlaatcv[,1],cutoffs)
auc(rocgrridgeSelV)

plot(rocridgeV[1,],rocridgeV[2,],type="l",lty=1,ann=F,col="black")
points(rocgrridgeV[1,],rocgrridgeV[2,],type="l",lty=1,col="grey")
points(rocgrridgeSelV[1,],rocgrridgeSelV[2,],type="l",lty=2,col="grey")

```

6 Preprocessing of methylation data

Here, we give a short description of the preprocessing steps used for both example data sets. An exact description of all probe and sample quality control steps is outside the scope of this paper, but can be provided on request. After the quality control steps, probes that contain SNPs at or near the target CpG-site are removed (because these may not measure methylation). Likewise, probes vulnerable to cross-hybridization are removed. Note that these two filtering steps remove only about 1% of the probes. Next, we applied **dasen** normalization, available in the R-package **watermelon** (Pidsley et al., 2013), which uses separate background corrections for the type I and type II probes, and separate type- and channel-specific (methylated and unmethylated) quantile normalizations. Finally, the normalized beta-values were transformed to a Gaussian scale by the commonly used arcsine square root transformation.

7 Explanation on Supplementary Figures

Figure 1 displays the iterative version of the `GRridge` algorithm, with optional post-hoc variable selection.

Figure 2 displays the stabilizing effect of applying isotonic regression to the initial estimates $(\hat{\tau}_g^{\text{init}})^2$ on the resulting penalty multipliers, $\lambda'_g \propto 1/(\hat{\tau}_g)^2$. We show the \log_{10} estimates of λ'_g for the first iteration in the second example, where 100 groups of variables are based on external p -values (see Section ‘Improved diagnostic classification by use of external data’ in main document). A group for which index g (x-axis) is small corresponds to relatively small p -values. The isotonic regression function f on g forces decreasing values of $(\hat{\tau}_g)^2 = \hat{f}(g)$ (and hence increasing values of λ'_g) for groups with increasing external p -values.

Figure 3 displays the ROC curves of the logistic ridge regression for the first example, applied to either standardized or unstandardized variables. In addition, it shows the `GRridge` ROC curve with automatic standardization using a sample variance based partition of the variables.

Figure 4 displays the histogram of the variable-wise penalty multipliers for the partition based on external p -values in the second example. It illustrates the relevance of the partition, because the range across variables is large.

Let $\beta^+ = (|\hat{\beta}_1|, \dots, |\hat{\beta}_p|)$ be the vector of absolute values of regression coefficient. Then define the scaled cumulative sorted absolute values by

$$c_\ell = \frac{\sum_{k=1}^{\ell} \beta_{(k)}^+}{\sum_{j=1}^p \beta_j^+}, \quad (3)$$

where $\beta_{(k)}^+$ is the k th smallest element of β^+ . Figure 5 displays $c_\ell, \ell = 1, \dots, p$ for the original ridge coefficient estimates and the group-regularized ones from the second example. We clearly observe that the curve increases much steeper at the end for group-regularized ridge than for ordinary ridge, indicating that the most extreme regression coefficients are relatively much larger for group-regularized ridge than for ordinary ridge.

Figure 6 illustrates the CVL-profile for the purpose of variable selection in the second data example. The number of selected variables, 42, is the smallest number of variables within a 1% margin of the maximum value of the CVL.

References

- Breiman, L. (2001). Random forests. *Machine Learning*, **45**, 5–32.
- Cule, E. et al. (2011). Significance testing in ridge regression for genetic data. *BMC Bioinf.*, **12**, 372.
- Doksum, K. et al. (2008). Nonparametric variable selection: the EARTH algorithm. *J. Amer. Statist. Assoc.*, **103**, 1609–1620.
- Gruber, M. (1998). *Improving Efficiency by Shrinkage: The James–Stein and Ridge Regression Estimators*. Statistics: A Series of Textbooks and Monographs.
- Hastie, T. et al. (2008). *The elements of statistical learning, 2nd ed.* Springer, New York.

- Hoerl, A.E. and Kennard, R.W. (1970). Ridge regression: Biased estimation for nonorthogonal problems. *Technometrics*, **12**, 55–67.
- Meier, L. et al. (2008). The group Lasso for logistic regression. *J. R. Stat. Soc. Ser. B Stat. Methodol.*, **70**, 53–71.
- Pidsley, R. et al. (2013). A data-driven approach to preprocessing Illumina 450K methylation array data. *BMC Genomics*, **14**, 293.
- Waldron, L. et al. (2011). Optimized application of penalized regression methods to diverse genomic data. *Bioinformatics*, **27**, 3399–3406.

8 Supplementary Figures and Tables

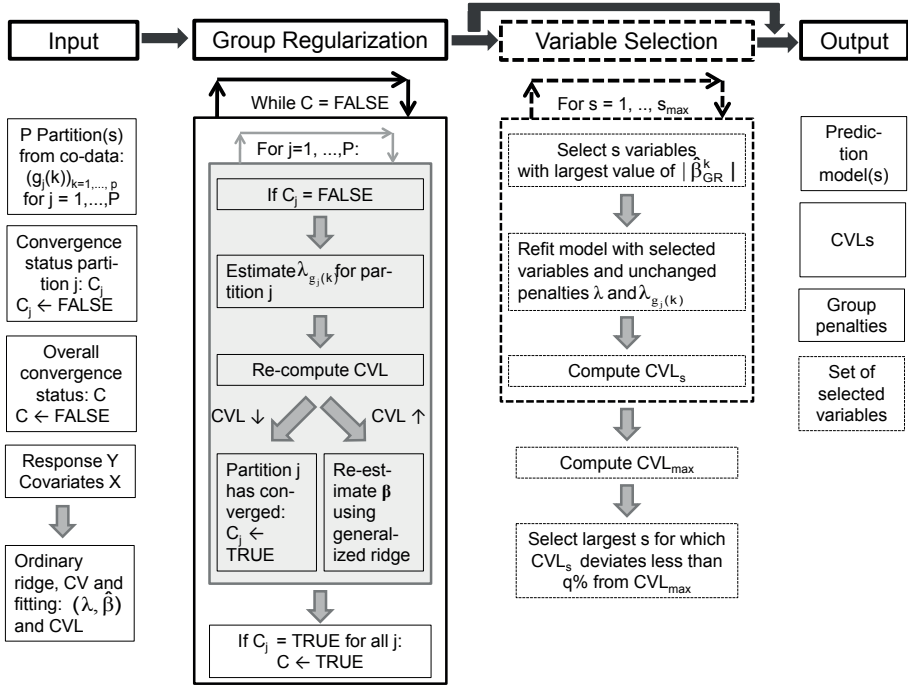


Figure 1: The iterative version of **GRridge** algorithm, including optional post-hoc variable selection. Within the Group Regularization module, the outer black box represents the re-penalization iteration, whereas the inner grey box represents the partition iteration. ‘CVL’ refers to cross-validated likelihood.

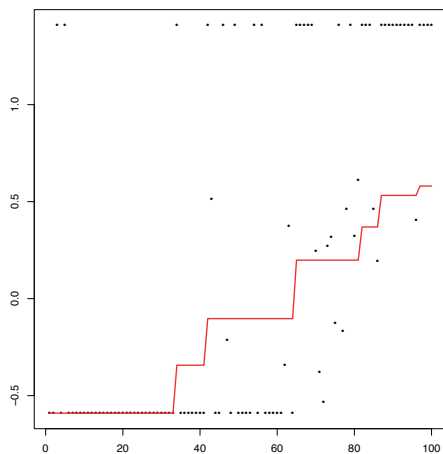


Figure 2: Effect of isotonic regression. X-axis: index of variable group, sorted in increasing order of external p -values. Y-axis: \log_{10} of the estimate of the group penalty multiplier λ'_g ; Black: initial estimates and Red: estimates based on isotonic regression.

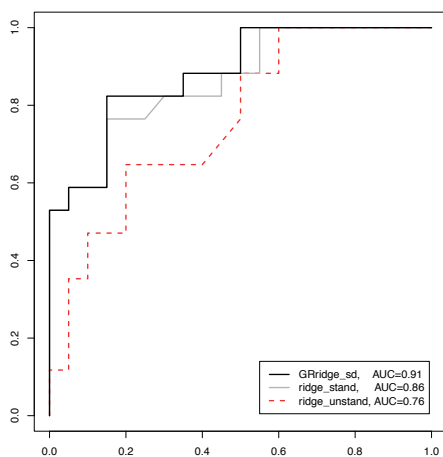


Figure 3: ROC curves for first example, automatic standardization. Group-regularized ridge (GRridge_sd), ridge on standardized variables ridge_stand and ridge on unstandardized variables ridge_unstand. X-axis: False Positive Rate, y-axis: True Positive Rate.

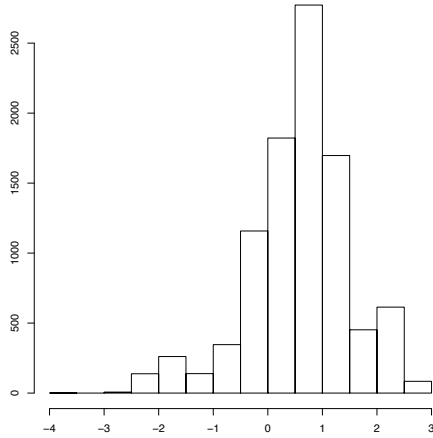


Figure 4: Histogram of p -value based penalty multipliers (\log_{10} scale) for second example.

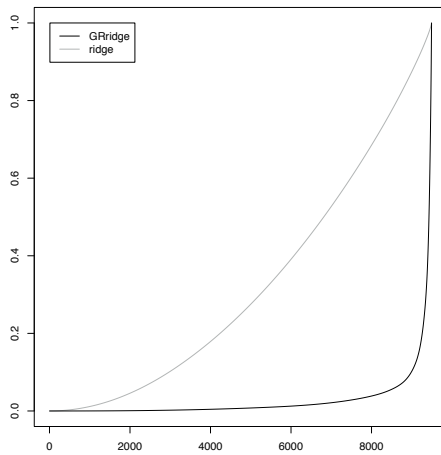


Figure 5: Order index ℓ (x-axis) against scaled cumulative sorted absolute values of regression coefficients $c_\ell(3)$ for ridge and group-regularized ridge

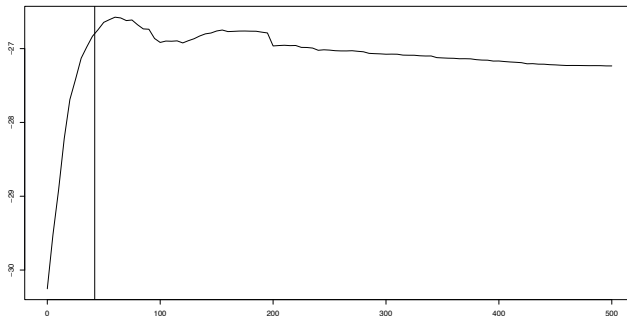


Figure 6: Number of variables included in the model(x-axis) versus cross-validated log-likelihood (CVL; y-axis). Vertical line denotes the number of selected variables, 42, allowing a 1% margin.

$(\bar{\beta}, \rho)$	f	q	Ridge	GRridge System	GRridge Iterative	Group- lasso
(0.1, 0.1)	1.3	0	0.69 (0.65)	0.76 (0.66)	0.76 (0.66)	0.59 (0.50)
		0.7	0.67 (0.67)	0.77 (0.72)	0.78 (0.71)	0.68 (0.64)
		0.9	0.67 (0.63)	0.86 (0.83)	0.86 (0.83)	0.79 (0.78)
	1.6	0	0.69 (0.69)	0.81 (0.74)	0.81 (0.74)	0.70 (0.62)
		0.7	0.69 (0.66)	0.81 (0.79)	0.80 (0.78)	0.67 (0.66)
		0.9	0.69 (0.66)	0.86 (0.83)	0.87 (0.84)	0.81 (0.78)
	2	0	0.70 (0.67)	0.81 (0.80)	0.81 (0.81)	0.74 (0.67)
		0.7	0.69 (0.67)	0.83 (0.74)	0.82 (0.74)	0.71 (0.65)
		0.9	0.69 (0.63)	0.90 (0.80)	0.90 (0.80)	0.81 (0.77)
(0.01, 0.5)	1.3	0	0.67 (0.61)	0.67 (0.62)	0.68 (0.62)	0.58 (0.50)
		0.7	0.65 (0.59)	0.69 (0.59)	0.68 (0.62)	0.53 (0.50)
		0.9	0.63 (0.61)	0.67 (0.61)	0.70 (0.65)	0.66 (0.50)
	1.6	0	0.71 (0.66)	0.77 (0.71)	0.77 (0.75)	0.67 (0.59)
		0.7	0.73 (0.68)	0.81 (0.66)	0.81 (0.75)	0.64 (0.61)
		0.9	0.70 (0.67)	0.81 (0.75)	0.81 (0.79)	0.68 (0.65)
	2	0	0.74 (0.72)	0.85 (0.72)	0.85 (0.80)	0.68 (0.67)
		0.7	0.72 (0.67)	0.72 (0.71)	0.82 (0.72)	0.70 (0.68)
		0.9	0.72 (0.67)	0.88 (0.69)	0.88 (0.84)	0.77 (0.70)

Table 4: Median (minimum) AUCs for four methods when $p_g = 200, G = 10, p = G * p_g = 2000$; $\bar{\beta}, \rho, f$ and q refer to: mean differential signal, correlation strength, relative differential signal for consecutive groups, and group-sparsity, respectively.

$(\bar{\beta}, \rho)$	f	q	Ridge	GRridge System	GRridge Iterative	Group- lasso
(0.1, 0.1)	1.3	0	0.73 (0.71)	0.77 (0.75)	0.76 (0.75)	0.67 (0.65)
		0.7	0.73 (0.72)	0.82 (0.79)	0.80 (0.79)	0.70 (0.69)
		0.9	0.73 (0.69)	0.91 (0.88)	0.91 (0.88)	0.88 (0.86)
	1.6	0	0.72 (0.70)	0.82 (0.81)	0.82 (0.82)	0.74 (0.74)
		0.7	0.70 (0.68)	0.82 (0.75)	0.83 (0.75)	0.76 (0.73)
		0.9	0.71 (0.71)	0.89 (0.88)	0.90 (0.89)	0.88 (0.87)
	2	0	0.74 (0.71)	0.87 (0.82)	0.87 (0.82)	0.80 (0.75)
		0.7	0.70 (0.69)	0.86 (0.84)	0.86 (0.83)	0.81 (0.77)
		0.9	0.72 (0.70)	0.90 (0.70)	0.90 (0.89)	0.88 (0.86)
(0.01, 0.5)	1.3	0	0.79 (0.75)	0.80 (0.78)	0.81 (0.79)	0.69 (0.64)
		0.7	0.77 (0.76)	0.85 (0.80)	0.85 (0.82)	0.72 (0.69)
		0.9	0.75 (0.71)	0.90 (0.78)	0.90 (0.87)	0.86 (0.84)
	1.6	0	0.78 (0.74)	0.88 (0.81)	0.88 (0.81)	0.76 (0.75)
		0.7	0.79 (0.77)	0.88 (0.81)	0.90 (0.87)	0.79 (0.75)
		0.9	0.78 (0.77)	0.94 (0.94)	0.94 (0.94)	0.90 (0.89)
	2	0	0.77 (0.75)	0.91 (0.85)	0.90 (0.85)	0.81 (0.79)
		0.7	0.80 (0.78)	0.90 (0.80)	0.92 (0.89)	0.81 (0.80)
		0.9	0.81 (0.77)	0.96 (0.94)	0.96 (0.95)	0.92 (0.91)

Table 5: Median (minimum) AUCs for four methods when $p_g = 500, G = 10, p = G * p_g = 5000$; $\bar{\beta}, \rho, f$ and q refer to: mean differential signal, correlation strength, relative differential signal for consecutive groups, and group-sparsity, respectively.

$(\bar{\beta}, \rho)$	f	q	Ridge	GRridge System	GRridge Iterative	Group- lasso
(0.1, 0.1)	1.3	0.88	0.64 (0.61)	0.73 (0.63)	0.81 (0.81)	0.71 (0.64)
		0.96	0.65 (0.63)	0.77 (0.75)	0.90 (0.85)	0.89 (0.84)
	1.6	0.88	0.64 (0.64)	0.75 (0.71)	0.83 (0.77)	0.74 (0.71)
		0.96	0.66 (0.65)	0.77 (0.76)	0.91 (0.89)	0.88 (0.87)
	2	0.88	0.66 (0.63)	0.74 (0.65)	0.85 (0.83)	0.81 (0.77)
		0.96	0.64 (0.62)	0.81 (0.76)	0.91 (0.91)	0.90 (0.87)
(0.01, 0.5)	1.3	0.88	0.70 (0.63)	0.83 (0.71)	0.87 (0.78)	0.76 (0.71)
		0.96	0.71 (0.65)	0.90 (0.70)	0.96 (0.94)	0.92 (0.90)
	1.6	0.88	0.72 (0.70)	0.81 (0.73)	0.93 (0.85)	0.79 (0.78)
		0.96	0.71 (0.70)	0.80 (0.79)	0.97 (0.94)	0.92 (0.90)
	2	0.88	0.71 (0.70)	0.89 (0.76)	0.91 (0.89)	0.84 (0.83)
		0.96	0.70 (0.68)	0.90 (0.70)	0.96 (0.96)	0.93 (0.91)

Table 6: Median (minimum) AUCs for four methods when $p_g = 500, G = 25, p = G * p_g = 12500$; $\bar{\beta}, \rho, f$ and q refer to: mean differential signal, correlation strength, relative differential signal for consecutive groups, and group-sparsity, respectively.

$(\bar{\beta}, \rho)$	f	q	Ridge	GRridge System	GRridge Iterative	Group- lasso
(0.1, 0.1)	1.3	0	0.184 (0.197)	0.163 (0.209)	0.163 (0.216)	0.205 (0.213)
		0.7	0.183 (0.183)	0.152 (0.170)	0.149 (0.177)	0.181 (0.194)
		0.9	0.173 (0.193)	0.100 (0.119)	0.098 (0.109)	0.129 (0.132)
	1.6	0	0.195 (0.201)	0.156 (0.177)	0.158 (0.177)	0.195 (0.220)
		0.7	0.188 (0.197)	0.149 (0.151)	0.150 (0.155)	0.196 (0.213)
		0.9	0.184 (0.192)	0.118 (0.129)	0.108 (0.115)	0.141 (0.143)
	2	0	0.198 (0.208)	0.157 (0.175)	0.151 (0.159)	0.185 (0.201)
		0.7	0.199 (0.203)	0.155 (0.182)	0.148 (0.178)	0.197 (0.204)
		0.9	0.206 (0.219)	0.158 (0.219)	0.108 (0.219)	0.149 (0.164)
(0.01, 0.5)	1.3	0	0.098 (0.120)	0.099 (0.120)	0.107 (0.124)	0.120 (0.173)
		0.7	0.095 (0.116)	0.084 (0.116)	0.095 (0.116)	0.108 (0.116)
		0.9	0.075 (0.079)	0.073 (0.079)	0.064 (0.088)	0.066 (0.078)
	1.6	0	0.112 (0.123)	0.090 (0.111)	0.092 (0.125)	0.122 (0.164)
		0.7	0.100 (0.125)	0.061 (0.123)	0.067 (0.094)	0.123 (0.131)
		0.9	0.108 (0.111)	0.058 (0.073)	0.055 (0.094)	0.103 (0.111)
	2	0	0.120 (0.161)	0.073 (0.161)	0.072 (0.121)	0.133 (0.160)
		0.7	0.118 (0.147)	0.118 (0.144)	0.072 (0.159)	0.129 (0.141)
		0.9	0.114 (0.127)	0.041 (0.127)	0.042 (0.059)	0.105 (0.121)

Table 7: Median (maximum) mean Brier residuals for four methods when $p_g = 200, G = 10, p = G * p_g = 2000$; $\bar{\beta}, \rho, f$ and q refer to: mean differential signal, correlation strength, relative differential signal for consecutive groups, and group-sparsity, respectively.

$(\bar{\beta}, \rho)$	f	q	Ridge	GRridge System	GRridge Iterative	Group- lasso
(0.1, 0.1)	1.3	0	0.203 (0.211)	0.183 (0.195)	0.188 (0.193)	0.219 (0.227)
		0.7	0.196 (0.202)	0.159 (0.173)	0.164 (0.179)	0.205 (0.214)
		0.9	0.187 (0.217)	0.103 (0.125)	0.105 (0.115)	0.115 (0.123)
	1.6	0	0.200 (0.215)	0.161 (0.181)	0.160 (0.171)	0.192 (0.196)
		0.7	0.213 (0.221)	0.170 (0.196)	0.166 (0.189)	0.191 (0.197)
		0.9	0.202 (0.215)	0.119 (0.157)	0.113 (0.125)	0.123 (0.126)
	2	0	0.208 (0.218)	0.143 (0.162)	0.142 (0.161)	0.173 (0.193)
		0.7	0.210 (0.226)	0.153 (0.165)	0.145 (0.163)	0.172 (0.182)
		0.9	0.214 (0.227)	0.132 (0.221)	0.116 (0.129)	0.129 (0.140)
(0.01, 0.5)	1.3	0	0.135 (0.160)	0.138 (0.150)	0.135 (0.147)	0.170 (0.190)
		0.7	0.129 (0.147)	0.083 (0.130)	0.086 (0.127)	0.159 (0.175)
		0.9	0.121 (0.129)	0.042 (0.102)	0.044 (0.062)	0.063 (0.073)
	1.6	0	0.152 (0.189)	0.108 (0.143)	0.104 (0.142)	0.154 (0.165)
		0.7	0.143 (0.152)	0.095 (0.135)	0.089 (0.100)	0.139 (0.153)
		0.9	0.141 (0.144)	0.037 (0.042)	0.036 (0.041)	0.066 (0.077)
	2	0	0.164 (0.170)	0.089 (0.126)	0.089 (0.125)	0.140 (0.149)
		0.7	0.155 (0.161)	0.093 (0.145)	0.076 (0.094)	0.144 (0.147)
		0.9	0.125 (0.145)	0.034 (0.048)	0.037 (0.048)	0.066 (0.083)

Table 8: Median (maximum) mean Brier residuals for four methods when $p_g = 500, G = 10, p = G * p_g = 5000$; $\bar{\beta}, \rho, f$ and q refer to: mean differential signal, correlation strength, relative differential signal for consecutive groups, and group-sparsity, respectively.

$(\bar{\beta}, \rho)$	f	q	Ridge	GRridge System	GRridge Iterative	Group- lasso
(0.1, 0.1)	1.3	0.88	0.235 (0.241)	0.230 (0.241)	0.175 (0.186)	0.212 (0.233)
		0.96	0.230 (0.246)	0.190 (0.231)	0.127 (0.146)	0.135 (0.156)
	1.6	0.88	0.239 (0.246)	0.225 (0.230)	0.174 (0.215)	0.214 (0.218)
		0.96	0.229 (0.239)	0.199 (0.217)	0.117 (0.134)	0.139 (0.141)
	2	0.88	0.234 (0.240)	0.215 (0.240)	0.156 (0.170)	0.177 (0.193)
		0.96	0.236 (0.238)	0.193 (0.210)	0.120 (0.128)	0.126 (0.145)
(0.01, 0.5)	1.3	0.88	0.210 (0.220)	0.151 (0.220)	0.120 (0.220)	0.183 (0.211)
		0.96	0.182 (0.199)	0.090 (0.196)	0.040 (0.058)	0.087 (0.095)
	1.6	0.88	0.195 (0.212)	0.163 (0.190)	0.093 (0.148)	0.172 (0.186)
		0.96	0.195 (0.209)	0.150 (0.183)	0.042 (0.086)	0.093 (0.102)
	2	0.88	0.209 (0.226)	0.123 (0.182)	0.119 (0.131)	0.154 (0.161)
		0.96	0.204 (0.216)	0.113 (0.201)	0.056 (0.067)	0.096 (0.100)

Table 9: Median (maximum) mean Brier residuals for four methods when $p_g = 500, G = 25, p = G * p_g = 12500$; $\bar{\beta}, \rho, f$ and q refer to: mean differential signal, correlation strength, relative differential signal for consecutive groups, and group-sparsity, respectively.