

Response-guided Principal Component Classification

Duncan Bennett
Adviser : Hao Helen Zhang

Nov 2020

Abstract

Response-guided principal component regression (RgPCR) is a generalization of ridge regression. This method improves upon principal component regression (PCR) by taking the response values into account during variable selection. In this paper, we will modify RgPCR for binary classification problems using ideas from logistic regression. This technique is called *response-guided principal component classification (RgPCC)*.

Contents

| | | |
|----------|---|-----------|
| 1 | Introduction | 1 |
| 2 | Literature Review | 2 |
| 2.1 | PCA | 2 |
| 2.2 | RgPCR | 2 |
| 2.3 | Logistic Regression | 3 |
| 3 | Response-guided Principal Component Classification | 4 |
| 3.1 | Method | 4 |
| 3.2 | Iterative Approximations | 5 |
| 3.3 | Algorithm for LASSO penalty | 6 |
| 4 | Numerical Studies | 7 |
| 4.1 | Creation of Simulated Data | 7 |
| 4.2 | Performance of RgPCC on Simulated Data | 8 |
| 4.3 | Performance of RgPCC on Realworld Data | 10 |
| 5 | Further Research | 11 |
| 5.1 | Speed | 11 |
| 5.2 | C_p -type Statistic for Parameter Tuning | 11 |
| 5.3 | Further Generalization | 11 |
| 6 | Figures | 13 |

1 Introduction

Principal component regression (PCR) was introduced (by Jeffers, 1967) to deal with multicollinearity. This method can achieve dimension reduction and improve prediction performance compared to

ordinary least squares. However, its downfall in regression is that the principal components depend solely on the design \mathbf{X} and in this sense the variable selection is "blind", as it does not take the response into account. In early 2020 Lang and Zou [2] introduced Response-guided Principal Component Regression (RgPCR) to remedy the "blind" selection of PCR. This is done by replacing the hard-thresholding of PCR with soft-thresholding via a penalty function. The result is that both the variance of the predictors and the association with the response of principal components is taken into account during thresholding.

In this paper, we will combine RgPCR with logistic regression for binary classification. To do this we optimized a penalized log likelihood function by quadratically approximating and taking advantage of the principal components of "psuedo data" in the resulting expression. The result is a principal component classification algorithm that takes the response into account during variable selection.

In section 2 we will give our motivation and cover the background knowledge necessary. In section 3 we will derive the RgPCC algorithm. In sections 4 and 5 we will compare the performance of RgPCC against other methods on simulated data and then on realworld data. Lastly, in section 6 and 7 we will propose further topics of study and summarize our results.

2 Literature Review

Here we describe two supervised learning methods that have motivated RgPCC. First, we summarized RgPCR, a method developed by Lang and Zou [2]. Then we summarized logistic regression for binary classification.

2.1 PCA

2.2 RgPCR

RgPCR is a generalization of ridge regression which we will briefly summarize. Let $\mathbf{X} = \mathbf{U}\mathbf{D}\mathbf{V}^T$ be the singular value decomposition of an $N \times p$ design matrix \mathbf{X} . Then with response y we can write the ridge regression solution as

$$\hat{\boldsymbol{\beta}}^{\text{ridge}} = \arg \min_{\boldsymbol{\beta}} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_2^2 + \lambda \sum_{j=1}^p \beta_j^2. \quad (1)$$

We may rewrite this as

$$\hat{\boldsymbol{\gamma}}^{\text{ridge}} = \arg \min_{\boldsymbol{\gamma}} \|\mathbf{y} - \mathbf{U}\boldsymbol{\gamma}\|_2^2 + \lambda \sum_{j=1}^p \frac{\gamma_j^2}{d_j^2}. \quad (2)$$

where $\boldsymbol{\gamma} = \mathbf{D}\mathbf{V}^T\boldsymbol{\beta}$ and d_j is the j^{th} diagonal entry of \mathbf{D} . Hence, ridge regression can be viewed as a weighted L_2 penalized regression in the space of principal components. The RgPCR generalization comes from replacing the L_2 penalization with an arbitrary non-decreasing function $p_\lambda(\cdot)$. Some examples of these are LASSO, SCAD and MCP. In general, we write this RgPCR solution as

$$\hat{\boldsymbol{\gamma}} = \arg \min_{\boldsymbol{\gamma}} \|\mathbf{y} - \mathbf{U}\boldsymbol{\gamma}\|_2^2 + \sum_{j=1}^p p_\lambda \left(\frac{\gamma_j}{d_j} \right), \quad (3)$$

There are three main advantages to (3). First, the regularization helps prevent overfitting. Secondly, the use of PCA and the orthogonality of U allows for this minimizer to be calculated component-wise. Lastly, the penalty allows the variable selection to be determined by large variance in the predictors and the association with the response of principal components (we can see this by the presence of $\hat{\gamma}_j^{\text{ols}}$ in the following).

The orthogonality of \mathbf{U} allows us to solve (3) componentwise. In particular,

$$\hat{\gamma}_j = \arg \min_{\gamma_j} (\hat{\gamma}_j^{\text{ols}} - \gamma_j)^2 + p_\lambda \left(\frac{\gamma_j}{d_j} \right). \quad (4)$$

where $\hat{\gamma}_j^{\text{ols}} = \mathbf{y}^T U_j$ (equation (4) follows from the orthogonality of \mathbf{U} and a calculation can be found in [2]). When $p_\lambda(t) = |t|$, the LASSO penalty, then we can find a closed form of (4),

$$\hat{\gamma}_j^{\text{lasso}} = \left(|\hat{\gamma}_j^{\text{ols}}| - \frac{\lambda}{2d_j} \right)^+ \cdot \text{sgn}(\hat{\gamma}_j^{\text{ols}}). \quad (5)$$

We can recover the solution in the original coordinates by

$$\hat{\boldsymbol{\beta}} = \sum_{j=1}^p \tilde{V}_j \frac{\hat{\gamma}_j}{d_j} \quad (6)$$

2.3 Logistic Regression

Logistic regression is a common technique for classification. For now we consider binary classification. Let \mathbf{X} be an $N \times p$ design matrix with \mathbf{y} the $N \times 1$ classifications. Then if we assume that

$$\log \left(\frac{\Pr(G = 1 | X = x)}{\Pr(G = 0 | X = x)} \right) = \boldsymbol{\beta}^T x \quad (7)$$

the log-odds are linear, then we may calculate the conditional probability

$$\Pr(G = 1 | X = x) = \frac{\exp(\boldsymbol{\beta}^T x)}{1 + \exp(\boldsymbol{\beta}^T x)}. \quad (8)$$

Note that we may compute all probabilities for the data simultaneously by

$$\mathbf{p} = \frac{\exp(\mathbf{X}\boldsymbol{\beta})}{1 + \exp(\mathbf{X}\boldsymbol{\beta})}. \quad (9)$$

where p_i is the probability that x_i is in class 1.

To fit such a line for the log-odds (7), we fit by maximizing the log-likelihood,

$$\ell(\mathbf{X}, \boldsymbol{\beta}) = \sum_{i=1}^N [y_i \log(p(x_i; \boldsymbol{\beta})) + (1 - y_i) \log(1 - p(x_i; \boldsymbol{\beta}))] \quad (10)$$

This maximization can be found in [1] and has a convenient interpretation. When maximizing by the Newton-Raphson method, each iteration can be interpreted as solving a weighted least squares problem. That is

$$\boldsymbol{\beta}^{\text{new}} = \boldsymbol{\beta}^{\text{old}} + (\mathbf{X}^T \mathbf{W} \mathbf{X})^{-1} \mathbf{X}^T (\mathbf{y} - \mathbf{p}) \quad (11)$$

$$= (\mathbf{X}^T \mathbf{W} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{W} \mathbf{z} \quad (12)$$

where

$$\mathbf{p} = \frac{\exp(\mathbf{X}\boldsymbol{\beta})}{1 + \exp(\mathbf{X}\boldsymbol{\beta})} \quad (13)$$

$$\mathbf{W} = \text{diag}[p_i(1 - p_i)] \quad (14)$$

$$\mathbf{z} = \mathbf{X}\boldsymbol{\beta} + \mathbf{W}^{-1}(\mathbf{y} - \mathbf{p}). \quad (15)$$

In particular, $\boldsymbol{\beta}^{\text{new}}$ is the solution to a weighted least squares problem with response \mathbf{z} , design \mathbf{X} and weights \mathbf{W} .

For our purposes, it will be more convenient to state this as the following equivalent problem, $\boldsymbol{\beta}^{\text{new}}$ estimates the minimizer of

$$\boldsymbol{\beta}^{\text{new}} = \arg \min_{\boldsymbol{\beta}} \|\mathbf{W}^{1/2}(\mathbf{z} - \mathbf{X}\boldsymbol{\beta})\|^2 \quad (16)$$

3 Response-guided Principal Component Classification

First, a note on notation. We let p denote the number of predictors in our data, p_λ to a non-decreasing penalty function and p_i to be the i^{th} component of a vector of probabilities \mathbf{p} .

Suppose we have a classification problem with design \mathbf{X} and classes \mathbf{y} . In logistic regression we fit by maximizing the log-likelihood, we take the opposite of this and penalize the coefficients as such,

$$\boldsymbol{\beta}^{\text{RgPCC}} = \arg \min_{\boldsymbol{\beta}} -\ell(\mathbf{X}, \boldsymbol{\beta}) + \sum_{j=1}^p p_\lambda(\beta_j) \quad (17)$$

3.1 Method

To solve (17), we can approximate the log-likelihood $\ell(\mathbf{X}, \boldsymbol{\beta})$ with a quadratic approximation centered at $\boldsymbol{\beta}^*$ (the minimizer of $-\ell(\mathbf{X}, \boldsymbol{\beta})$). This quadratic approximation can be interpreted as the least squares error

$$-\ell(\mathbf{X}, \boldsymbol{\beta}) \approx \|\mathbf{W}^{1/2}(\mathbf{z} - \mathbf{X}\boldsymbol{\beta})\|^2 = \|\tilde{\mathbf{y}} - \tilde{\mathbf{X}}\boldsymbol{\beta}\|^2 \quad (18)$$

where \mathbf{W} , and \mathbf{z} are as in (26 - 28) with (26) evaluated at $\boldsymbol{\beta} = \boldsymbol{\beta}^*$ and $\mathbf{W}^{1/2}\mathbf{X} = \tilde{\mathbf{X}}$, $\tilde{\mathbf{y}} = \mathbf{W}^{1/2}\mathbf{z}$. The derivation of this approximation can be found in [3]. Note that, we may often refer to $\tilde{\mathbf{X}}$ and $\tilde{\mathbf{y}}$ as the *pseudo data* and *pseudo response* respectively.

Equation (17) can be rewritten as an approximate RgPCR problem

$$\arg \min_{\boldsymbol{\beta}} \|\tilde{\mathbf{y}} - \tilde{\mathbf{X}}\boldsymbol{\beta}\|^2 + \sum_{j=1}^p p_{\lambda}(\beta_j) \quad (19)$$

For particular choices of p_{λ} , this may be solved using techniques from [2]. That is, by viewing the above in terms of principal components we may solve the equivalent problem

$$\arg \min_{\boldsymbol{\gamma}} \|\tilde{\mathbf{y}} - \tilde{\mathbf{U}}\boldsymbol{\gamma}\|^2 + \sum_{j=1}^p p_{\lambda}\left(\frac{\gamma_j}{d_j}\right) \quad (20)$$

where $\tilde{\mathbf{X}} = \tilde{\mathbf{U}}\tilde{\mathbf{D}}\tilde{\mathbf{V}}^T$ is the SVD decomposition of $\tilde{\mathbf{X}}$ and $\boldsymbol{\gamma} = \tilde{\mathbf{D}}\tilde{\mathbf{V}}^T\boldsymbol{\beta}$. The matrices $\tilde{\mathbf{U}}$ and $\tilde{\mathbf{V}}$ are orthogonal matrices and are $N \times p$ and $p \times p$ matrices respectively. Let d_j denotes the diagonal entries of $\tilde{\mathbf{D}}$ for $1 \leq j \leq p$. For our description here, we assume that $\tilde{\mathbf{X}}$ has full rank, but the following still holds with a thin SVD decomposition.

As with RgPCR, we may take advantage of the orthogonality of $\tilde{\mathbf{U}}$ and optimize componentwise. That is, with $\hat{\gamma}_j^{\text{ols}} = \mathbf{y}^T \tilde{\mathbf{U}}_j$,

$$\arg \min_{\gamma_j} (\hat{\gamma}_j^{\text{ols}} - \gamma_j)^2 + p_{\lambda}\left(\frac{\gamma_j}{d_j}\right). \quad (21)$$

From [2], when $p_{\lambda}(t) = |t|$ is the LASSO penalty, we have the closed form solution of (21)

$$\hat{\gamma}_j = \left(|\hat{\gamma}_j^{\text{ols}}| - \frac{\lambda}{d_j} \right)^+ \cdot \text{sgn}(\hat{\gamma}_j^{\text{ols}}), \quad (22)$$

where a^+ denotes the positive real part of a real number a . We can then rewrite this in the original coordinates using

$$\hat{\boldsymbol{\beta}} = \sum_{j=1}^p \tilde{\mathbf{V}}_j \frac{\hat{\gamma}_j}{d_j}. \quad (23)$$

3.2 Iterative Approximations

In practice, a single quadratic approximation of $-\ell(\mathbf{X}, \boldsymbol{\beta})$ may not produce a minimum close to the true minimum. Just as in logistic regression, we take the quadratic approximations iteratively to improve our approximation of the minimum. That is, for some approximate minimizer $\boldsymbol{\beta}^{(n)}$, we solve

$$\arg \min_{\boldsymbol{\beta}} -\ell(\mathbf{X}, \boldsymbol{\beta}) + \sum_{j=1}^p p_{\lambda}(\beta_j) \quad (24)$$

which can be approximated by

$$-\ell(\mathbf{X}, \boldsymbol{\beta}) \approx \|\tilde{\mathbf{y}}^{(n)} - \tilde{\mathbf{X}}^{(n)}\boldsymbol{\beta}\|^2 \quad (25)$$

where

$$\mathbf{p}^{(n)} = \frac{\exp(\mathbf{X}\boldsymbol{\beta}^{(n)})}{1 + \exp(\mathbf{X}\boldsymbol{\beta}^{(n)})} \quad (26)$$

$$\mathbf{W}^{(n)} = \text{diag}[p_i^{(n)}(1 - p_i^{(n)})] \quad (27)$$

$$\mathbf{z}^{(n)} = \mathbf{X}\boldsymbol{\beta}^{(n)} + (\mathbf{W}^{(n)})^{-1}(\mathbf{y} - \mathbf{p}^{(n)}). \quad (28)$$

$$\tilde{\mathbf{X}}^{(n)} = (\mathbf{W}^{(n)})^{1/2}\mathbf{X} \quad (29)$$

$$\tilde{\mathbf{y}}^{(n)} = (\mathbf{W}^{(n)})^{1/2}\mathbf{z}^{(n)} \quad (30)$$

As with the first iteration, we can solve

$$\arg \min_{\boldsymbol{\beta}} \|\tilde{\mathbf{y}}^{(n)} - \tilde{\mathbf{X}}^{(n)}\boldsymbol{\beta}\|^2 + \sum_{j=1}^p p_{\lambda}(\beta_j) \quad (31)$$

to estimate the solution to (25). Then, as with the first iteration, we can express the above using principal components

$$\arg \min_{\boldsymbol{\gamma}} \|\tilde{\mathbf{y}}^{(n)} - \tilde{\mathbf{U}}^{(n)}\boldsymbol{\gamma}\|^2 + \sum_{j=1}^p p_{\lambda}\left(\frac{\gamma_j}{d_j^{(n)}}\right) \quad (32)$$

and take advantage of the orthogonality of $\tilde{\mathbf{U}}$ so solve componentwise

$$\arg \min_{\gamma_j} (\hat{\gamma}_j^{\text{ols},(n)} - \gamma_j)^2 + p_{\lambda}\left(\frac{\gamma_j}{d_j^{(n)}}\right). \quad (33)$$

3.3 Algorithm for LASSO penalty

Here, we will briefly describe the algorithm used to approximate a solution to (17) when we use the LASSO penalty. This will generalize to any penalty that has an RgPCR closed form solution.

(Step 1) Guess an initial $\hat{\boldsymbol{\beta}}^{(0)}$

(Step 2) With $\hat{\boldsymbol{\beta}}^{(0)}$ calculate the following

$$\begin{aligned} \mathbf{p}^{(0)} &= \frac{\exp(\mathbf{X}\boldsymbol{\beta}^{(0)})}{1 + \exp(\mathbf{X}\boldsymbol{\beta}^{(0)})} = \frac{\exp(\mathbf{U}\boldsymbol{\gamma}^{(0)})}{1 + \exp(\mathbf{U}\boldsymbol{\gamma}^{(0)})} \\ \mathbf{W}^{(0)} &= \text{diag}[p_i^{(0)}(1 - p_i^{(0)})] \\ \mathbf{z}^{(0)} &= \mathbf{X}\boldsymbol{\beta}^{(0)} + (\mathbf{W}^{(0)})^{-1}(\mathbf{y} - \mathbf{p}^{(0)}). \end{aligned}$$

(Step 3) Compute the solution to the RgPCR problem

$$\arg \min_{\boldsymbol{\gamma}} \|\tilde{\mathbf{y}}^{(0)} - \tilde{\mathbf{U}}^{(0)}\boldsymbol{\gamma}\|^2 + \lambda \sum_{j=1}^p \left| \frac{\gamma_j}{d_j^{(0)}} \right| \quad (34)$$

using the closed form solution

$$\hat{\gamma}_j^{(1)} = \left(|\hat{\gamma}_j^{\text{ols}}| - \frac{\lambda}{d_j} \right)^+ \cdot \text{sgn}(\hat{\gamma}_j^{\text{ols}}), \quad (35)$$

where $\hat{\gamma}_j^{\text{ols}} = \tilde{y}^{(0)} \tilde{U}_j^{(0)}$

(Step 4) repeat steps 2 and 3 using $\hat{\gamma}^{(1)}$ and similarly using $\hat{\gamma}^{(k)}$ until $\varepsilon^{(k)} = \frac{\|\beta^{(k)} - \beta^{(k+1)}\|}{\|\beta^{(k)}\|} < \varepsilon$ for some desired tolerance ε .

(Step 5) Recover the solution in the original coordinates with

$$\hat{\beta}^{(n)} = \sum_{j=1}^p \tilde{V}_j^{(n)} \frac{\hat{\gamma}_j^{(n)}}{d_j^{(n)}}.$$

4 Numerical Studies

In this section we compare RgPCC with LASSO penalty to conventional logistic regression. We first summarize how our simulated data was created and then discuss the results of the methods mentioned previously. We will consider a variety of sample sizes, dimensions and sparsities in γ and compare error rates in $\mathbf{p}, \beta, \gamma$ and classification. We will also vary which components of γ are nonzero. Again, please note that p is the number of predictors while \mathbf{p} is a vector of probabilities.

| p | N | γ sparsity |
|-----|----------|-------------------|
| 12 | 100, 200 | 1, 3, 5 |
| 50 | 100, 200 | 1, 3 |
| 80 | 100, 200 | 1 |

Figure 1: Parameters of data creation.

The convergence of our algorithm is based on the tolerance of $\varepsilon = 0.1$, which means we will terminate our Newton-Raphson method when

$$\varepsilon^{(k)} = \frac{\|\beta^{(k)} - \beta^{(k+1)}\|}{\|\beta^{(k)}\|} < 0.1.$$

4.1 Creation of Simulated Data

Let \mathbf{X} be our design matrix where each row of \mathbf{X} is independently generated from $N(0, \Sigma)$ where $\Sigma_{ij} = \rho^{|i-j|}$ with $\rho = 0.8$.

To generate our response data, let $\mu = \mathbf{X}\beta^* = \mathbf{U}\gamma^*$ where $\gamma^* = \mathbf{D}\mathbf{V}^T\beta^*$. Then the classes \mathbf{y} come from a Bernoulli distribution with parameter $\mathbf{p} = \frac{\exp(\mathbf{U}\gamma^*)}{1 + \exp(\mathbf{U}\gamma^*)} = \frac{\exp(\mu)}{1 + \exp(\mu)}$. Therefore the response data is dependent on γ^* and the data \mathbf{X} . We make a variety of choices for γ^* (denoted just as γ below) which have varying amounts of sparsity and created varying levels of linear separability.

4.2 Performance of RgPCC on Simulated Data

Here we fit RgPCC, logistic regression and PC logistic regression to our simulated data. To test the quality of the fit, we see how well our models predict the true probability vector \mathbf{p} (the parameter of the Bernoulli distribution above) over a test set of size $5N$. We measure the accuracy of our predicted $\hat{\mathbf{p}}$ using the 1-norm, 2-norm and EGKL. These are defined as follows:

$$\begin{aligned} \|v\|_1 &= \\ \|v\|_2 &= \\ EGKL(v) &= \end{aligned}$$

We fit the models using samples of size N for dimension p from Figure 1. We also vary the nonzero components of the true γ . We neglect the case where $N < p$, although there are no theoretical issues in this case and will be covered in further research.

We use the following values for the true γ :

$$\begin{aligned} \gamma_1 &= (25, 0, \dots, 0) \\ \gamma_2 &= (20, 10, 10, 0, \dots, 0) \\ \gamma_3 &= (15, 10, 5, 5, 3, 0, \dots, 0) \\ \gamma'_1 &= (0, 0, 0, 0, 0, 25, 0, \dots, 0) \\ \gamma'_2 &= (0, 0, 0, 0, 0, 20, 10, 10, 0, \dots, 0) \\ \gamma'_3 &= (0, 0, 0, 0, 0, 15, 10, 5, 5, 3, 0, \dots, 0) \end{aligned}$$

Below we show only a few cases of all the tests. The rest of the results are present in the last section but have the same conclusion as the one presented here.

| | method | sample_size | one.norm | two.norm | EGKL | class.error | gamma.size |
|----|-------------|-------------|----------|----------|----------|-------------|------------|
| 1 | log | 100 | 128.9002 | 48.2978 | 135.2297 | 0.3776 | 12 |
| 2 | pcalog | 100 | 99.0854 | 28.8548 | 38.2041 | 0.4807 | 5.14 |
| 3 | RgPCC.AIC | 100 | 117.0917 | 41.3108 | 92.3249 | 0.3775 | 10.78 |
| 4 | RgPCC.BIC | 100 | 77.461 | 20.1746 | 28.3024 | 0.3784 | 6.19 |
| 5 | RgPCC.MSE | 100 | 80.1013 | 18.9059 | 25.2497 | 0.3936 | 1.97 |
| 6 | RgPCC.pMSE | 100 | 76.0027 | 17.8124 | 23.8169 | 0.3849 | 2.81 |
| 7 | RgPCC.MSECV | 100 | 128.7374 | 48.1967 | 134.5861 | 0.3776 | 12 |
| 8 | log | 200 | 189.2017 | 49.8042 | 87.5984 | 0.389 | 12 |
| 9 | pcalog | 200 | 151.9253 | 34.7084 | 42.483 | 0.4834 | 5.16 |
| 10 | RgPCC.AIC | 200 | 177.856 | 44.8999 | 73.7718 | 0.3891 | 11.17 |
| 11 | RgPCC.BIC | 200 | 101.4748 | 17.4053 | 19.2933 | 0.3887 | 5.53 |
| 12 | RgPCC.MSE | 200 | 118.2799 | 22.4926 | 26.6875 | 0.386 | 7.63 |
| 13 | RgPCC.pMSE | 200 | 94.8438 | 15.3939 | 16.9056 | 0.3907 | 4.25 |
| 14 | RgPCC.MSECV | 200 | 155.4774 | 35.7876 | 51.6239 | 0.3893 | 10.2 |

Figure 2: For γ'_2 and $p = 12$.

| | method | sample_size | one.norm | two.norm | EGKL | class.error | gamma.size |
|----|-------------|-------------|----------|----------|----------|-------------|------------|
| 1 | log | 100 | 219.7315 | 119.8399 | 3337.828 | 0.4469 | 50 |
| 2 | pcalog | 100 | 147.4434 | 62.2117 | 241.1183 | 0.412 | 16.61 |
| 3 | RgPCC.AIC | 100 | 168.0875 | 77.802 | 288.5035 | 0.4282 | 32.88 |
| 4 | RgPCC.BIC | 100 | 88.2784 | 25.8719 | 37.0038 | 0.4077 | 5.7 |
| 5 | RgPCC.MSE | 100 | 86.9706 | 25.1717 | 35.639 | 0.4087 | 5.55 |
| 6 | RgPCC.pMSE | 100 | 86.0024 | 24.6228 | 34.8048 | 0.4066 | 5.33 |
| 7 | RgPCC.MSECV | 100 | 95.7481 | 30.1532 | 46.0895 | 0.4088 | 8.06 |
| 8 | log | 200 | 295.6583 | 112.8611 | 392.142 | 0.4326 | 50 |
| 9 | pcalog | 200 | 204.7182 | 58.4586 | 107.592 | 0.4175 | 18.21 |
| 10 | RgPCC.AIC | 200 | 148.4187 | 33.7531 | 45.4301 | 0.4129 | 10.64 |
| 11 | RgPCC.BIC | 200 | 127.0452 | 25.9331 | 31.4919 | 0.4123 | 7.18 |
| 12 | RgPCC.MSE | 200 | 174.8136 | 44.7253 | 68.5684 | 0.4161 | 19.38 |
| 13 | RgPCC.pMSE | 200 | 127.0452 | 25.9331 | 31.4919 | 0.4123 | 7.18 |
| 14 | RgPCC.MSECV | 200 | 127.0452 | 25.9331 | 31.4919 | 0.4123 | 7.18 |

Figure 3: For γ'_2 and $p = 50$.

Overall, we can see that RgPCC provides slight improvements to logistic regression and principal component logistic regression when the true γ is nonzero in the leading principal components. However, when the nonzero components of γ are not the leading principal components we see a significant different between principal component logistic regression and our other methods.

We can see that the testing error for \mathbf{p} is improved with RgPCC by at least a factor of 4 in each case. RgPCC avoids overfitting the data as much as logistic regression does. In particular, there is drastic improvement in classification when γ^* is very sparse (that is for γ_1). Similar results occur when $N = 200$ but are left in the table and results section.

We also run a few examples for higher dimensional data. Our method should excel in this via its use of principal components and sparsity in γ .

where

$$\begin{aligned}\gamma_1 &= (25, 0, 0, 0, 0, 0, \dots, 0) \\ \gamma_2 &= (10, 5, 5, 0, 0, 0, \dots, 0)\end{aligned}$$

Here we see similar results as in the low dimensional case. At this point, we know that RgPCC performs better, if no just as well, as logistic regression does. However, there are still many aspect of the method we wish to study.

4.3 Performance of RgPCC on Realworld Data

In this section we apply four real datasets to compare the prediction performance of RgPCC with conventional binary classification methods. We look at the following data sets.

- Divorce Predictors (170 instances, 54 predictors, binary response)
- Cryotherapy (90 instances, 7 predictors, binary response)

- Audit (777 instances, 18 predictors, binary response)
- Ecoli (336 instances, 8 predictors, binary response)

We will compare the methods RgPCC, logistic, PCA logistic and ridge regression. We will tune RgPCC by AIC and BIC. Thus we will compare a total of five different methods.

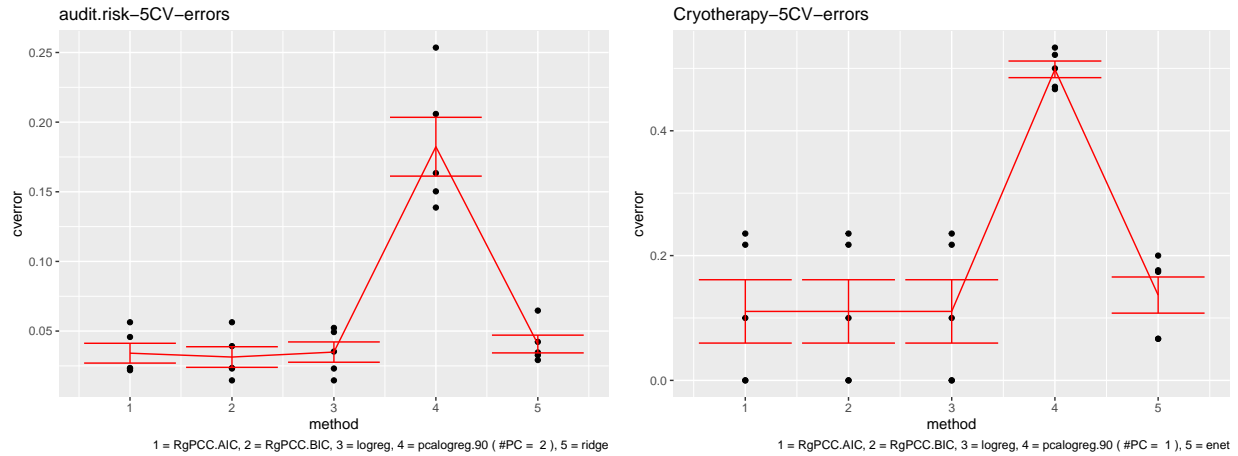
To compare the performance of our binary classification methods based on the error of each fold of 5-fold cross validation. We then perform Tukey’s multiple pairwise comparison to test if the methods differ significantly or not.

The following is a summary of the error results, more can be found in the section 6.

| | RgPCC AIC | RgPCC BIC | Logistic | PC Logistic | Ridge |
|---------|-----------|-----------|----------|-------------|----------|
| Divorce | 0.023529 | 0.023529 | 0.058823 | 0.035294 | 0.023529 |
| Cryo | 0.12222 | 0.12222 | 0.12222 | 0.5 | 0.14444 |
| Audit | 0.033548 | 0.030967 | 0.034838 | 0.18193 | 0.04129 |
| Ecoli | 0.041666 | 0.038690 | 0.041666 | 0.053571 | 0.044642 |

Figure 4: The mean classification errors over each of the 5 folds.

When running the Tukey’s multiple pairwise comparisons we noticed that in most cases the difference in the errors was insignificantly different with p -values around $0.99 \geq \alpha = 0.05$. However, there were a few cases in which PC logistic was significantly different than all the other methods (Audit data and Cryotherapy). This data is given below.



Thus, in terms of classification error RgPCC is comparable to logistic and ridge regression in application settings.

5 Further Research

As we continue work on RgPCC, we plan to investigate and improve a in a few key areas.

5.1 Speed

This method of RgPCC improves upon testing error in binary classification compared to logistic regression. We also believe that this use of principal component analysis should improve the speed

of the computation, especially when there is sparsity in the solution with respect to the principal components. We wish to further explore this, especially in higher dimensions and when $p > N$.

5.2 C_p -type Statistic for Parameter Tuning

Using cross-validation for parameter selection is data and time intensive. In Zou and Lang's paper [2], the regularization parameter in RgPCR can be selected using Stein's unbiased risk estimation. In later works, we will also investigate if this method can be generalized to RgPCC.

5.3 Further Generalization

In many cases, data is not presented in a manner such that classes are linearly separable. Kernel methods embed data in higher dimensional spaces in a way that makes classes linearly separable, without have to work in the higher dimensional space directly. We would like to be able to apply RgPCC in these nonlinear cases as well.

We would also like to extend this to multiclass classification through ideas from logistic regression for multiclass classification. We can also take the ideas from [3] to generalize this method to other classifiers such as SVM.

Bibliography

- [1] T. HASTIE, R. TIBSHIRANI, AND J. FRIEDMAN (2009) — *The Elements of Statistical Learning : Data Mining, Inference, and Prediction (second edition)*, New York, NY: Springer.
- [2] W. LANG, AND H. ZOU (2020) — *A simple method to improve principal components regression*, WILEY, <https://onlinelibrary.wiley.com/doi/abs/10.1002/sta4.288>
- [3] H. WANG AND C. LENG (2007) — *Unified LASSO Estimation by Least Squares Approximation*, Journal of the American Statistical Association, 102:479, 1039-1048, DOI; 10.1198/016214507000000509, <https://doi.org/10.1198/016214507000000509>

6 Figures