

Fitting High-dimensional Regressions Models

G. de los Campos

11/29/2020

Up to now, in this course we have discussed methods that estimate parameters by either maximizing the likelihood function (ML) or by minimizing the residual sum of squares (OLS). These methods have reasonably good properties (e.g., OLS is unbiased and has minimum variance among the class of linear unbiased estimators, ML is asymptotically unbiased and asymptotically efficient). However, the performance of these methods can be sub-optimal when the number of parameters is large relative to sample size.

Therefore, in this note we consider methods that are tailored for problems involving a large number of predictors. We will cover the following approaches:

- **1) Variable selection using screening methods:** This approach selects predictors using a marginal association test (i.e., testing the association of the response and each predictor, one predictor at a time) and builds models using the top- q ($q=1,2,\dots$). An optimal model DF can be chosen by evaluating the ability of the model to predict testing data for models with 1DF, 2DF,...
- **2) Penalized Regressions:** This approach estimates effects using a penalized sum of squares. We will consider three methods: Ridge Regression, Lasso, and Elastic Net. The extent of regularization will be controlled by a parameter (λ) which will be chosen to maximize prediction accuracy in testing data.
- **3) Bayesian Shrinkage and Variable Selection methods:** In Bayesian regression the choice of prior will determine whether the model performs shrinkage, variable selection, or a combination of the two. We will present examples using shrinkage and variable selection priors.

Data

To illustrate the application of the methods listed above, we will use a data set available in the [BGLR](#) R-package. This data set provides four phenotypes (see object `wheat.Y`) for 599 wheat lines that were genotyped at 1,279 genetic markers (see object `wheat.X`).

Loading the data

This code below loads the data, center, and scales the the genotypes. While centering and scaling is not strictly needed, it is often a good practice when using penalized (e.g., Lasso) or Bayesian regressions.

```
library(BGLR)
data(wheat)
head(wheat.Y)
```

```
##           1           2           4           5
## 775    1.6716295 -1.72746986 -1.89028479  0.0509159
## 2166 -0.2527028  0.40952243  0.30938553 -1.7387588
## 2167  0.3418151 -0.64862633 -0.79955921 -1.0535691
## 2465  0.7854395  0.09394919  0.57046773  0.5517574
## 3881  0.9983176 -0.28248062  1.61868192 -0.1142848
## 3889  2.3360969  0.62647587  0.07353311  0.7195856
```

```
dim(wheat.X)
```

```
## [1] 599 1279
```

```
X=scale(wheat.X,center=TRUE,scale=TRUE)
y=wheat.Y[,2] # picks one phenotype

N<-nrow(X) ; p<-ncol(X)
```

We will compare models based on their ability to predict data that was not used to fit the models. The following code produces a training-testing partition that we will use for all methods.

Creating a Training-Testing partition

```
set.seed(12345)
tst<-sample(1:N,size=150,replace=FALSE)
XTRN<-X[-tst,]
yTRN<-y[-tst]
XTST<-X[tst,]
yTST<-y[tst]
```

1) Variable screening

We first rank predictors based on a marginal association test. Note that screening is done within the training data only.

```
pValues<-numeric()
for(i in 1:p){
  fm<-lsfit(y=yTRN,x=XTRN[,i])
  pValues[i]<-ls.print(fm,print.it=F)$coef[[1]][2,4] # extracts p-value, similar to lm() but a bit fa
}

```

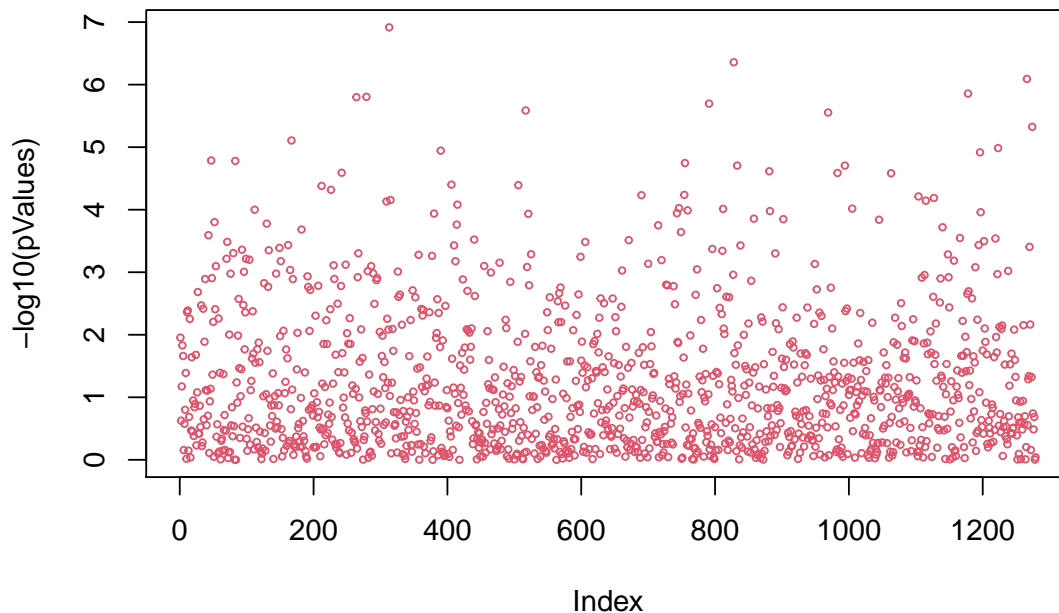


Figure 1: Figure 1: Marginal association p-value

Let's now build models using the top- q ($q=1, \dots, 300$) markers. The script: - Ranks markers based on p-values (from smallest to largest). - Fit models using the top 1, top 2, \dots , top- q markers using data from the training set. - For each of the fitted model the script computes the correlation between phenotype and predictions within the training data and in the testing data.

Building prediction models using the top-q markers

```
##### VARIABLE SELECTION #####
mrk_rank<-order(pValues); corTRN<-numeric(); corTST<-numeric()
for(i in 1:300){
  tmpIndex<- mrk_rank[1:i]
  ZTRN=XTRN[,tmpIndex,drop=F]
  ZTST=XTST[,tmpIndex,drop=F]

  fm<-lm(yTRN~ZTRN)
  bHat=coef(fm)[-1]
  bHat<-ifelse(is.na(bHat),0,bHat)

  yHatTRN=ZTRN%%bHat
  corTRN[i]<-cor(yTRN,yHatTRN)

  yHatTST=ZTST%%bHat
  corTST[i]<-cor(yTST,yHatTST)
}
```

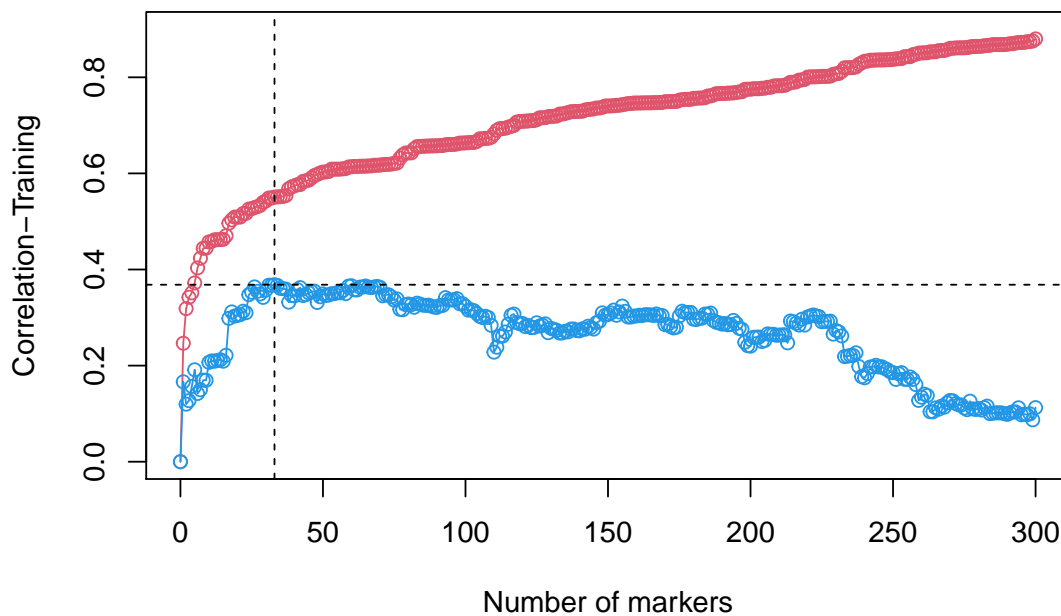


Figure 2: Figure 2: Correlation between predictions and phenotypes in training and testing set by model DF.

Remarks

- Goodness of fit in the training data set (`corTRN`, blue) increases with DF
- However prediction accuracy in testing data (`corTST`, red) increases, reaches a plateau, and then decreases.

The curves presented in the previous figure are estimates subject to sampling variability. To quantify this and to reduce the variance of estimates we can conduct many training-testing partitions and average across them. This is illustrated in the following figure, each of the skyblue lines is an estimate derived from a training-testing partition. The solid red line is an average across partitions.

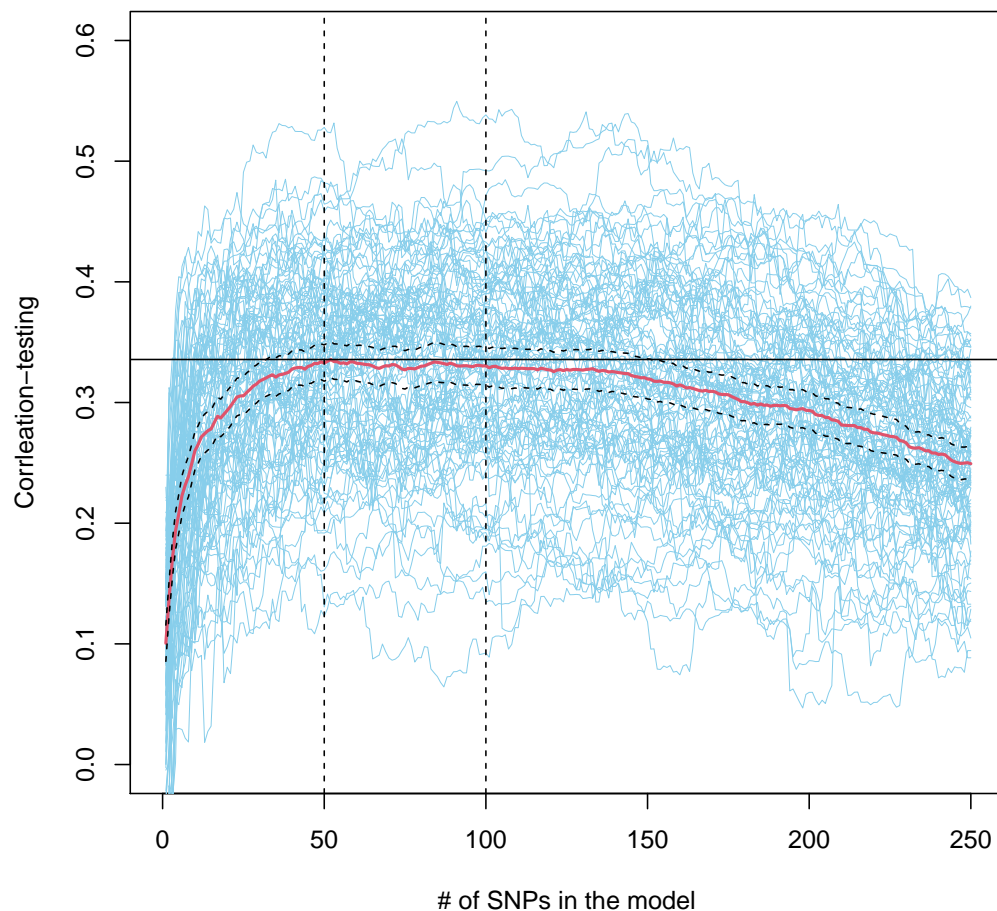


Figure 3: Figure 2b: Correlation between predictions and phenotypes in testing data by model DF. (100 training-testing partitions)

Regularized Regression

A groundbreaking paper by [James & Stein \(1961\)](#) showed that in some settings the least squares estimator (also ML in their example) could be inadmissible; that is, they showed that in some circumstances there was another estimator that had lower Mean-Squared Error (MSE) over the entire parameter space. Their estimator shrinks the least square estimates towards zero; thus reducing the variance of estimates.

Recall that the MSE of an estimator can be decomposed as the sum of the variance plus the squared of the bias of the estimator $MSE = E[(\hat{\theta} - \theta)^2] = \text{Variance} + \text{Bias}^2$. Shrinkage reduces the variance of estimates at the expense of some bias. However, when the number of parameters to be estimated is large, the reduction in variance overcomes the increase in bias; thus leading to a reduction in MSE.

There are many ways to obtain regularized estimates; two commonly used approaches are penalized and Bayesian methods. We discussed each of them in the next two sections. In most cases there is a duality between the two approaches by which a penalized estimator can be viewed as the posterior mode from a Bayesian model.

2) Penalized regressions using glmnet

In a penalized regression, estimates are obtained by minimizing a penalized log-likelihood or, in the case of linear models, a penalized residual sum of squares:

$$\hat{\beta} = \operatorname{argmin}\{(y - X\beta)'(y - X\beta) + \lambda J(\beta)\}$$

where $J(\beta)$ is a penalty function. Common choices for the penalty function are the

- L2-norm, $J(\beta) = \sum_j \beta_j^2$ (aka Ridge Regression, [Hoerl and Kennard 1970](#)),
- L-1 norm $J(\beta) = \sum_j |\beta_j|$ (aka Lasso, [Tibshirani, 1996](#)), and,
- Linear combinations of the two $J(\beta) = (1 - \alpha) \sum_j \beta_j^2 + \alpha \sum_j |\beta_j|$ (aka [Elastic Net, Zhou and Hastie, 2005](#)), for some $\alpha \in [0, 1]$.

Choosing $\lambda = 0$ leads Ordinary Least Squares estimates. Ridge regression shrunk OLS estimates towards zero, without making variable selection. Lasso and Elastic Net combine variable selection and shrinkage.

Commonly, these models are fitted over a grid of values of the regularization parameter (λ); an optimal value for that parameter is often chosen by evaluating the ability of the fitted models to predict data that was not used to train the models (i.e., testing data).

2.1: Ridge Regression (RR)

In the RR ([Hoerl and Kennard 1970](#)) $J(\beta) = \sum_j \beta_j^2 = \beta' \beta$; thus, the objective function becomes

$$\hat{\beta} = \operatorname{argmin}\{(y - X\beta)'(y - X\beta) + \lambda \beta' \beta\} = \operatorname{argmin}\{y'y + \beta'(X'X + I\lambda)\beta - 2\beta'X'y\}$$

The solution can be shown to be

$$\hat{\beta} = (X'X + I\lambda)^{-1}X'y$$

Adding λ to the diagonal entries of $X'X$ shrinks estimates towards zero. This is illustrated in the following simplified example.

```
set.seed(195021)
# Toy simulation
n=50
p=3 # number of predictors
W=matrix(nrow=n,ncol=p,rnorm(n*p))
b=c(-1,1,2) # true effects
signal=W%*%b
error=rnorm(sd=sd(signal),n=length(signal))
```

```

y=signal+error

# centering to avoid the need of including an intercept
W=scale(W,center=T,scale=F)
y=y-mean(y)

# OLS
WW=crossprod(W)
Wy=crossprod(W,y)
bOLS=solve(WW,Wy)

# Ridge regression
lambda=3
C=WW
diag(C)=diag(C)+lambda
bRR_3=solve(C,Wy)

lambda=10
C=WW
diag(C)=diag(C)+lambda
bRR_10=solve(C,Wy)

lambda=100
C=WW
diag(C)=diag(C)+lambda
bRR_100=solve(C,Wy)
round(cbind('true_effect'=b, 'ols'=bOLS, 'RR_3'=bRR_3, 'RR_10'=bRR_10, 'RR_100'=bRR_100), 3)

```

```

##      true_effect
## [1,]          -1 -1.073 -1.040 -0.965 -0.461
## [2,]           1  0.964  0.871  0.707  0.192
## [3,]           2  2.509  2.350  2.051  0.797

```

Let's compare estimates in just this MC replicate.

```
sum((b-bOLS)^2)
```

```
## [1] 0.2656491
```

```
sum((b-bRR_3)^2)
```

```
## [1] 0.140734
```

```
sum((b-bRR_10)^2)
```

```
## [1] 0.08950725
```

```
sum((b-bRR_100)^2)
```

```
## [1] 2.389746
```

In this example using $\lambda = 3$ or $\lambda = 10$ improved the estimates (smaller distance to the true parameter values), but using $\lambda = 100$ induced too much shrinkage towards zero, thus increasing the squared-difference between estimates and true parameter values.

Unlike the Ridge Regression, Lasso and Elastic Net estimates do not have a closed form; however, estimates can be derived using iterative algorithms (e.g., a coordinate-descent gradient) such as the ones implemented

in the `glmnet` R-package.

2.2: Fitting Penalized Regressions using the `glmnet` R-package

The following code shows how to implement Ridge Regression, Lasso, and Elastic Net using the `glmnet` package. By default, `glmnet` fits models over a grid of 100 values of the regularization parameter λ . The plots produced at the end of the script display, for each of the models, the correlation between predictions and observations in testing data, by value of λ .

```
library(glmnet)

# alpha 0 gives Ridge Regression
fmRR=glmnet(y=yTRN,x=XTRN,alpha=0)
dim(fmRR$beta)

## [1] 1279 100

length(fmRR$lambda)

## [1] 100

range(fmRR$lambda)

## [1] 2.477649 247.764898

# alpha 1 gives Lasso
fmL=glmnet(y=yTRN,x=XTRN,alpha=1)

# alpha between 0 and 1 gives elastic net
fmEN=glmnet(y=yTRN,x=XTRN, alpha=0.5)

COR.RR=rep(NA,100)
COR.L=rep(NA,100)
COR.ENet=rep(NA,100)

# evaluating correlation in TST set
for(i in 1:100){
  COR.RR[i]=cor(yTST,XTST%%fmRR$beta[,i])
  COR.L[i]=cor(yTST,XTST%%fmL$beta[,i])
  COR.ENet[i]=cor(yTST,XTST%%fmEN$beta[,i])
}

## Warning in cor(yTST, XTST %% fmL$beta[, i]): the standard deviation is zero
## Warning in cor(yTST, XTST %% fmEN$beta[, i]): the standard deviation is zero
```

Remarks

- The `glmnet` function fits each of the models over a grid of values of λ , the rules used to choose those values are described in [Friedman, Hastie, and Tibshirani \(2010\)](#).
- The matrix `$beta` has the solutions (estimated effects) obtained for each value of λ .
- After fitting the model we evaluate prediction accuracy by correlating the testing phenotypes `yTST` with predictions (see loop above).
- In the case of Lasso and Elastic Net the default values of `lambda` led to an internal maxima, i.e., an internal region with maximum correlation. This is not the case for the Ridge Regression, the grid of values of λ used in that case may need to use smaller values of λ . The following example produce new fits using a user-provided grid of values for the regularization parameter.

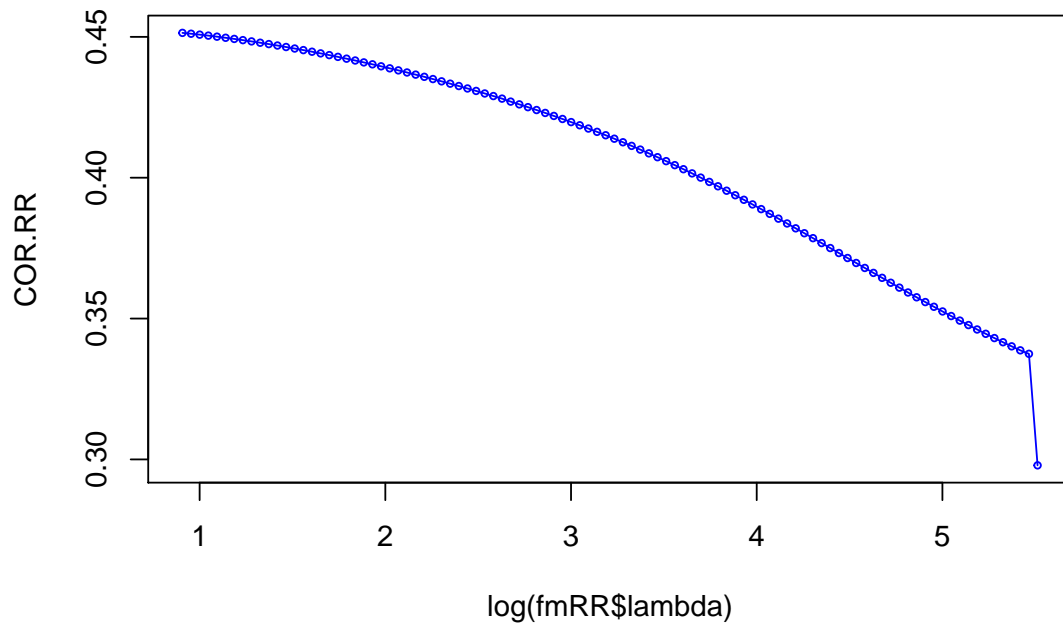


Figure 4: Figure 3a: Correlation between predictions and phenotypes in testing, Rdige Regression

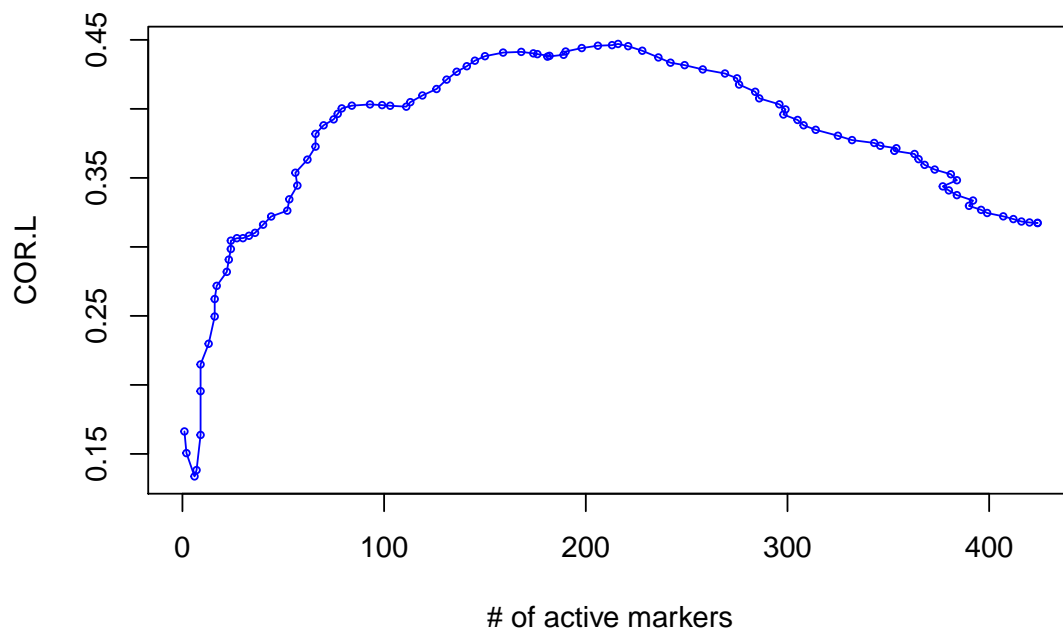


Figure 5: Figure 3b: Correlation between predictions and phenotypes in testing, Lasso

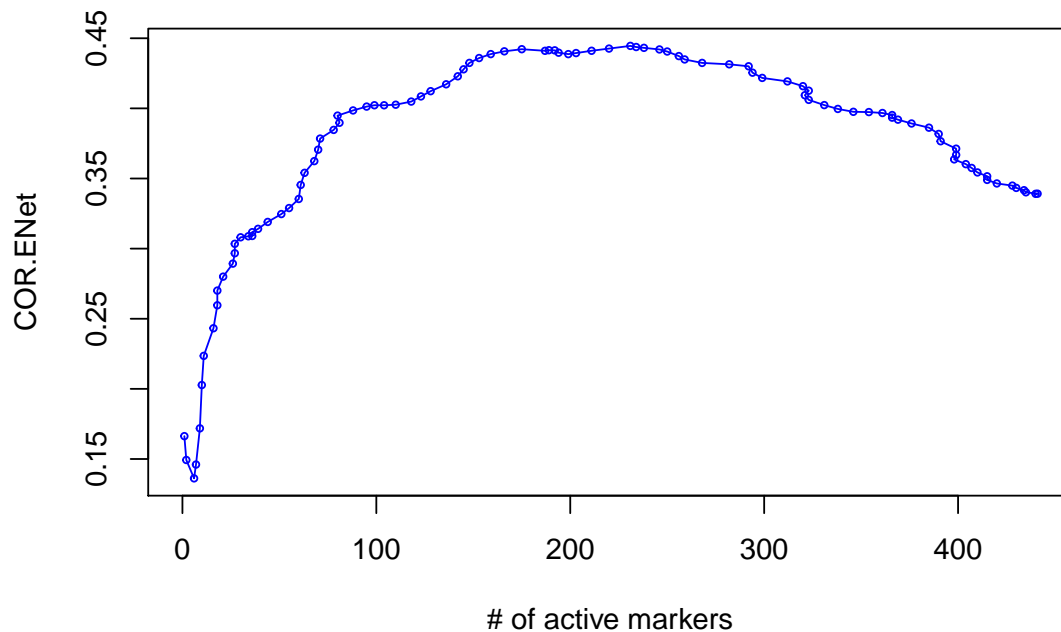


Figure 6: Figure 4: Correlation between predictions and phenotypes in testing, Elastic Net

```
lambda=c(fmRR$lambda,min(fmRR$lambda)*seq(from=0.9,to=0.01,length=50))
fmRR=glmnet(y=yTRN,x=XTRN,alpha=0,lambda=lambda)

# evaluating correlation in TST set
COR.RR=rep(NA,length(fmRR$lambda))
for(i in 1:length(fmRR$lambda)){
  COR.RR[i]=cor(yTST,XTST*fmRR$beta[,i])
}
```

- The three models achieve correlations of about 0.45 (Ridge Regression does slightly better)
- This is much better than what we obtained selecting markers based on their marginal association (correlation ~ 0.37 , Example 1).
- Remember that estimates of accuracy, such as the ones discussed above, are point estimates subject to sampling variability; there is sampling variance emerging from the sampling of training and testing data. In the In-class assignment you will be asked to repeat the examples using many training-testing partitions; we will use those results to assess sampling variability on these estimates and also to get a more precise estimate.

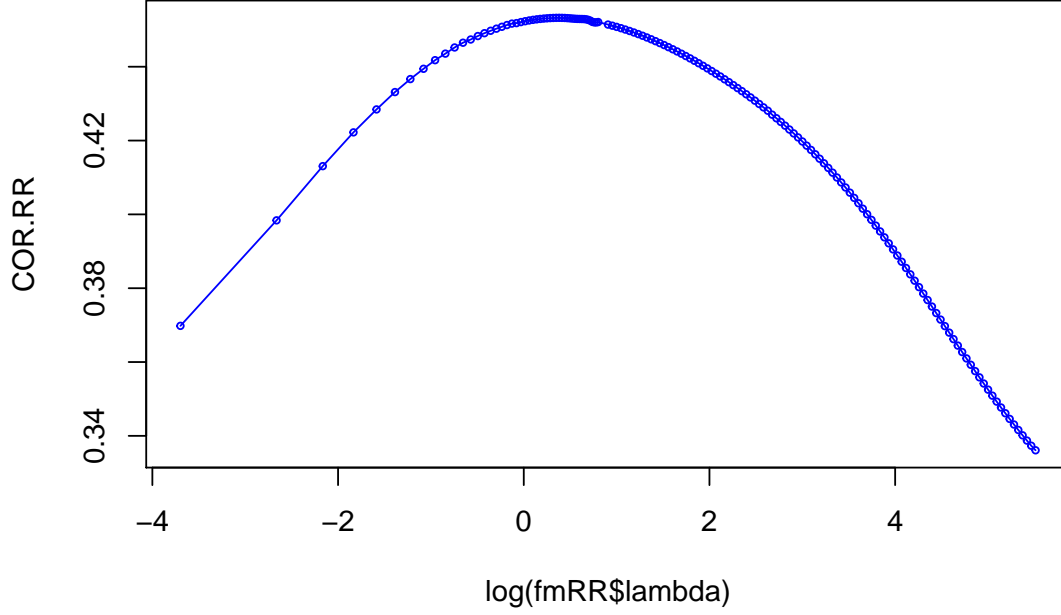


Figure 7: Figure 3c: Correlation between predictions and phenotypes in testing, Ridge Regression

3) Bayesian Regressions

In a Bayesian model, estimates are obtained from the posterior distribution of the model unknowns (e.g., regression coefficients).

Recall that from Baye's rule we have that $p(A|B) = \frac{p(A,B)}{p(B)} = \frac{p(B|A)P(A)}{p(B)}$. Taking A to be the model parameters (e.g., β in a regression model), and B to be the data (y), we have that

$$p(\beta|y) = \frac{p(y|\beta)p(\beta)}{p(y)}$$

Above,

- $p(\beta|y)$ is the posterior distribution of the parameters given the data (the object we use to summarize knowledge and uncertainty about model parameters),
- $p(y|\beta)$ is the conditional distribution of the data given the parameters, the likelihood function when viewed as a function of β ,
- $p(\beta)$ is the prior distribution of the model unknowns, the object we use to summarize *prior knowledge*, and
- $p(y) = \int p(y|\beta)p(\beta) d\beta$.

The last object, $p(y)$ is the marginal distribution of the data and does not involve the unknown parameters. It is an integration constant that normalizes the posterior density but does not affect the shape of the posterior density is not affected by it. Therefore, we see that the posterior distribution is proportional to the product of the likelihood times the prior distribution

$$p(\beta|y) \propto p(y|\beta)p(\beta)$$

This makes evident how the posterior distribution (and inferences from it, e.g., the posterior mean, or the posterior mode) depends on both evidence provided by the data, quantified via the *likelihood function*, and *prior knowledge* summarized by the prior distribution.

3.1) A Bayesian model with a Gaussian likelihood and a Gaussian prior

Let's consider a linear model $y = X\beta + \varepsilon$ with a Gaussian likelihood,

$$p(y|X, \beta) = N(X\beta, I\sigma_\varepsilon^2)$$

The ML estimator can be shown to be the OLS estimator

$$\hat{\beta} = (X'X)^{-1}X'y$$

Consider now using a Gaussian IID prior with zero-mean and variance σ_β^2 , that is

$$p(\beta) = N(0, I\sigma_\beta^2)$$

The posterior distribution becomes

$$p(\beta|y, \sigma_\varepsilon^2, \sigma_\beta^2) \propto N(X\beta, I\sigma_\varepsilon^2) \times N(0, I\sigma_\beta^2)$$

This can be shown to be a Multivariate Normal distribution with mean $\tilde{\beta} = (X'X + I\lambda)^{-1}X'y$ and variance-covariance matrix $V = (X'X + I\lambda)^{-1}\sigma^2$, where $\lambda = \sigma^2/\sigma_\beta^2$. $\tilde{\beta} = (X'X + I\lambda)^{-1}X'y$ is the Ridge-regression estimator. Thus, we see that Ridge Regression estimates can be seen as the posterior mean (also the posterior mode) of the vector of effects in a Gaussian regression model with IID Gaussian prior.

The regularization parameter, $\lambda = \sigma^2/\sigma_\beta^2$, is a noise-to-signal ratio. In penalize regressions, this parameter is often chosen using cross-validation (see section on penalized regressions, above). In a Bayesian context, the variances can be treated as unknown (e.g., by assigning them a scaled-inverse chi-square prior); thus inferring effects and variances jointly from the training data. This is illustrated in the following example which uses the [BGLR R-package](#).

```
library(BGLR)
nIter=6000 # I set this to small value that way it will run quickly, for more serious analyses use lo
burnIn=1000 # and longer burnin

# Gaussian prior ("Bayesian Ridge-Regression")
LP=list( list(X=XTRN,model='BRR') ) # 2-level list, allows specifying different types of random and f
fmBRR=BGLR(y=yTRN,ETA=LP,nIter=nIter,burnIn=burnIn,saveAt='BRR_',verbose=FALSE)
# Retriving samples from the variance parameters

vE=scan('BRR_varE.dat',quiet=TRUE)
vB=scan('BRR_ETA_1_varB.dat',quiet=TRUE)
```

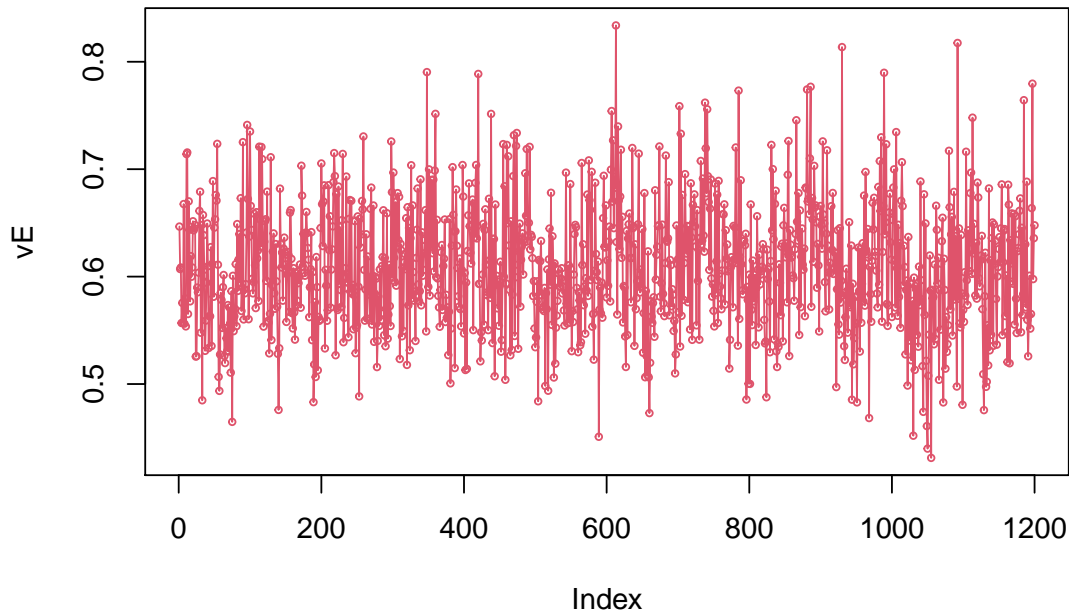


Figure 8: Figure 4: Trace plot of the error variance.

```
## Prediction accuracy in the testing set
cor( yTST, XTST%*%fmBRR$ETA[[1]]$b)

##           [,1]
## [1,] 0.4502181

max(COR.RR)

## [1] 0.4532439
```

While the RR appears to outperform slightly the Bayesian model in prediction accuracy, the estimate of prediction accuracy for the RR is likely upwardly biased because, below, in the comparison we choose lambda based on the same testing data that is used to evaluate accuracy. On the other hand, predictions from the Bayesian model were derived using the training data only.

3.2) Various shrinkage and variable selection priors

The Gaussian prior used in the previous section induces shrinkage without performing any variable selection. In the last decade a plethora of Bayesian models using different priors have been developed. These priors can be classified in three main groups: (i) Gaussian, (ii) Thick-tailed priors, this group includes the double-exponential (used in the [Bayesian Lasso](#)) and the scaled-t prior, and (ii) finite mixture priors with a point

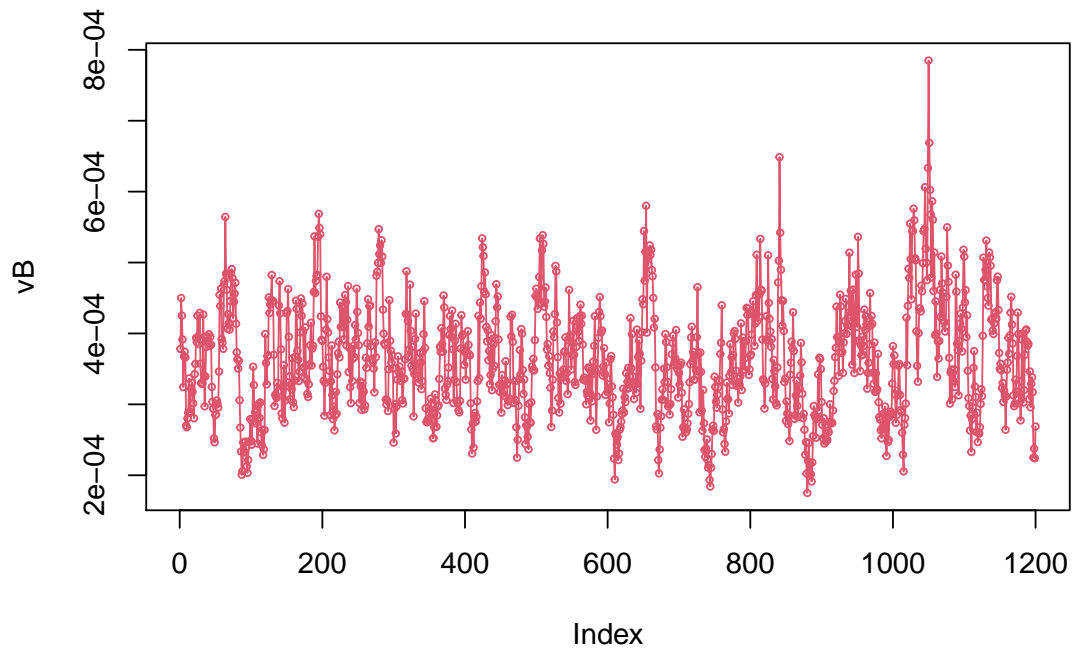


Figure 9: Figure 5: Trace plot of the variance of effects.

of mass at zero, these models perform both variable selection and shrinkage. The following Figure displays these priors.

These priors are scaled by multiple *hyper-parameters* (e.g., degree of freedom, scale, prior probability of non-null effects). Fortunately, some of these parameters can be treated as unknown and can be inferred by data.

The [BGLR R-package](#) implements some of these priors. The following example illustrates how to fit Bayesian models with different prior distributions of effects. For additional examples and a description of the methods implemented you can check the following [GitHub repository](#) (include multiple examples) and the following [manuscript](#).

2a) Model fitting

```
# Scaled-t
LP[[1]]$model='BayesA'
fmBA=BGLR(y=yTRN,ETA=LP,nIter=nIter,burnIn=burnIn,saveAt='BA_',verbose=FALSE)

# Double-Exponential
LP[[1]]$model='BL'
fmBL=BGLR(y=yTRN,ETA=LP,nIter=nIter,burnIn=burnIn,saveAt='BL_',verbose=FALSE)

# Spike-slab (Gaussian)
LP[[1]]$model='BayesC'
fmBC=BGLR(y=yTRN,ETA=LP,nIter=nIter,burnIn=burnIn,saveAt='BC_',verbose=FALSE)

# Spike-slab (Scaled-t)
LP[[1]]$model='BayesB'
fmBB=BGLR(y=yTRN,ETA=LP,nIter=nIter,burnIn=burnIn,saveAt='BB_',verbose=FALSE)
```

2b) Retrieving samples and estimates

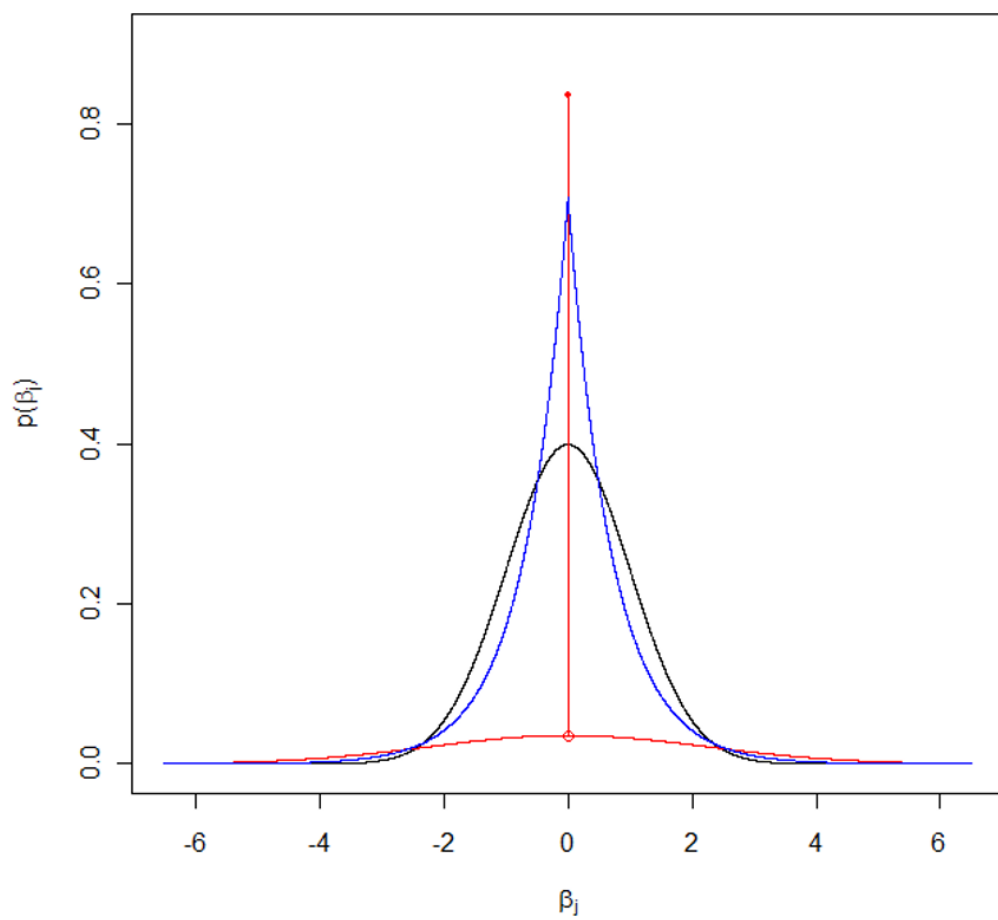


Figure 10: Figure 4: Prior distributions of effects commonly used in Bayesian models

```

bayes=c(
  'BRR'   =cor( yTST, XTST%%fmBRR$ETA[[1]]$b),
  'BL'    =cor( yTST, XTST%%fmBL$ETA[[1]]$b),
  'BayesA'=cor( yTST, XTST%%fmBA$ETA[[1]]$b),
  'BayesB'=cor( yTST, XTST%%fmBB$ETA[[1]]$b),
  'BayesC'=cor( yTST, XTST%%fmBC$ETA[[1]]$b)
)
round(bayes,3)

```

```

##      BRR      BL BayesA BayesB BayesC
## 0.450 0.450 0.451 0.455 0.448

```

In general we see no big difference in prediction accuracy, all the Bayesian models achieved in this example a prediction correlation close than the one achieved by penalized regressions.

These models offer more than just predictions. This code illustrates how to extract other parameters (for more details follow the links provided above).

2c) Retrieving samples and estimates

```

##          775          2166          2465          3881          3889          4248
## -0.76327204 0.40523461 0.09598448 -0.43601817 0.39317259 -0.05501539
## [1] -0.007336212
## [1] 0.6099191
## $logLikAtPostMean
## [1] -468.7343
##
## $postMeanLogLik
## [1] -524.2074
##
## $pD
## [1] 110.9462
##
## $DIC
## [1] 1159.361
## [1] 6000

```

```

# Gaussian prior
head(fmBRR$ETA[[1]]$b) # posterior means of effects

```

```

##          wPt.0538          wPt.8463          wPt.6348          wPt.9992          wPt.2838
## 6.897426e-04 -3.197053e-03 -1.624006e-03 -5.027481e-05 -1.239100e-03
##          wPt.8266
## 1.392390e-03

```

```

head(fmBRR$ETA[[1]]$SD.b) # posterior SDs

```

```

##          wPt.0538          wPt.8463          wPt.6348          wPt.9992          wPt.2838          wPt.8266
## 0.01763717 0.01804330 0.01781049 0.01855480 0.01793263 0.01848906

```

```

# Bayes A
head(fmBA$ETA[[1]]$b)

```

```

##          wPt.0538          wPt.8463          wPt.6348          wPt.9992          wPt.2838
## 4.787594e-04 -1.945393e-03 9.506763e-05 1.217393e-04 -1.006084e-03
##          wPt.8266

```

```
## 4.190577e-04

# Bayesian Lasso
head(fmBL$ETA[[1]]$b)

##      wPt.0538      wPt.8463      wPt.6348      wPt.9992      wPt.2838
## 0.0007948838 -0.0024648168 -0.0003986224 -0.0009162976 -0.0008811480
##      wPt.8266
## 0.0010594661

# BayesC
head(fmBC$ETA[[1]]$b)

##      wPt.0538      wPt.8463      wPt.6348      wPt.9992      wPt.2838
## 0.0002725590 -0.0014128833 -0.0011997475 0.0001290981 -0.0010514157
##      wPt.8266
## -0.0007484630

head(fmBC$ETA[[1]]$d) # posterior probability of inclusion

## [1] 0.475 0.497 0.495 0.486 0.489 0.472

# BayesB
head(fmBB$ETA[[1]]$b)

##      wPt.0538      wPt.8463      wPt.6348      wPt.9992      wPt.2838
## -0.0003210662 -0.0026599218 -0.0017579226 0.0002400352 -0.0012135738
##      wPt.8266
## -0.0002733823

head(fmBB$ETA[[1]]$d) # posterior probability of inclusion

## [1] 0.426 0.416 0.414 0.415 0.388 0.417
```