

Assessed Practical Statistical Programming: Designing a greedy grid algorithm for the Elastic-Net regression

Practical Number: P535

May 1, 2017

Abstract

A greedy grid algorithm is implemented to estimate the optimal tuning parameters of an elastic-net regression fitted to a data set. The best pair of parameters is chosen to be the one that maximizes the predictive correlation, computed using a 10-fold cross-validation. The cross-validation is implemented both sequentially and using parallel computing, and the performances are compared. The same thing is done in order to compute the cross-validation 10 times. It is found that using 2 slave processes speeds up the computation but that using 4 slaves only leads to marginal gains. The architecture of the cross-validation function and the greedy grid algorithm is discussed and it is argued that the code implemented has conceptual and practical limitations

Keywords: Elastic-Net regression; Greedy grid algorithm; Parallel computing; Statistical programming;

1 Introduction

In this report, we will comment on the implementation of a greedy grid algorithm (also known as hill climbing algorithm) used to find the optimal combination of the tuning parameters λ_1 and λ_2 in the *elastic net* regression. The *elastic net* regression is a penalized linear regression model that combines an L1 penalty with an L2 penalty to produce regularized estimates of the regression vector of coefficients β and can be written as

$$(\mathbf{y} - \mathbf{X}\beta)^T(\mathbf{y} - \mathbf{X}\beta) + \lambda_1 \sum_{i=0}^p |\beta_i| + \lambda_2 \sum_{i=0}^p \beta_i^2, \quad \lambda_1, \lambda_2 \geq 0, \quad (1)$$

where λ_1 is the tuning parameter for the L1 norm and λ_2 for the squared L2 norm. The greedy grid algorithm will simultaneously propose new candidates for λ_1 and λ_2 for which the predictive correlation between the observed responses and fitted responses will be computed using 10-fold cross-validation. If one of the new proposed pairs (λ_1, λ_2) gives a lower predictive correlation than the current one, then the pair is accepted as the new optimal point, and we repeat the process until there is no significant improvement (for a certain threshold). The algorithm is tested on the nki70 data set available in R. This report will focus on the R code and the parallelization and not on the statistical results and interpretation when applied on the nki70 data set.

Firstly, the time complexity of solving the penalized least squares problem for a given (λ_1, λ_2) will be estimated. Then, the cross-validation function and the algorithm will be displayed to facilitate the explanation of the parallel implementation. The architecture and choices made in writing the algorithm will also be discussed and the limitations will be highlighted. The full R code is available in Appendix A.

2 Time complexity

Estimating the time complexity of solving the penalized least squares problem for a given pair (λ_1, λ_2) is important since this operation will be repeated many times in the algorithm. For each iteration of the algorithm, it will be solved 10 times in the cross-validation, which is repeated 10 times, and there are 9 new proposed pairs (λ_1, λ_2) (in fact there are 8, but we compute 9), which gives 900 computations for each iteration. Rewriting equation (1) in matrix form leads to

$$(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T(\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) + \lambda_1\boldsymbol{\beta}^T\mathbf{e} + \lambda_2\boldsymbol{\beta}^T\boldsymbol{\beta}, \quad \lambda_1, \lambda_2 \geq 0, \quad (2)$$

where \mathbf{e} is a $p \times 1$ vector with $e_i = 1$ if $\beta_i \geq 0$ and -1 if $\beta_i < 0$, $i = 1, \dots, p$. Minimizing (2) with respect to $\boldsymbol{\beta}$ gives

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}^T\mathbf{X} + \mathbf{I}_p\lambda_2)^{-1}(\mathbf{X}^T\mathbf{y} - \frac{1}{2}\lambda_1\mathbf{e}), \quad (3)$$

where \mathbf{I}_p is the $p \times p$ identity matrix. Taking that multiplying a $r \times s$ matrix and a $s \times t$ matrix takes $O(rst)$ operations, adding matrices has quadratic complexity and computing the inverse of a $r \times r$ matrix is $O(r^3)$ using a Cholesky decomposition or Gram-Schmidt, we can then estimate the time complexity of computing $\hat{\boldsymbol{\beta}}$ as follows

1. Compute $\mathbf{X}^T\mathbf{X}$ and add it to $\mathbf{I}_p\lambda_2$
2. Take the result and compute its inverse
3. Compute $\mathbf{X}^T\mathbf{y}$ and subtract $\frac{1}{2}\lambda_1\mathbf{e}$
4. Multiply the results of steps 2 and 3

Step 1 has complexity $O(np^2 + 2p^2)$ where the second term comes from multiplying λ_2 with the identity matrix and then adding it (hence the factor of 2). In step 2, we take the inverse of a $p \times p$ which has complexity $O(p^3)$. Similar to step 1, step 3 has complexity $O(np + 2p)$. Finally, step 4 consists in multiplying a $p \times p$ by $p \times 1$ and has complexity of $O(p^2)$. Adding all the steps we get that the total time complexity is $O(p^3 + p^2(n + 1) + p(n + 2))$

3 R code

3.1 Cross-validation function

For each new proposed pair of tuning parameters in the greedy grid algorithm, we need a function to compute the predictive correlation using 10-fold cross-validation as explained in the introduction. The code in Source Code 1 is the cross-validation function. We will highlight the part where parallel computing is implemented and briefly compare the performances of the parallelized and non parallelized functions.

Source Code 1: Parallelized cross-correlation function

```
parallel.xval = function(data, cluster, k, lambda1, lambda2)
{
  n = nrow(data)
```

```

folds = split(sample(n), seq_len(k))

xval.fold = function(fold, lambda1, lambda2) {

  dtrain = data[-fold, ]
  dtest = data[fold, ]
  ElasticNet.fold = penalized(time, penalized = ~ . - time, data = dtrain,
    lambda1 = lambda1, lambda2 = lambda2, trace = FALSE)

  pred = predict(ElasticNet.fold, data = dtest)
  return(data.frame(PRED = pred[, "mu"], OBS = dtest$time))
}

# Added code to implement parallelisation

clusterExport(cluster, list("data"))
clusterEvalQ(cluster, library(penalized))

# Here, ParLapply, the parallel version of lapply() is used
pairs = parLapply(cluster, folds, xval.fold, lambda1 = lambda1,
  lambda2 = lambda2)

return(do.call("rbind", pairs))
}

```

The function *parallel.xval* performs the 10-fold cross-validation. Note that once we have splitted the dat into the 10 folds, this is an embarrassingly parallel problem since the function *xval.fold* is applied 10 times independently and then the results are combined to compute the predictive correlation in the post processing phase. This is thus done using parallel computing with the *parLapply* function (the parallel equivalent of the *lapply* function). Table 1 below displays the elapsed time for 1 run of the function (the average running time of 10 runs) without parallel computing, and with parallel computing (2 slaves and then 4 slaves, which is the maximum that can be run on this computer). The inputs used here are $k = 10$, $\lambda_1 = \lambda_2 = 10$.

Implementation	Elapsed time (s)	Overhead time (s)
Sequential	1.795	Not applicable
Parallel (2 slaves)	1.061	0.164
Parallel (4 slaves)	1.035	0.586

Table 1: Elapsed time and overhead time to run the cross-validation function.

One can see from Table 1 that there is a considerable gain in running time from using 2 slaves but then the gain from 2 slaves to 4 slaves is negligible as the overhead time is about 3 times as high. Figure 1 shows an history of the parallel process for 2 slaves.



Figure 1: Organization of the workload shared between the two slave processes.

From Figure 1, we can see that the 2 slave processes share the work on the different folds of the cross-validation (green rectangles). And we can see that there is little overhead time and so the actual elapsed time is not too far from the ideal time. One could plot the same graph for 4 slaves but it becomes messy because of the many red lines indicating each time a slave process starts and finishes a task.

Furthermore, *parallel.xval()* needs to be performed 10 times, independently, and then take the average of the 10 predictive correlations. This is also embarrassingly parallel. We would then like to parallelize the 10 repetitions, the R code of both the sequential and parallelized version is shown in the Source Code 2 below.

Source Code 2: Repetition of the cross-validation function

```
# Sequential repetition using replicate
replicate(10, xval(data = markers.df, cluster = cl, k = 10,
lambda1 = 10, lambda2 = 10)[1,2])

# Parallelised implementation of the repetition using parSapply

cl <- makeCluster(2)
clusterEvalQ(cl, library(penalized))
clusterExport(cl, list("nki70", "xval", "markers.df"))
cur.predcor <- parSapply(cl, 1:10, xval)[2,]
stopCluster(cl)
```

In the above code, the same input were used as previously. Note that *xval* (see Appendix A) is the sequential version of the function *parallel.xval* in Source Code 1. In a similar way when we compared the sequential and parallelized implementations of Source Code 1, Table 2 below summarizes the performances of the two versions in Source Code 2.

Implementation	Elapsed time (s)	Overhead time (s)
Sequential	2.681	Not applicable
Parallel (2 slaves)	1.141	Not applicable
Parallel (4 slaves)	0.804	Not applicable

Table 2: Elapsed time and overhead time for 10 runs of the cross-validation function.

In this case, there is again a big gain from parallelizing the computation using 2 slaves, but the time gain starts fading when using 4 slaves. Note, that the actual time for 2 slaves is 1.141, while the sequential time is 2.681 which is more than twice as much as the actual time. This is because the sequential and parallel methods in Source Code 2 are not equivalent and they use a completely different implementation. Therefore, the sequential time and the parallel time are not linked like they were in Source Code 1 (where *parLapply* was the parallel equivalent of *lapply*). Figure 2 shows how the workload is shared when using 4 slaves. Slave 4 and 1 performs the function more times than 2 and 3 since 4 is not a factor of 10.

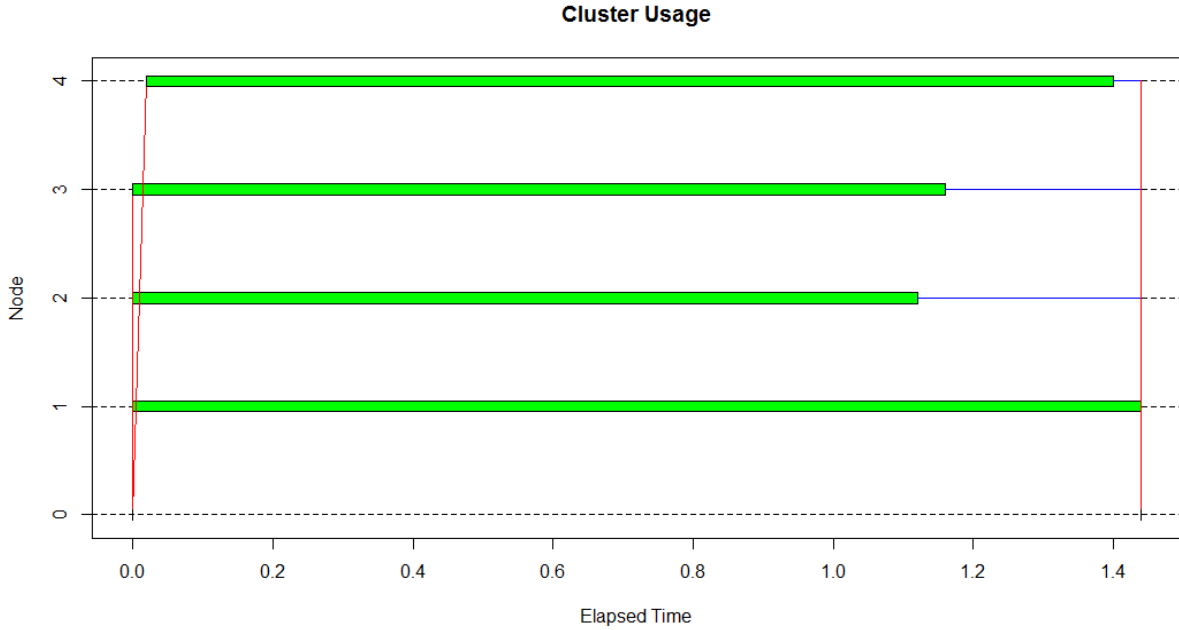


Figure 2: Organization of the workload shared between the 4 slave processes.

One can see that there is some overhead time for the fourth slave (it starts working a bit later). Since the data are exported once at the beginning, the small overhead is certainly due to the fact that the indices in *parSapply* are generated sequentially. They could be used all at once since they are independent and doing so would lead to a slight gain in time.

What we would like to do is to combine the parallel implementation in Source Code 1 with the one in Source Code 2. However, this does not work since it is not possible to fork slave processes (i.e, use slaves within slaves). This is a limitation in the implementation of our code and would need to be investigated further. Let us now look at the

greedy grid algorithm where we use the parallelization of Source Code 1 and the sequential implementation in Source Code 2.

3.2 Greedy grid algorithm

The full algorithm is displayed in Source code 3 below. We will firstly comment on the architecture (data structure, code organization) and then on the main flaws of the greedy grid algorithm.

Source Code 3: Greedy grid algorithm

```
tune.elasticNet = function(data, cluster, k, initial.lambda1
,initial.lambda2, stepping1, stepping2, epsilon)
{

  cur.lambda1 = initial.lambda1
  cur.lambda2 = initial.lambda2

  cur.predcor <- mean(replicate(10,parallel.xval(data, cluster = cluster,
k = k, lambda1 = cur.lambda1, lambda2 = cur.lambda2)[1,2]))

  # repeat the body until the condition in if is not met for the 9
# combinations of lambdas at each iteration
  new.lambda1 <- 0
  new.lambda2 <- 0

  repeat {
    new.lambda1 <- cur.lambda1 + seq(-stepping1, stepping1, stepping1)
    new.lambda2 <- cur.lambda2 + seq(-stepping2, stepping2, stepping2)
    correl <- matrix(0, nrow=3, ncol=3)

    for (j in 1:3)
    {
      for (m in 1:3)
      {
        correl[j,m] <- mean(replicate(10,cor(parallel.xval(data,
cluster=cluster, k=k, lambda1 = new.lambda1[j]
, lambda2 = new.lambda2[m]))[1,2]))
      }
    }

    if (max(correl) - cur.predcor >= epsilon)
    {
      indice <- which(correl == max(correl), arr.ind = TRUE)
      lambda1_ind <- indice[1]
      lambda2_ind <- indice[2]
      cur.lambda1 <- new.lambda1[lambda1_ind]
      cur.lambda2 <- new.lambda2[lambda2_ind]
      cur.predcor <- max(correl)
    }
  }
```

```

else {
  break
}
if(cur.lambda1 < stepping1) {stepping1 <- cur.lambda1/2}
if(cur.lambda2 < stepping2) {stepping2 <- cur.lambda2/2}
}
return(c(cur.predcor = cur.predcor, cur.lambda1 = cur.lambda1,
cur.lambda2 = cur.lambda2))
}

```

Note that while the computation in Source Code 1 and 2 are embarrassingly parallel, the greedy grid algorithm is now fine-grained parallel. This is because it first performs the cross-correlation function 10 times then obtain the resulting predictive correlation and iterate until the stopping criteria is reached. Let us comment on the important features of the greedy grid algorithm. We first need to have an initial predictive correlation which is assigned to the variable *cur.predcor*. This could have been done outside the function and just add the current predictive correlation as an extra input in the *tune.elasticNet* function. Then, in the repeat loop, the new proposed tuning parameters are assigned. Note that one could have used parallel computing to calculate the proposed values but it would not lead to a significant increase in performance since it is a trivial operation. Here, the structure used to store these values is a 3 by 3 matrix because it makes it easy to store and access the elements using a double loop on the matrix indices. One could have omitted to compute and store the pair $(\lambda_1 + 0, \lambda_2 + 0)$ since the predictive correlation was already computed for this point. However, it is not as straightforward to then assign and retrieve the pairs (λ_1, λ_2) of interest. Finally, we look for the maximum value in that matrix and retrieve its indices to know which pair of lambdas was responsible for it.

The solution given by the greedy algorithm is not guaranteed to be the optimal solution (regardless of the stepping size) as it may get stuck in a local optima. It is thus a good idea to try a grid of starting values which could be implemented using parallel computing. For instance, one could put a piece of code before the repeat loop in Source Code 3 which would look at several initial pairs, compute their predictive correlations and then pick the pairing giving the best correlation as the starting point. Furthermore, the stepping size is an hyper parameter that also affects the final result and must be taken into account. However, if enough starting points are considered across a domain of realistic values, the stepping size should not have such an impact.

4 Conclusions

While the focus of this report is more about the implementation of the algorithm rather than the application of it to the data set, the results will still be presented. After, looking at a few different starting values and step sizes, the best predictive correlation we obtained is 0.3615 at $\lambda_1 = 1$ and $\lambda_2 = 2$. Figure 3 shows the regression line of the elastic net with $\lambda_1 = 1$ and $\lambda_2 = 2$.

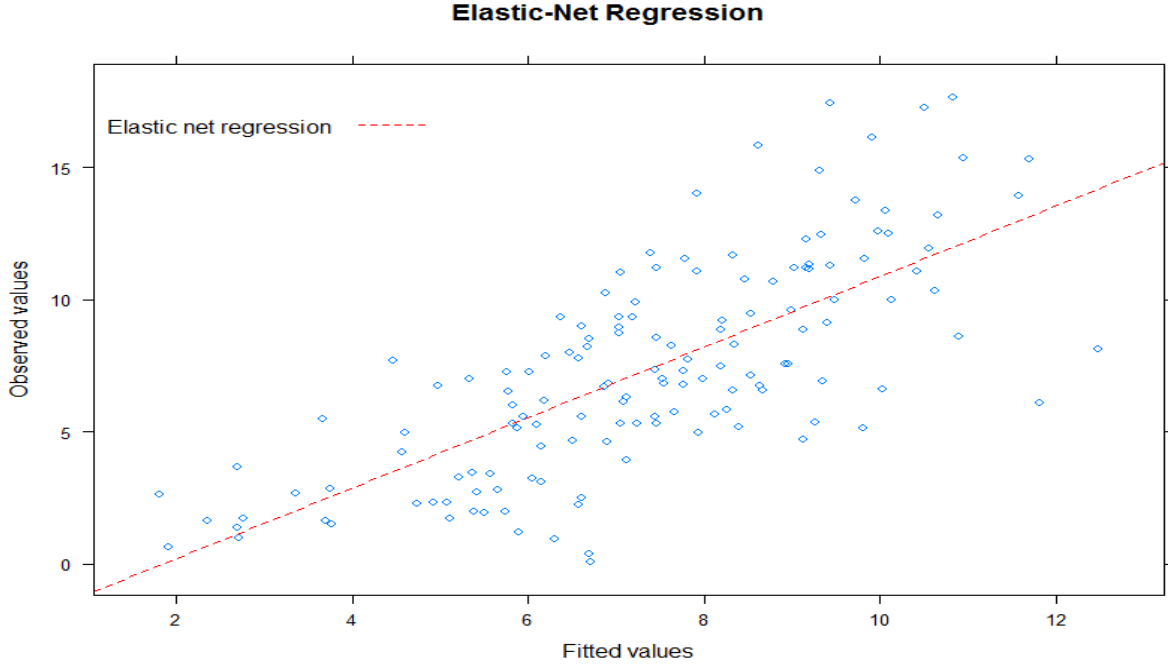


Figure 3: Elastic net regression line on the tested data.

One can see that the variance seems to increase and thus the correlation is better for the lower values than the higher values.

We have seen that there are two main parts that can be implemented using parallel computing: the cross-validation function and computing it 10 times. We have seen that in both cases, using 2 slave processes reduced consequently the amount of time needed to perform these tasks. However, using 4 slave processes only increased the performance marginally, largely due to a more important overhead time. One of the main issue with our code architecture is that it requires to nest (fork) slave processes, which cannot be done. It would thus be of interest to investigate different ways of implementing the cross-validation function and the algorithm in order to be able to parallelize both parts at once. We also discussed that we could add a piece of code which would look beforehand at an ideal starting point for the algorithm since guessing a sensible (but arbitrary nonetheless) starting point can lead to suboptimal solution as we could get stuck in a local optimum. This part of code could be implemented in parallel. Overall, the inclusion of parallel computing sped up the process significantly and improving the architecture of the code could lead to a better performance.

Appendix A

Source Code 4: Complete R code

```
# Load the libraries

library(lattice)
library(Rcpp)
library(grid)
library(ggplot2)
library(penalized)
library(parallel)
library(latticeExtra)
library(snow)
library(snowfall)
library(doSNOW)
library(doParallel)

# Get the data to test the algorithm

data("nki70")
nki70
dim(nki70)
markers.df <- as.data.frame(nki70[, -c(2:7)])
dim(markers.df)

# First step. Function not parallelised to obtain the initial
# predictive correlation.

xval = function(i,data, k, lambda1, lambda2,...) {

  n = nrow(data)

  folds = split(sample(n), seq_len(k))

  xval.fold = function(fold) {

    dtrain = data[-fold, ]
    dtest = data[fold, ]

    ElasticNet.fold = penalized(time, penalized = ~ . - time,
    data = dtrain, lambda1 = lambda1, lambda2 = lambda2, trace = FALSE)

    pred = predict(ElasticNet.fold, data = dtest)

    return(data.frame(PRED = pred[, "mu"], OBS = dtest$time))

  }
}
```

```

pairs = lapply(folds, xval.fold)
pairs2 = do.call("rbind", pairs)
return(cor(pairs2))
}

xval(i=1,markers.df, k=10, lambda1 = 20, lambda2 = 20)[1,2]
mean(replicate(10,xval(i=1,markers.df, k=10, lambda1 = 20,
lambda2 = 20)[1,2]))

# Make 100 runs of the 10 fold cross validation function and
# take the average time

system.time( for (i in 1:10) mean(replicate(10,xval(markers.df, k=10,
lambda1 = 20, lambda2 = 20)[1,2])))/10

# Parallel function to implement the algorithm

parallel.xval = function(data, cluster, k, lambda1, lambda2)
{
  n = nrow(data)

  folds = split(sample(n), seq_len(k))

  xval.fold = function(fold, lambda1, lambda2) {

    dtrain = data[-fold, ]
    dtest = data[fold, ]

    ElasticNet.fold = penalized(time, penalized = ~ . - time,
data = dtrain, lambda1 = lambda1, lambda2 = lambda2, trace = FALSE)

    pred = predict(ElasticNet.fold, data = dtest)

    return(data.frame(PRED = pred[, "mu"], OBS = dtest$time))
  }

  # Added code to implement parallelisation

  clusterExport(cluster, list("data"))
  clusterEvalQ(cluster, library(penalized))

  # Here, parLapply
  pairs = parLapply(cluster, folds, xval.fold, lambda1 = lambda1,
lambda2 = lambda2)

  return(do.call("rbind", pairs))
}

```

```

}

cl = makeCluster(2)
system.time( for (i in 1:10) mean(replicate(10,cor(parallel.xval(markers.df,
cluster=cl, k=10, lambda1 = 20, lambda2 = 20))[1,2])))/10
stopCluster(cl)

replicate(10,cor(parallel.xval(markers.df,cluster=cl, k=10,
lambda1 = 20, lambda2 = 20))[1,2])

# Plot the process of the slaves during the parallel computing

par(mfrow=c(1,1))

cl = snow::makeCluster(2)
plot(snow.time(replicate(10,parallel.xval(markers.df,
cluster=cl, k=10, lambda1 = 20, lambda2 = 20))))
snow::stopCluster(cl)

cl = snow::makeCluster(4)
plot(snow.time(replicate(10,parallel.xval(markers.df,cluster=cl,
k=10, lambda1 = 20, lambda2 = 20))))
snow::stopCluster(cl)

# Make the 10 repetitions of the cross-validation function parallel

xval2 <- function(i,data=markers.df, k=10, lambda1=10, lambda2=10,...)
{n = nrow(data)

folds = split(sample(n), seq_len(k))

xval.fold = function(fold) {

  dtrain = data[-fold, ]
  dtest = data[fold, ]

  ElasticNet.fold = penalized(time, penalized = ~ . - time,
data = dtrain, lambda1 = lambda1, lambda2 = lambda2, trace = FALSE)

  pred = predict(ElasticNet.fold, data = dtest)

  return(data.frame(PRED = pred[, "mu"], OBS = dtest$time))

}

pairs = lapply(folds, xval.fold)
pairs2 = do.call("rbind", pairs)
return(cor(pairs2))}

```

```

cl <- makeCluster(2)
clusterEvalQ(cl,library(penalized))
clusterExport(cl,list("nki70","xval","markers.df"))
cur.predcor <- parSapply(cl, 1:10, xval2)[2,]
stopCluster(cl)

# Compute the time for the parallelised repetition and the
# non parallelised

system.time( for (i in 1:10) replicate(10,xval(data = markers.df,
cluster = cl, k = 10, lambda1 = 10, lambda2 = 10)[1,2]))/10

cl <- makeCluster(4)
clusterEvalQ(cl,library(penalized))
clusterExport(cl,list("nki70","xval","markers.df"))
system.time( for (i in 1:10) cur.predcor <-
parSapply(cl, 1:10, xval2)[2,])/10
stopCluster(cl)

# Plot the process of the slaves during the parallel computing

par(mfrow=c(1,1))

cl = snow::makeCluster(4)
clusterEvalQ(cl,library(penalized))
clusterExport(cl,list("nki70","xval","markers.df"))
plot(snow.time(parSapply(cl, 1:10, xval2)[2,]))
snow::stopCluster(cl)

# Find the right parameters by now implementing the grid-search algorithm

tune.elasticNet = function(data, cluster, k, initial.lambda1,
initial.lambda2, stepping1, stepping2, epsilon)
{

  cur.lambda1 = initial.lambda1
  cur.lambda2 = initial.lambda2

  predcorVec <- 0

  cur.predcor <- mean(replicate(10,cor(parallel.xval(data,
cluster = cluster, k = k, lambda1 = cur.lambda1, lambda2 = cur.lambda2))[1,2]))
  new.lambda1 <- 0
  new.lambda2 <- 0
  # repeat the body until the condition in if is not met for

```

```

# the 9 combinations of lambdas at each iteration
repeat {

  new.lambda1 <- cur.lambda1 + seq(-stepping1, stepping1, stepping1)
  new.lambda2 <- cur.lambda2 + seq(-stepping2, stepping2, stepping2)
  correl <- matrix(0, nrow=3, ncol=3)

  for (j in 1:3)
  {
    for (m in 1:3)
    {
      correl[j,m] <- mean(replicate(10, cor(parallel.xval(data,
        cluster=cluster, k=k, lambda1 = new.lambda1[j],
        lambda2 = new.lambda2[m]))[1,2]))
    }
  }

  if (max(correl) - cur.predcor >= epsilon)
  {

    indices <- which(correl == max(correl), arr.ind = TRUE)
    lambda1_ind <- indices[1]
    lambda2_ind <- indices[2]
    cur.lambda1 <- new.lambda1[lambda1_ind]
    cur.lambda2 <- new.lambda2[lambda2_ind]
    cur.predcor <- max(correl)

  } #THEN

  else {

    break
  }

  if(cur.lambda1 < stepping1) {stepping1 <- cur.lambda1/2}

  if(cur.lambda2 < stepping2) {stepping2 <- cur.lambda2/2}

}

return(c(cur.predcor = cur.predcor, cur.lambda1 = cur.lambda1,
  cur.lambda2 = cur.lambda2))

}

# Compute the Elastic net function

cl = makeCluster(2)

```

```
tune.elasticNet(data = markers.df, cluster = cl, k = 10,
initial.lambda1 = 20, initial.lambda2 = 20, stepping1=6,
stepping2=6, epsilon=1e-8)
stopCluster(cl)
```

Plot fitted values of the Elastic net vs the observed values

```
ElasticNet2 = penalized(time, penalized = markers.df[,2:71] ,
data = nki70, lambda1 = 2, lambda2 = 1, trace = FALSE)
ElasticNet3 = penalized(time, penalized = markers.df[,2:71] ,
data = nki70, lambda1 = 10, lambda2 = 10, trace = FALSE)
```

Plot the regression line

```
par(mfrow=c(2,2))
```

```
xyplot(nki70[,1] ~fitted(ElasticNet2), main = "Elastic-Net Regression",
xlab = "Fitted values", ylab = "Observed values",
key=list(text=list(lab=c("Elastic net regression")),
corner=c(0,0.9),
lines=list(lty=2, lwd=1, col="red")),
panel = function(...) {
  panel.xyplot(...)
  panel.lmline(..., col = "red", lty = 2)
})
```

```
xyplot(nki70[,1] ~fitted(ElasticNet3), main = "Elastic-Net Regression",
xlab = "Fitted values", ylab = "Observed values",
key=list(text=list(lab=c("Elastic net regression")),
corner=c(0,0.9),
lines=list(lty=2, lwd=1, col="red")),
panel = function(...) {
  panel.xyplot(...)
  panel.lmline(..., col = "red", lty = 2)
})
```

Plot the residuals

```
xyplot(residuals(ElasticNet2) ~ fitted(ElasticNet2), main = "Residuals
of the elastic net regression",
xlab = "fitted values", ylab = "residuals",
panel = function(...) {
  panel.abline(h = 0, col = "grey", lty = 2)
  panel.xyplot(...)
  panel.loess(..., col = "darkred", lty = 2)
})
```