# Forward regression in R: from the extreme slow to the extreme fast

Michail Tsagris and Manos Papadakis

Department of Computer Science, University of Crete, Heraklion, Creece

mtsagris@yahoo.gr, papadakm95@gmail.com

**Abstract**

Forward regression has been criticised heavily and one of the many reasons is regarding its speed and its stopping criteria. The main focus of this paper is on demonstrating how to make it efficient, using R. Our method works for continuous predictor variables only, as the use of the partial correlation plays the most important role.

**Keywords**: forward regression, partial correlation coefficient, computational efficiency

## 1 Introduction

It is common truth that efficient code is not a matter of computer, operating system, or programming language, but a matter of implementation. Despite this ordinal common reasoning argument, not many programmers, scientific developers and mainly researchers put too much emphasis on this aspect.

Researchers looking for speed in their calculation tend to use Matlab, Python, C++, etc. In general, the aforementioned languages or software are faster than R. What most people are not aware of though is that R has many "powerful" functions whose efficiency is comparable to matlab for example.

A drawback of forward regression is speed. It has been long recognised that it is computationally prohibitive in many cases. Nevertheless it is taught in the undergraduate departments of statistics across the globe.

This paper deals exactly with this issue, right at the heart of the problem, how to make forward regression efficient. To this end, we will make use of R's built-in command *cor*; a very efficiently written function whose abilities we will take advantage here.

## 2 Forward selection method for linear models

Despite the title of this Section, the general idea of the forward selection method is applicable to all kinds of models, linear or not, regardless of normally distributed error or not. Suppose we have a univariate response variable $\mathbf{y} \in \mathbb{R}^n$ and a matrix of $p$ predictor variables $\mathbf{X} \in \mathbb{R}^{\mathbf{n} \times \mathbf{p}}$. The algorithm below summarizes the method.

<u>Forward selection method</u>

1. Standardise all predictor variables, (make them have zero mean and unit variance) and center the values of the response variable.

2. Examine all univariate associations between $\mathbf{y}$ and $\mathbf{x}_i$, for $i = 1, \ldots, p$. The $j$-th most significant among the significant ones enters the model.

3. Examine all conditional associations $(\mathbf{y}, \mathbf{x}_i) | CS$. The most significant among the significant ones enters the model. $CS$ stands for conditioning set.

4. Repeat Step 3, increasing the $CS$ every time by 1 (one variable is added) until no significant association exists.

In the next 2 subsections we will see how to perform Steps 2-4. That is, how to assess the significance of the (conditional) associations.

## 2.1   Significance of a univariate association

In order to assess the univariate associations, Pearson's correlation coefficient is calculated. The significance of the sample correlation $r$ will be assessed via Fisher's z-transform (Fisher, 1915)

$$z = \frac{1}{2} \log \frac{1+r}{1-r} = \tanh^{-1}(r)$$

Under $H_0$,

$$z \sim N\left(\frac{1}{2} \log \frac{1+\rho}{1-\rho}, \frac{1}{n-3}\right),$$

where $\rho$ is the true correlation coefficient. Alternatively, a $t_{n-3}$ distribution can be used. Note that this is different from assessing the significance of the slope coefficient of a simple linear regression in which case we would use the following test statistic

$$T = \frac{r\sqrt{n-2}}{\sqrt{1-r^2}}$$

Under $H_0$, $T \sim t_{n-2}$. Note that in the simple linear regression

$$y_i = \hat{\alpha} + \hat{\beta} z_i + e_i \tag{1}$$

$T$ coincides with the Wald test statistic for the $H_0 : b = 0$ which is given by

$$W = \frac{\hat{\beta}}{\text{s.e.}(\hat{\beta})} \tag{2}$$

and its square is equal to the usual F-test

$$T^2 = F = \frac{SSE_0 - SSE_1}{SSE_1 / (n - 2)},$$

where $SSE_1$ is the sum of squares of the errors of the regression model (1) and $SSE_0$ is simply $(n - 1)\sigma_y^2$, where $\sigma_y^2$ is the variance of $\mathbf{y}$. The z or t test based upon Fisher's z-transformation is asymptotically equivalent to the Wald (and hence the $T^2$ and the F tests).

The convenience of the z-transformation lies in the fact that its standard error has a known formula and requires no further calculation and this is taken into advantage when the computational cost is examined.

## 2.2 Significance of a conditional association

In the second and latter steps of the forward regression, the significance of the conditional correlations is to be examined. The two models are

$$
\begin{aligned}
H_0 : y_i &= \hat{\alpha}_1 + \mathbf{z}_i\hat{\boldsymbol{\beta}}_1 && +e_i \\
H_1 : y_i &= \hat{\alpha}_2 + \mathbf{z}_i\hat{\boldsymbol{\beta}}_2 && +\gamma x_i && +e_i.
\end{aligned}
\tag{3}
$$

Hence, under $H_0$ $\gamma = 0$. The partial F-test is

$$F = \frac{SSE_0 - SSE_1}{SSE_1 / (n - q - 1 - 1)}.$$

The conditional correlation must be used with the cardinality of the conditional set being $|\mathbf{z}| = q$, i.e there are $q$ variables. When there is only one conditioning variable, the conditional correlation has a closed form, but we provide the general way of calculating it (Fisher et al., 1924)

$$r(y, x | z) = r(e_{1i}, e_{2i}), \tag{4}$$

where $e_{1i}$ and $e_{2i}$ are the residuals of the two regression models

$$
\begin{aligned}
e_{1i} &= y_i - \hat{\alpha}_1 + \mathbf{z}_i\hat{\boldsymbol{\beta}}_1 && \tag{5} \\
e_{2i} &= x_i - \hat{\delta}_0 + \mathbf{z}_i\hat{\boldsymbol{\delta}}_1 && \tag{6}
\end{aligned}
$$

The variance of $r(y, x | z)$ is equal to $1/(n - 3 - q)$.

The $\gamma$ coefficient of the second model (3) is related to the partial correlation (4) via the following relationship

$$r(y, x | z) = \frac{W}{\sqrt{W^2 + n - q - 1}},$$

3

where $W$ is the Wald test statistic (2) for the $\gamma$ coefficient.

Note that, both for the partial correlation and the partial F test, two linear regression models are fitted. Hence, one can argue that the computational cost is the same. We will see however later on that this is not true when using R. Alternatively, the relationship equation can be used with one fitted model.

# 3   Forward selection with linear models in R

The R code for the forward selection method using the partial correlation can be found in the package **Rfast** (Papadakis et al., 2017). The first two steps of the method are rather straightforward. We remind the reader that all predictor variables have been centred and scaled to have unit variance and the response variable is only centred. In the first step, the command **cor(y, x)**, where **y** is a vector and **x** is a matrix is used. In the second step, with one conditioning variable, the partial correlation is applied.

The next (if necessary) steps are the most important ones regarding the speed gains. For this reason we take into advantage a very powerful base command in R, the **.lm.fit**, which is 10 times faster than **lm.fit**.

In the next two lines of R code, **$z$ is the selected variables at step $k$ (conditioning set)**, $x$ is the whole set of variables, including $z$ and $y$ is the response variable. The vector of partial correlations (4) between the response and each of the predictor variables conditioning on the selected variables is then calculated in the third line.

```
z <- x[, sela]
e1 <- .lm.fit(z, y)$residuals
e2 <- .lm.fit(z, x)$residuals
yx.z <- cor(e2, e1)
```

By taking a close look in the above R code segment, you will see that the calculation of the partial correlation coefficient via the residuals (5) is what makes the forward regression fast. According the algorithm of forward search, one must "scan" all non selected predictor variables at each step. So, with 1000 predictor variables, if 10 variables were to be selected, one would have to create roughly $10,000$ regression models. Using a *for* loop in R, this has already become slow.

We will now bridge the mathematics with the R code. The estimates of the linear regression coefficients are given by

$$\hat{\mathbf{B}} = \left(\mathbf{X}^T\mathbf{X}\right)^{-1}\mathbf{X}^T\mathbf{Y}. \tag{7}$$

Let us translate the two residual objects of the above R code into mathematics using (7)

$$\mathbf{e}_1 = \mathbf{Y} - \left(\mathbf{Z}^T\mathbf{Z}\right)^{-1}\mathbf{Z}^T\mathbf{Y} \ \ \text{(vector of residuals)}$$
$$\mathbf{E}_2 = \mathbf{X} - \left(\mathbf{Z}^T\mathbf{Z}\right)^{-1}\mathbf{Z}^T\mathbf{X} \ \ \text{(matrix of residuals)}.$$

Hence, instead of having to "scan" the design matrix of predictor variables at every step, all we need to do is two multiplications, one by the left and one by the right. The using the **cor** command, we get the vector of partial correlation coefficients at almost no time.

## 3.1 Two examples comparing time

We compared our implementation against the more general implementation of forward regression available in **MXM** (Lagani et al., 2017). **MXM** is a variable selection oriented package, offering many different methods for many types of data. The comparison might not look fair, mainly because the **MXM**'s command **lm.fsreg** for linear models is generic, it is designed to accept any type of predictor variables, continuous and/or categorical and uses command **lm** for this reason. In addition, the algorithm has not been optimized, using *.lm.fit* for example and incrementally updating the design matrix. That implementation is based on the algorithm of the forward regression. At each step, all the remaining variables are "scanned" and their significance is examined.

We did two single examples, with only 500 predictor variables. In the first case no variable is associated with the response variable. In the second case, the response is linearly related with three predictor variables. The relevant code to produce the data and the output using the package **microbenchmark** are given below.

```
x <- matrix(runif(500 * 500, 1, 100), ncol = 500)
## first case scenario
y <- rnorm(500)
mb <- microbenchmark(cor.fsreg(y, x), lm.fsreg(y, x) )

## second case scenario
y <- 3 * x[, 10] + 2 * x[, 100] + 3 * x[, 20] + rnorm(100, 0, 5)
mb2 <- microbenchmark(cor.fsreg(y, x), lm.fsreg(y, x) )
```

In the first example, **Rfast** required 22 miliseconds (0.02 seconds), whereas **MXM** required 2,204 miliseconds (2.2 seconds). In the second example **Rfast** required 64 miliseconds (0.064 seconds), whereas **MXM** required 6,590 miliseconds (6.59 seconds).

We can see that our implementation is 100 times faster than **MXM**'s generic implementation when no predictor variable is truly significant. It is 274 times faster when there are only three predictor variables truly significant. In the first case, two variables were selected, whereas in the second case, the three significant variables were selected.

When using the standard algorithm, and a new candidate variable is tried, the crossproduct of the design matrix may not be invertible, because the candidate variable is highly correlated with another one, and thus a check should be made at every step. However, this implementation of the forward regression will work just fine, because the QR decomposition implementation behind handles these cases.
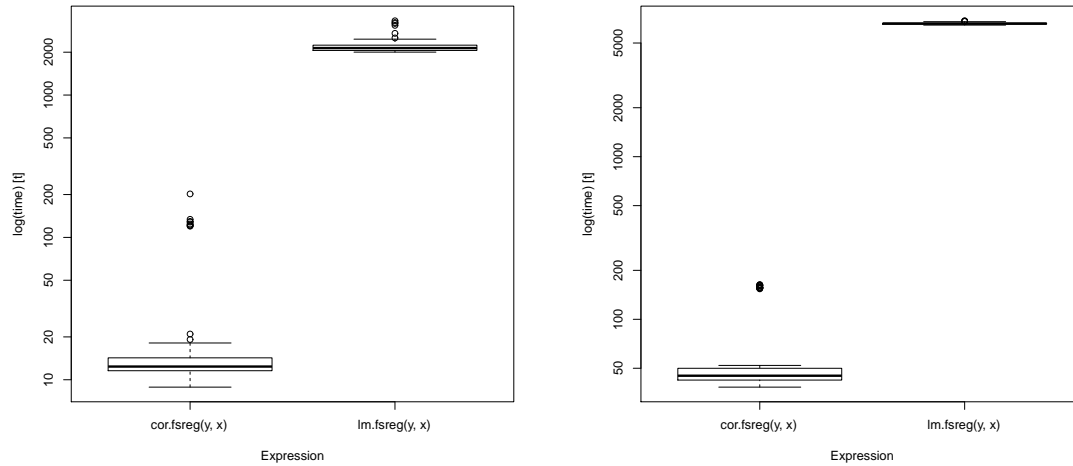
Figure 1: Box plots of execution time for both packages. The left graph refers to the first example, whereas the right graph to the second graph.

# 4  Conclusions

We have showed how to speed up dramatically the forward regression in the linear models case with continuous predictor variables only. A comparison with the generic (not fully optimised) forward regression implementation in **MXM** showed dramatic speed-up factors.

We demonstrated the fact that it is not a matter of computer, the computational power, or even of the algorithm many times, but a matter of the implementation. This sentence is very common, yet not everybody writes functions and packages without taking this into account. The second message here is that in order to have really fast functions, one must take cases and optimise each case separately, rather than have a generic function. It becomes clear though, that the optimisation of the generic function is the optimal solution.

# References

Fisher, R. A. (1915). Frequency distribution of the values of the correlation coefficient in samples from an indefinitely large population. *Biometrika*, 10(4):507–521.

Fisher, R. A. et al. (1924). The distribution of the partial correlation coefficient. *Metron*, 3(3-4):329–332.

Lagani, V., Athineou, G., Farcomeni, A., Tsagris, M., and Tsamardinos, I. (2017). Feature Selection with the R Package MXM: Discovering Statistically-Equivalent Feature Subsets. *Journal of Statistical Software*, 80(7).

Papadakis, M., Tsagris, M., Dimitriadis, M., Fafalios, S., Tsamardinos, I., Fasiolo, M.,

Borboudakis, G., Burkardt, J., Zou, C., and Lakiotaki, K. (2017). *Rfast: Fast R Functions*. R package version 1.8.6.