

# Using stepreg

Walter K. Kremers, Mayo Clinic, Rochester MN

10 May 2024

## The Package

The `stepreg()` and `cv.stepreg()` functions in the *glmnet* package were written for convenience and stability as opposed to speed or broad applicability. When fitting lasso models we wanted to compare these to standard stepwise regression models. Keeping a more modern approach we tune by either number of terms included in the model (James, Witten, Hastie and Tibshirani, *An Introduction to Statistical Learning with applications in R*, 2nd ed., Springer, New York, 2021) or by the  $p$  critical value for model inclusion, as this too is a common tuning parameter when fitting stepwise models.

When fitting lasso models we often use one-hot coding for predictor factors when setting up the design matrix. This allows lasso to identify and add to the model a term for any one group that might be particularly different from the others. By the penalty lasso stabilizes the model coefficients and keeps them from going to infinity, while ridge will generally uniquely identify coefficients despite any strict collinearities.

Before writing this program we tried different available packages to fit stepwise models for the Cox regression framework but all we tried had difficulties with numerical stability for the large and wide clinical datasets we were working with, and which involved one-hot coding. There may well be a package that would be stable for the data we were analyzing but we decided to write this small function to be able to tune for stability.

This program is slow but our goal was not for routine usage but to use the stepwise procedure on occasion as a reference for the lasso models. For many clinical datasets the lasso clearly outperformed the stepwise procedure, and ran much faster. For many simulated data sets with simplified covariance structures, i.e. independence of the underlying predictors, the lasso did not appear to do much better than the stepwise procedure tuned by number of model terms or  $p$ .

## Data requirements

The data requirements for `stepreg()` and `cv.stepreg()` are similar to those of `cv.glmnet()` and we refer to the *Using glmnet* vignette for a description.

## An example dataset

To demonstrate usage of `cv.stepreg` we first generate a data set for analysis, run an analysis and evaluate. Following the *Using glmnet* vignette, the code

```
# Simulate data for use in an example survival model fit  
# first, optionally, assign a seed for random number generation to get applicable results  
set.seed(116291950)  
simdata=glmnet.simdata(nrows=1000, ncols=100, beta=NULL)
```

generates simulated data for analysis. We extract data in the format required for input to the *cv.stepreg* (and *glmnet*) programs.

```
# Extract simulated survival data
xs = simdata$xs           # matrix of predictors
y_ = simdata$yt          # vector of survival times
event = simdata$event    # indicator of event vs. censoring
```

Inspecting the predictor matrix we see

```
# Check the sample size and number of predictors
print(dim(xs))
```

```
## [1] 1000 100
```

```
# Check the rank of the design matrix, i.e. the degrees of freedom in the predictors
Matrix::rankMatrix(xs)[[1]]
```

```
## [1] 94
```

```
# Inspect the first few rows and some select columns
print(round(xs[1:10,c(1:12,18:20)],digits=6))
```

```
##           X1 X2 X3 X4 X5 X6 X7 X8 X9 X10 X11 X12           X18           X19           X20
## [1,]  1  0  0  0  1  0  1  0  0  0  0  0  0 -1.208898  0.056971 -0.565631
## [2,]  1  1  0  0  0  0  0  0  1  0  0  0  0  0.395354  0.427313  0.185235
## [3,]  1  0  0  1  0  1  0  0  0  0  0  0  0  1.044608 -0.746960  0.964274
## [4,]  1  1  0  0  0  0  0  1  0  0  0  0  0  0.028859 -1.277651  0.203243
## [5,]  1  0  0  1  0  1  0  0  0  0  0  0  0 -1.205172 -1.287454 -1.698229
## [6,]  1  0  0  0  1  0  1  0  0  0  0  0  0 -1.158210 -0.068841  1.458800
## [7,]  1  0  0  0  1  0  0  0  1  0  0  0  0  0.151713  1.095396  1.476831
## [8,]  1  0  0  1  0  0  1  0  0  0  0  0  0 -0.139246 -0.424550  0.073340
## [9,]  1  1  0  0  0  0  0  1  0  0  0  0  0 -0.069326  0.172792  1.039656
## [10,] 1  0  0  1  0  0  1  0  0  0  0  0  0  0.677420  1.185946 -1.473551
```

## Cross validation (CV) informed stepwise model fit

To fit stepwise regression models where the number of model terms are informed by cross validation to select *df*, the number of model terms, and *p*, the entry threshold, we can use the function *cv.stepreg()* function.

```
# Fit a relaxed lasso model informed by cross validation
cv.stepwise.fit = cv.stepreg(xs,NULL,y_,event,family="cox",folds_n=5,steps_n=30,track=0)
```

Note, in the derivation of the stepwise regression models, individual coefficients may be unstable even when the model may be stable which elicits warning messages. Thus we “wrapped” the call to *cv.stepreg()* within the *suppressWarnings()* function to suppress excessive warning messages in this vignette. The first term in the call to *cv.stepreg()*, *xs*, is the design matrix for predictors. The second input term, here *NULL*, is for the start time in case (start, stop) time data setup is used in a Cox survival model. The third term is the outcome variable for the linear regression or logistic regression model and the time of event or censoring in case of the Cox model, and finally the fourth term is the event indicator variable for the Cox model taking

the value 1 in case of an event or 0 in case of censoring at time  $y_*$ . The forth term would be NULL for either linear or logistic regression. Currently the options for family are “gaussian” for linear regression, “binomial” for logistic regression (both using the *stats* glm() function) and “cox” for the Cox proportional hazards regression model using the coxph() function of the R *survival* package. If one sets track=1 the program will update progress in the R console. For track=0 it will not. To summarize the model fit and inspect the coefficient estimates we use the summary() function.

```
# summarize model fit ...
summary(cv.stepwise.fit)
```

```
##
## CV best df = 16, CV best p enter = 0.01 for 16 predictors
## in the full data model, from 100 candidate predictors
##
## df loglik.null loglik pvalue concordance std X2
## 1 16 -3709.825 -3705.723 0.004178366 0.8796415 0.005219351 -2.544254
## 2 16 -3709.825 -3705.723 0.004178366 0.8796415 0.005219351 -2.544254
## X3 X7 X10 X11 X12 X14 X16
## 1 -0.4123862 -0.5812514 0.6538633 -0.4939628 0.4246715 -1.387424 -1.647604
## 2 -0.4123862 -0.5812514 0.6538633 -0.4939628 0.4246715 -1.387424 -1.647604
## X18 X19 X20 X21 X23 X24 X25
## 1 0.7966722 -1.150425 -0.4928893 -0.1818494 1.075441 0.7174526 -0.4877742
## 2 0.7966722 -1.150425 -0.4928893 -0.1818494 1.075441 0.7174526 -0.4877742
## X62
## 1 -0.1259569
## 2 -0.1259569
```

To extract beta’s or calculate predicted we use the predict() function.

```
# get betas ...
betas = predict(cv.stepwise.fit)
t( betas[1:20,] )
```

```
## X1 X2 X3 X4 X5 X6 X7 X8 X9 X10 X11
## df 0 -2.544254 -0.4123862 0 0 0 -0.5812514 0 0 0.6538633 -0.4939628
## p 0 -2.544254 -0.4123862 0 0 0 -0.5812514 0 0 0.6538633 -0.4939628
## X12 X13 X14 X15 X16 X17 X18 X19 X20
## df 0.4246715 0 -1.387424 0 -1.647604 0 0.7966722 -1.150425 -0.4928893
## p 0.4246715 0 -1.387424 0 -1.647604 0 0.7966722 -1.150425 -0.4928893
```

```
# predicted ...
preds = predict(cv.stepwise.fit, xs)
t( preds[1:14,] )
```

```
## [,1] [,2] [,3] [,4] [,5] [,6] [,7]
## df -4.652185 -2.777916 -1.515435 -0.979273 0.3337369 -5.318352 -1.121909
## p -4.652185 -2.777916 -1.515435 -0.979273 0.3337369 -5.318352 -1.121909
## [,8] [,9] [,10] [,11] [,12] [,13] [,14]
## df -2.543347 -2.617922 -4.385983 -0.4020953 -4.200559 5.43046 -3.462096
## p -2.543347 -2.617922 -4.385983 -0.4020953 -4.200559 5.43046 -3.462096
```

## Nested cross validation

Because the values for lambda and gamma informed by CV are specifically chosen to give a best fit, model fit statistics for the CV derived model will be biased. To address this one can perform a CV on the CV derived estimates, that is a nested cross validation as argued for in SRDM (Simon R, Radmacher MD, Dobbin K, McShane LM. Pitfalls in the Use of DNA Microarray Data for Diagnostic and Prognostic Classification. J Natl Cancer Inst (2003) 95 (1): 14-18. <https://academic.oup.com/jnci/article/95/1/14/2520188>). This is done here by the `nested.glmnetr()` function.

```
# A nested cross validation to evaluate a cross validation informed stepwise fit
#nested.cox.fit = nested.glmnetr(xs,NULL,y_,event,family="cox",
#                               dostep=1,doaic=1,folds_n=5,steps_n=30,track=0)
y_ = simdata$y_
nested.gau.fit = nested.glmnetr(xs,NULL,y_,NULL,family="gaussian",
                               dostep=1,doaic=1,folds_n=3,steps_n=30,track=1)
```

```
#names(nested.gau.fit)
summary(nested.gau.fit)
```

```
## Sample information including number of records, number of columns in
## design (predictor, X) matrix, and df (rank) of design matrix:
## family n xs.columns xs.df null.dev/n
## gaussian 1000 100 94 7.96
##
## For LASSO, Stepwise regression tuned by df and p, and AIC, average (Ave) model
## performance measures from the 3-fold (nested) cross validation are given together
## with naive summaries calculated using all data without cross validation
##
## Ave DevRat Ave Int Ave Slope Ave R-square Ave Non Zero
## LASSO min 0.8705 -0.0407 1.0317 0.8715 51.3333
## LASSO minR 0.8687 -0.0107 1.0157 0.8696 29.3333
## LASSO minR.GO 0.8704 0.0051 0.9989 0.8707 13.3333
## Ridge 0.8539 -0.1241 1.0979 0.8608 99.0000
## Naive DevRat Naive R-square Non Zero
## LASSO min 0.8863 0.9420 51
## LASSO minR 0.8769 0.9364 13
## LASSO minR.GO 0.8769 0.9364 13
## Ridge 0.8919 0.9448 99
##
## Ave DevRat Ave Int Ave Slope Ave R-square Ave Non Zero
## Stepwise df tuned 0.8671 0.0052 0.9940 0.8674 16.6667
## Stepwise p tuned 0.8627 0.0282 0.9839 0.8631 24.0000
## Stepwise AIC 0.8629 0.0293 0.9803 0.8632 30.0000
## Naive DevRat Naive R-square Non Zero
## Stepwise df tuned 0.8827 0.9395 19
## Stepwise p tuned 0.8833 0.9399 20
## Stepwise AIC 0.8878 0.9422 30
```

For this example we use 3 folds. We would generally use between 5 or 10 folds in practice, to get reasonable run times and to better allow variability in variable selection.

Before providing analysis results the output first reports sample size and since this is for a Cox regression, the number of events, followed by the number of predictors and the df (degrees of freedom) of the design matrix,

as well as some information on “Tuning parameters” to compare the lasso model with stepwise procedures as described in JWHT (James, Witten, Hastie and Tibshirani, An Introduction to Statistical Learning with applications in R, Springer, New York, 2021). In general we have found in practice that the lasso performs better.

Next are the nested cross validation results. First are the per record (or per event in case of the Cox model) log-likelihoods which reflect the amount of information in each observation. Since we are not using large sample theory to base inferences we feel the per record are more intuitive, and they allow comparisons between datasets with unequal sample sizes. Next are the average number of model terms which reflect the complexity of the different models, even if in a naive sense, followed by the agreement statistics, concordance or r-square. These nested cross validated concordances should be essentially unbiased for the given design, unlike the naive concordances where the same data are used to derive the model and calculate the concordances (see SRDM).

In addition to evaluating the CV informed model fits using another layer of CV, the `nested.glmnet()` function does the CV fits based upon the whole data set. Here we see, not unexpectedly, that the concordances estimated from the nested CV are slightly smaller than the concordances naively calculated using the original dataset. Depending on the data the nested CV and naive agreement measures can be very similar or disparate.

Fit information for the CV fit can be gotten by extracting the `object$cv.stepreg.fit` object and calling the `summary()` and `predict()` functions.

```
# Summary of a CV model fit from a nested CV output object
summary(nested.gau.fit$cv.stepreg.fit)
```

```
##
## CV best df = 19, CV best p enter = 0.03 for 20 predictors
##      in the full data model, from 100 candidate predictors
##
##   df loglik.null   loglik    pvalue   rsquare rsquareadj      Int      X2
## 1 19  -2456.327 -1384.781 0.01532871 0.8827085 0.8804345 0.7934083 -2.374707
## 2 20  -2456.327 -1382.101 0.02060680 0.8833355 0.8809522 0.7951200 -2.373546
##      X3      X4      X6      X8      X10      X12      X14
## 1 -0.2896507 0.3961826 0.5041838 0.2439999 0.7493202 0.4179345 -1.623983
## 2 -0.2901617 0.3946920 0.5111328 0.2397952 0.7451425 0.4111942 -1.624389
##      X17      X18      X19      X20      X21      X23      X24
## 1 1.747628 0.8906858 -1.102188 -0.5406626 -0.1252505 1.090738 0.6988531
## 2 1.746315 0.8919730 -1.105330 -0.5425721 -0.1270239 1.090795 0.6985010
##      X25      X28      X43      X62      X79
## 1 -0.4341470 0.07509635 0.00000000 -0.08267686 -0.07654123
## 2 -0.4287512 0.07668351 0.07283277 -0.08163437 -0.07438254
```

```
# get betas ...
betas = predict(nested.gau.fit$cv.stepreg.fit)
t( betas[1:10,] )
```

```
##      Int X1      X2      X3      X4 X5      X6 X7      X8 X9
## df 0.7934083 0 -2.374707 -0.2896507 0.3961826 0 0.5041838 0 0.2439999 0
## p 0.7951200 0 -2.373546 -0.2901617 0.3946920 0 0.5111328 0 0.2397952 0
```

```
# get predicted ...
preds = predict(nested.gau.fit$cv.stepreg.fit,xs)
t( preds[1:8,] )
```

```
##      [,1]      [,2]      [,3]      [,4]      [,5]      [,6]      [,7]      [,8]
## df -2.063766 -0.2322597 1.603771 1.830213 3.245500 -2.49770 1.462672 0.7644801
## p  -2.130016 -0.1919559 1.643084 1.777001 3.244594 -2.48281 1.389711 0.6710150
```