

Variable Selection for Health Care Demand in Germany

Zhu Wang
Connecticut Children's Medical Center
University of Connecticut School of Medicine
zwang@connecticutchildrens.org

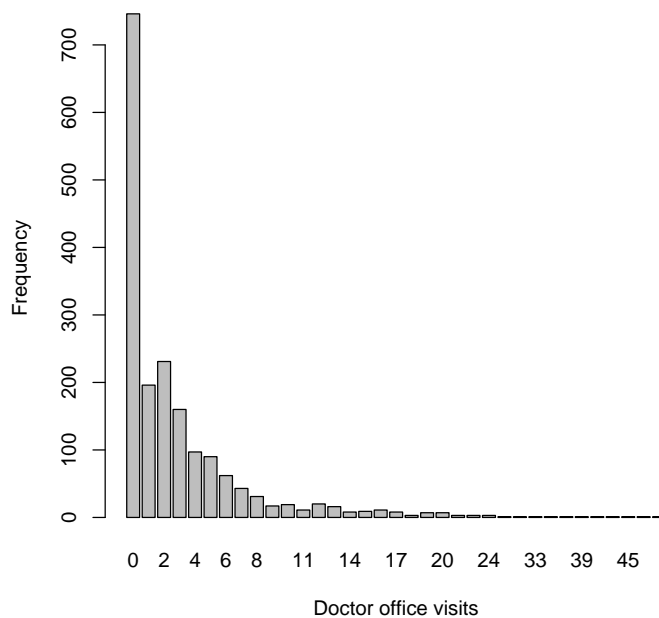
July 21, 2018

This document reproduces the data analysis presented in Wang et al. (2015). In an effort to optimizing the computing algorithms, the penalized regression can be slightly different. For a description of the theory behind application illustrated here we refer to the original manuscript.

Riphahn et al. (2003) utilized a part of the German Socioeconomic Panel (GSOEP) data set to analyze the number of doctor visits. The original data have twelve annual waves from 1984 to 1995 for a representative sample of German households, which provide broad information on the health care utilization, current employment status, and the insurance arrangements under which subjects are protected. The data set contains number of doctor office visits for 1,812 West German men aged 25 to 65 years in the last three months of 1994. As shown in the figure, many doctor office visits are zeros, which can be difficult to fit with a Poisson or negative binomial model. Therefore, zero-inflated negative binomial (ZINB) model is considered.

```
R> library("mpath")
R> library("zic")
R> library("pscl")
R> data(docvisits)

R> barplot(with(docvisits, table(docvisits)), ylab = "Frequency",
           xlab = "Doctor office visits")
```



We include the linear spline variables *age30* to *age60* and their interaction terms with the health satisfaction *health*.

```
R> dt <- docvisits[, -(2:3)]
R> tmp <- model.matrix(~age30 * health + age35 * health +
  age40 * health + age45 * health + age50 * health +
  age55 * health + age60 * health, data = dt)[, -(1:9)]
R> dat <- cbind(dt, tmp)
```

Full ZINB model with all predictor variables.

```
R> m1 <- zeroinfl(docvisits ~ . | ., data = dat, dist = "negbin")
R> summary(m1)
R> cat("loglik of zero-inflated model", logLik(m1))
R> cat("BIC of zero-inflated model", AIC(m1, k = log(dim(dat)[1])))
R> cat("AIC of zero-inflated model", AIC(m1))
```

Backward stepwise variable selection with significance level $\alpha=0.01$.

```
R> fitbe <- be.zeroinfl(m1, data = dat, dist = "negbin",
  alpha = 0.01, trace = FALSE)
R> summary(fitbe)
R> cat("loglik of zero-inflated model with backward selection",
  logLik(fitbe))
R> cat("BIC of zero-inflated model with backward selection",
  AIC(fitbe, k = log(dim(dat)[1])))
```

Compute LASSO estimates.

```
R> fit.lasso <- zipath(docvisits ~ . | ., data = dat, family = "negbin",
  nlambda = 100, lambda.zero.min.ratio = 0.001, maxit.em = 300,
  maxit.theta = 25, theta.fixed = FALSE, trace = FALSE,
  penalty = "enet", rescale = FALSE)
```

Estimated coefficient parameters with smallest BIC value.

```
R> minBic <- which.min(BIC(fit.lasso))
R> coef(fit.lasso, minBic)
R> cat("theta estimate", fit.lasso$theta[minBic])
```

Compute standard errors of coefficients and theta (the last one for theta).

```
R> se(fit.lasso, minBic, log = FALSE)
```

Compute AIC, BIC, log-likelihood values of the selected model.

```
R> AIC(fit.lasso)[minBic]
R> BIC(fit.lasso)[minBic]
R> logLik(fit.lasso)[minBic]
```

Compute log-likelihood value via 10-fold cross-validation using 2 CPU cores.

```
R> n <- dim(dat)[1]
R> K <- 10
R> set.seed(197)
R> foldid <- split(sample(1:n), rep(1:K, length = n))
R> fitcv <- cv.zipath(docvisits ~ . | ., data = dat, family = "negbin",
  nlambda = 100, lambda.count = fit.lasso$lambda.count[1:30],
  lambda.zero = fit.lasso$lambda.zero[1:30], maxit.em = 300,
  maxit.theta = 1, theta.fixed = FALSE, penalty = "enet",
  rescale = FALSE, foldid = foldid, n.cores = 2)
R> cat("cross-validated loglik", max(fitcv$cv))
```

Compute MCP estimates. We compute solution paths for the first 30 pairs of shrinkage parameters (the EM algorithm can be slow), and then evaluate results as for the LASSO estimates. For cross-validation, set maximum number of iterations in estimating scaling parameter 1 (maxit.theta=1) to reduce computation costs.

```
R> tmp <- zipath(docvisits ~ . | ., data = dat, family = "negbin",
  gamma.count = 2.7, gamma.zero = 2.7, lambda.zero.min.ratio = 0.1,
  maxit = 1, maxit.em = 1, maxit.theta = 2, theta.fixed = FALSE,
  penalty = "mnet")
R> fit.mcp <- zipath(docvisits ~ . | ., data = dat, family = "negbin",
  gamma.count = 2.7, gamma.zero = 2.7, lambda.count = tmp$lambda.count[1:30],
  lambda.zero = tmp$lambda.zero[1:30], maxit.em = 300,
  maxit.theta = 25, theta.fixed = FALSE, penalty = "mnet")
```

Estimated coefficient parameters with smallest BIC value.

```
R> minBic <- which.min(BIC(fit.mcp))
R> coef(fit.mcp, minBic)
R> cat("theta estimate", fit.mcp$theta[minBic])
```

Compute standard errors of coefficients and theta (the last one for theta).

```
R> se(fit.mcp, minBic, log = FALSE)
```

Compute AIC, BIC, log-likelihood values of the selected model.

```
R> AIC(fit.mcp)[minBic]
R> BIC(fit.mcp)[minBic]
R> logLik(fit.mcp)[minBic]
```

Compute log-likelihood value via 10-fold cross-validation using 2 CPU cores.

```
R> fitcv <- cv.zipath(docvisits ~ . | ., data = dat, family = "negbin",
  gamma.count = 2.7, gamma.zero = 2.7, lambda.count = tmp$lambda.count[1:30],
  lambda.zero = tmp$lambda.zero[1:30], maxit.em = 300,
  maxit.theta = 1, theta.fixed = FALSE, penalty = "mnet",
  rescale = FALSE, foldid = foldid, n.cores = 2)
R> cat("cross-validated loglik", max(fitcv$cv))
```

Compute SCAD estimates.

```
R> tmp <- zipath(docvisits ~ . | ., data = dat, family = "negbin",
  gamma.count = 2.5, gamma.zero = 2.5, lambda.zero.min.ratio = 0.01,
  maxit = 1, maxit.em = 1, maxit.theta = 2, theta.fixed = FALSE,
  penalty = "snet")
R> fit.scad <- zipath(docvisits ~ . | ., data = dat, family = "negbin",
  gamma.count = 2.5, gamma.zero = 2.5, lambda.count = tmp$lambda.count[1:30],
  lambda.zero = tmp$lambda.zero[1:30], maxit.em = 300,
  maxit.theta = 25, theta.fixed = FALSE, penalty = "snet")
```

Estimated coefficient parameters with smallest BIC value.

```
R> minBic <- which.min(BIC(fit.scad))
R> coef(fit.scad, minBic)
R> cat("theta estimate", fit.scad$theta[minBic])
```

Compute standard errors of coefficients and theta (the last one for theta).

```
R> se(fit.scad, minBic, log = FALSE)
```

Compute AIC, BIC, log-likelihood values of the selected model.

```
R> AIC(fit.scad)[minBic]
R> BIC(fit.scad)[minBic]
R> logLik(fit.scad)[minBic]
```

Compute log-likelihood value via 10-fold cross-validation using 2 CPU cores.

```
R> fitcv <- cv.zipath(docvisits ~ . | ., data = dat, family = "negbin",
  gamma.count = 2.5, gamma.zero = 2.5, lambda.count = tmp$lambda.count[1:30],
  lambda.zero = tmp$lambda.zero[1:30], maxit.em = 300,
  maxit.theta = 1, theta.fixed = FALSE, penalty = "snet",
  rescale = FALSE, foldid = foldid, n.cores = 2)
R> cat("cross-validated loglik", max(fitcv$cv))
```

Running time for the entire analysis.

```

R> print(proc.time() - ptm)

R> sessionInfo()

R version 3.4.4 (2018-03-15)
Platform: x86_64-pc-linux-gnu (64-bit)
Running under: Ubuntu 14.04.5 LTS

Matrix products: default
BLAS: /usr/lib/libblas/libblas.so.3.0
LAPACK: /usr/lib/lapack/liblapack.so.3.0

locale:
 [1] LC_CTYPE=en_US.UTF-8      LC_NUMERIC=C
 [3] LC_TIME=en_US.UTF-8      LC_COLLATE=C
 [5] LC_MONETARY=en_US.UTF-8  LC_MESSAGES=en_US.UTF-8
 [7] LC_PAPER=en_US.UTF-8     LC_NAME=C
 [9] LC_ADDRESS=C             LC_TELEPHONE=C
[11] LC_MEASUREMENT=en_US.UTF-8 LC_IDENTIFICATION=C

attached base packages:
[1] stats      graphics  grDevices  utils      datasets  methods
[7] base

other attached packages:
[1] pscl_1.4.6           MASS_7.3-50
[3] zic_0.8.1            coda_0.16-1
[5] lattice_0.20-33      RcppArmadillo_0.4.450.1.0
[7] Rcpp_0.12.7          mpath_0.3-5
[9] knitr_1.14

loaded via a namespace (and not attached):
 [1] magrittr_1.5          splines_3.4.4         doParallel_1.0.8
 [4] bst_0.3-15            gbm_2.1.3             R.cache_0.10.0
 [7] foreach_1.4.4        stringr_1.0.0         tools_3.4.4
[10] parallel_3.4.4       grid_3.4.4            glmnet_1.9-8
[13] R.oo_1.18.0           iterators_1.0.7       survival_2.41-3
[16] numDeriv_2012.9-1    R.rsp_0.20.0         Matrix_1.2-5
[19] formatR_1.2.1        R.utils_1.34.0        codetools_0.2-15
[22] rpart_4.1-13         evaluate_0.8          stringi_0.4-1
[25] compiler_3.4.4       R.methodsS3_1.7.0

```

References

- Regina T Riphahn, Achim Wambach, and Andreas Million. Incentive effects in the demand for health care: a bivariate panel count data estimation. *Journal of Applied Econometrics*, 18(4):387–405, 2003.
- Zhu Wang, Shuangge Ma, and Ching-Yun Wang. Variable selection for zero-inflated and overdispersed data with application to health care demand in

germany. *Biometrical Journal*, 33(29):5192–208, 2015. Article first published online: 8 JUN 2015 DOI: 10.1002/bimj.201400143.