

Generalized Linear Model with Elastic Net

Alina

25 11 2020

```
library(glmnet)
library(clusterSim)
setwd("~/GitHub/fallstudien_2_projekt_1/datasets")
data <- read.csv("dataset_ul.csv", header = TRUE, sep = ",")

#data <- data[c(data$provider == "o2", data$scenario == "highway"),]
```

Data

read the table and select the covariates for the model

```
data <- subset(data, select = c(scenario,
                                provider,
                                rsrp_dbm,
                                rsrq_db,
                                rssnr_db,
                                payload_mb,
                                f_mhz,
                                throughput_mbits,
                                time_s,
                                ci))
```

now eliminate the rows with NA's; "complete.cases" returns a logical vector indicating which cases are complete, i.e., have no missing values

```
data <- data[complete.cases(data),]
```

separate the full dataset in train and test data, whereby 75% of the data get to be the training set and the rest will be the test set

```
#set.seed(101)
sample <- sample.int(n = nrow(data), size = floor(.75*nrow(data)), replace = F)
train <- data[sample, ]
test <- data[-sample, ]
```

in our case there are two variables that contain factors, therefore we have to encode them, so that the model can handle them; one-hot ENCODING is used

```

X <- makeX(train, test = test)

train <- X[["x"]]
test <- X[["xtest"]]

x.train <- subset(train, select = -throughput_mbits)
y.train <- subset(train, select = throughput_mbits)

x.test <- subset(test, select = -throughput_mbits)
y.test <- subset(test, select = throughput_mbits)

```

scale the variables for feature importance, whereby the test set has to be scaled with mean and standard deviation from train set

```

scaled.train <- scale(train)
scaled.test <- scale(test, center = attr(scaled.train, "scaled:center"),
                                scale = attr(scaled.train, "scaled:scale"))

```

Hyperparameter Tuning

while the parameter lambda is determined via cross validation, the best choice for alpha has to be detected by the user himself; “a” is the number of values of alphas that are examined and “alpha.opt” is the optimal value of the input alphas (“alpha.opt” = 0 LASSO, “alpha.opt” = 1 RIDGE regression)

```

a <- 20

fits <- list()
for (i in 0:a){
  fit <- paste0("alpha", i)
  fits[[fit]] <- cv.glmnet(x.train, y.train, type.measure = "mse",
                          nfolds = 30, alpha = i/a, family = "gaussian")
}

results <- data.frame()
for (i in 0:a){
  fit <- paste0("alpha", i)
  predicted <- predict(fits[[fit]], s = fits[[fit]]$lambda.min,
                      newx = x.test)
  mse <- mean((y.test - predicted)^2)
  res <- data.frame(alpha = i/a, mse = mse)
  results <- rbind(results, res)
}

alpha.opt <- results$alpha[results$mse == min(results$mse)]
cat("optimal alpha:", alpha.opt)

```

```
## optimal alpha: 0
```

Modelfitting

fit the glmnet model with cross validation for the penalty parameter lambda the parameter alpha for the elastic net model has to be set by user

```
fit.cv.scaled <- cv.glmnet(subset(scaled.train, select = -throughput_mbits),
                          subset(scaled.train, select = throughput_mbits),
                          type.measure = "mse", nfolds = 30, alpha = alpha.opt)

fit.cv <- cv.glmnet(x.train, y.train, type.measure = "mse", nfolds = 30, alpha = alpha.opt)
```

Prediction

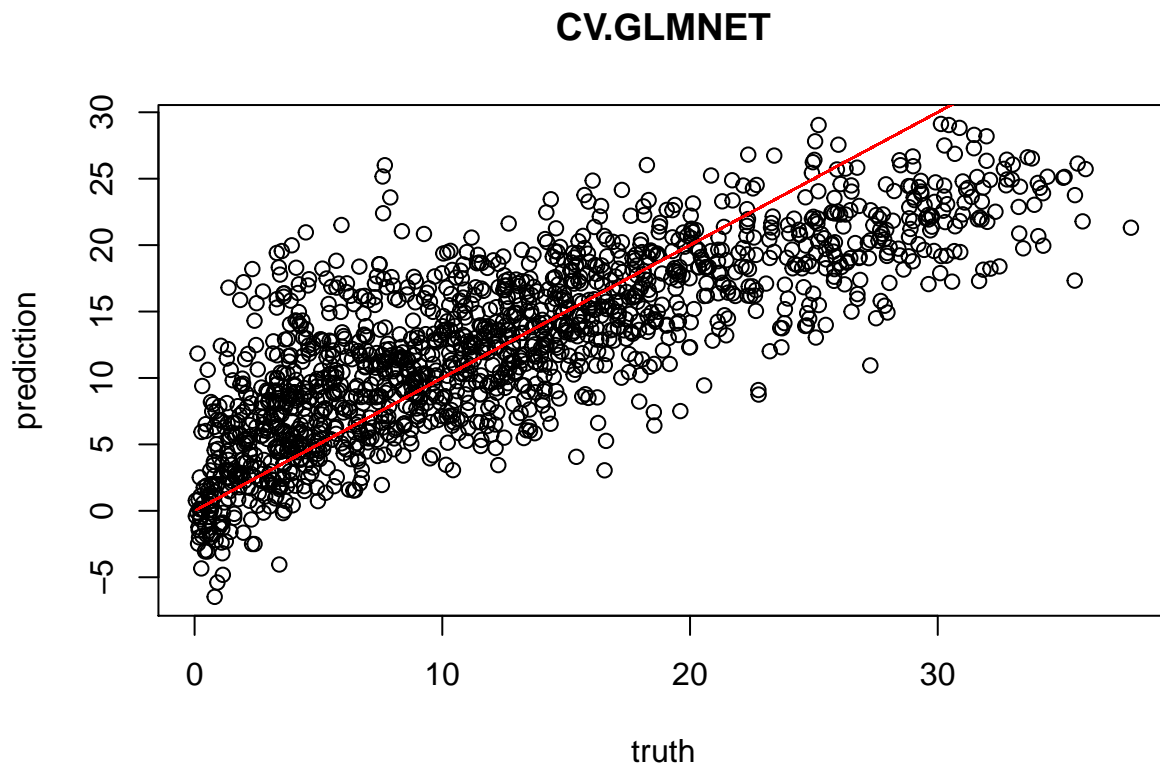
from the fitted `cv.glmnet` model we now generate the predictions with the covariates from the test set, we hereby use the penalty parameter `lambda` that generates the lowest error in the cv process

```
pred.cv <- predict(object = fit.cv, newx = x.test, s = "lambda.min", type = "response")
```

Results

plot the predictions from the `cv.glmnet` model against the truth values from our test set

```
plot(y.test, pred.cv, main = "CV.GLMNET", xlab = "truth", ylab = "prediction")
lines(y.test, y.test, col = "red")
```



we want to know the model rating, therefore we calculate the R-squared, MSE and MAE

```
yq <- mean(y.test)
R2 <- sum((pred.cv-yq)^2)/sum((y.test-yq)^2)
cat("R2", R2)
```

```
## R2 0.5961709
```

```
cat("\n")
```

```
n <- 1/length(pred.cv)
mse <- n*sum((y.test-pred.cv)^2)
cat("MSE", mse)
```

```
## MSE 28.82485
```

```
cat("\n")
```

```
mae <- n*sum(abs(y.test-pred.cv))
cat("MAE", mae)
```

```
## MAE 4.159306
```

Feature Importance

compare the absolute coefficients of the model, the larger the value the more information does the corresponding covariate brings

```
coef <- abs(coef(fit.cv.scaled))
#coef
```

```
coef[order(coef)]
```

```
## <sparse>[ <logic> ] : .M.sub.i.logical() maybe inefficient
```

```
## [1] 2.383183e-16 2.963195e-03 3.990639e-03 2.270952e-02 4.174826e-02
## [6] 5.526033e-02 6.480626e-02 1.174793e-01 1.286586e-01 1.336686e-01
## [11] 1.348341e-01 1.656257e-01 1.916363e-01 3.247893e-01
```

single models

run the model with the following data

```
#data <- data[data$provider == "vodafone",] R2 ~ 64
#data <- data[data$provider == "o2",] R2 ~ 54
#data <- data[data$provider == "tmobile",] R2 ~ 53
```