

Statistics 452: Statistical Learning and Prediction

Review Part 2: Predicting a HUI score

Brad McNeney

2018-11-26

Data

- ▶ One of the datasets given to the Stat 652 class includes 590 explanatory variables and a health utilities index score, called HUIDHSI.

This derived variable is a Health Utilities Index which provides a description of an individual's overall functional health, based on eight attributes: vision, hearing, speech, ambulation (ability to get around), dexterity (use of hands and fingers), emotion (feelings), cognition (memory and thinking) and pain (in HUP module). The version of the index used in CCHS is adapted from the HUI Mark 3 (HUI3). The index is designed to produce an overall health utility score. This multi-attribute utility index produces a score ranging from 1.000 (perfect health), through 0.000 (health status equal to death) to -0.360 (health status worse than death).

- ▶ I have unzipped the data file in my copy of the Project652 directory on github.

```
hs <- read.csv("../..../Project652/HStrain.csv")
```

- ▶ 591 variables, grouped into 38 categories, indicated by the first three letters of the variable names:

```
cn <- colnames(hs)
table(substr(cn,start=1,stop=3))
```

```
##
## ADL ADM ALC CAG CCC CGE CIH CR1 CR2 DHH DPS DS2 EDU FAL GEN GEO HC2 HUI
##  4  4  5 45 31  8 27 34 19  9 33  4  2 15 10  2 19  1
## HUP HWT IAL IN2 LBF LON MED NUR OH3 OWN PA2 RET RPL SDC SLP SLS SMK SPA
##  1  5  6  8 19  4 33 12 27  2 49 32 19  6  1  6 21 24
## SSA TRA
## 25 19
```

```
head(sort(cn),n=4) # Activities of Daily Living
```

```
## [1] "ADL_04A" "ADL_06A" "ADLDCLS" "ADLDT0I"
```

Survey information not useful for prediction

- ▶ The variables that start with ADM are to do with administering the survey and are not useful for prediction.
- ▶ For example, ADM_RNO is a sequential record number, ADM_N09 indicates whether the interview was by phone, in-person, etc.
- ▶ I will remove these.

```
library(dplyr)
hs <- select(hs, -starts_with("ADM"))
```

Summary variables

- ▶ Several categories of variable have an overall summary score, or classification, developed by survey experts.
- ▶ For example, the Activities of Daily Living (ADL) variables are summarized by ADLDCLS (page 158 of data dictionary):

ADLDCLS - Instrumental & Basic Activities of Daily Living Class. - Based on ADLDCLST and ADLDMEA. This variable is an overall summary measure of ratings of the ADL capacity-instrumental and physical dimensions. The instrument and the derived variable classification are developed from the activities of daily living component of the OARS Multidimensional Functional Assessment Questionnaire (OMFAQ). See documentation on derived variables.

```
summary(hs$ADLDCLS)
```

##	MILD IMPAIRMENT	MOD IMPAIRMENT	NO FUNC IMPAIR	SEV IMPAIRMENT
##	1098	258	8567	56
##	TOTAL IMPAIRMENT			
##	21			

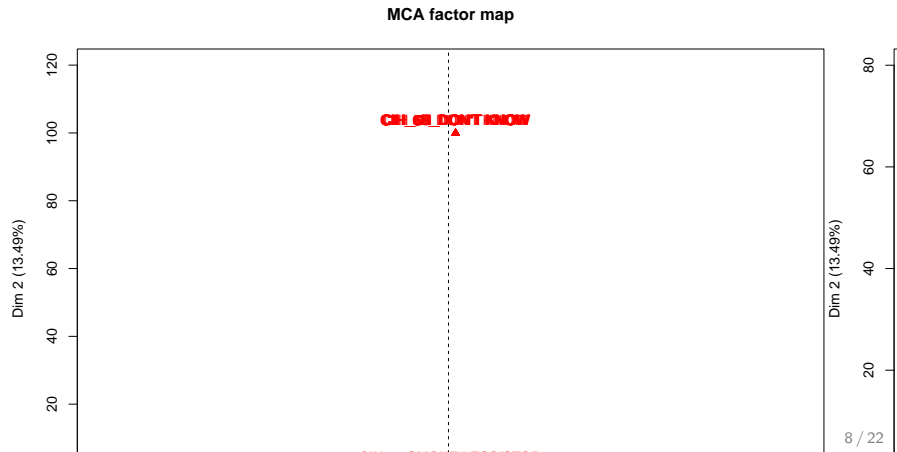
Choice of summary variable

- ▶ Some sets of variables do not have a single score, but may have several that could be useful.
- ▶ For example the caregive variables CAG tell us about care responsibilities (parent, spouse, neighbor, etc.)
 - ▶ CAGDFAP records the frequency of care (rarely, monthly, weekly, daily, etc.)
 - ▶ CAGDIAP records frequency and number of hours (daily - 1 hour, daily - 3 hours, etc.)

Making our own score

- ▶ Some groups of variables have no summary score.
- ▶ Can compute a few PCs from a group of survey questions.
- ▶ Example, CIH variables (questions about improving health).

```
library(FactoMineR)
res.mca <- MCA(select(hs, starts_with("CIH")))
```




```
CIHPCs <- res.mca$ind$coord[,1:4] # first 4 explain 50%  
colnames(CIHPCs) <- paste("CIH",colnames(CIHPCs))
```

My choices

```
hsred <- select(hs,  
  ADLDCLS,ALCDTTM,CAGDFAP,CCCF1,CCDCPD,  
  CR1FRHC,CR2DTHC,CR2DFAR,DPSDSF,EDUDR04,  
  FALGO2,GENDHDI,GENDMHI,HC2FCOP,  
  HUIDHSI, # response  
  HUPDPAD,HWTGBMI,IN2GHH,LONDSCR,MEDF1,  
  NURDHNR,PA2DSCR,SLP_02,SLSDCLS,  
  SMKDSTY,SPAFFAR,starts_with("SSAD"))  
hsred <- data.frame(hsred,CIHPCs)
```

Training and test sets

- ▶ To compare methods we'll divide our 10000 observations into training and test sets.
- ▶ I'll go with a 70% training set

```
set.seed(123)
n.train <- 7000
train <- sample(1:nrow(hs),replace=FALSE,size=n.train)
hsred.train <- hsred[train,]
hsred.test <- hsred[-train,]
```

Subset selection

```
library(leaps)
rr <- regsubsets(HUIDHSI ~ ., data=hsred.train, nvmax=40,
                 method="forward")
```

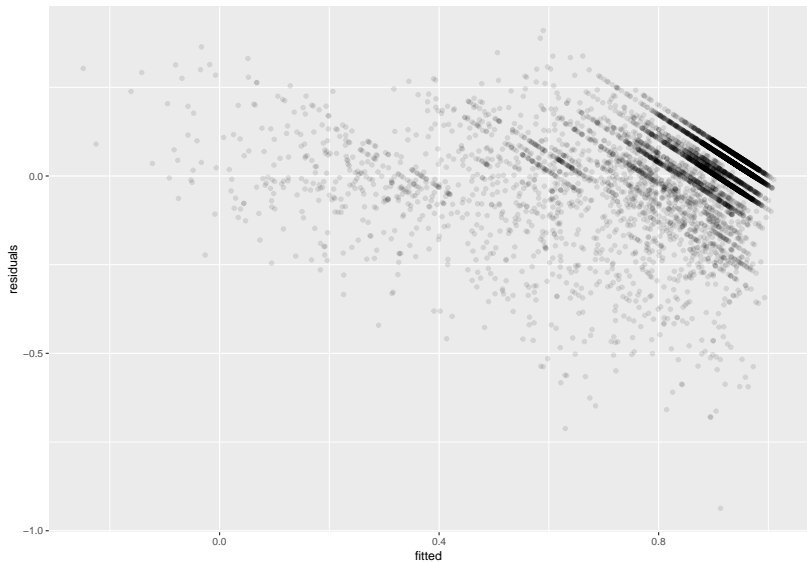
Reordering variables and trying again:

```
ss <- summary(rr)
pbest <- which.min(ss$bic)
# coef(rr, id=30) Important variables are
# Activities of daily living, general health,
# satisfaction with life, smoking status, two
# of our "improve health" PCs, and some others
```

```
Xfull.train <- model.matrix(HUIDHSI ~ .,data=hsred.train)
Xred <- Xfull.train[,ss$which[pbest,]]
Y.train <- hsred.train$HUIDHSI
ll <- lm.fit(Xred,Y.train)
Xfull.test <- model.matrix(HUIDHSI ~ .,data=hsred.test)
Xred <- Xfull.test[,ss$which[pbest,]]
pred.test <- Xred %*% ll$coef
Y.test <- hsred.test$HUIDHSI
mean((Y.test - pred.test)^2)
```

```
## [1] 0.01356958
```

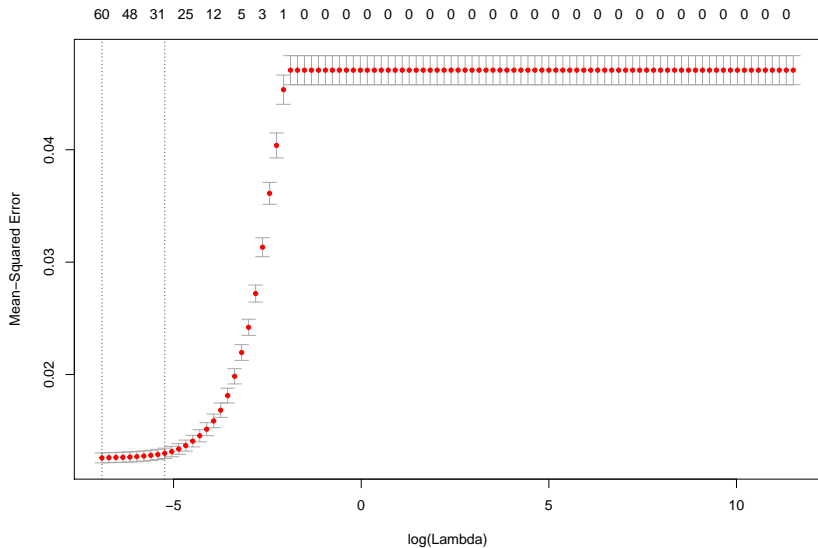
```
dd <- data.frame(fitted=ll$fitted,residuals=ll$residuals)
library(ggplot2)
ggplot(dd,aes(x=fitted,y=residuals)) + geom_point(alpha=.1)
```



Lasso

```
library(glmnet)
lambdas <- 10^{seq(from=-3,to=5,length=100)}
cv.lafit <- cv.glmnet(Xfull.train,Y.train,alpha=1,lambda=lambdas)
```

```
plot(cv.lafit)
```



```
la.best.lam <- cv.lafit$lambda.1se
```


Lasso coefficients

- ▶ A slightly larger set of non-zero coefficients than subset selection, but very similar.

```
ll <- glmnet(Xfull.train,Y.train,alpha=1,lambda=la.best.lam)
coef(ll)
```

```
## 82 x 1 sparse Matrix of class "dgCMatrix"
##                                     s0
## (Intercept)                7.750715e-01
## (Intercept)                .
## ADLDCLSMOD IMPAIRMENT      -9.602517e-02
## ADLDCLSNO FUNC IMPAIR      9.230568e-02
## ADLDCLSSEV IMPAIRMENT     -1.579962e-01
## ADLDCLSTOTAL IMPAIRMENT   -9.771975e-02
## ALCDDTMOCCASIO. DRINKER    .
## ALCDDTMREGULAR DRINKER     .
## CAGDFAPOCC OR RARELY       .
## CAGDFAPREG BASIS DLY       .
## CAGDFAPREG BASIS LESS      .
## CAGDFAPREG BASIS MNTH      .
## CAGDFAPREG BASIS WK        .
## CCCF1HAS NO CHRON CON      1.108552e-02
## CCCDCPDNOT HAVE COPD       .
## CR1FRHCREC FORMAL H C     -6.268507e-03
## CR2DTHCDID NOT REC H C     3.336383e-02
## CR2DTHCFORMAL H C ONLY     .
## CR2DTHCINFORMAL H C ONL    .
## CR2DFAROCC OR RARELY       .
## CR2DFARREG BASIS DLY       -9.683494e-03
## CR2DFARREG BASIS LESS      .
## CR2DFARREG BASIS MNTH      .
## CR2DFARREG BASIS WK        .
```

```
pred.test <- predict(l1,Xfull.test)
mean((Y.test-pred.test)^2)
```

```
## [1] 0.01378689
```

Random forests

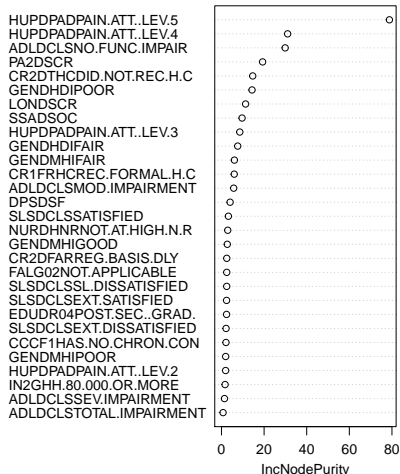
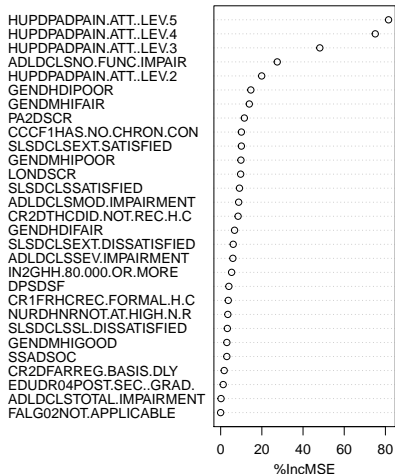
- ▶ Computationally intensive, and computation grows with the number of features and observations.
- ▶ With the HS data we have 80 features – need to filter some out, or need to reduce the sample size while we explore which features are important.

Use features found by lasso or subset selection

```
nonz <- (as.numeric(coef(l1))!=0)[-1] # rm intercept  
hsred2.train <- data.frame(HUIDHSI=Y.train,Xfull.train[,nonz])
```

```
library(randomForest)
set.seed(1)
bb <- randomForest(HUIDHSI ~ ., data=hsred2.train, ntree=200,
                    mtry=sqrt(ncol(hsred2.train)), importance=TRUE)
varImpPlot(bb)
```

bb



```
hsred2.test <- data.frame(HUIDHSI=Y.test,Xfull.test[,nonz])  
pred.test <- predict(bb,newdata=hsred2.test)  
mean((Y.test - pred.test)^2)
```

```
## [1] 0.01394767
```