

# Report code

Hannah Sansford, Ettore Fincato, Harry Tata

This markdown file contains the code used to produce all plots and conduct all statistical analysis contained within the report.

## Setup

First load the important packages.

```
library(tidyverse)
library(reshape2)
library(ggplot2)
library(glmnet)
library(ROCR)
library(caret)
library(cvms)
library(tibble)
library(pROC)
library(checkmate)
library(BBmisc)
library(testit)
library(devtools)
library(mlr)
```

Next, install and load the `higgsboson` package.

```
install_github("https://github.com/hsansford1/higgsboson")
library(higgsboson)
```

Access the data included in the package. First load the data used to train the models.

```
train <- higgsboson::training
```

Put the data in the format required for training:

```
df_train <- train[,2:33] #remove eventid
df_train <- df_train[,-31] #remove weights

df_train$Label <- ifelse(df_train$Label=="s",1,0) #encode "s" and "b" to 1 - 0 (resp.)
df_train$Label <- as.factor(df_train$Label) #need this as factor for caret package
```

Use the `reweight` function to normalise the weights (use `??reweight` to see function help). `Ns()` and `Nb()` are hardcoded values of  $N_s$  and  $N_b$  for the higgsboson dataset (see ‘Problem Formulation’ section of report for explanation of these values).

```
weights <- reweight(train$Weight, df_train$Label, Ns(), Nb())
```

Set all missing values equal to zero.

```
df_train[df_train==999] <- 0
```

To get a standardised data set `st_train`, run the following.

```
st_train <- df_train
st_train <- as.data.frame(scale(df_train[,1:30]))
st_train["Label"] <- df_train$Label
```

## Weighted Logistic Regression with custom AMS metric

The `mlr` package allows the use of custom metrics. First, define the classification task and make the logistic learner

```
trainTask <- makeClassifTask(data = st_train,
                             target = "Label",
                             positive = 1,
                             weights = weights
                             )

logistic.learner <- makeLearner("classif.logreg", predict.type = "response")
```

Use the custom AMS measure created in the `higgsboson` package `AMS_measure()` to conduct cross-validation on fitting the weighted logistic regression (WLR). This gives us an idea of the AMS we can expect to get from a test set.

```
set.seed(22)
AMS_mlr <- AMS_measure()
cv.logistic <- crossval(learner = logistic.learner, task = trainTask, iters = 5,
                       stratify = TRUE,
                       measures = AMS_mlr,
                       show.info = F, models=TRUE)

cv.logistic$aggr
```

```
## AMS_weighted.test.mean
##                0.2326592
```

```
cv.logistic$measures.test
```

```
##   iter AMS_weighted
## 1    1    0.4093149
## 2    2    0.1659211
## 3    3    0.2190233
## 4    4    0.2178583
## 5    5    0.1511781
```

Now, we fit a WLR model using the whole training data set and check how it performs on the test data.

```
train_control <- trainControl(method = "cv", number = 10)

fmodel <- caret::train(Label ~ .,
                       data = st_train,
                       trControl = train_control,
                       method = "glm",
                       weights = weights,
                       family=binomial()
                       )
```

```

test <- higgsboson::test

df_test <- test[,2:33] #remove eventid
df_test <- df_test[,-31] #remove weights
df_test$Label <- ifelse(df_test$Label=="s",1,0) #encode "s" and "b" to 1 - 0 (resp.) for logistic regression
df_test$Label <- as.factor(df_test$Label) #need this as factor for caret
# set missing values to 0 and standardise data
df_test[df_test==999] <- 0
st_test <- as.data.frame(scale(df_test[,1:30]))
st_test$Label <- df_test$Label

weights_test <- reweight(test$Weight, st_test$Label, Ns(), Nb()) # extract weights

truth <- st_test$Label
response <- predict(fmodel, newdata = st_test[, -length(st_test)])

AMS_weighted(truth, response, weights_test)

## [1] 0.1818793

```

## Two-stage maximisation

Below is the code for conducting the two-stage maximisation of the AMS.

First we create a training/validation split of the training data that preserves the overall class distribution.

```

trainIndex <- createDataPartition(st_train$Label, p = .8, list = FALSE, times = 1)

Train <- st_train[ trainIndex,]
Valid <- st_train[-trainIndex,]

```

We then fit a WLR model using the new training set `Train`, before visualising how the AMS varies with threshold `theta`. This informs the users decision of how to set the parameters `theta_0` and `theta_1` in the function `threshold_CV` in the `higgsboson` package (use `??threshold_CV` for help on this function). We can see that the peak definitely occurs in the range 0 to 0.1.

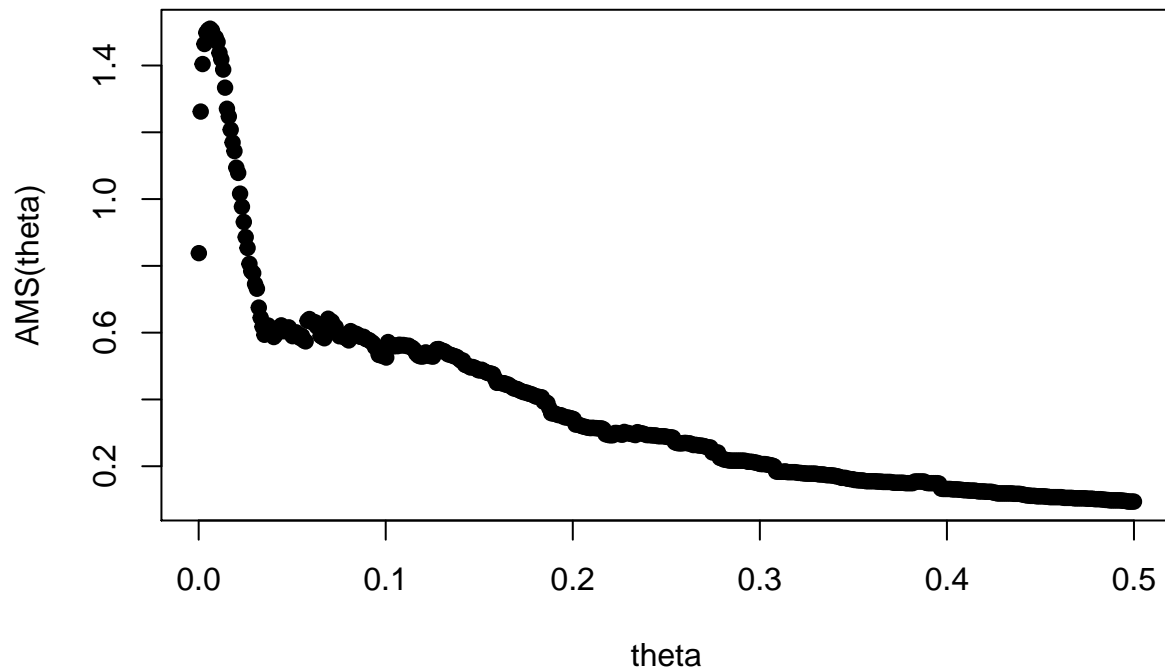
```

weights_Train <- reweight(weights[trainIndex], Train$Label, Ns(), Nb())
weights_Valid <- reweight(weights[-trainIndex], Valid$Label, Ns(), Nb())

logreg_weighted2 <- caret::train(Label ~ .,
                                data = Train,
                                method = "glm",
                                weights = weights_Train,
                                family=binomial()
)

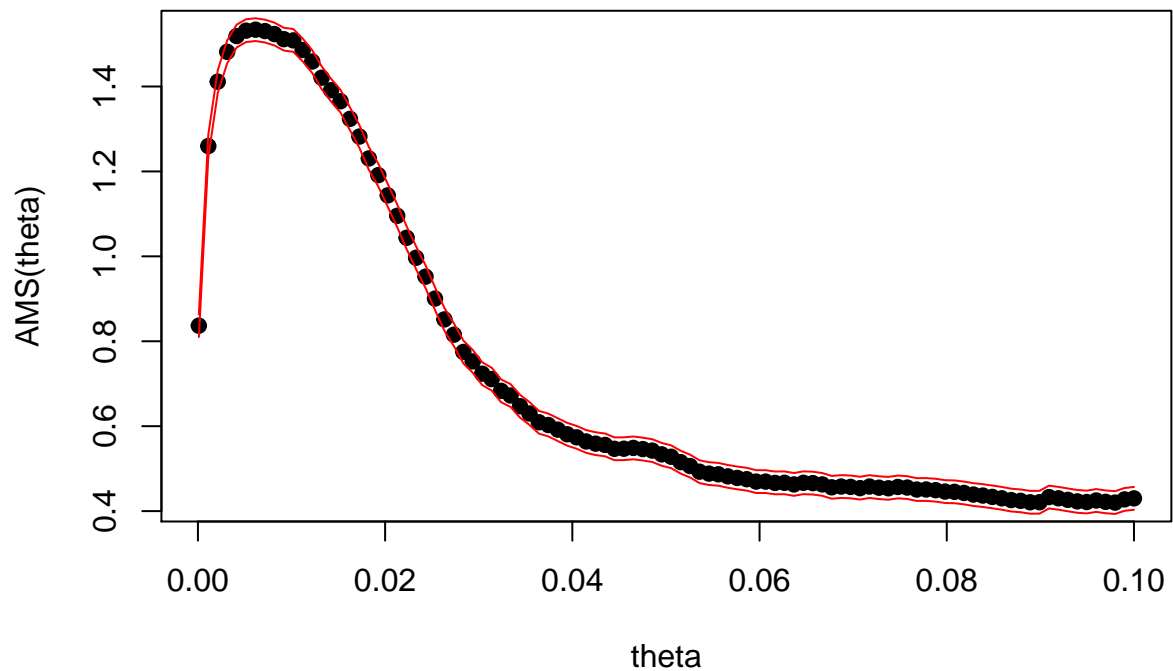
#Plot AMS for different values of threshold theta
theta_vals <- as.data.frame(seq(0.0001, 0.5, length.out=500)) # generate small sample thresholds theta
AMS_vals <- apply(theta_vals, 1, AMS(logreg_weighted2, Valid[,1:30], Valid[31], weights_Valid)) #compute AMS
plot(as.array(unlist(theta_vals)), AMS_vals, xlab="theta", ylab="AMS(theta)", pch=19) #plot it

```



Now, using the plot above, we can set `theta_0 = 0.0001` and `theta_1=0.1` in the `threshold_CV` function.

```
theta_CV <- threshold_CV(st_train, st_train$Label, weights, theta_0=0.0001, theta_1=0.1, n=100)
```



```
theta_CV
```

```
## $max_theta
## [1] 0.006154545
##
## $max_AMS
## [1] 1.533845
##
```

```
## $AMS_sd
## [1] 0.0009010219 0.0081645519 0.0063437419 0.0091000053 0.0087860769
## [6] 0.0153004775 0.0191553973 0.0292494018 0.0243173612 0.0200236573
## [11] 0.0185361417 0.0222922858 0.0229103826 0.0123141979 0.0172119711
## [16] 0.0230585444 0.0128590198 0.0164859944 0.0184639406 0.0175362141
## [21] 0.0297883054 0.0206818117 0.0218645263 0.0256676748 0.0258060322
## [26] 0.0413358862 0.0371291751 0.0437032356 0.0402525605 0.0342755910
## [31] 0.0252147162 0.0324946217 0.0389079786 0.0312523558 0.0285351710
## [36] 0.0264681853 0.0324136982 0.0295005308 0.0308692859 0.0276963752
## [41] 0.0244938514 0.0263578155 0.0272951238 0.0330144774 0.0356349456
## [46] 0.0302733815 0.0304448219 0.0309394441 0.0310298614 0.0316408636
## [51] 0.0290351472 0.0279686988 0.0296152078 0.0294406889 0.0322525304
## [56] 0.0378886477 0.0336396061 0.0310738896 0.0316473786 0.0338247218
## [61] 0.0288254517 0.0327440199 0.0279686435 0.0279303778 0.0219023953
## [66] 0.0279367954 0.0252695669 0.0238164194 0.0251349573 0.0330309283
## [71] 0.0300993661 0.0245543368 0.0210921148 0.0232300219 0.0275479964
## [76] 0.0275489099 0.0262726356 0.0296809337 0.0290542283 0.0304715761
## [81] 0.0378346326 0.0360171603 0.0365626233 0.0392169890 0.0363205482
## [86] 0.0342974294 0.0271761943 0.0264979206 0.0250127671 0.0281748662
## [91] 0.0228397551 0.0241149352 0.0265903816 0.0278718288 0.0240810086
## [96] 0.0227182821 0.0232010237 0.0230121965 0.0285160068 0.0241127326
##
## $max_thetas
## [1] 0.007163636 0.007163636 0.005145455 0.006154545 0.005145455
```

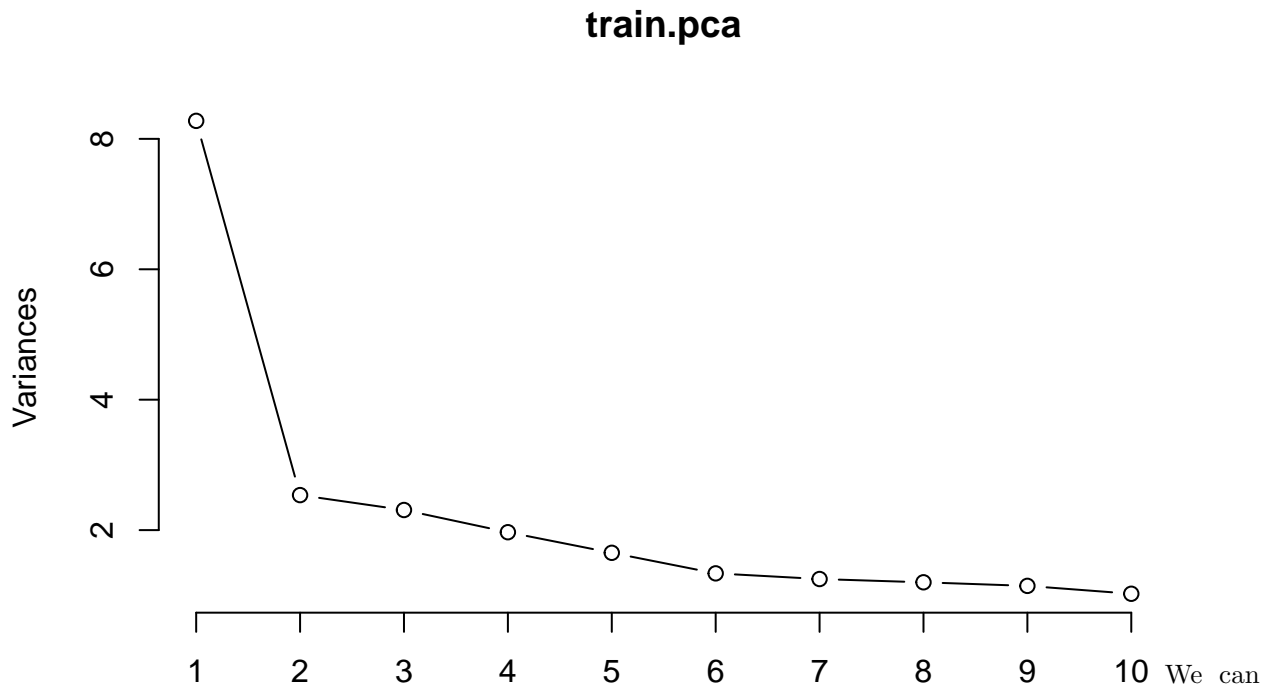
## Principal-Component Analysis

Alternatively, we can perform PCA on the data set before fitting a model and tuning the threshold. This will reduce the dimension of the data set, improving the speed of the function and, hopefully, reduce any overfitting of the more complicated model.

```
train.pca <- prcomp(st_train[,1:30])
summary(train.pca)
```

```
## Importance of components:
##
##          PC1      PC2      PC3      PC4      PC5      PC6      PC7
## Standard deviation  2.8768 1.59275 1.51937 1.40279 1.28531 1.15624 1.11806
## Proportion of Variance 0.2759 0.08456 0.07695 0.06559 0.05507 0.04456 0.04167
## Cumulative Proportion 0.2759 0.36042 0.43737 0.50296 0.55803 0.60259 0.64426
##
##          PC8      PC9      PC10     PC11     PC12     PC13     PC14
## Standard deviation  1.09531 1.07001 1.01201 1.00306 0.93818 0.87197 0.86820
## Proportion of Variance 0.03999 0.03816 0.03414 0.03354 0.02934 0.02534 0.02513
## Cumulative Proportion 0.68425 0.72242 0.75656 0.79009 0.81943 0.84478 0.86990
##
##          PC15     PC16     PC17     PC18     PC19     PC20     PC21
## Standard deviation  0.80557 0.74824 0.71231 0.66525 0.60098 0.5198 0.48412
## Proportion of Variance 0.02163 0.01866 0.01691 0.01475 0.01204 0.0090 0.00781
## Cumulative Proportion 0.89153 0.91020 0.92711 0.94186 0.95390 0.9629 0.97072
##
##          PC22     PC23     PC24     PC25     PC26     PC27     PC28
## Standard deviation  0.46176 0.39368 0.37755 0.35851 0.30716 0.25633 0.24541
## Proportion of Variance 0.00711 0.00517 0.00475 0.00428 0.00314 0.00219 0.00201
## Cumulative Proportion 0.97782 0.98299 0.98774 0.99203 0.99517 0.99736 0.99937
##
##          PC29     PC30
## Standard deviation  0.13754 3.534e-06
## Proportion of Variance 0.00063 0.000e+00
## Cumulative Proportion 1.00000 1.000e+00
```

```
screeplot(train.pca, type="lines")
```



We can see that the first two PCA's account for the majority of the variance, but we shall keep the first three as a precaution.

```
st_train_dimred <- data.frame("PCA1"=train.pca$x[,1], "PCA2"=train.pca$x[,2], "PCA3"=train.pca$x[,3])
st_train_dimred["Label"] <- st_train$Label
```

Now, we can perform the first stage of the two-stage procedure as above.

```
trainIndex <- createDataPartition(st_train_dimred$Label, p = .8,
                                  list = FALSE,
                                  times = 1)

Train <- st_train_dimred[ trainIndex,]
Valid  <- st_train_dimred[-trainIndex,]

weights_Train <- reweight(weights[trainIndex], Train$Label, Ns(), Nb())
weights_Valid <- reweight(weights[-trainIndex], Valid$Label, Ns(), Nb())

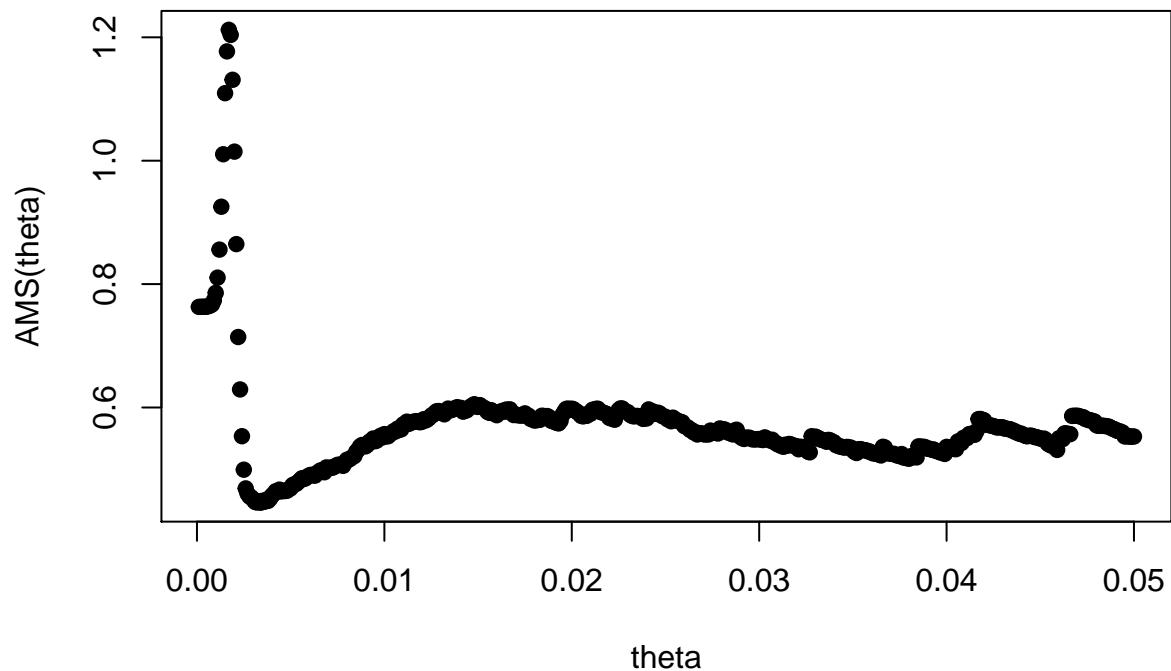
logreg_weighted_pca <- caret::train(Label ~ .,
                                     data = Train,
                                     method = "glm",
                                     weights = weights_Train,
                                     family=binomial()
)
print(logreg_weighted_pca)

## Generalized Linear Model
##
## 200001 samples
##      3 predictor
##      2 classes: '0', '1'
```

```
##
## No pre-processing
## Resampling: Bootstrapped (25 reps)
## Summary of sample sizes: 200001, 200001, 200001, 200001, 200001, 200001, ...
## Resampling results:
##
##   Accuracy   Kappa
## 0.6580688 0.0006967816
```

Now we can plot the AMS for varying values of the threshold  $\theta$ .

```
theta_vals <- as.data.frame(seq(0.0001, 0.05, length.out=500)) # generate small sample thresholds theta
AMS_vals <- apply(theta_vals, 1, AMS(logreg_weighted_pca, Valid[,1:3], Valid[4], weights_Valid)) #compute
plot(as.array(unlist(theta_vals)), AMS_vals, xlab="theta", ylab="AMS(theta)", pch=19) #plot it
```



```
max_theta <- theta_vals[which.max(AMS_vals),1]
max_theta
```

```
## [1] 0.0017
```

```
max_AMS <- AMS_vals[which.max(AMS_vals)]
max_AMS
```

```
## [1] 1.212122
```

We can see that the peak occurs between very small values of  $\theta$ , hence we try  $\theta_0=0.0001$  and  $\theta_1=0.006$  in the `threshold_CV` function.

```
theta_CV <- threshold_CV(st_train_dimred[,1:3], st_train_dimred$Label, weights=weights, theta_0=0.0001,
```

```
## Warning in eval(family$initialize): non-integer #successes in a binomial glm!
```

```
## Warning in eval(family$initialize): non-integer #successes in a binomial glm!
```

```
## Warning in eval(family$initialize): non-integer #successes in a binomial glm!
```

[illegible]

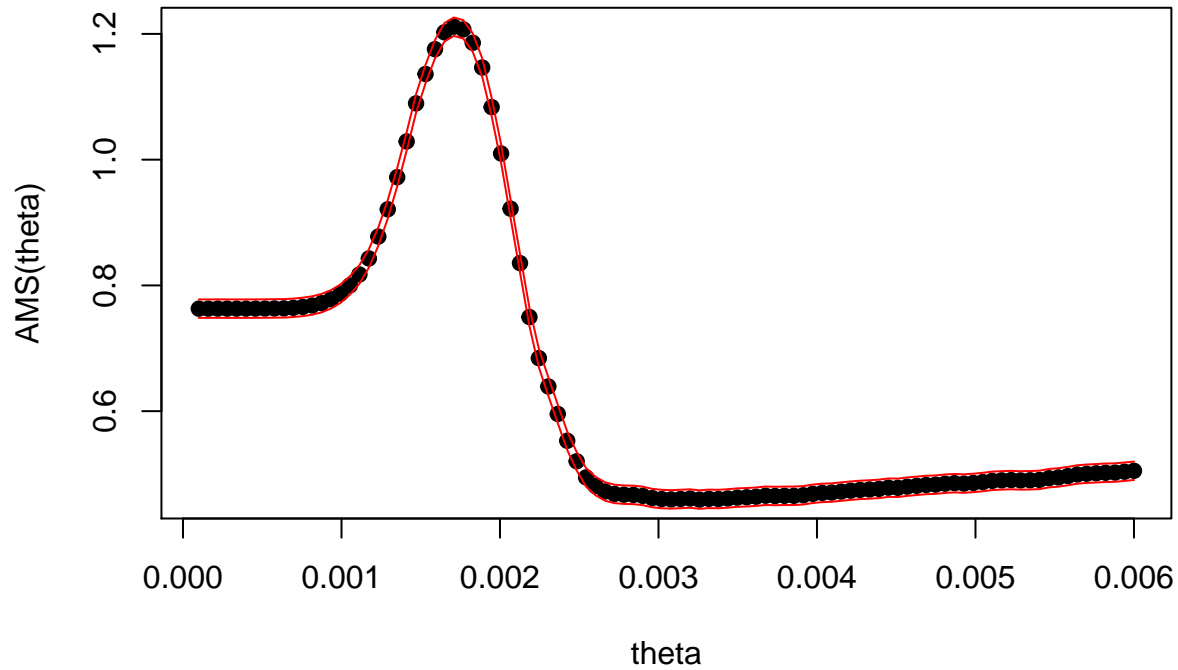




[illegible]

[illegible]

[illegible]



```
theta_CV  #less variance, both of theta and of AMS
```

```
## $max_theta
## [1] 0.001709091
##
## $max_AMS
## [1] 1.211464
##
## $AMS_sd
## [1] 3.113443e-05 6.098031e-05 6.840868e-05 1.170284e-04 1.749527e-04
## [6] 2.209178e-04 3.883600e-04 4.564580e-04 5.170606e-04 3.776670e-04
## [11] 5.213827e-04 3.488913e-04 3.110193e-04 3.193567e-04 6.264979e-04
## [16] 5.104123e-04 1.029586e-03 1.038027e-03 2.333071e-03 3.326344e-03
## [21] 4.266756e-03 4.851101e-03 3.822262e-03 5.139048e-03 5.138905e-03
## [26] 6.212007e-03 6.120419e-03 6.890763e-03 8.564148e-03 1.375585e-02
## [31] 1.648589e-02 2.528649e-02 3.215153e-02 3.508819e-02 3.125052e-02
## [36] 2.741701e-02 1.802814e-02 1.752584e-02 1.804633e-02 2.051882e-02
## [41] 2.108876e-02 1.987291e-02 1.806499e-02 1.665281e-02 1.749482e-02
## [46] 1.882053e-02 1.967156e-02 1.870699e-02 1.928601e-02 2.099750e-02
## [51] 2.160873e-02 2.150023e-02 2.142512e-02 2.032449e-02 2.101275e-02
## [56] 2.226225e-02 2.246597e-02 2.163525e-02 2.010057e-02 2.062877e-02
## [61] 2.036974e-02 2.019866e-02 2.090797e-02 2.077318e-02 2.117374e-02
## [66] 2.104925e-02 2.077360e-02 2.024101e-02 2.038898e-02 1.998980e-02
## [71] 2.091814e-02 2.154387e-02 2.159600e-02 2.152450e-02 2.228930e-02
## [76] 2.107314e-02 2.062601e-02 2.058820e-02 1.956955e-02 1.975093e-02
## [81] 1.990464e-02 1.867180e-02 2.008927e-02 2.111247e-02 1.938265e-02
## [86] 1.896941e-02 1.770108e-02 1.669148e-02 1.616101e-02 1.496573e-02
## [91] 1.461834e-02 1.512790e-02 1.473607e-02 1.449848e-02 1.455825e-02
## [96] 1.445515e-02 1.463020e-02 1.374000e-02 1.363554e-02 1.418338e-02
##
## $max_thetas
## [1] 0.001709091 0.001709091 0.001709091 0.001709091 0.001709091
```