

Call Center Data

Syed Rahman

04/24/2019

Introduction

For this study, the data was obtained from <http://iew3.technion.ac.il/serveng/callcenterdata>, which contains daily data for a call center for an Israeli bank for the year 1999. Four days are missing. But for all the other 361 days we have information for time of call, length of call, etc.

```
months = format(ISOdate(1999, 1:12, 1), "%B")
months = tolower(months)
data = read.table("january.txt", header = TRUE)
for (month in months[-1]) {
  data = rbind(data, read.table(paste(month, ".txt", sep = ""),
    header = TRUE))
}
head(data, n = 6)
```

```
##   vru.line call_id customer_id priority type   date vru_entry vru_exit
## 1  AA0101   33116     9664491         2   PS 990101   0:00:31  0:00:36
## 2  AA0101   33117         0         0   PS 990101   0:34:12  0:34:23
## 3  AA0101   33118     27997683        2   PS 990101   6:55:20  6:55:26
## 4  AA0101   33119         0         0   PS 990101   7:41:16  7:41:26
## 5  AA0101   33120         0         0   PS 990101   8:03:14  8:03:24
## 6  AA0101   33121         0         0   PS 990101   8:18:42  8:18:51
##   vru_time q_start  q_exit q_time outcome ser_start ser_exit ser_time
## 1         5 0:00:36 0:03:09    153   HANG   0:00:00   0:00:00         0
## 2        11 0:00:00 0:00:00         0   HANG   0:00:00   0:00:00         0
## 3         6 6:55:26 6:55:43     17  AGENT   6:55:43   6:56:37         54
## 4        10 0:00:00 0:00:00         0  AGENT   7:41:25   7:44:53        208
## 5        10 0:00:00 0:00:00         0  AGENT   8:03:23   8:05:10        107
## 6         9 0:00:00 0:00:00         0  AGENT   8:18:50   8:23:25        275
##       server
## 1 NO_SERVER
## 2 NO_SERVER
## 3   MICHAL
## 4   BASCH
## 5   MICHAL
## 6   KAZAV
```

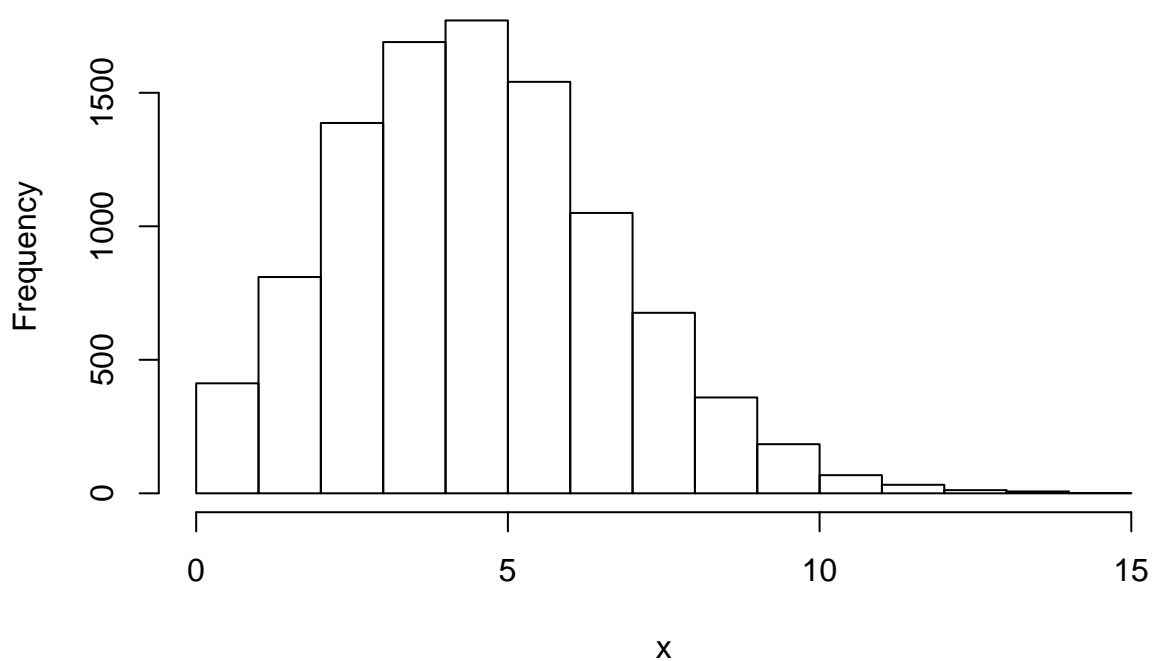
There were 444448 observations (calls) and 17 columns (features). As the call center was only open from 7AM to midnight, I removed all the data points outside of this interval.

The goal here is to predict the number of calls coming in during the second half of the day using the number of calls recieved during the first half of the day. For this purpose, a lot of the data needs to be pre-processed. The parameters of interest are the mean, which is 102 dimensional vector, and the covariance matrix which has 5151 unique elements (I have worked with much larger datasets in general. Some of the datasets I have used at work have upto a billion rows and hundreds of columns).

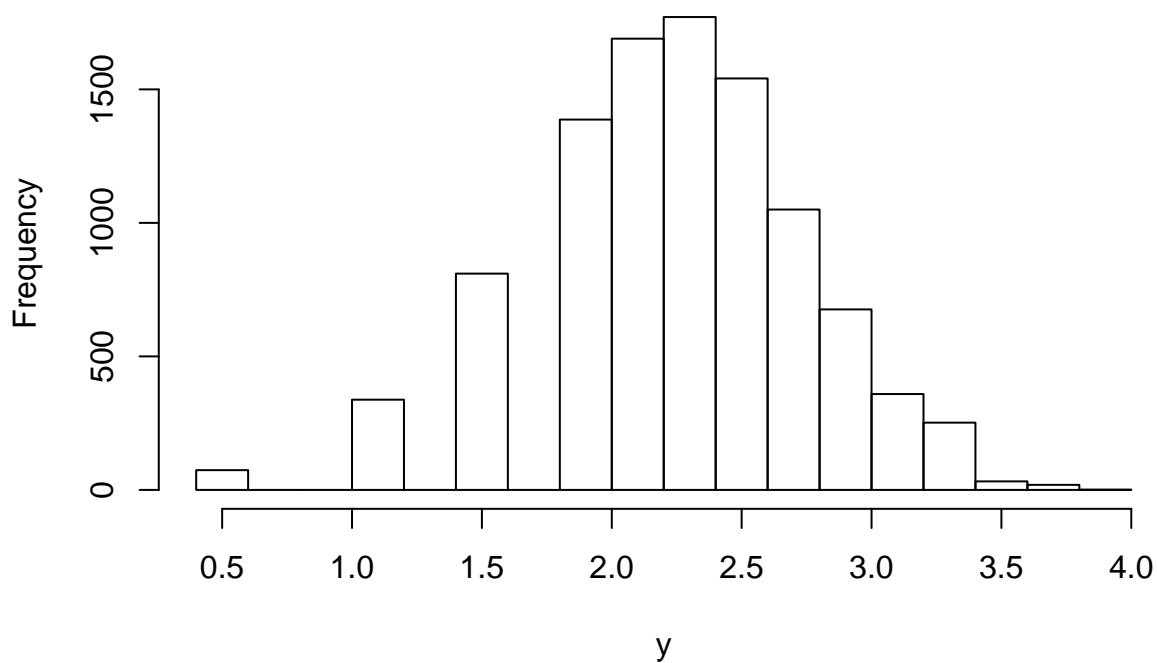
Model

Getting back to the problem at hand, suppose x_1 and x_2 is the number of calls received during the first half of the day and the second half of the day, respectively. Let $x = (x_1^t, x_2^t)^t$ and $y = \sqrt{x + \frac{1}{4}}$. It can be shown that if $x \sim \text{Poisson}(\lambda)$, then $y \sim \mathcal{N}(\mu, \Sigma)$.

Poisson data



Poisson data with square root transformation



Suppose we partition the covariance matrix of y as follows:

$$\Sigma = \begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{pmatrix}$$

In such a case, $y_2|y_1 \sim \mathcal{N}(\mu_2 + \Sigma_{21}\Sigma_{11}^{-1}(y_1 - \mu_1), \Sigma_{22} - \Sigma_{21}\Sigma_{11}^{-1}\Sigma_{12})$. Hence the best mean squared predictor of y_2 given y_1 (i.e. the predictor that minimizes mean squared error) is

$$\mathbb{E}[y_2|y_1] = \mu_2 + \Sigma_{21}\Sigma_{11}^{-1}(y_1 - \mu_1).$$

This is done using the follow code:

```
predict.mean <- function(x1, mu, Sigma) {
  p1 <- length(x1)
  p <- length(mu)
  p2 <- p - p1
  mu1 <- mu[1:p1]
  mu2 <- mu[(p1 + 1):p]
  Sigma11 <- Sigma[1:p1, 1:p1]
  # Sigma12 <- Sigma[1:p1, (p1+1):p]
  Sigma21 <- Sigma[(p1 + 1):p, 1:p1]
  # Sigma22 <- Sigma[(p1+1):p, (p1+1):p]
  x2 <- mu2 + Sigma21 %*% solve(Sigma11, x1 - mu1)
  return(x2)
}
```

Hence we must estimate μ_2 and Σ from the data. To this end, we divide the dataset into a training dataset and a test dataset. Typically we do this in a random fashion, but since this is time-series data, I take the first 70% of the data to be the training dataset and the remaining 30% of the dataset to be the test dataset. Once we have the estimates, we make predictions for the calls coming in for the second half of the day using the calls that came in for the first half of the day for the test data set and then compare with true values.

I divide the day up into 10 minute chunks and the first goal is to reformat the data so that there is an easy way to view the number of calls received at the call center on each day for each time period.

```
data = filter(data, outcome == "AGENT") ### keeping only agent outcome, rest are hangups or transforms
data = select(data, date, vru_entry) ### only interested in date and vru entry. discarding all other
data$timedate = paste(data$date, data$vru_entry)
data$timedate = strptime(data$timedate, format = "%y%m%d %H:%M:%S") ### converting to R date format
tt = seq(from = ISOdate(1999, 1, 1, 0, 0, 0, tz = "EST"), to = ISOdate(1999,
  12, 31, 0, 0, 0, tz = "EST"), by = "10 min") ### creating 10 minute chunks
data$timeperiods = cut(data$timedate, breaks = tt) ### dividing up calls accordin to 10 minute chunks

data$periods <- sapply(strsplit(as.character(data$timeperiods),
  " "), "[", 2) ### keeping only the time period. date discarded

data$timeperiods <- as.character(data$timeperiods)
data$count <- as.numeric(ave(data$timeperiods, data$timeperiods,
  FUN = length)) ### creating a column with counts of how many calls were received in a certain time

data2 = select(data, date, periods, count) ### keeping only relevant columns such as date, periods and

data_wide <- reshape(data2, timevar = "periods", idvar = c("date"),
  direction = "wide")

data_wide[is.na(data_wide)] = 0
```

```

### reshaping long format to a wide format and setting all NA
### values to 0

#### reordering column by time
data_wide = data_wide[, order(names(data_wide))]
data_wide$date2 = strptime(data_wide$date, format = "%y%m%d")
data_wide$days = weekdays(data_wide$date2) ### adding column for days of the week
data_wide$holidays = 1 ### adding a holidays column. 1 means not a holiday. 2 means a holiday. mostly :

# downloading list of holidays for 1999 for Israel
holidays = readLines("http://world.std.com/~reinhold/jewishholidays.txt")
holidays = holidays[grepl("1999", holidays)]
holidays = unlist(strsplit(holidays, "\t"))[seq(1, 28, 2)]
holidays = strptime(holidays, "%a. %e %B %Y ")

for (i in 1:length(data_wide$date2)) {
  for (count in 1:length(holidays)) {
    if (sum(data_wide$date2[i] == holidays[count]) > 0) {
      data_wide$holidays[i] = 2
    }
  }
}

data_wide$Colour = ifelse(data_wide$days == "Friday", 1, ifelse(data_wide$days ==
  "Monday", 2, ifelse(data_wide$days == "Saturday", 3, ifelse(data_wide$days ==
  "Sunday", 4, ifelse(data_wide$days == "Thursday", 5, ifelse(data_wide$days ==
  "Tuesday", 6, 7)))))) ### adding a colour column based on day of week

#### deleting all times before 7AM when the center opens and
#### date column and convert data to matrix format
X = as.matrix(data_wide[, -c(1:29, 132:137)])
Xsvd = svd(X)
data_wide$date2 <- as.POSIXct(data_wide$date2)
data_fri = filter(data_wide, days == "Friday")
dim(data_fri)

## [1] 52 137

data_sat = filter(data_wide, days == "Saturday")
data_rest = filter(data_wide, days != "Saturday" & days != "Friday")
X = as.matrix(data_wide[, -c(1:29, 132:137)])
X_fri = as.matrix(data_fri[, -c(1:29, 132:137)])
X_sat = as.matrix(data_sat[, -c(1:29, 132:137)])
X_rest = as.matrix(data_rest[, -c(1:29, 132:137)])
#### for poisson data apply transformation sqrt(x+1/4) to data
#### to make more normal
Y = sqrt(X + 1/4)
Y_fri = sqrt(X_fri + 1/4)
Y_sat = sqrt(X_sat + 1/4)
Y_rest = sqrt(X_rest + 1/4)

```

Scree Plots

The scree plot indicates that about 25% is explained by the 1st eigenvector, about 5% by the 2nd one and so on.

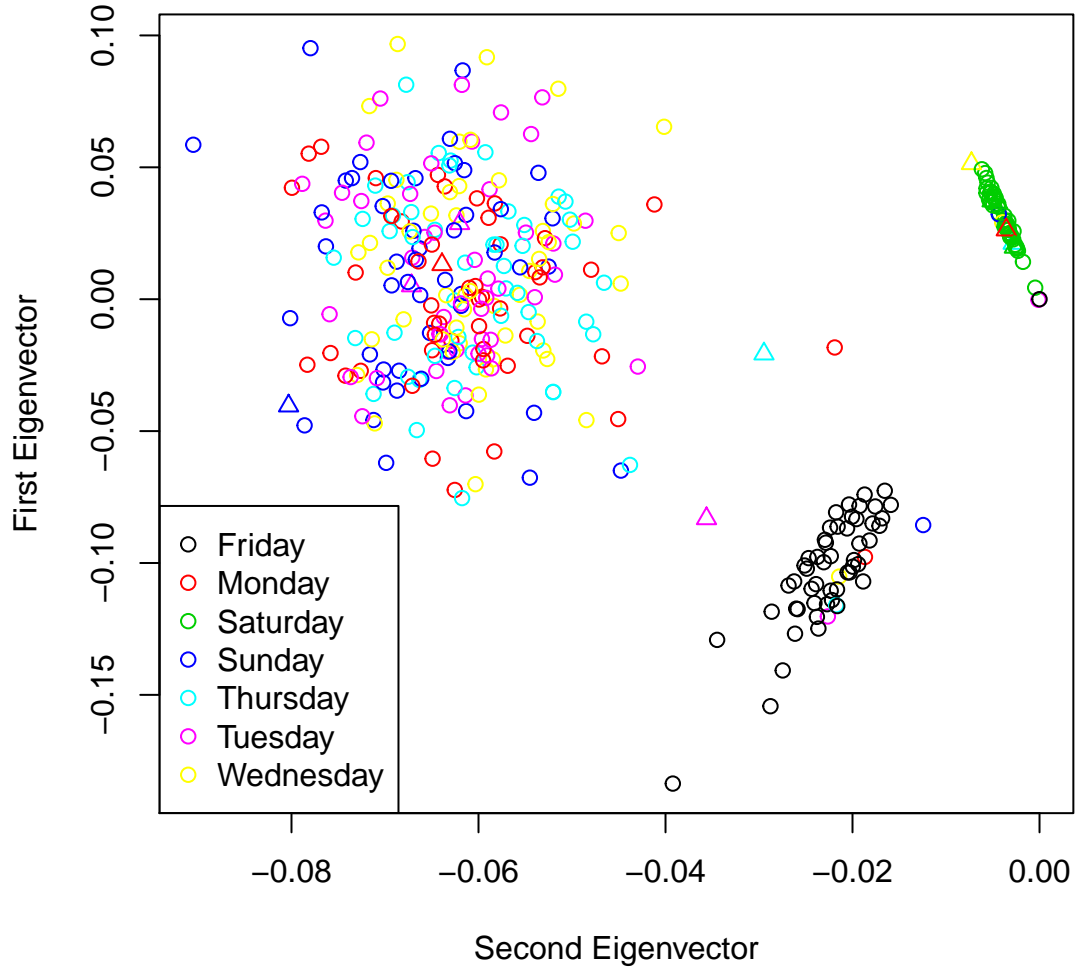
Scaled Eigenvalues for Call Data



Bi-Plots

The bi-plot indicates that the weekends (Fridays and Saturdays in Israel) are separate from the rest of the data. This indicates that for the prediction phase, it might be best to work on all of these separately. The triangles indicate ± 1 day from Holidays. Almost all the outlier seem to be either the weekends or Holidays.

Biplot for Call Data



Model Evaluation

To evaluate the models we use Absolute Error(AE) where

$$AE_t = \frac{1}{T} \sum_{i=T+1}^{364} |\hat{y}_{it} - y_{it}|$$

and T is the size of the training dataset.

Full dataset

Initially we do the evaluations on the full dataset (including weekends). We estimate μ_2 by using the sample mean and estimate Σ using the sample covariance method and CSCS. For CSCS, we need to pick a penalty parameter, which we do using cross-validation of the likelihood.

```
CSCS2 <- function(Y, lambda, L = NULL, maxitr = 100, tol = 1e-04,
  warmstart = FALSE) {
```

```

#### inputs Y: n by p matrix of data lambda: l1-penalty
#### parameter maxitr: maximum number of iterations allowed
#### (diagonal/off-diagonal optimization) tol: if maximum
#### difference of iterates is less than tol, consider converged
#### warmstart: warmstarting actually made the runs slower
#### (overhead may be too expensive)

#### outputs L: lower triangular matrix of cholesky factor
#### computed with CSCS algorithm itr_log: (p-1)-vector of
#### number of iterations eps_log: (p-1)-vector of number
#### maximum difference for considering convergence

n = nrow(Y)
p = ncol(Y)

if (is.null(L))
  L = diag(p)

S = (t(Y) %*% Y)/n

itr_log = eps_log = NULL

L[1, 1] = 1/sqrt(S[1, 1])
for (k in 2:p) {
  ## Loop in Algorithm 2

  ## nu_k vector (equation 2.5)
  nuk_old = nuk_new = c(rep(0, k - 1), 1)
  r = 0

  repeat {
    ## Loop in Algorithm 1

    r = r + 1

    km1_ = 1:(k - 1) ## 1, ..., k-1 is off-diagonal elements indices

    ## Update off-diagonal terms
    hk = lassoshooting(XtX = S[km1_, km1_, drop = FALSE],
      Xty = -nuk_old[k] * S[km1_, k], lambda = 0.5 *
        lambda)
    nuk_new[km1_] = hk$coefficients

    ## Update diagonal term
    sumterm = sum(nuk_new[km1_] * S[k, km1_])
    nuk_new[k] = (-sumterm + sqrt(sumterm^2 + 4 * S[k,
      k]))/(2 * S[k, k])

    ## Check convergence
    maxdiff = max(abs(nuk_new - nuk_old))
    if (maxdiff < tol || r >= maxitr) {
      L[k, 1:k] = nuk_new
    }
  }
}

```

```

        eps_log = c(eps_log, maxdiff)
        itr_log = c(itr_log, r)
        break
    } else {
        nuk_old = nuk_new
    }
}
}

list(L = L, itr = itr_log, eps = eps_log)
}

cv = function(lambda, K, train) {
  cvec = rep(0, K)
  set.seed(12345)
  index = sample(1:dim(train)[1], replace = FALSE)
  foldsize = ceiling(dim(train)[1]/K)
  inTrain <- list()
  for (i in 1:(K - 1)) {
    inTrain[[i]] <- index[((i - 1) * foldsize + 1):((i -
      1) * foldsize + foldsize)]
  }
  inTrain[[K]] <- index[((K - 1) * foldsize + 1):length(index)]
  for (k in 1:K) {
    Ytrain = train[-c(inTrain[[k]]), ]
    Ytest = train[c(inTrain[[k]]), ]
    E = Ytrain
    out = CSCS2(as.matrix(E), lambda)$L
    Sigma_v = t(out) %*% out
    s_v = dim(Ytest)[1]
    sumterm = sum(diag(Ytest %*% (Sigma_v) %*% t(Ytest)))
    cvec[k] = s_v * log(det(solve(Sigma_v))) + sumterm
  }
  cv = (1/K) * sum(cvec)
  return(cv)
}

training = 1:252
train = scale(Y[training, ], center = TRUE, scale = FALSE)
test = Y[-training, ]
mu = colMeans(train)
S = (t(train) %*% train)/(dim(train)[1])
SError <- matrix(0, nrow = nrow(test), ncol = 51)
for (k in 1:nrow(test)) {
  S.mu2 <- predict.mean(test[k, 1:51], mu = mu, Sigma = S)
  SError[k, ] <- (S.mu2 - as.numeric(test[k, 52:102]))
}
E.S <- colMeans(abs(SError))
Time = 52:102
AE = E.S
Method = rep("S", length(AE))
Sdata = data.frame(Method, AE, Time)

```

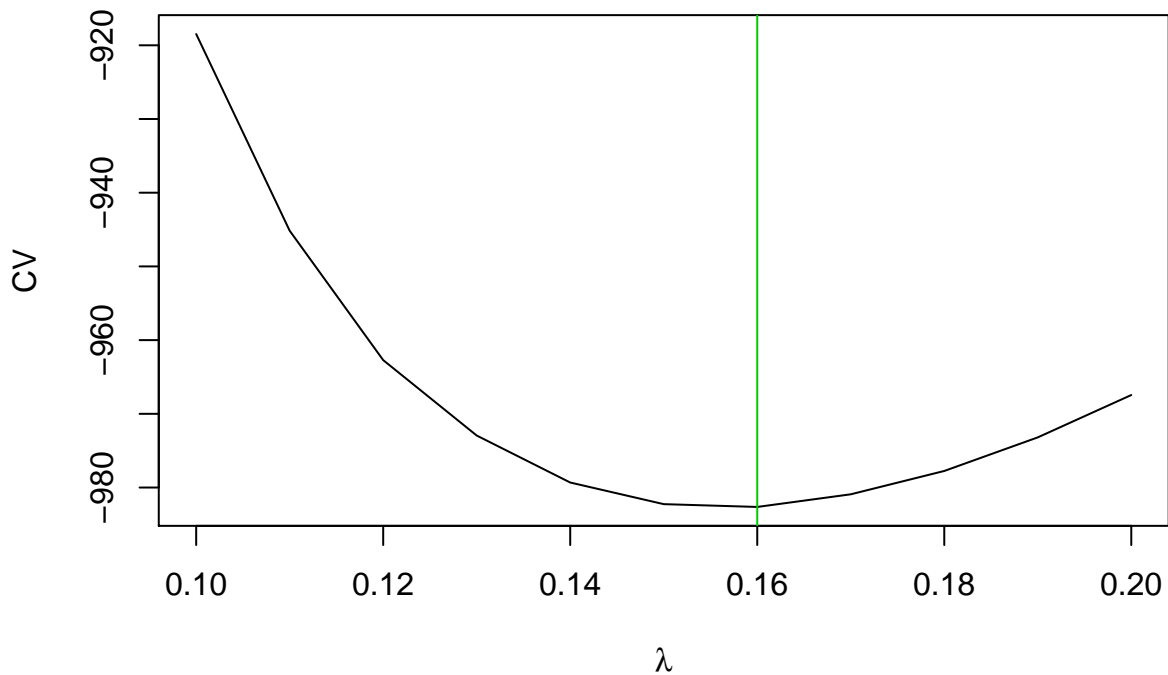


```

lambdal = 0.01
lambda = seq(0.1, 0.2, by = lambdal)
out = rep(0, length(lambda))
for (k in 1:length(lambda)) {
  out[k] = cv(lambda[k], 5, train)
}

```

CV for CSCS

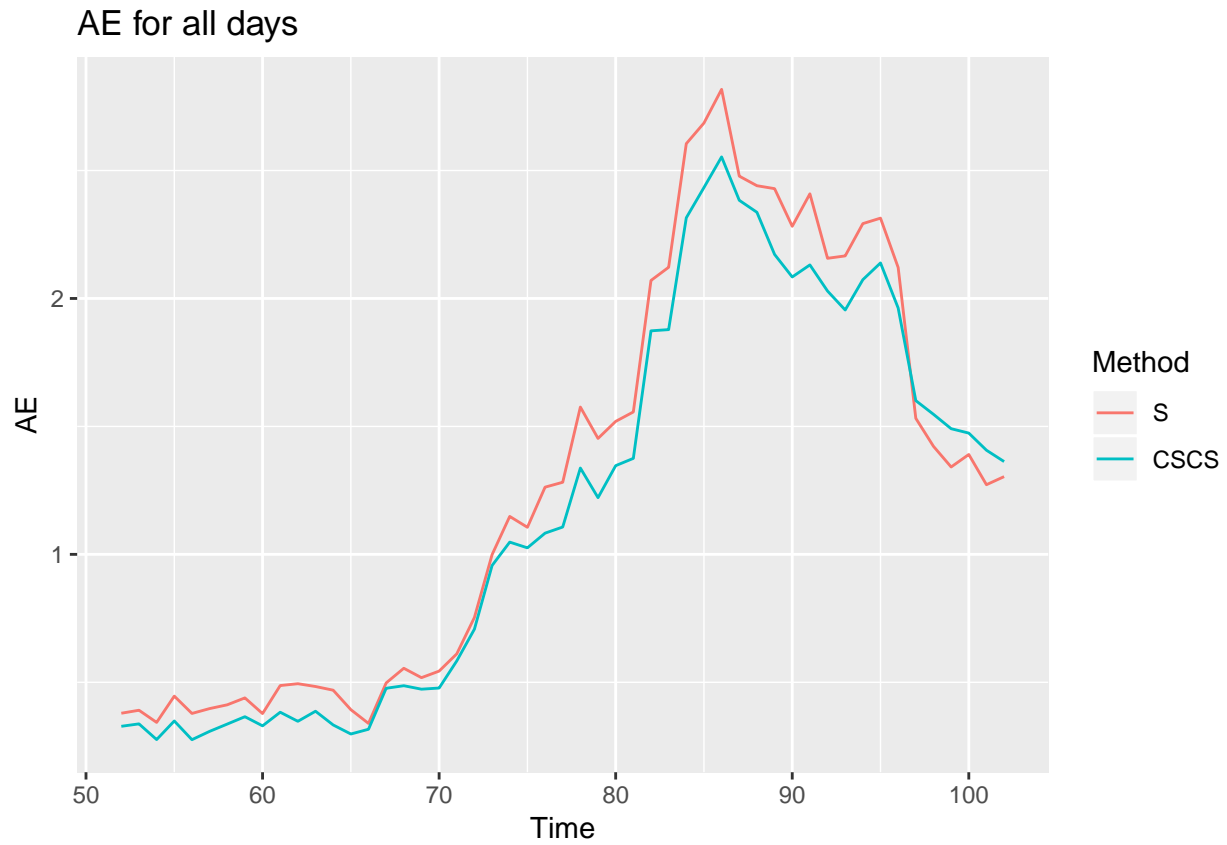


```

lambdastar = lambda[which.min(out)]
E = train
out = CSCS2(as.matrix(E), lambdastar)$L
CSCS = solve(t(out) %*% out)
CSCSError <- matrix(0, nrow = nrow(test), ncol = 51)
for (k in 1:nrow(test)) {
  CSCS.mu2 <- predict.mean(test[k, 1:51], mu = mu, Sigma = CSCS)
  CSCSError[k, ] <- (CSCS.mu2 - as.numeric(test[k, 52:102]))
}
E.CSCS <- colMeans(abs(CSCSError))
Time = 52:102
AE = E.CSCS
Method = rep("CSCS", length(AE))
CSCSdata = data.frame(Method, AE, Time)

```

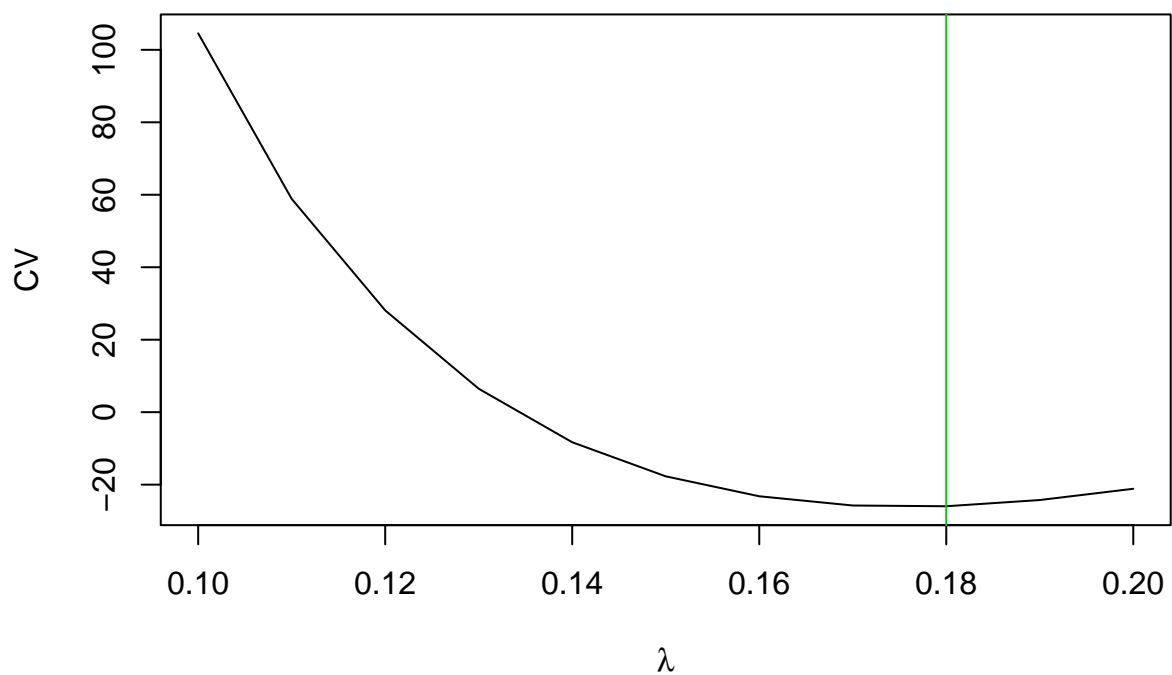
The following plot shows the comparison of using the sample covariance matrix and CSCS on the full dataset. It's quite clear that CSCS does better than the sample covariance matrix in most cases.



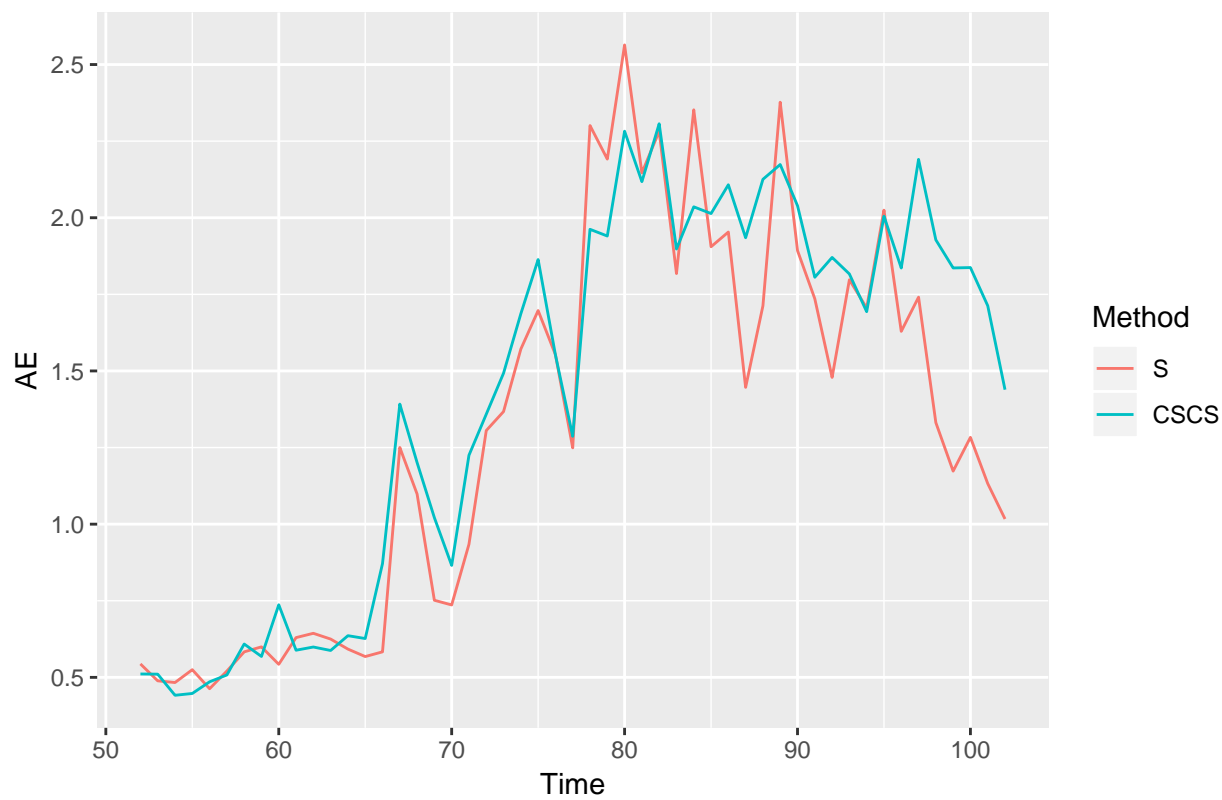
Weekdays

This is the comparison for just the weekdays. Here, there isn't a clear winner. This is probably our dataset is more homogenous. CSCS is a more robust covariance estimator than the sample covariance matrix. As we focussed on dataset with fewer outliers, the sample covariance matrix performed better comparatively.

CV for CSCS



AE for weekdays



Connections to graphical models

Finally, below is a graph showing how the time periods are related. Here, Node i corresponds to Time Period i .

