

# STAT 206 Homework 8

Due Thursday, December 6, 5:00 PM

**General instructions for homework:** Homework must be completed as an R Markdown file. Be sure to include your name in the file. Give the commands to answer each question in its own code block, which will also produce plots that will be automatically embedded in the output file. Each answer must be supported by written statements as well as any code used. (Examining your various objects in the “Environment” section of RStudio is insufficient – you must use scripted commands.)

## Part I - Metropolis-Hasting algorithm

Suppose  $f \sim \Gamma(2, 1)$ .

1. Write an independence MH sampler with  $g \sim \Gamma(2, \theta)$ .

```
ind.chain <- function(x, n, theta = 1) {  
  ## if theta = 1, then this is an iid sampler  
  m <- length(x)  
  x <- append(x, double(n))  
  for(i in (m+1):length(x)){  
    x.prime <- rgamma(1, shape=2, scale=theta)  
    u <- exp((x[(i-1)]-x.prime)*(1-theta))  
    if(runif(1) < u)  
      x[i] <- x.prime  
    else  
      x[i] <- x[(i-1)]  
  }  
  return(x)  
}
```

2. What is  $R(x_t, X^*)$  for this sampler?

$$\exp[(x_t - x^*)(1 - \theta)]$$

3. Generate 10000 draws from  $f$  with  $\theta \in \{1/2, 1, 2\}$ .

```
trial1 <- ind.chain(1, 10000, 1/2)  
trial2 <- ind.chain(1, 10000, 1)  
trial3 <- ind.chain(1, 10000, 2)
```

4. Write a random walk MH sampler with  $h \sim N(0, \sigma^2)$ .

```
rw.chain <- function(x, n, sigma = 1) {  
  m <- length(x)  
  x <- append(x, double(n))  
  for(i in (m+1):length(x)){  
    x.prime <- x[(i-1)] + rnorm(1, sd = sigma)  
    u <- exp((x[(i-1)]-x.prime)*(x.prime/x[(i-1)]))  
    if(runif(1) < u && x.prime > 0)  
      x[i] <- x.prime  
    else  
      x[i] <- x[(i-1)]  
  }  
}
```

```
return(x)
}
```

5. What is  $R(x_t, X^*)$  for this sampler?

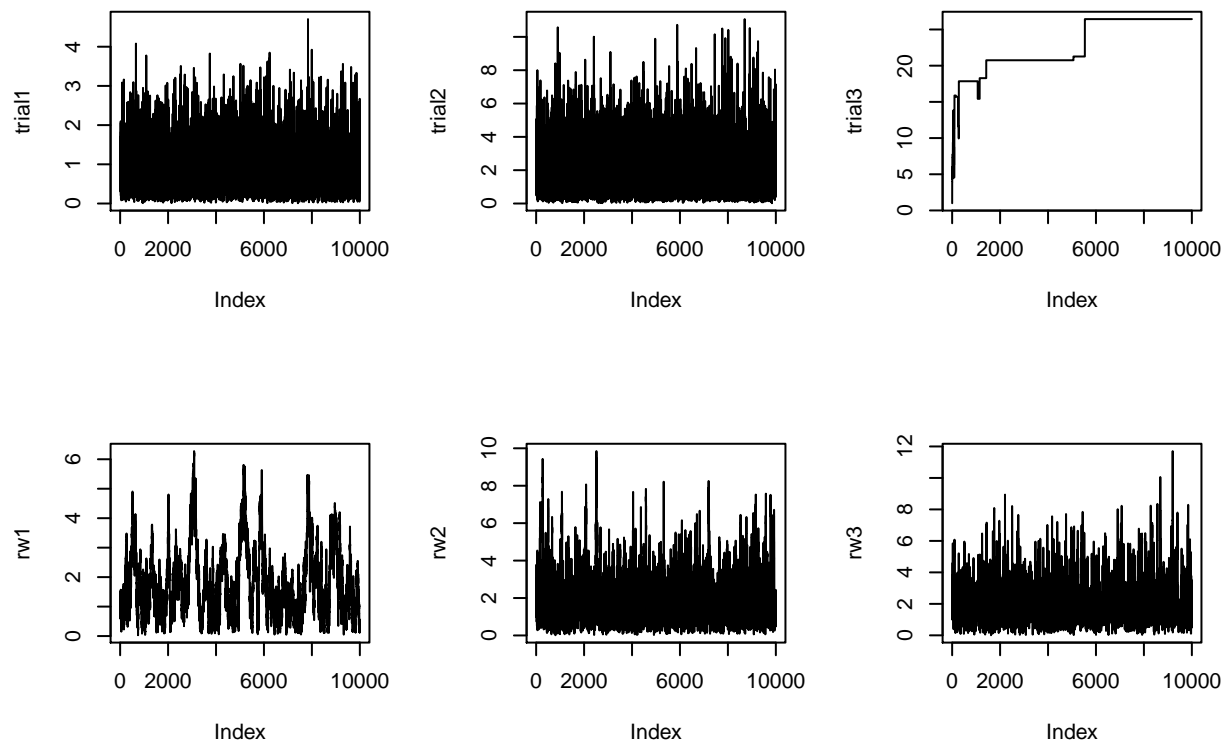
$$\exp[(x_t - x^*)^2] * (x^2 / x_t)$$

6. Generate 10000 draws from  $f$  with  $\sigma \in \{.2, 1, 5\}$ .

```
rw1 <- rw.chain(1, 10000, .2)
rw2 <- rw.chain(1, 10000, 1)
rw3 <- rw.chain(1, 10000, 5)
```

7. In general, do you prefer an independence chain or a random walk MH sampler? Why?

```
par(mfrow=c(2,3))
plot(trial1, type="l")
plot(trial2, type="l")
plot(trial3, type="l")
plot(rw1, type="l")
plot(rw2, type="l")
plot(rw3, type="l")
```



According to the plot above, a random walk MH sampler is better.

8. Implement the fixed-width stopping rule for you preferred chain.

```
library(mcmcse)
```

```
## mcmcse: Monte Carlo Standard Errors for MCMC
## Version 1.3-2 created on 2017-07-03.
## copyright (c) 2012, James M. Flegal, University of California, Riverside
## John Hughes, University of Colorado, Denver
```

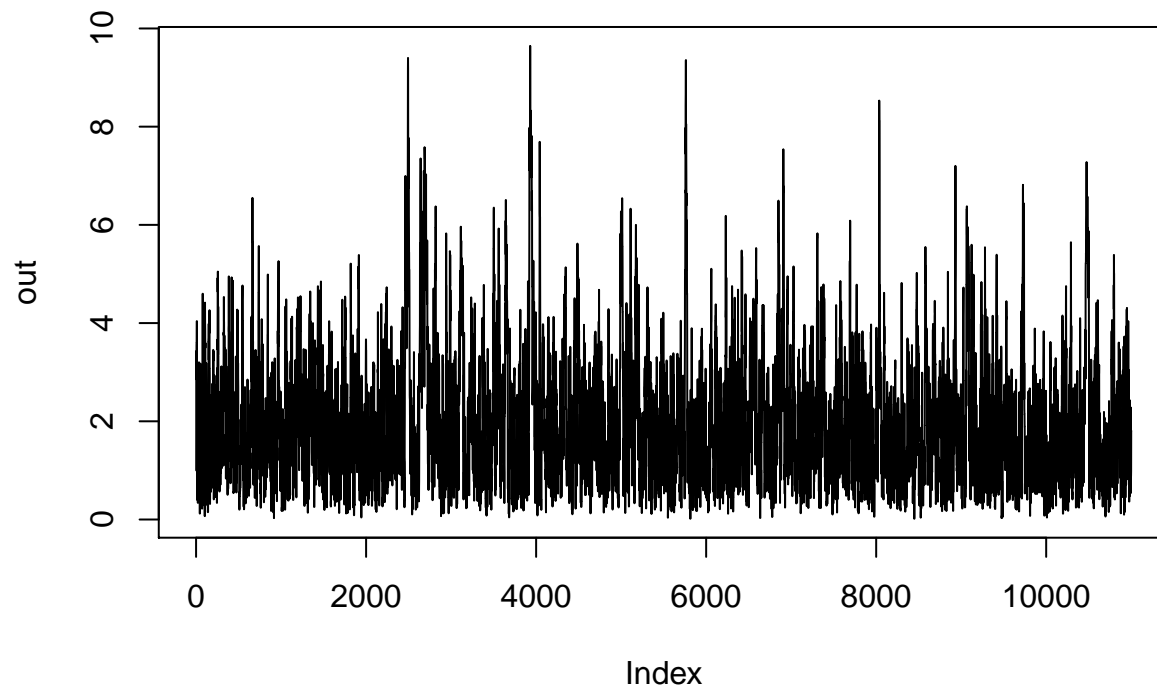
```
##                               Dootika Vats, University of Warwick
##                               Ning Dai, University of Minnesota
## For citation information, type citation("mcmcse").
## Type help("mcmcse-package") to get started.
```

```
out<-1
start<-1000
r<-1000
eps<-0.1
sigma<-1

out <- rw.chain(out, start, sigma)
MCSE <- mcse(out)$se
N <- length(out)
t <- qt(.975, (floor(sqrt(N) - 1)))
muhat <- mean(out)
check <- MCSE * t

while(eps < check) {
  out <- rw.chain(out, r, sigma)
  MCSE <- append(MCSE, mcse(out)$se)
  N <- length(out)
  t <- qt(.975, (floor(sqrt(N) - 1)))
  muhat <- append(muhat, mean(out))
  check <- MCSE[length(MCSE)] * t
}

plot(out, type = "l")
```



## Part II - Anguilla eel data

Consider the **Anguilla** eel data provided in the **dismo** R package. The data consists of 1,000 observations from a New Zealand survey of site-level presence or absence for the short-finned eel (*Anguilla australis*). We will use six out of twelve covariates. Five are continuous variables: **SegSumT**, **DSDist**, **USNative**, **DMaxSlope** and **DSSlope**; one is a categorical variable: **Method**, with five levels **Electric**, **Spo**, **Trap**, **Net** and **Mixture**.

Let  $x_i$  be the regression vector of covariates for the  $i$ th observation of length  $k$  and  $\beta = (\beta_0, \dots, \beta_9)$  be the vector regression coefficients. For the  $i$ th observation, suppose  $Y_i = 1$  denotes presence and  $Y_i = 0$  denotes absence of *Anguilla australis*. Then the Bayesian logistic regression model is given by

$$\begin{aligned} Y_i &\sim \text{Bernoulli}(p_i) , \\ p_i &\sim \frac{\exp(x_i^T \beta)}{1 + \exp(x_i^T \beta)} \text{ and,} \\ \beta &\sim N(\mathbf{0}, \sigma_\beta^2 \mathbf{I}_k) , \end{aligned}$$

where  $\mathbf{I}_k$  is the  $k \times k$  identity matrix. For the analysis,  $\sigma_\beta^2 = 100$  was chosen to represent a diffuse prior distribution on  $\beta$ .

9. Implement an MCMC sampler for the target distribution using the **MCMClogit** function in the **MCMCpack** package.

```
library(dismo)

## Loading required package: raster
## Loading required package: sp
library(dplyr)

##
## Attaching package: 'dplyr'

## The following objects are masked from 'package:raster':
##
##   intersect, select, union

## The following objects are masked from 'package:stats':
##
##   filter, lag

## The following objects are masked from 'package:base':
##
##   intersect, setdiff, setequal, union

library(MCMCpack)

## Loading required package: coda
## Loading required package: MASS

##
## Attaching package: 'MASS'

## The following object is masked from 'package:dplyr':
##
##   select

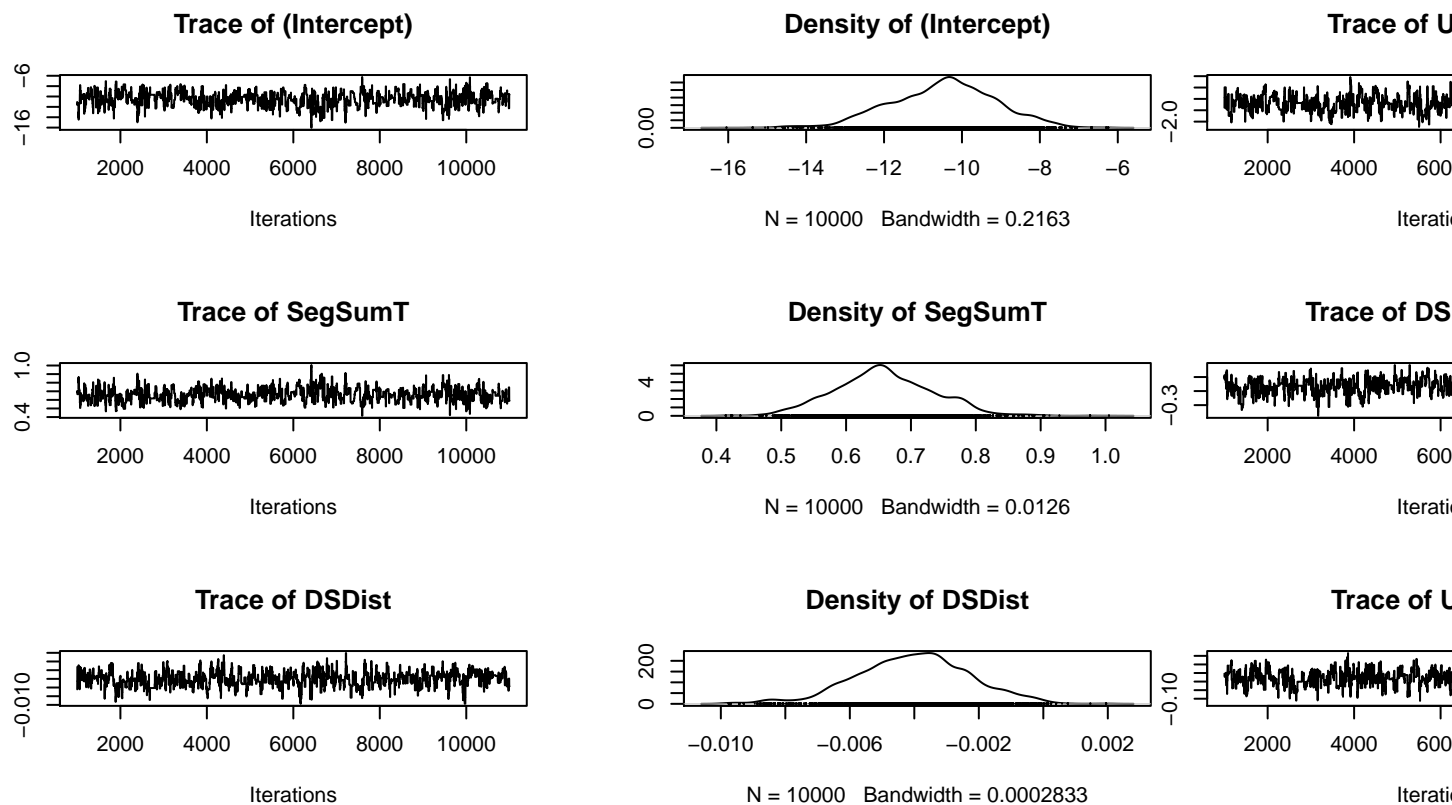
## The following objects are masked from 'package:raster':
##
##   area, select
```

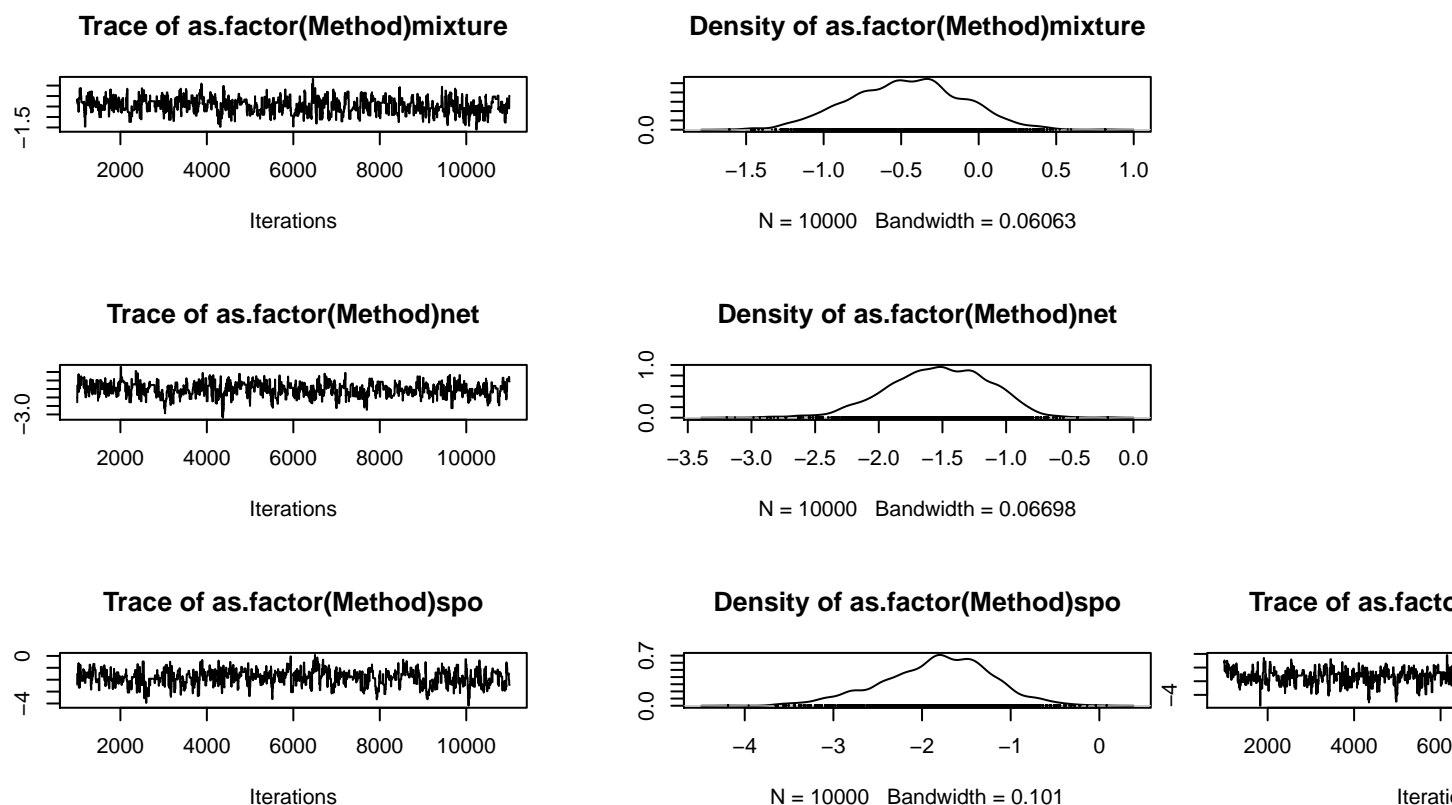
```
## ##
## ## Markov Chain Monte Carlo Package (MCMCpack)
## ## Copyright (C) 2003-2018 Andrew D. Martin, Kevin M. Quinn, and Jong Hee Park
## ##
## ## Support provided by the U.S. National Science Foundation
## ## (Grants SES-0350646 and SES-0350613)
## ##
```

```
library(MASS)
data(Anguilla_train)
aedata<-subset(Anguilla_train, select=c("Angaus", "SegSumT", "DSDist", "USNative", "DSMaxSlope", "USSlope"))
posterior <- MCMClogit(Angaus~SegSumT+DSDist+USNative+DSMaxSlope+USSlope+as.factor(Method), b0=0, B0=0.0)
```

10. Comment on the mixing properties for your sampler. Include at least one plot in support of your comments.

```
plot(posterior)
```





All the samplers seem to mix well after 1000 burn-in period.

11. Run your sampler for 100,000 iterations. Estimate the posterior mean along with an 80% Bayesian credible interval for each regression coefficient in the model. Be sure to include uncertainty estimates.

```
library(mcmcse)
posterior.it <- MCMClogit(Angaus~SegSumT+DSDist+USNative+DSMaxSlope+USSlope+as.factor(Method), b0=0, BO=1)
for (i in 1:10){
  left.est<-mcse.q(posterior.it[,i], .1)$est
  left.se<-mcse.q(posterior.it[,i], .1)$se
  right.est<-mcse.q(posterior.it[,i], .9)$est
  right.se<-mcse.q(posterior.it[,i], .9)$se
  print(paste("The 80% Bayesian credible interval for", expression(beta), i-1, "is (", left.est, "+/-", right.se, ")"))
}
```

```
## [1] "The 80% Bayesian credible interval for beta 0 is ( -12.2426282589807 +/- 0.0419237941246685 , -12.2426282589807 +/- 0.0419237941246685 )"
## [1] "The 80% Bayesian credible interval for beta 1 is ( 0.555750838735138 +/- 0.00182502112873983 , 0.555750838735138 +/- 0.00182502112873983 )"
## [1] "The 80% Bayesian credible interval for beta 2 is ( -0.00626101683144014 +/- 5.24479144272683e-05 , -0.00626101683144014 +/- 5.24479144272683e-05 )"
## [1] "The 80% Bayesian credible interval for beta 3 is ( -1.63392655774915 +/- 0.0103410398062143 , -1.63392655774915 +/- 0.0103410398062143 )"
## [1] "The 80% Bayesian credible interval for beta 4 is ( -0.244644298854542 +/- 0.00159683859939209 , -0.244644298854542 +/- 0.00159683859939209 )"
## [1] "The 80% Bayesian credible interval for beta 5 is ( -0.0758672979408923 +/- 0.000482667771468207 , -0.0758672979408923 +/- 0.000482667771468207 )"
## [1] "The 80% Bayesian credible interval for beta 6 is ( -0.905797523720252 +/- 0.00890730532111897 , -0.905797523720252 +/- 0.00890730532111897 )"
## [1] "The 80% Bayesian credible interval for beta 7 is ( -2.0226585412564 +/- 0.0103706502657955 , -2.0226585412564 +/- 0.0103706502657955 )"
## [1] "The 80% Bayesian credible interval for beta 8 is ( -2.63374203880239 +/- 0.0194911353470413 , -2.63374203880239 +/- 0.0194911353470413 )"
## [1] "The 80% Bayesian credible interval for beta 9 is ( -3.29305746017398 +/- 0.0169016328667814 , -3.29305746017398 +/- 0.0169016328667814 )"
```

12. Compare your Bayesian estimates to those obtained via maximum likelihood estimation.

```

beta.bayes.est<-c()
for (i in 1:10){
  beta.bayes.est[i]<-mcse(posterior.it[,i])$est
}
beta.bayes.est

## [1] -10.435352111  0.656157347 -0.004026127 -1.171925814 -0.171456890
## [6] -0.052225223 -0.467525699 -1.520005399 -1.828439743 -2.591878413

```

```
library(Rlab)
```

```
## Rlab 2.15.1 attached.
```

```
##
```

```
## Attaching package: 'Rlab'
```

```
## The following object is masked from 'package:MASS':
```

```
##
```

```
## michelson
```

```
## The following object is masked from 'package:dplyr':
```

```
##
```

```
## count
```

```
## The following objects are masked from 'package:stats':
```

```
##
```

```
## dexp, dgamma, dweibull, pexp, pgamma, pweibull, qexp, qgamma,
```

```
## qweibull, rexp, rgamma, rweibull
```

```
## The following object is masked from 'package:datasets':
```

```
##
```

```
## precip
```

```
aedata<-subset(Anguilla_train, select=c("Angaus", "SegSumT", "DSDist", "USNative", "DSMaxSlope", "USSlop
```

```
aedata$mixture<-ifelse(aedata$Method=="mixture", 1, 0)
```

```
aedata$net<-ifelse(aedata$Method=="net", 1, 0)
```

```
aedata$spo<-ifelse(aedata$Method=="spo", 1, 0)
```

```
aedata$trap<-ifelse(aedata$Method=="trap", 1, 0)
```

```
aedata<-aedata[, -7]
```

```
aedata.x<-cbind(c(rep(1, times=1000)), aedata[, -1])
```

```
nllh<-function(p=c(rep(0, times=10))){
```

```
  p<-as.matrix(p)
```

```
  aedata.x<-as.matrix(aedata.x)
```

```
  num<-exp(aedata.x %*% p)
```

```
  p<-num/(1+num)
```

```
  sum<-0
```

```
  for (i in 1:1000){
```

```
    sum<-dbern(aedata[,1][i], p[i], log=TRUE)
```

```
  }
```

```
  return(-sum)
```

```
}
```

```
nlm(nllh, p=c(rep(0, times=10)))$estimate
```

```
## [1] -0.5000001 -7.3500270 -18.5751725 -0.2150000 -0.1450000
```

```
## [6] -7.2500263 0.0000000 0.0000000 0.0000000 0.0000000
```

## Part II - Permutation tests

The Cram'ér von Mises statistic estimates the integrated square distance between distributions. It can be computed using the following formula

$$W = \frac{mn}{(m+n)^2} \left[ \sum_{i=1}^n (F_n(x_i) - G_m(x_i))^2 + \sum_{j=1}^m (F_n(y_j) - G_m(y_j))^2 \right]$$

where  $F_n$  and  $G_m$  are the corresponding empirical cdfs.

13. Implement the two sample Cram'ér von Mises test for equal distributions as a permutation test. Apply it to the `chickwts` data.

```
data(chickwts)
attach(chickwts)
library(RVAideMemoire)

## *** Package RVAideMemoire v 0.9-70 ***

##
## Attaching package: 'RVAideMemoire'

## The following object is masked from 'package:raster':
##
##      cv

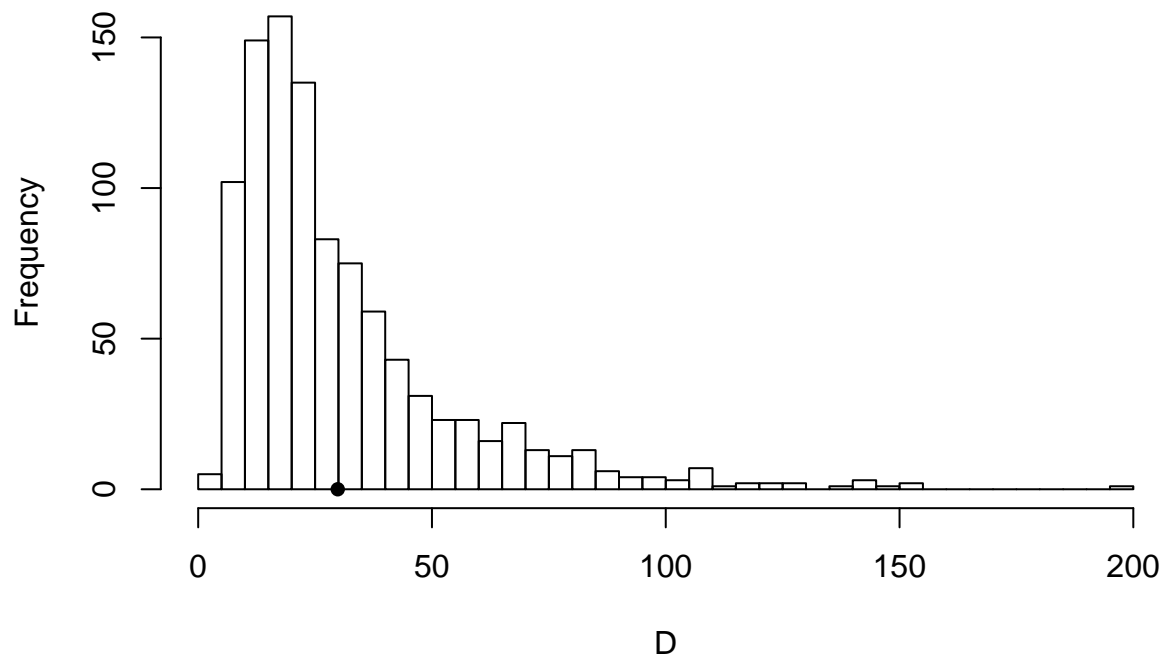
X <- as.vector(chickwts$weight[chickwts$feed=="soybean"])
Y <- as.vector(chickwts$weight[chickwts$feed=="linseed"])
B <- 999
Z <- c(X,Y)
K <- 1:26
D0<-CvM.test(X,Y)$statistic
D<-numeric(B)
for( i in 1:B){
  k <- sample(K, size=14, replace=F)
  x1 <- Z[k]
  y1 <- Z[-k]
  D[i] <- CvM.test(x1, y1)$statistic
}
p <- mean(c(D0,D) >= D0)
p

## [1] 0.374

hist(D, breaks=50, main="Permuation Distribution")
points(D0, 0, cex=1, pch=16)
```



## Permutation Distribution



14. How would you implement the bivariate Spearman rank correlation test for independence as a permutation test? The Spearman rank correlation test statistic can be obtained from the function `cor` with `method="spearman"`. Compare the achieved significance level of the permutation test with the p-value reported by `cor.test` on the same samples.

```
Score <- c(58, 48, 48, 41, 34, 43, 38, 53, 41, 60, 55, 44,
          43, 49, 47, 33, 47, 40, 46, 53, 40, 45, 39, 47,
          50, 53, 46, 53)
SAT <- c(590, 590, 580, 490, 550, 580, 550, 700, 560, 690, 800, 600,
        650, 580, 660, 590, 600, 540, 610, 580, 620, 600, 560, 560,
        570, 630, 510, 620)
r.obt <- cor(Score, SAT, method="spearman")
nreps <- 5000
r.random <- numeric(nreps)
for (i in 1:nreps) {
  Y <- Score
  X <- sample(SAT, 28, replace = FALSE)
  r.random[i] <- cor(X, Y, method="spearman")
}
prob <- length(r.random[r.random >= r.obt])/nreps
prob
```

```
## [1] 0.0016
```

```
cor.test(Score, SAT, method="spearman")$p.value
```

```
## Warning in cor.test.default(Score, SAT, method = "spearman"): Cannot
## compute exact p-value with ties
```

```
## [1] 0.0046757
```

The achieved significance level of the permutation test is slightly smaller than the p-value reported by `cor.test` on the same samples.