

# SimulationStudySummary

Xiang Ji

9/7/2017

This R markdown file summarizes Simulation Study results.

```
rm(list=ls()) # clean up workspace
setwd("/Users/xji3/GitFolders/YeastIGCTract/SimulationStudy/")

Tract.list <- c(3.0, 10.0, 50.0, 100.0, 200.0, 300.0, 400.0, 500.0)
# First read in HMM results
# from summary file
for(tract in Tract.list){
  hmm.tract.summary <- NULL
  for(sim in 1:100){
    hmm.summary <- paste("./summary/Tract_", toString(tract), '.0/sim_',
                        toString(sim), '/HMM_YDR418W_YELO54C_MG94_nonclock_sim_',
                        toString(sim), '_1D_summary.txt', sep = "")
    if (file.exists(hmm.summary)){
      all <- readLines(hmm.summary, n = -1)
      col.names <- paste("sim_", toString(sim), sep = "")
      row.names <- strsplit(all[length(all)], ' ')[[1]][-1]
      summary_mat <- as.matrix(read.table(hmm.summary,
                                          row.names = row.names,
                                          col.names = col.names))
      hmm.tract.summary <- cbind(hmm.tract.summary, summary_mat)
    }
  }
  assign(paste("HMM_Tract_", toString(tract), "_summary", sep = ""), hmm.tract.summary)
}

# from plots
for(tract in Tract.list){
  hmm.tract.plots <- NULL
  for(sim in 1:100){
    hmm.plot <- paste("./plot/Tract_", toString(tract), '.0/sim_',
                    toString(sim), '/HMM_YDR418W_YELO54C_lnL_sim_',
                    toString(sim), '_1D_surface.txt', sep = "")
    if (file.exists(hmm.plot)){
      lnL.surface <- read.table(hmm.plot)
      max.idx <- which.max(lnL.surface[, 2])
      new.summary <- matrix(c(3.0*exp(-lnL.surface[max.idx, 1]), lnL.surface[max.idx, 2]), 2, 1)
      rownames(new.summary) <- c("tract in nt", "lnL")
      colnames(new.summary) <- paste("sim_", toString(sim), sep = "")
      hmm.tract.plots <- cbind(hmm.tract.plots, new.summary)
    }
  }
  assign(paste("HMM_Tract_", toString(tract), "_plot", sep = ""), hmm.tract.plots)
}
```

```

# Now read in PSJS summary results
for(tract in Tract.list){
  PSJS.tract.summary <- NULL
  for(sim in 1:100){
    PSJS.summary <- paste("./summary/Tract_", toString(tract), '.0/sim_',
                          toString(sim), '/PSJS_HKY_rv_sim_',
                          toString(sim), "_Tract_", toString(tract), '.0_summary.txt', sep = "")
    if (file.exists(PSJS.summary)){
      all <- readLines(PSJS.summary, n = -1)
      col.names <- paste("sim_", toString(sim), sep = "")
      row.names <- strsplit(all[length(all)], ' ')[[1]][-1]
      summary_mat <- as.matrix(read.table(PSJS.summary,
                                          row.names = row.names,
                                          col.names = col.names))
      PSJS.tract.summary <- cbind(PSJS.tract.summary, summary_mat)
    }
  }
  assign(paste("PSJS_Tract_", toString(tract), "_summary", sep = ""), PSJS.tract.summary)
}

# Now read in actual mean tract length in each simulated dataset
for (tract in Tract.list){
  sim.tract <- NULL
  for(sim in 1:100){
    sim_log <- paste("./Tract_", toString(tract), ".0/sim_", toString(sim),
                    "/YDR418W_YELO54C_sim_", toString(sim), "_IGC.log", sep = "")
    # now read in log file
    log_info <- read.table(sim_log, header = TRUE)
    #tract.length <- log_info[, "stop_pos"] - log_info[, "start_pos"] + 1
    tract.length <- log_info[, "tract_length"]
    new.info <- matrix(c(mean(tract.length), sd(tract.length)), 2, 1)
    rownames(new.info) <- c("mean tract length", "sd tract length")
    colnames(new.info) <- paste("sim_", toString(sim), sep = "")
    sim.tract <- cbind(sim.tract, new.info)
  }
  assign(paste("sim.tract.", toString(tract), sep = ""), sim.tract)
}

```

OK, Now show the performance summary

```

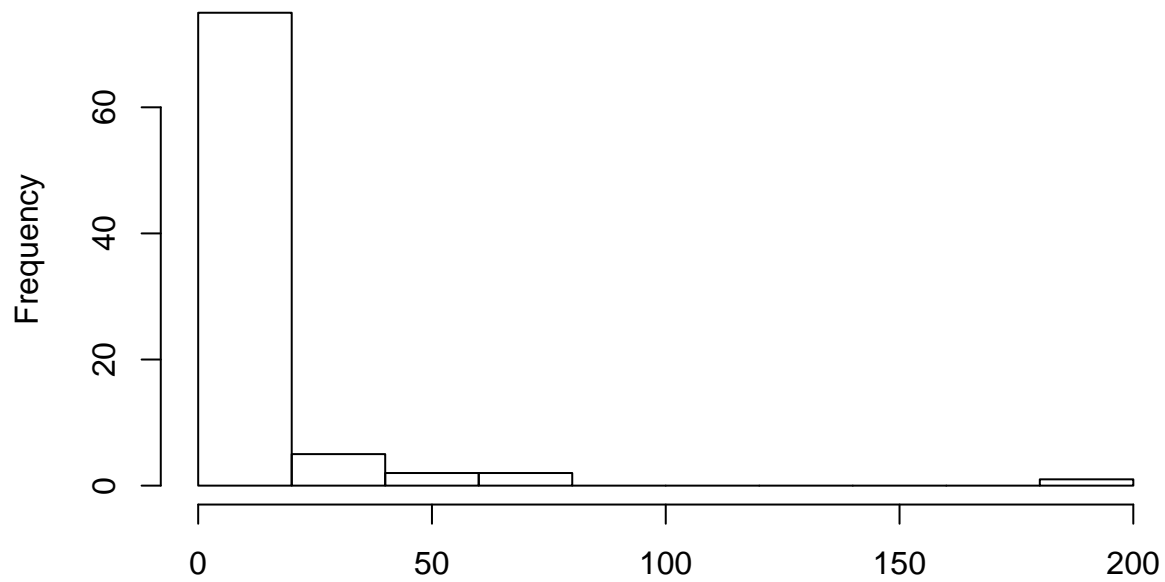
# HMM results
for (tract in Tract.list){
  # show how many stuck at boundary 1000 nt first
  print(paste("Tract = ", toString(tract), sep = ""))
  summary_mat <- get(paste("HMM_Tract_", toString(tract), "_plot", sep = ""))

  # histogram of inferred tract length
  hist(summary_mat[1, summary_mat[1, ] < 999.], main = paste("Tract = ", toString(tract), sep = ""))
  print(paste("Among total 100 simulated data sets, ", toString(sum(summary_mat[1, ] > 999)),
              " datasets stuck at 1000", sep = ""))
  print(c("mean", mean(summary_mat[1, summary_mat[1, ] < 999.]),
          "sd", sd(summary_mat[1, summary_mat[1, ] < 999.])))
}

```

```
## [1] "Tract = 3"
```

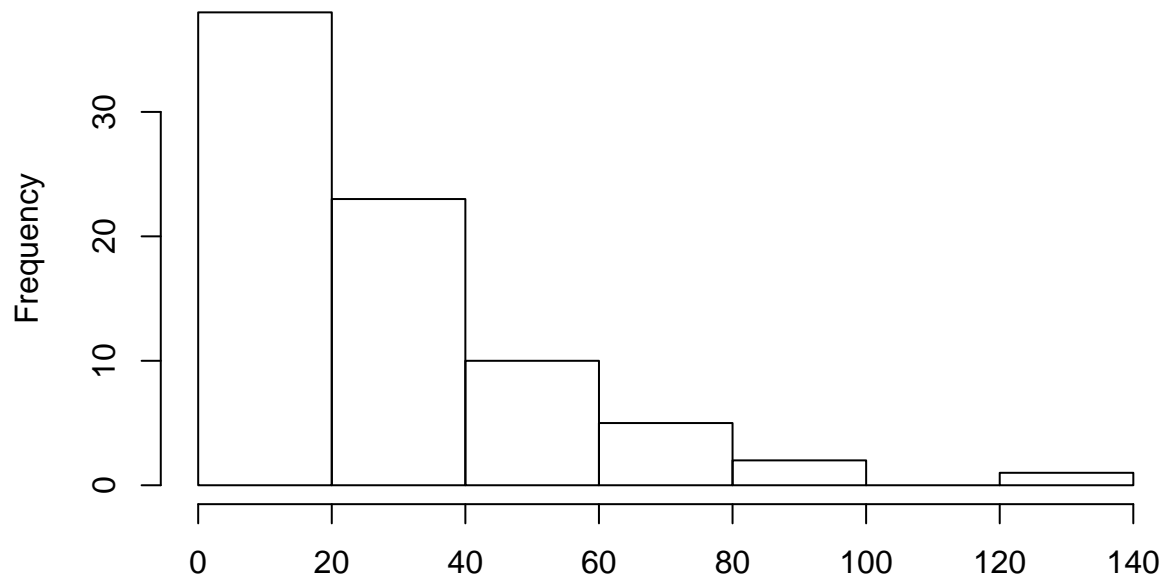
### Tract = 3



`summary_mat[1, summary_mat[1, ] < 999]`

```
## [1] "Among total 100 simulated data sets, 15 datasets stuck at 1000"  
## [1] "mean"          "11.6941176470561" "sd"  
## [4] "23.409920327626"  
## [1] "Tract = 10"
```

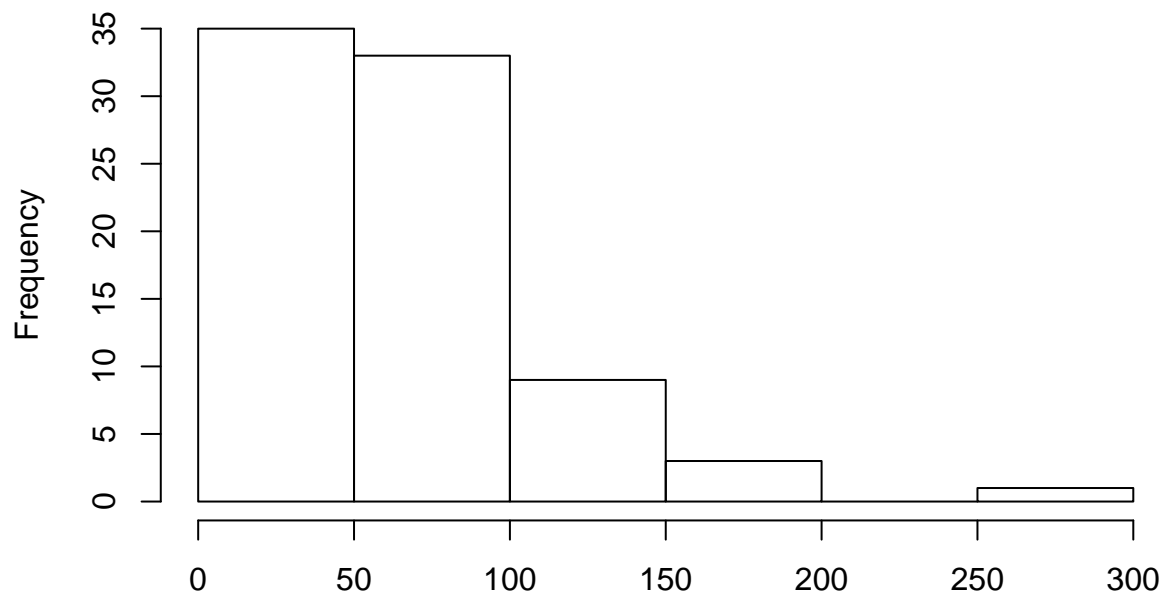
### Tract = 10



`summary_mat[1, summary_mat[1, ] < 999]`

```
## [1] "Among total 100 simulated data sets, 21 datasets stuck at 1000"  
## [1] "mean"          "28.5316455696352" "sd"  
## [4] "24.4245753935622"  
## [1] "Tract = 50"
```

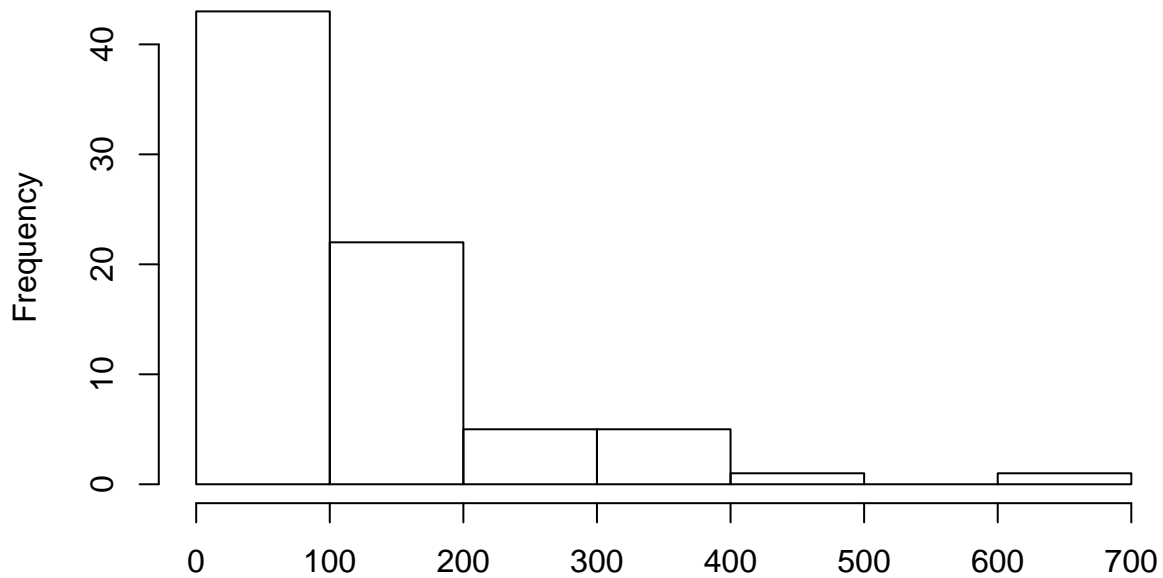
### Tract = 50



`summary_mat[1, summary_mat[1, ] < 999]`

```
## [1] "Among total 100 simulated data sets, 19 datasets stuck at 1000"
## [1] "mean"          "67.0864197531772" "sd"
## [4] "50.0167965615616"
## [1] "Tract = 100"
```

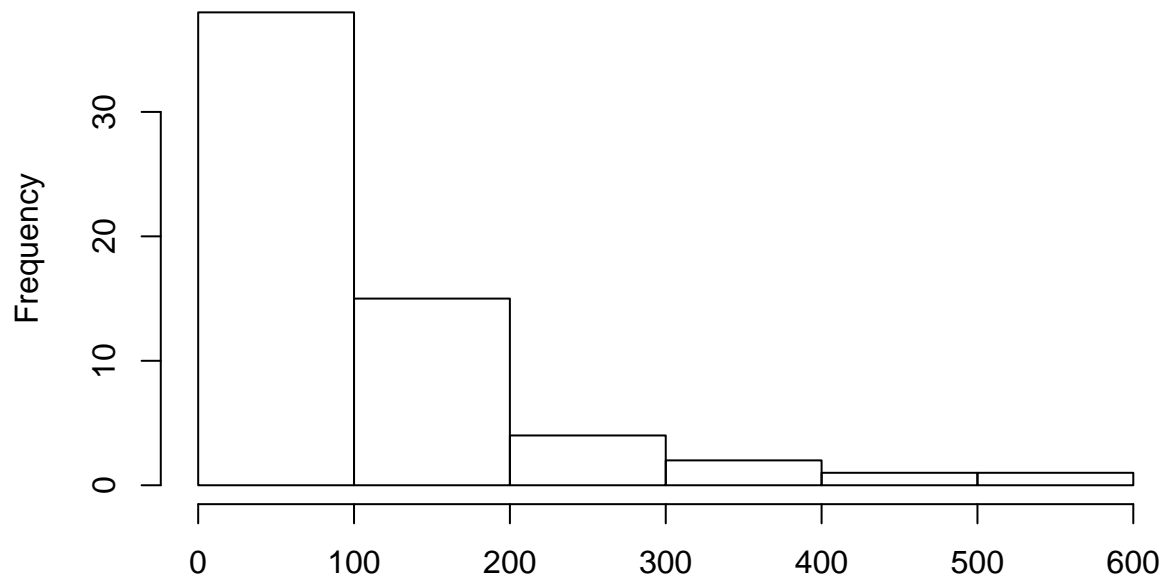
### Tract = 100



`summary_mat[1, summary_mat[1, ] < 999]`

```
## [1] "Among total 100 simulated data sets, 23 datasets stuck at 1000"
## [1] "mean"          "121.597402597374" "sd"
## [4] "111.189768998766"
## [1] "Tract = 200"
```

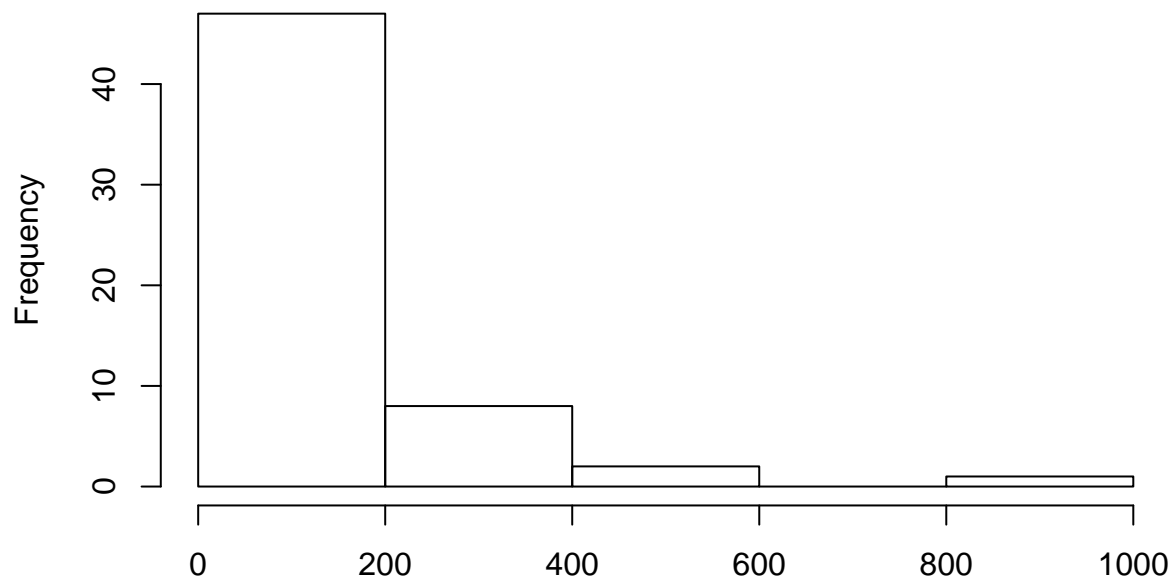
### Tract = 200



summary\_mat[1, summary\_mat[1, ] < 999]

```
## [1] "Among total 100 simulated data sets, 38 datasets stuck at 1000"  
## [1] "mean"          "111.590163934384" "sd"  
## [4] "109.950955285812"  
## [1] "Tract = 300"
```

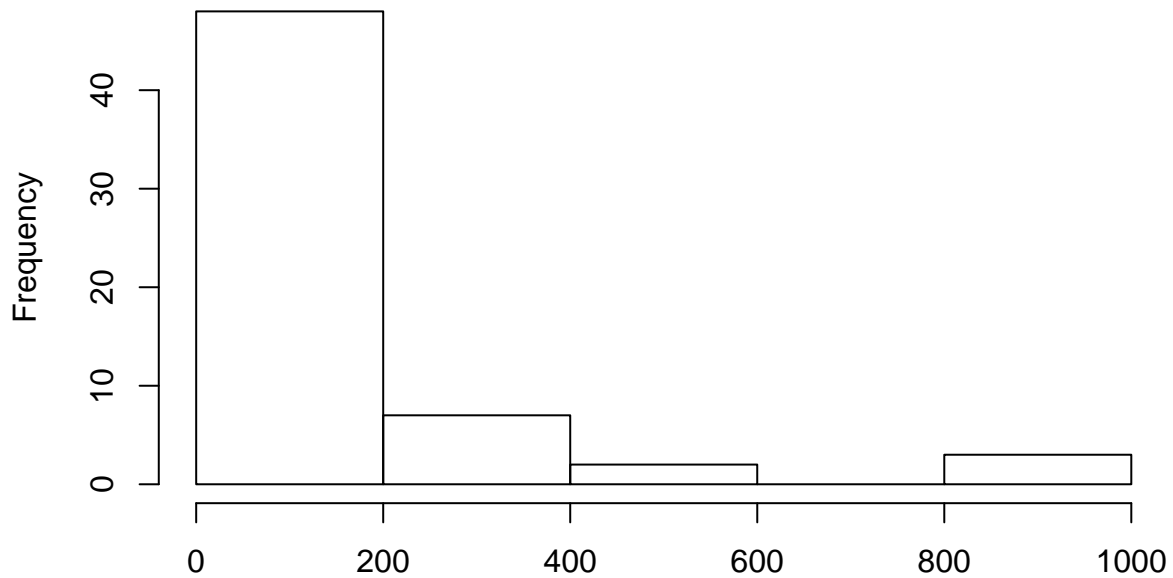
### Tract = 300



summary\_mat[1, summary\_mat[1, ] < 999]

```
## [1] "Among total 100 simulated data sets, 42 datasets stuck at 1000"
## [1] "mean"          "124.293103448246" "sd"
## [4] "157.010001761835"
## [1] "Tract = 400"
```

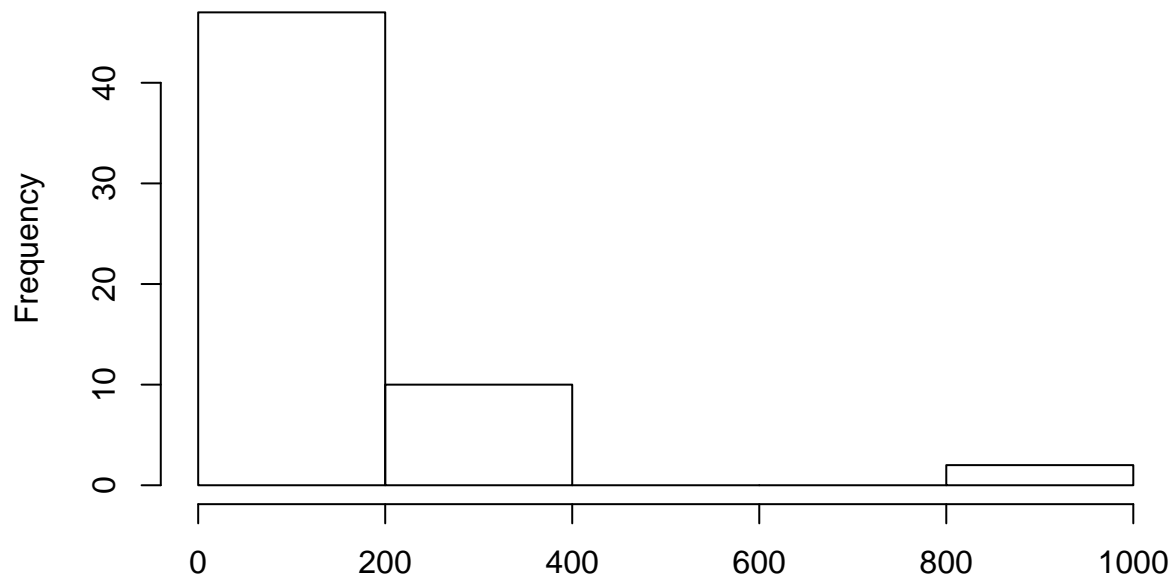
### Tract = 400



`summary_mat[1, summary_mat[1, ] < 999]`

```
## [1] "Among total 100 simulated data sets, 40 datasets stuck at 1000"
## [1] "mean"          "153.333333333315" "sd"
## [4] "219.73201449035"
## [1] "Tract = 500"
```

## Tract = 500



summary\_mat[1, summary\_mat[1, ] < 999]

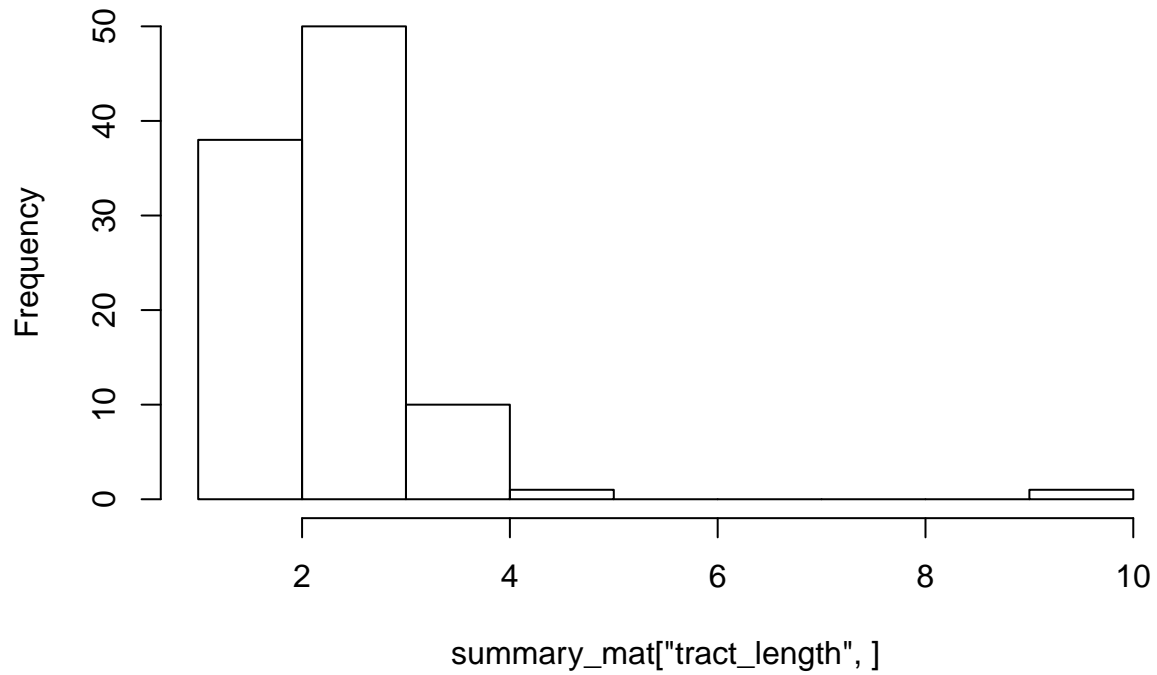
```
## [1] "Among total 100 simulated data sets, 41 datasets stuck at 1000"
## [1] "mean"          "136.30508474575" "sd"
## [4] "187.058284687408"
```

*# PSJS results*

```
for(tract in Tract.list){
  summary_mat <- get(paste("PSJS_Tract_", toString(tract), "_summary", sep = ""))
  # histogram of inferred tract length
  hist(summary_mat["tract_length", ], main = paste("Tract = ", toString(tract), sep = ""))
  print(c("mean", mean(summary_mat["tract_length", ]),
          "sd", sd(summary_mat["tract_length", ])))
}
```

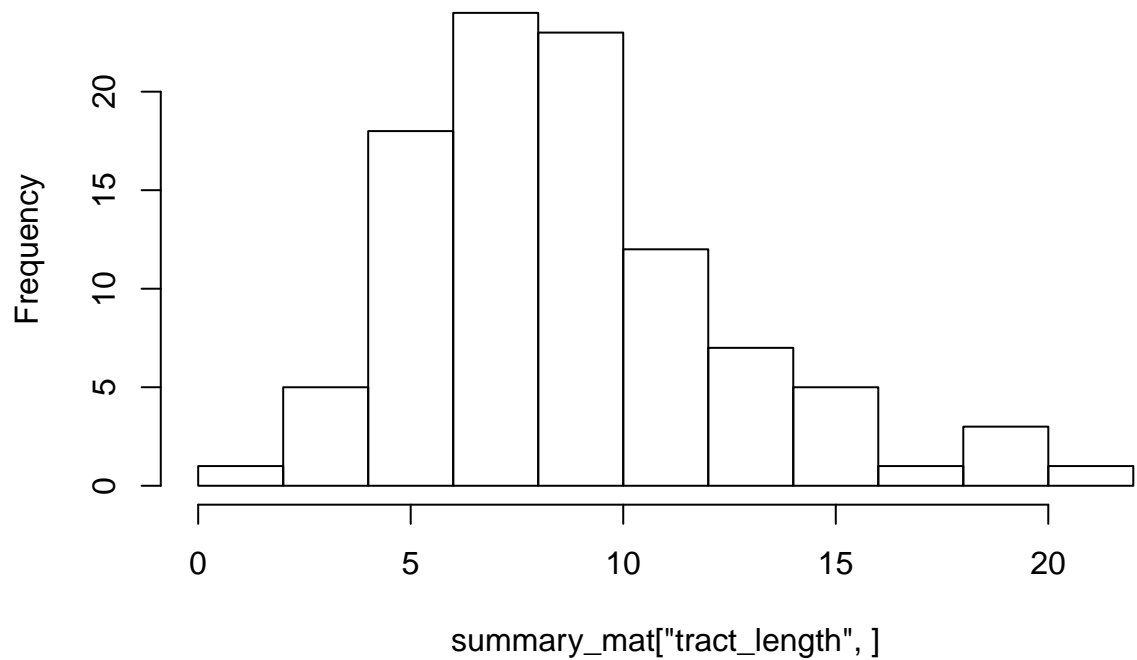


### Tract = 3



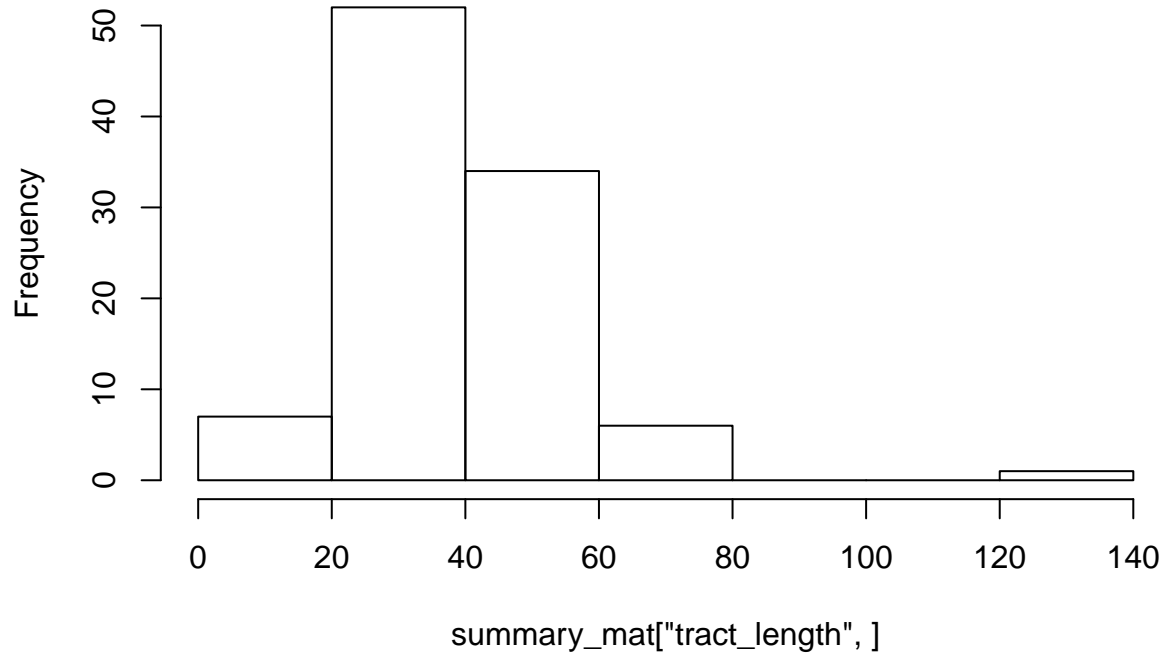
```
## [1] "mean"          "2.29718013607812" "sd"  
## [4] "1.02874535788003"
```

### Tract = 10



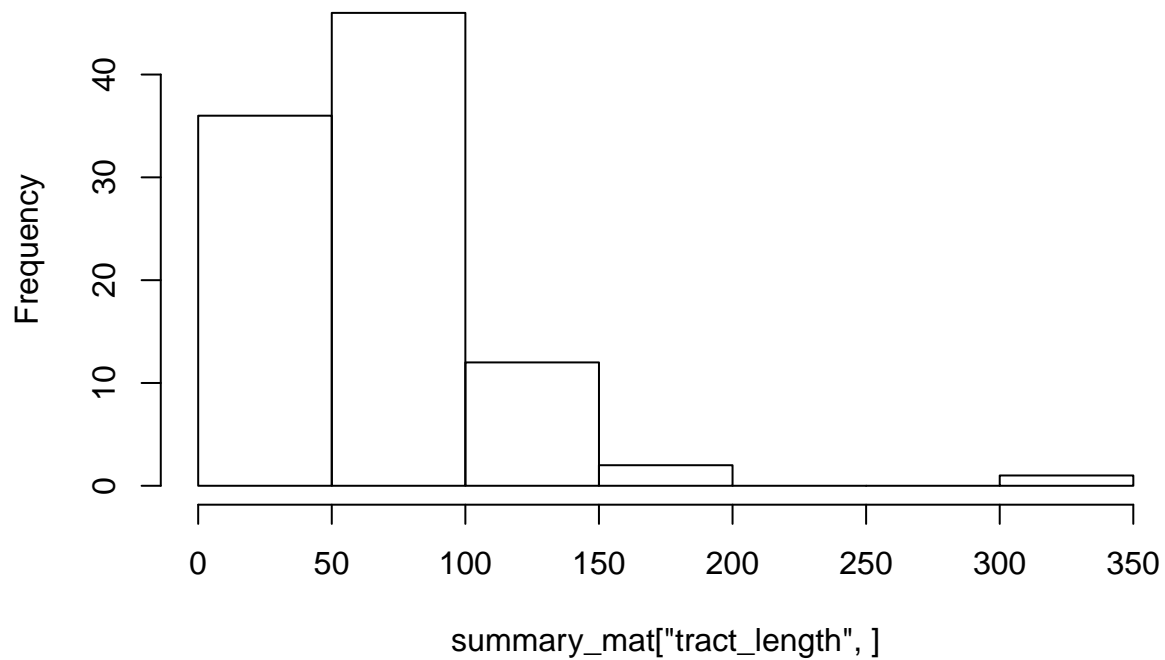
```
## [1] "mean"          "8.72620746076138" "sd"  
## [4] "3.72714198660433"
```

### Tract = 50



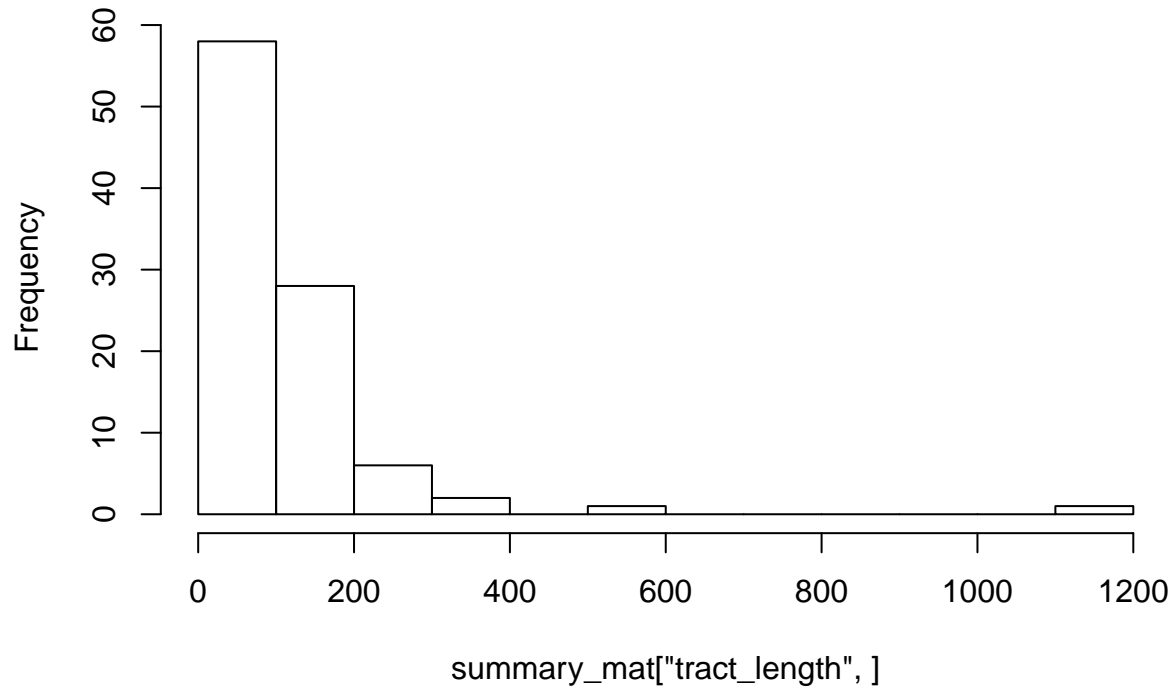
```
## [1] "mean"          "38.1001820929933" "sd"  
## [4] "16.207953647501"
```

### Tract = 100



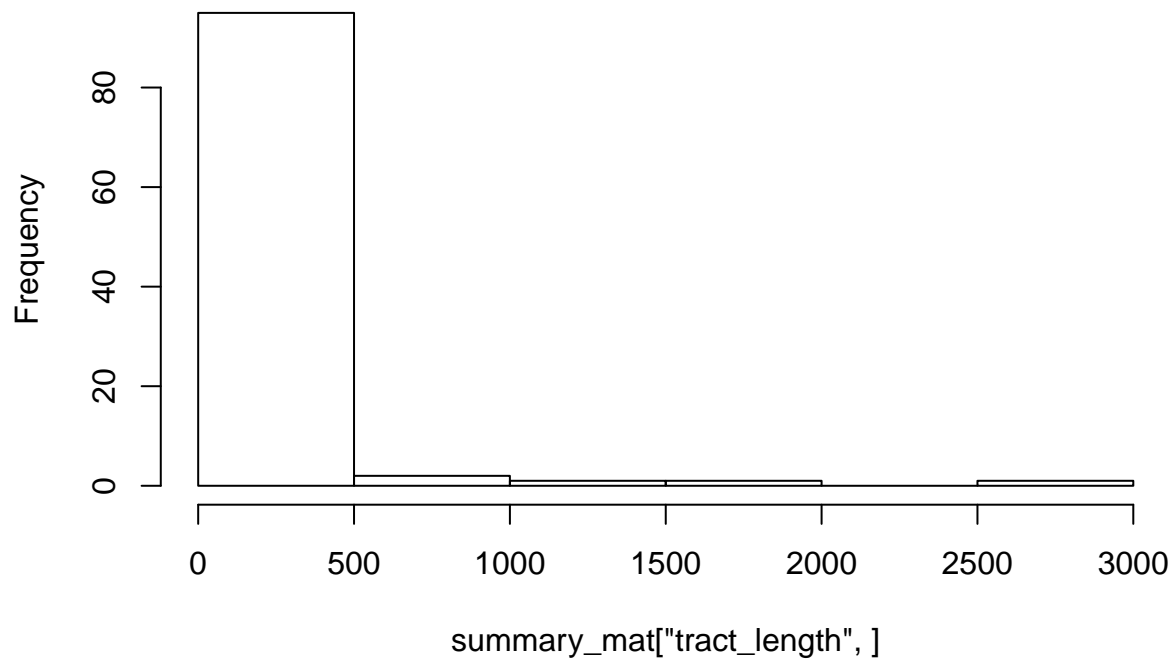
```
## [1] "mean"          "66.9946146602955" "sd"  
## [4] "40.527169109434"
```

### Tract = 200



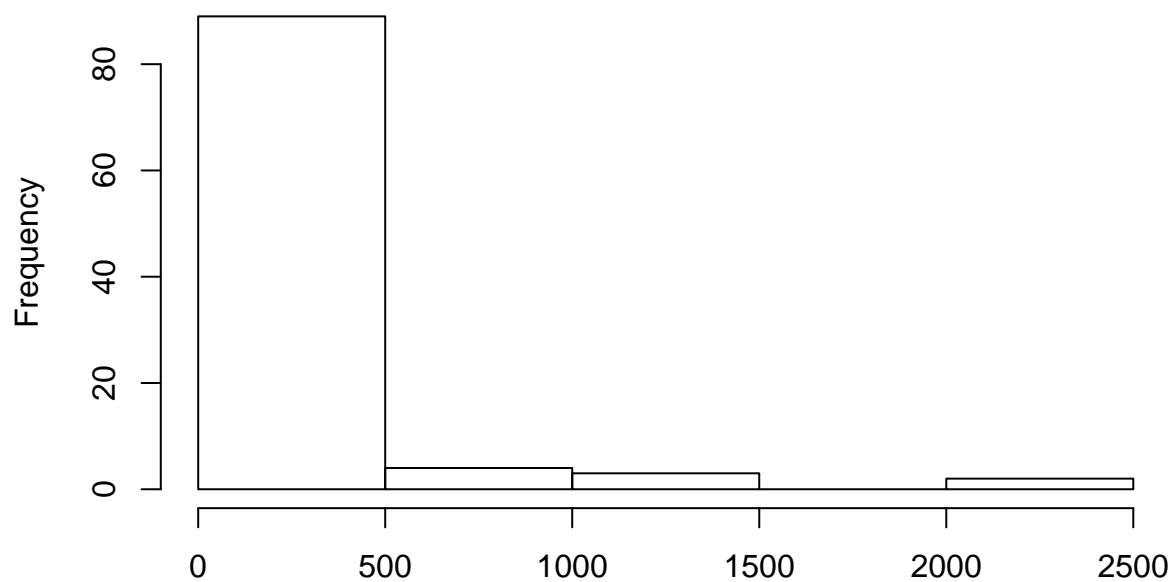
```
## [1] "mean"          "115.573878487051" "sd"  
## [4] "132.756671098129"
```

### Tract = 300



```
## [1] "mean"          "189.995520518941" "sd"  
## [4] "356.357994056025"
```

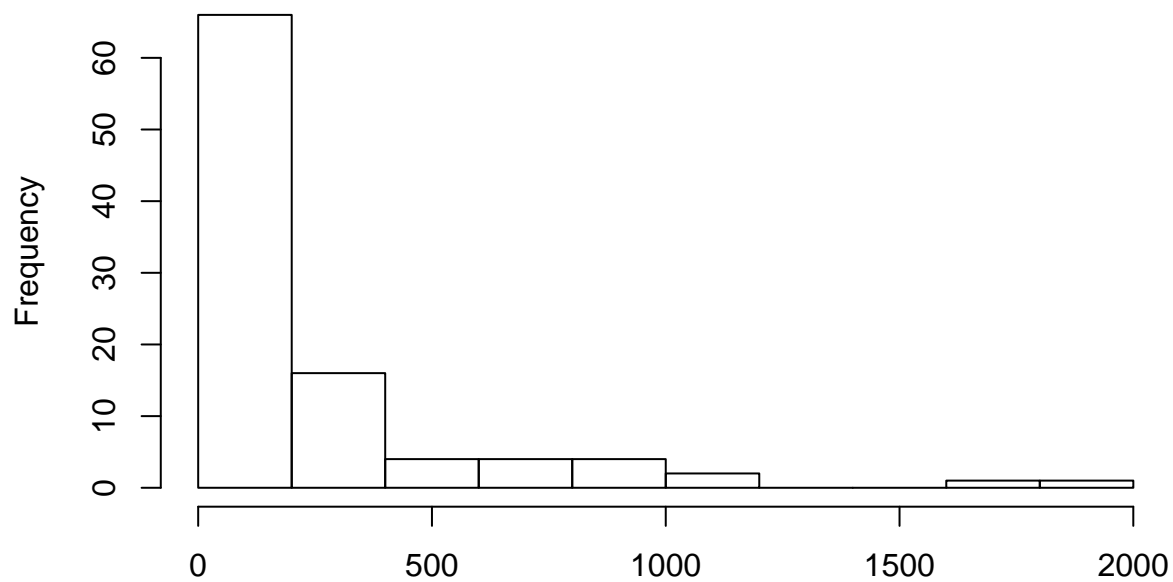
### Tract = 400



summary\_mat["tract\_length", ]

```
## [1] "mean"          "217.14487788415" "sd"  
## [4] "377.699962735285"
```

### Tract = 500



summary\_mat["tract\_length", ]

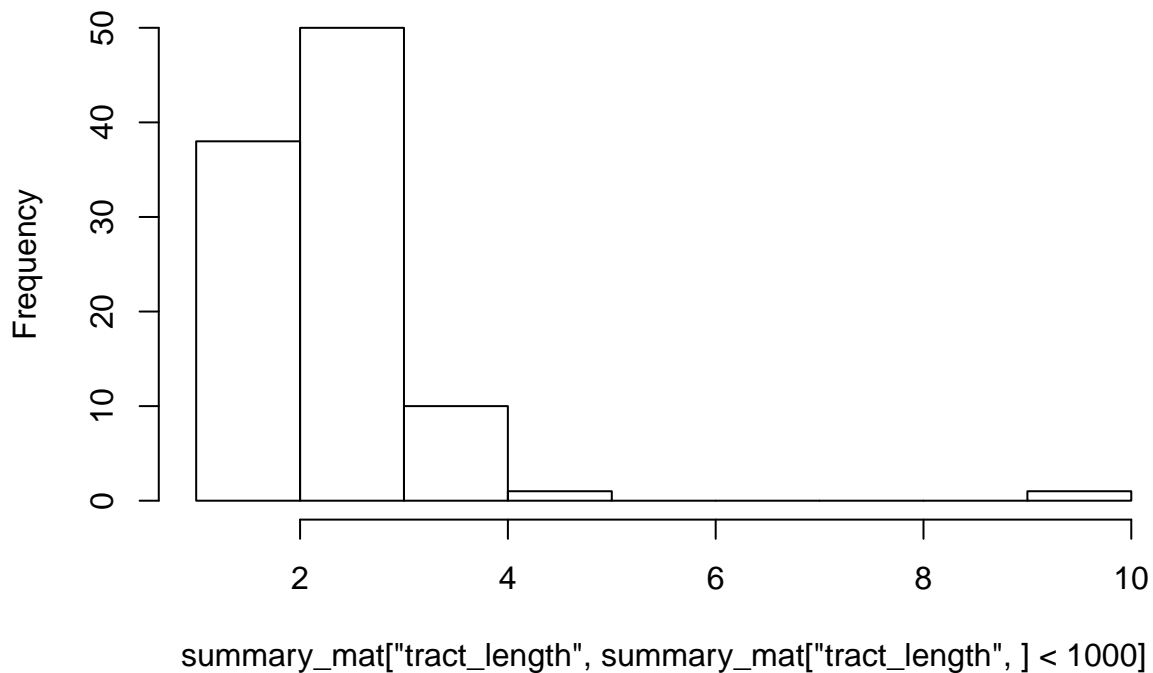
```
## [1] "mean"          "240.064448783542" "sd"  
## [4] "348.396527453917"
```

```

# exclude suspiciously long inferred tract length, plot again
for(tract in Tract.list){
  summary_mat <- get(paste("PSJS_Tract_", toString(tract), "_summary", sep = ""))
  # histogram of inferred tract length
  hist(summary_mat["tract_length", summary_mat["tract_length", ] < 1000.0], main = paste("Tract = ", toString(tract)),
  print(c("mean", mean(summary_mat["tract_length", summary_mat["tract_length", ] < 1000.0]),
          "sd", sd(summary_mat["tract_length", summary_mat["tract_length", ] < 1000.0)))))
}

```

### Tract = 3

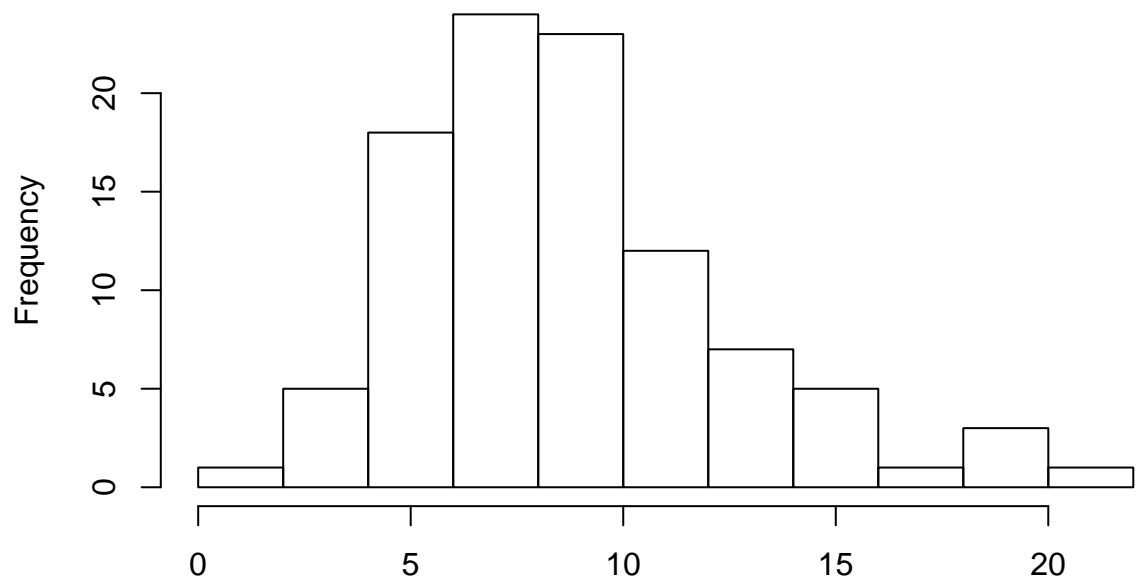


```

## [1] "mean"          "2.29718013607812" "sd"
## [4] "1.02874535788003"

```

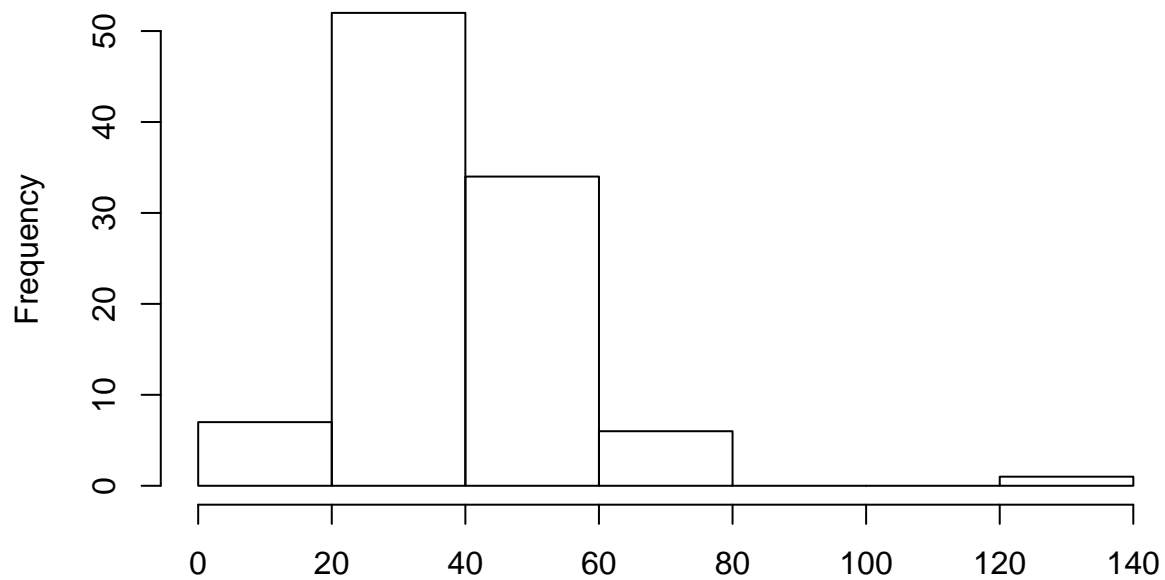
### Tract = 10



```
summary_mat["tract_length", summary_mat["tract_length", ] < 1000]
```

```
## [1] "mean"          "8.72620746076138" "sd"  
## [4] "3.72714198660433"
```

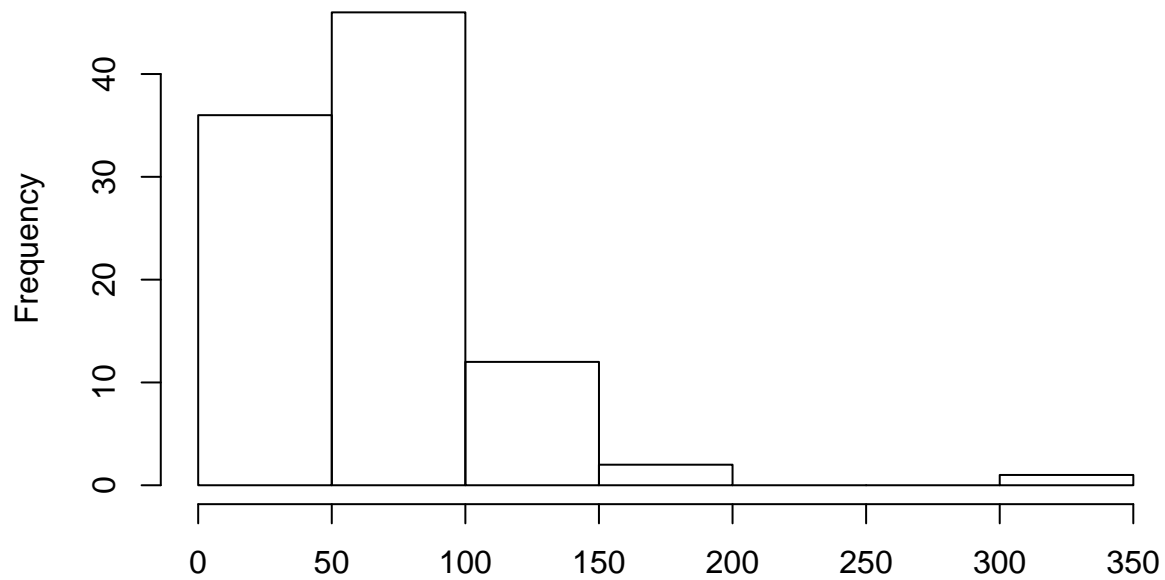
### Tract = 50



```
summary_mat["tract_length", summary_mat["tract_length", ] < 1000]
```

```
## [1] "mean"          "38.1001820929933" "sd"  
## [4] "16.207953647501"
```

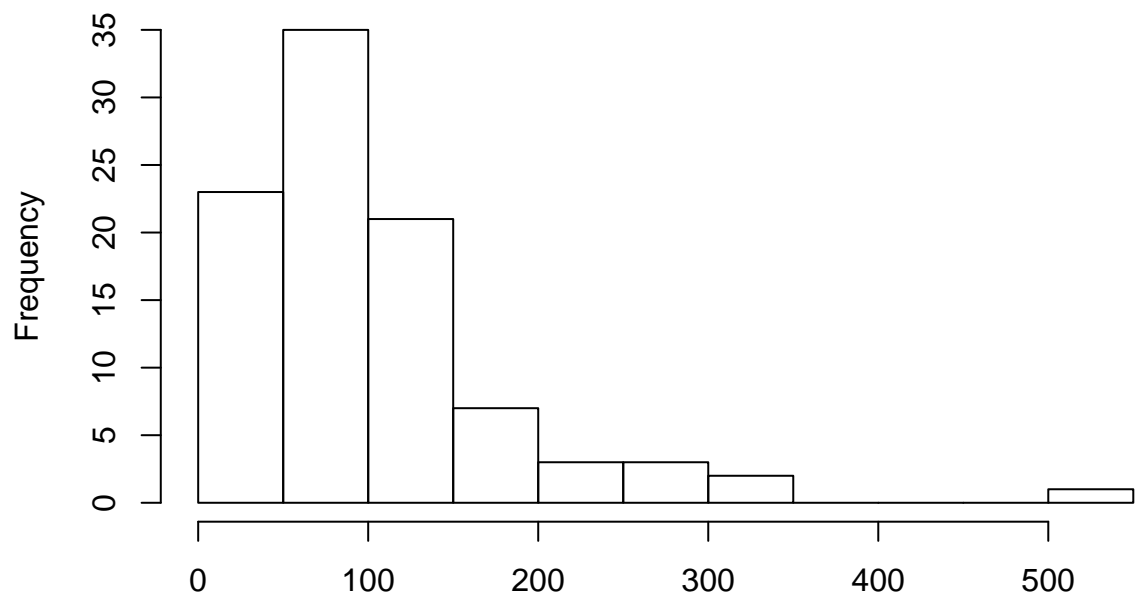
### Tract = 100



```
summary_mat["tract_length", summary_mat["tract_length", ] < 1000]
```

```
## [1] "mean"          "66.9946146602955" "sd"  
## [4] "40.527169109434"
```

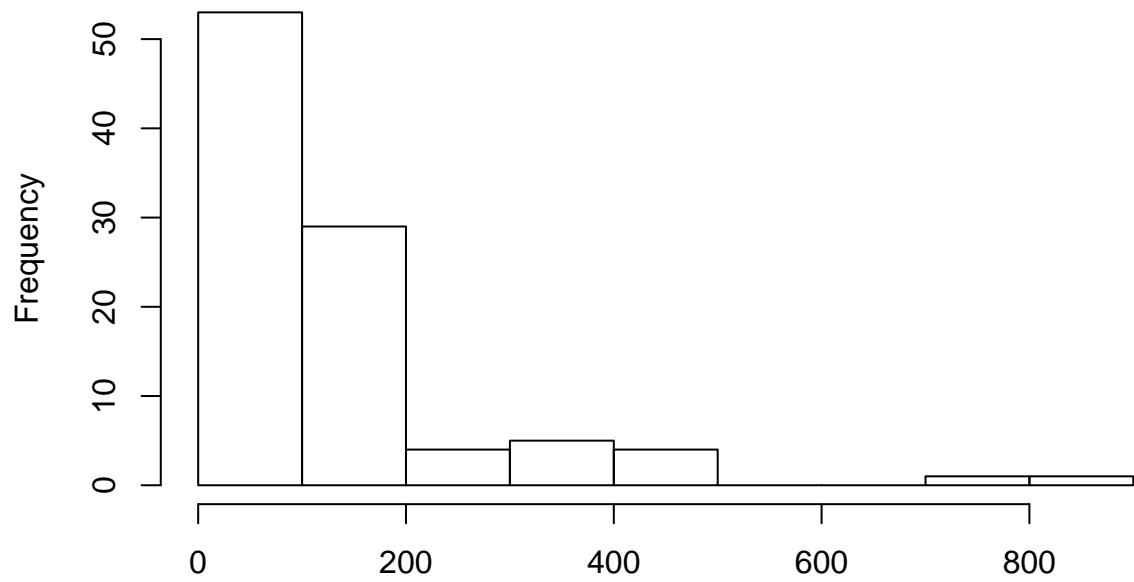
### Tract = 200



```
summary_mat["tract_length", summary_mat["tract_length", ] < 1000]
```

```
## [1] "mean"          "105.044108496939" "sd"  
## [4] "83.9910122074207"
```

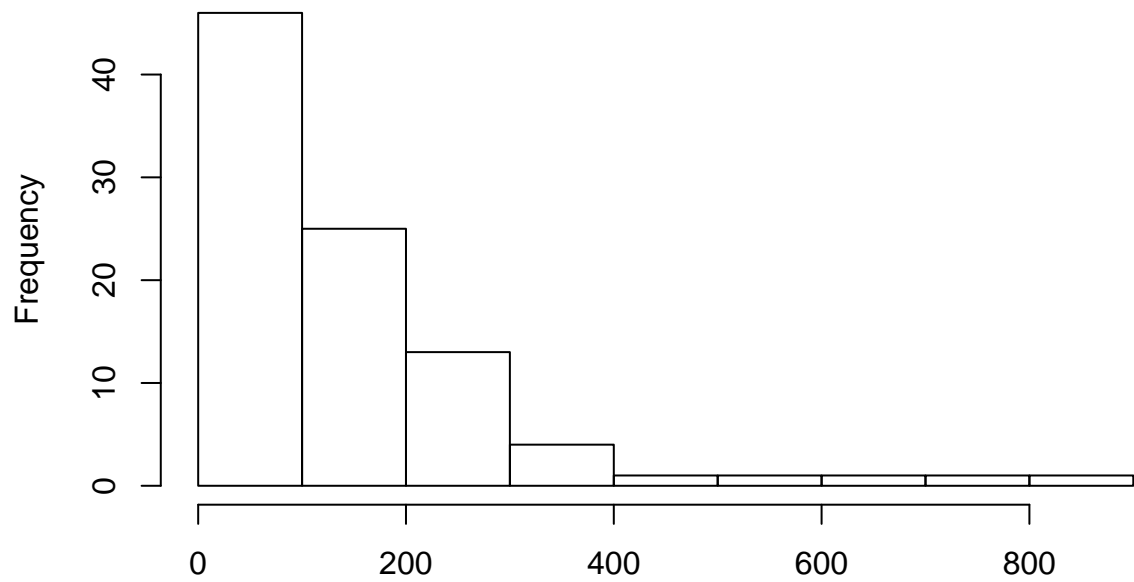
### Tract = 300



```
summary_mat["tract_length", summary_mat["tract_length", ] < 1000]
```

```
## [1] "mean"          "135.843591947529" "sd"  
## [4] "146.132608739573"
```

### Tract = 400

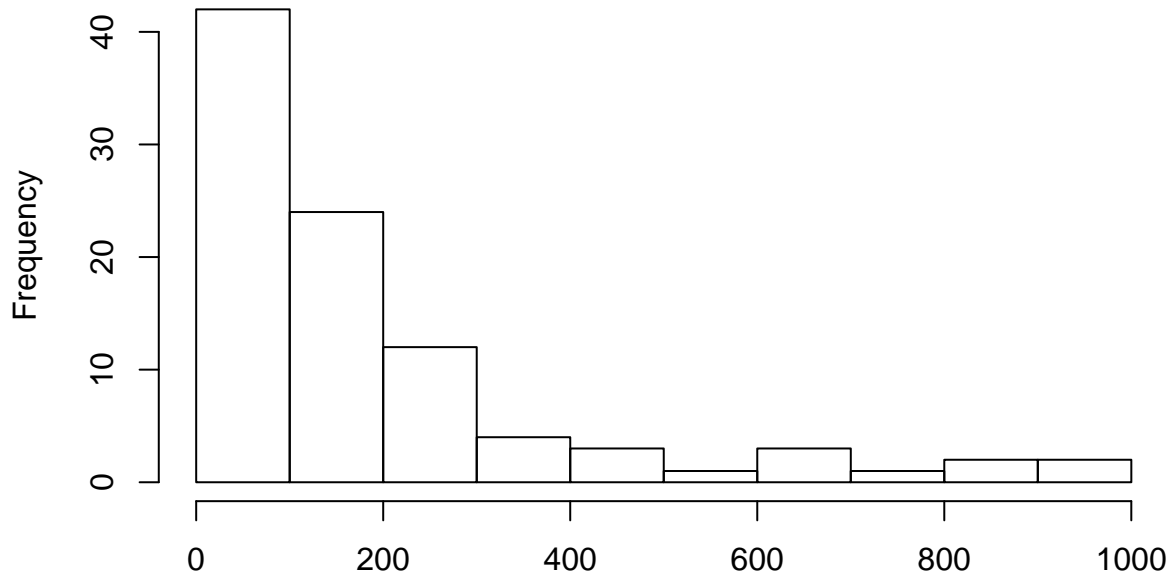


```
summary_mat["tract_length", summary_mat["tract_length", ] < 1000]
```

```
## [1] "mean"          "142.57022420328" "sd"  
## [4] "154.570518796298"
```



## Tract = 500



summary\_mat["tract\_length", summary\_mat["tract\_length", ] < 1000]

```
## [1] "mean"          "187.222502177302" "sd"
## [4] "223.924381481332"
```

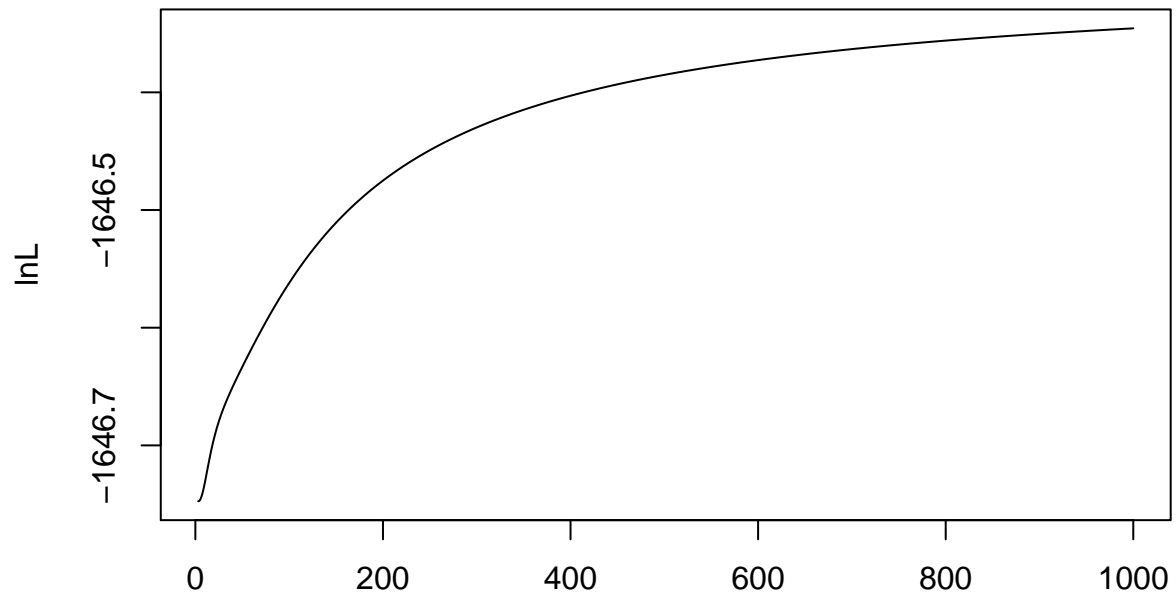
A plot of HMM surface that infer tract length at boundary from each tract length condition

*# plot the first two datasets that HMM inferred tract length stuck at boundary of 1000.0*

```
for (tract in Tract.list){
  print(paste("Tract = ", toString(tract), sep = ""))
  summary_mat <- get(paste("HMM_Tract_", toString(tract), "_plot", sep = ""))
  count = 0
  for(iter in 1:99){
    if(summary_mat[1, iter] > 999. & count < 2){
      to.plot <- read.table(paste("./plot/Tract_", toString(tract), ".0/", colnames(summary_mat)[iter],
                                "/HMM_YDR418W_YEL054C_lnL_", colnames(summary_mat)[iter],
                                "_1D_surface.txt", sep = ""))
      plot(3.0*exp(-to.plot[, 1]), to.plot[, 2], main = paste("Tract = ", toString(tract), " ", colnames(summary_mat)[iter]),
           ylab = "lnL")
      count = count + 1
    }
  }
}
```

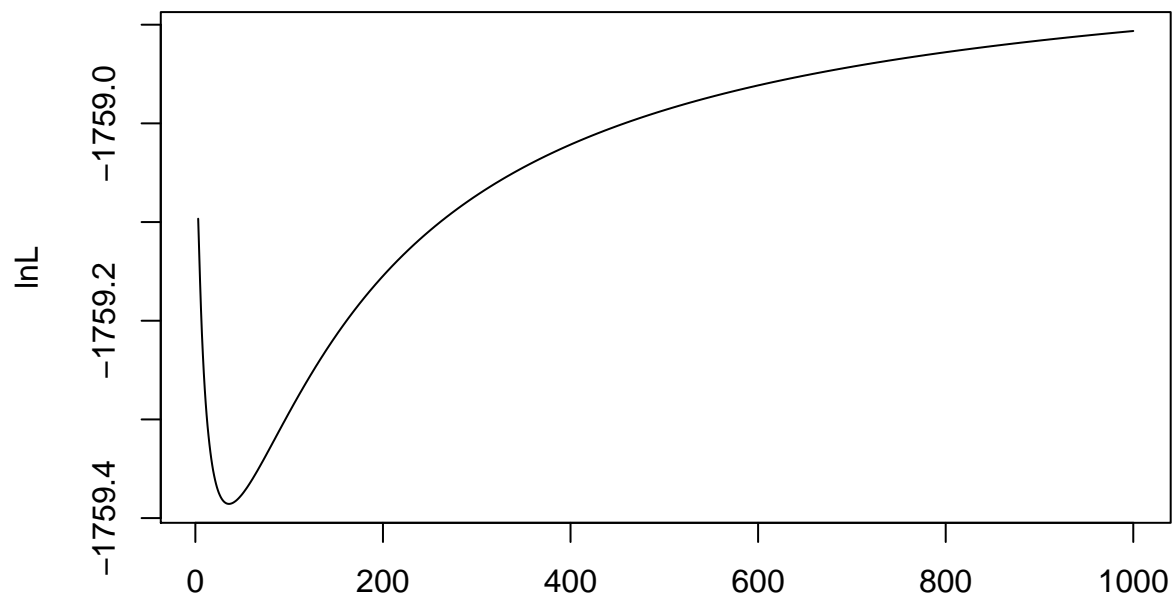
```
## [1] "Tract = 3"
```

**Tract = 3 sim\_2**



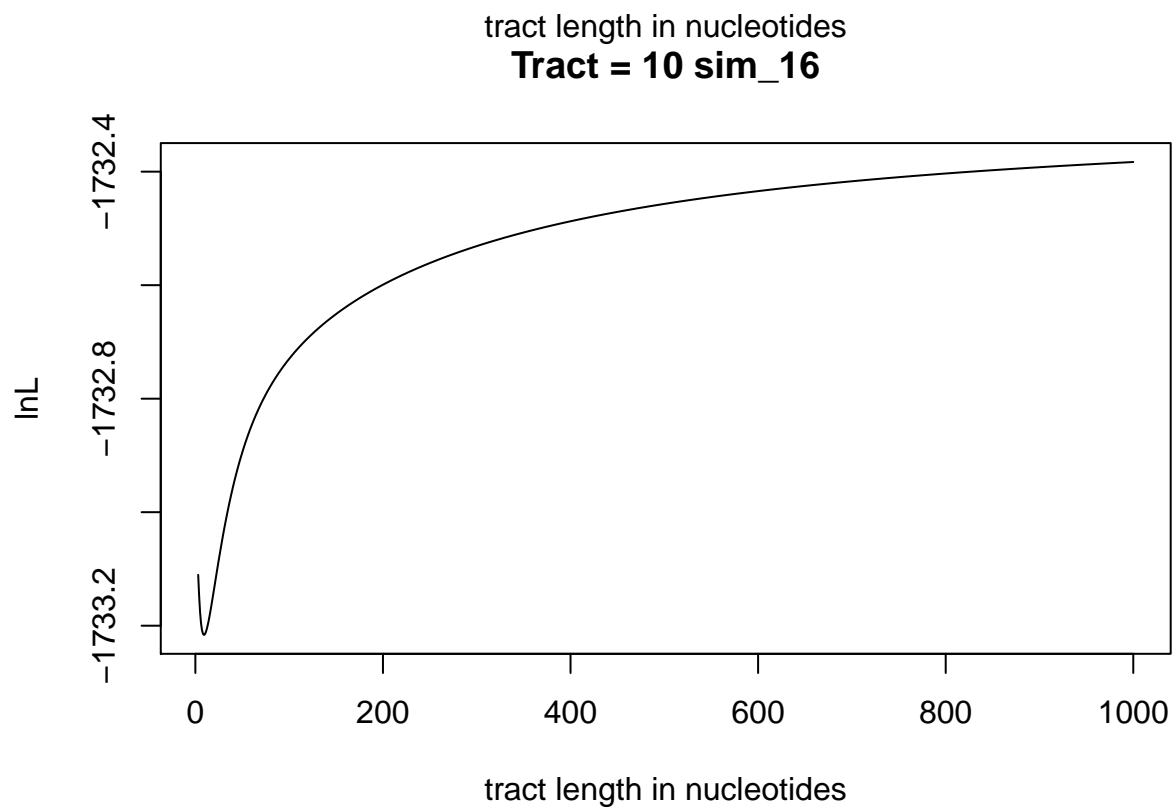
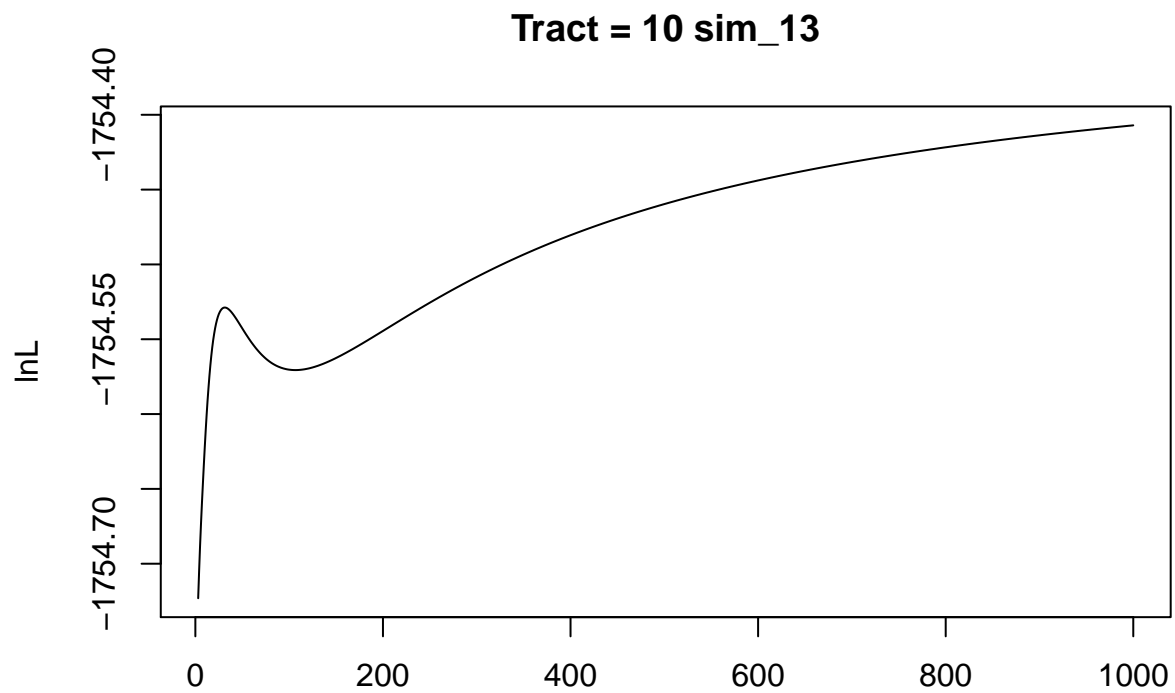
tract length in nucleotides

**Tract = 3 sim\_15**



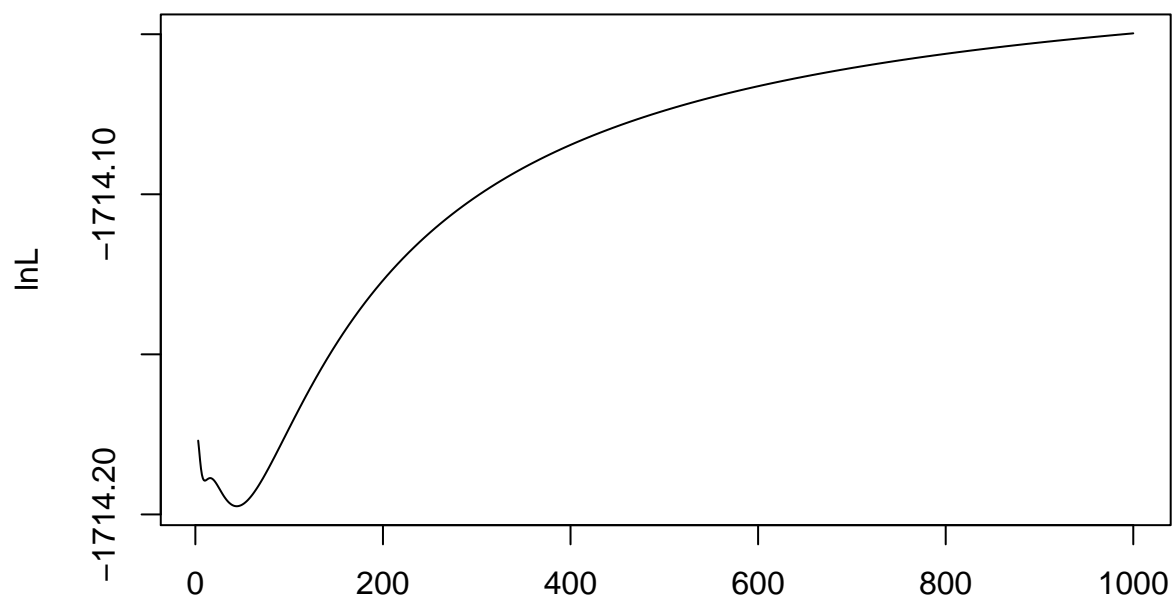
tract length in nucleotides

```
## [1] "Tract = 10"
```



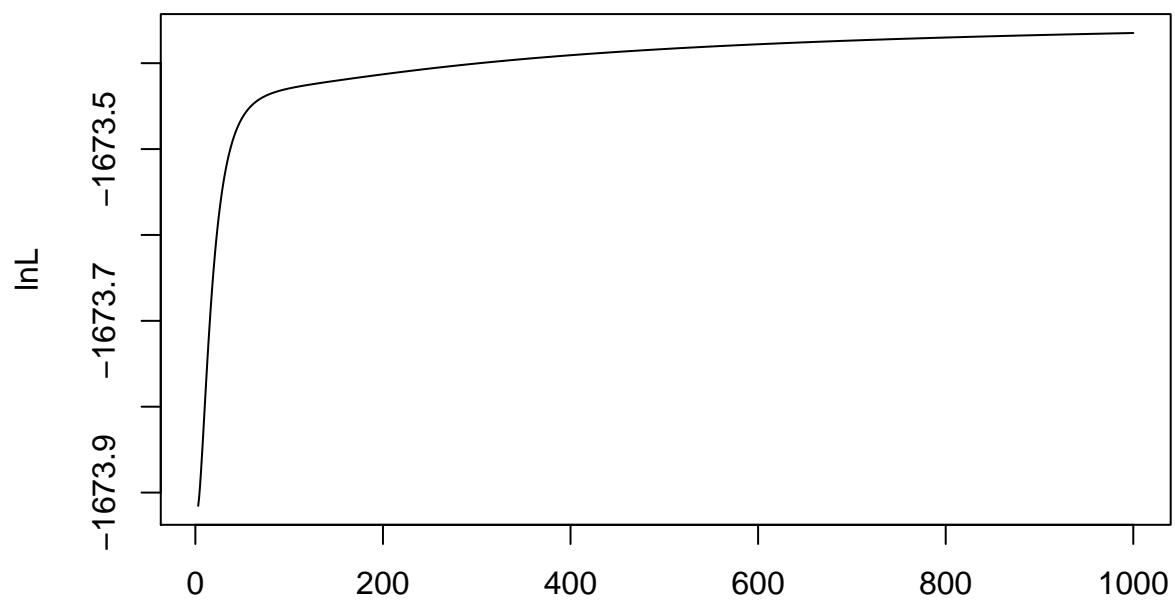
```
## [1] "Tract = 50"
```

**Tract = 50 sim\_2**



tract length in nucleotides

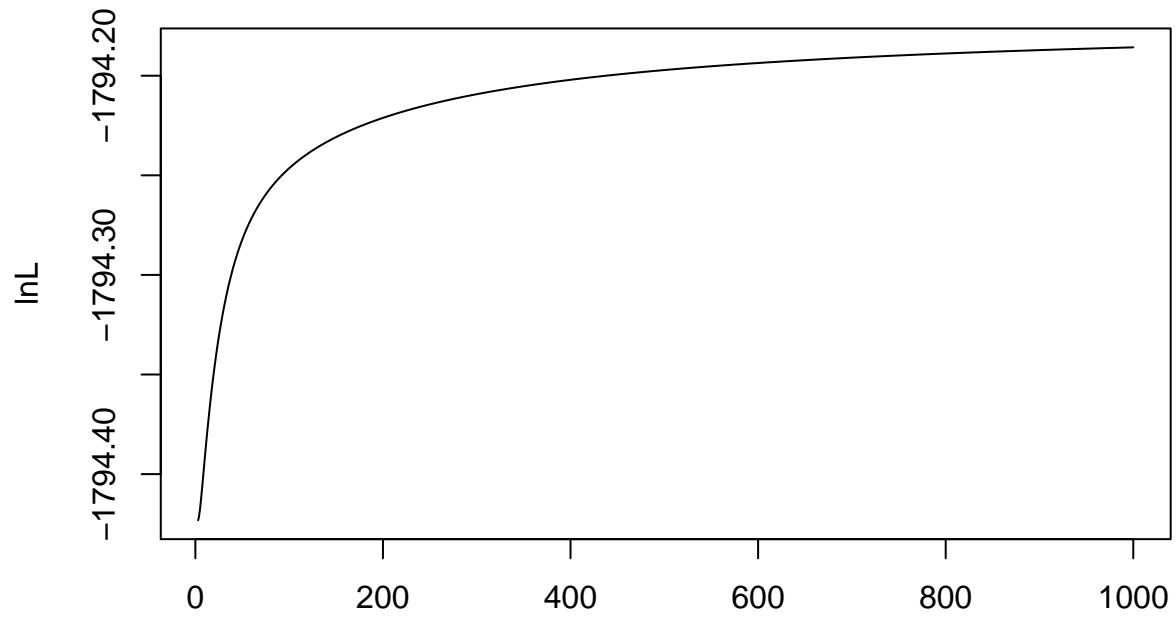
**Tract = 50 sim\_17**



tract length in nucleotides

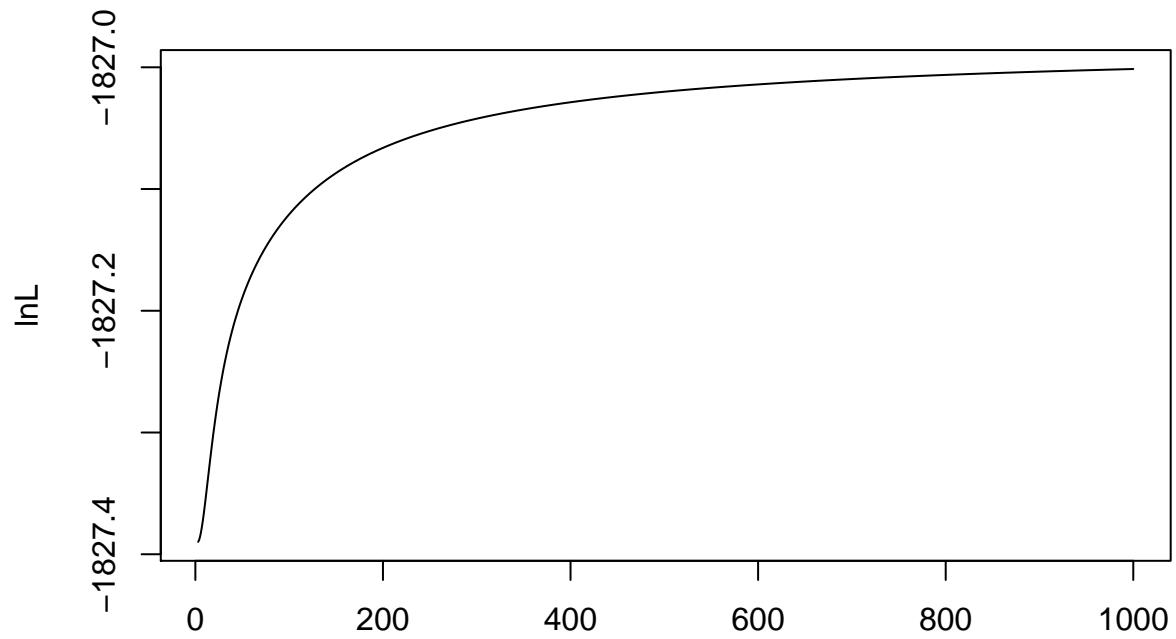
```
## [1] "Tract = 100"
```

**Tract = 100 sim\_13**



tract length in nucleotides

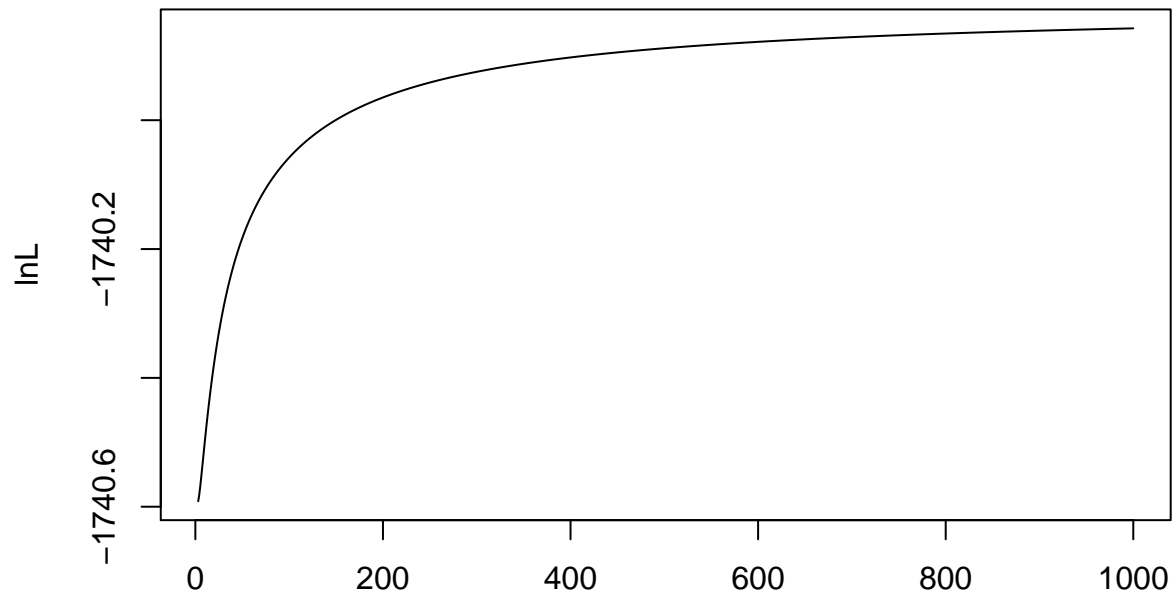
**Tract = 100 sim\_18**



tract length in nucleotides

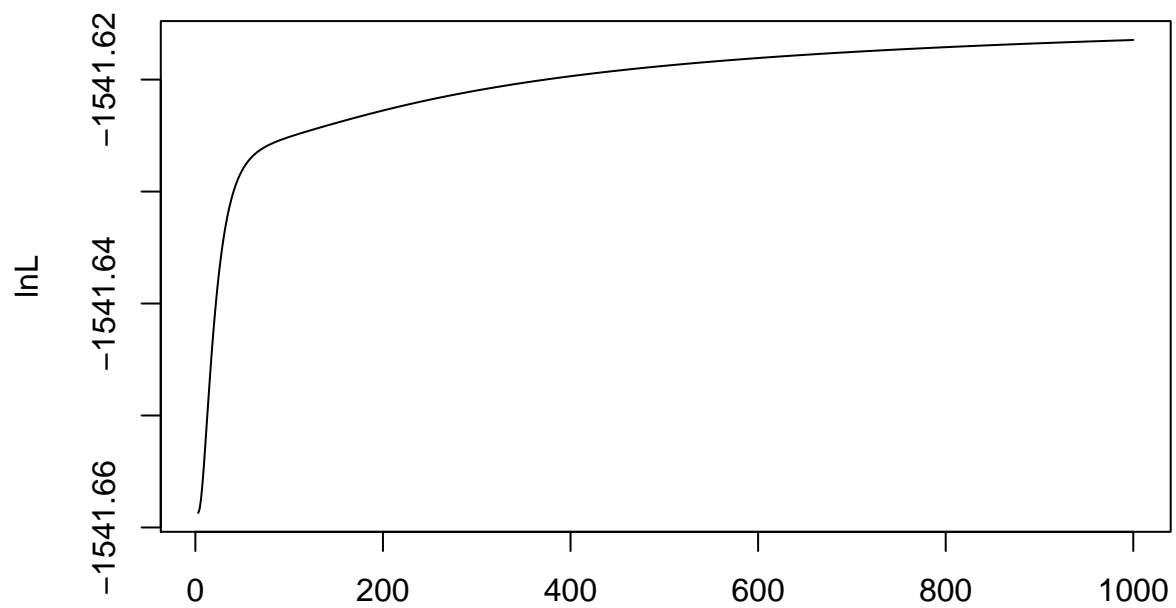
```
## [1] "Tract = 200"
```

**Tract = 200 sim\_1**



tract length in nucleotides

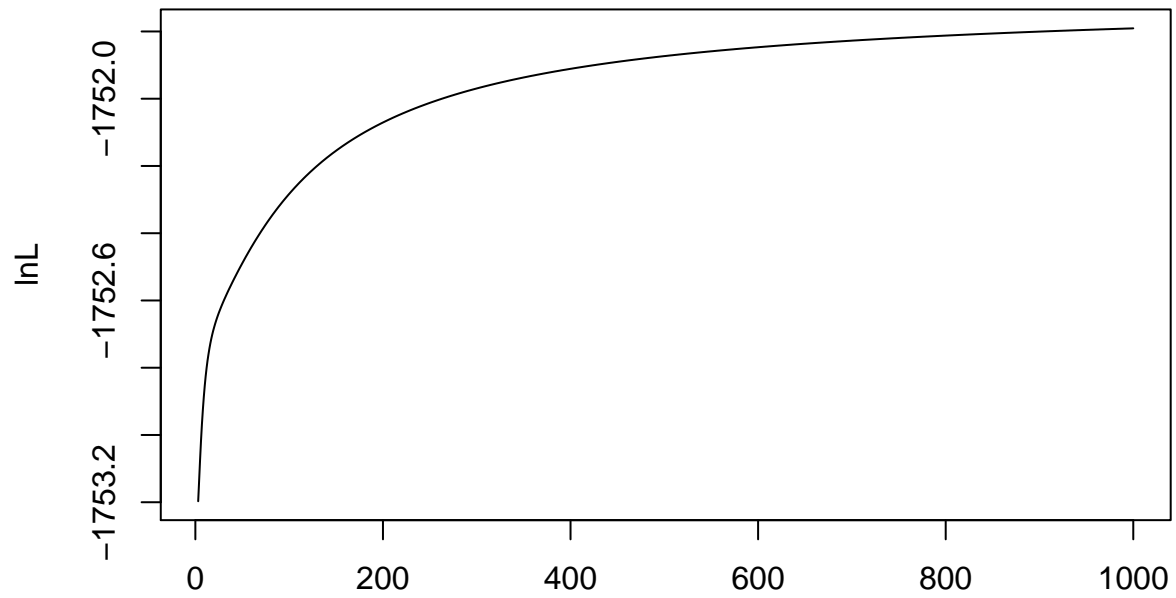
**Tract = 200 sim\_3**



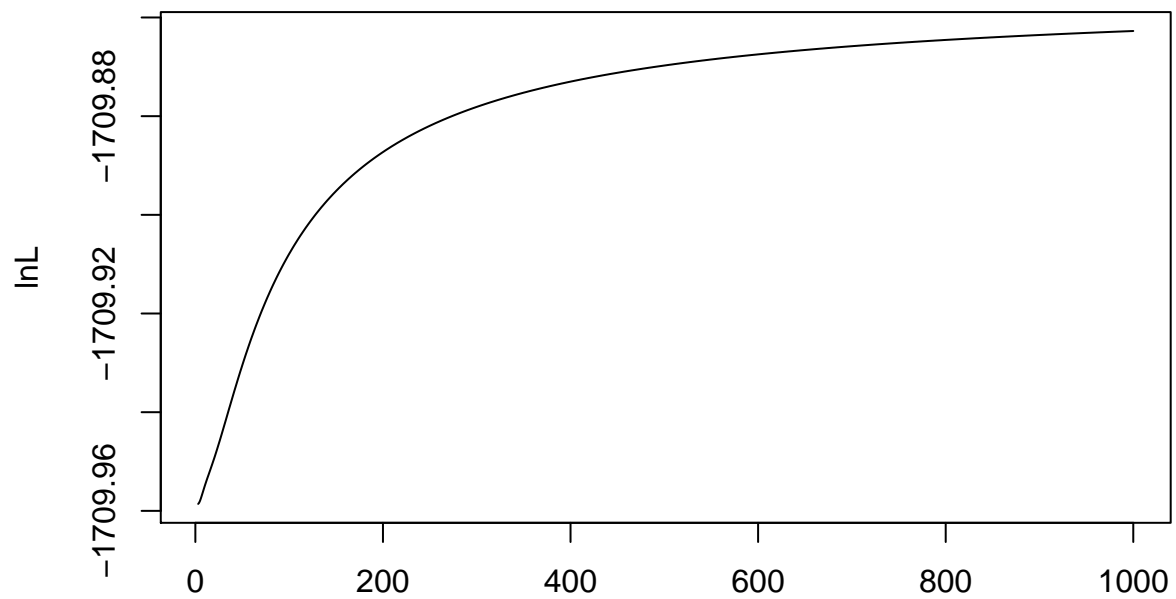
tract length in nucleotides

## [1] "Tract = 300"

**Tract = 300 sim\_6**



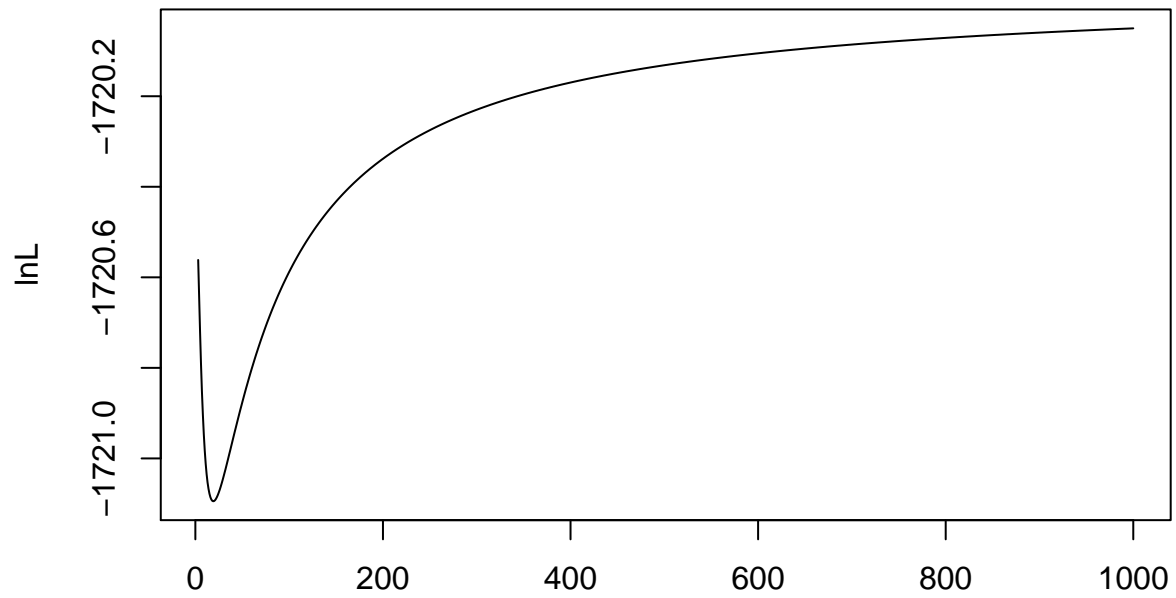
tract length in nucleotides  
**Tract = 300 sim\_13**



tract length in nucleotides

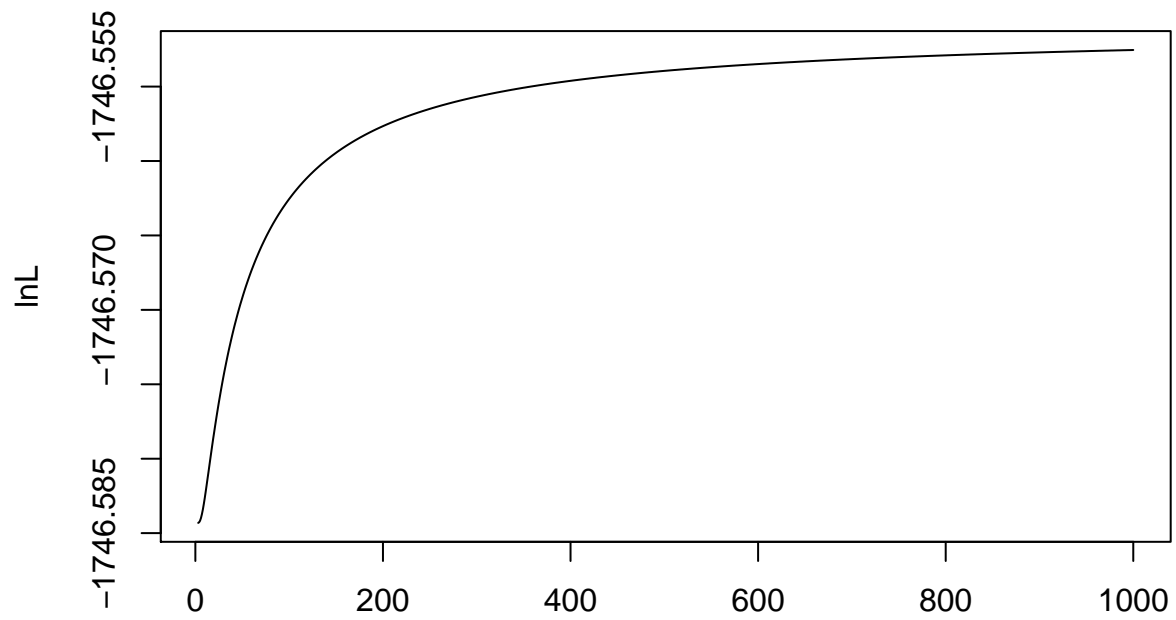
```
## [1] "Tract = 400"
```

**Tract = 400 sim\_1**



tract length in nucleotides

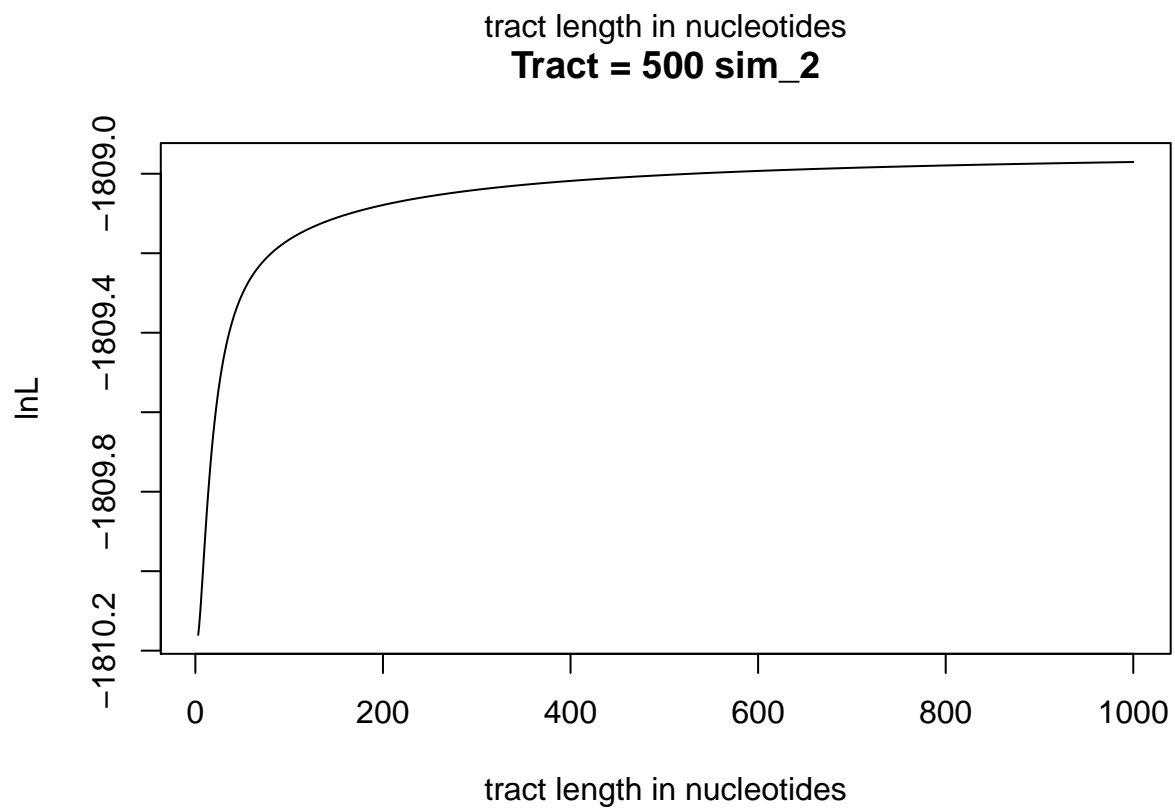
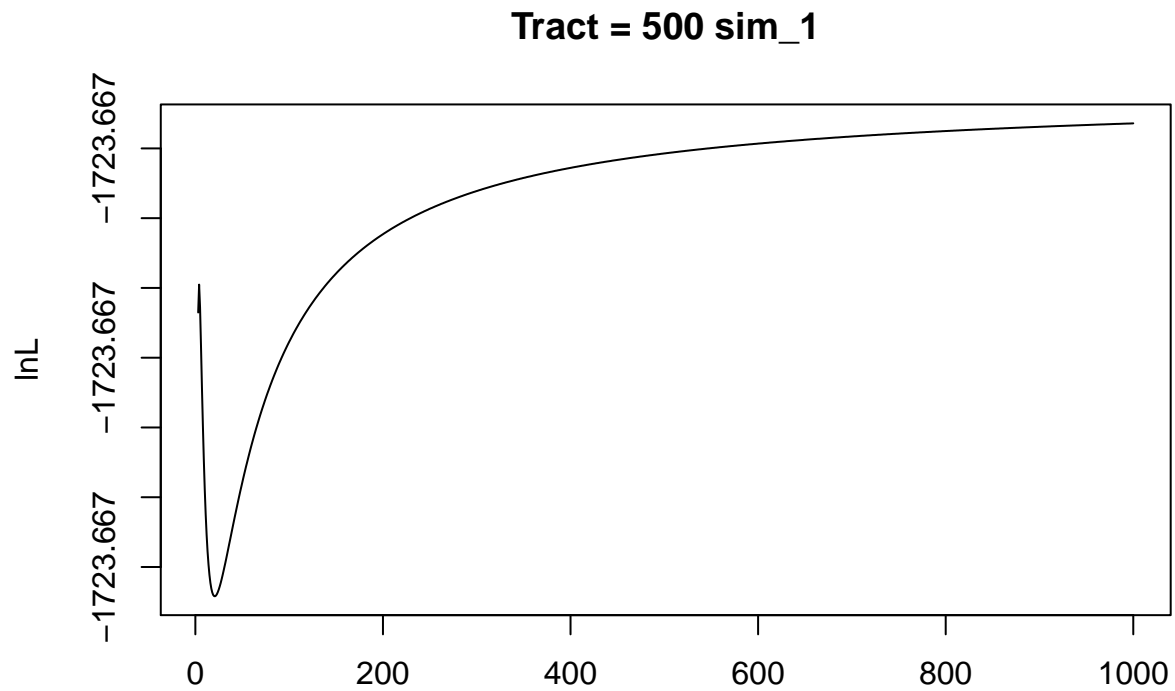
**Tract = 400 sim\_3**



tract length in nucleotides

## [1] "Tract = 500"

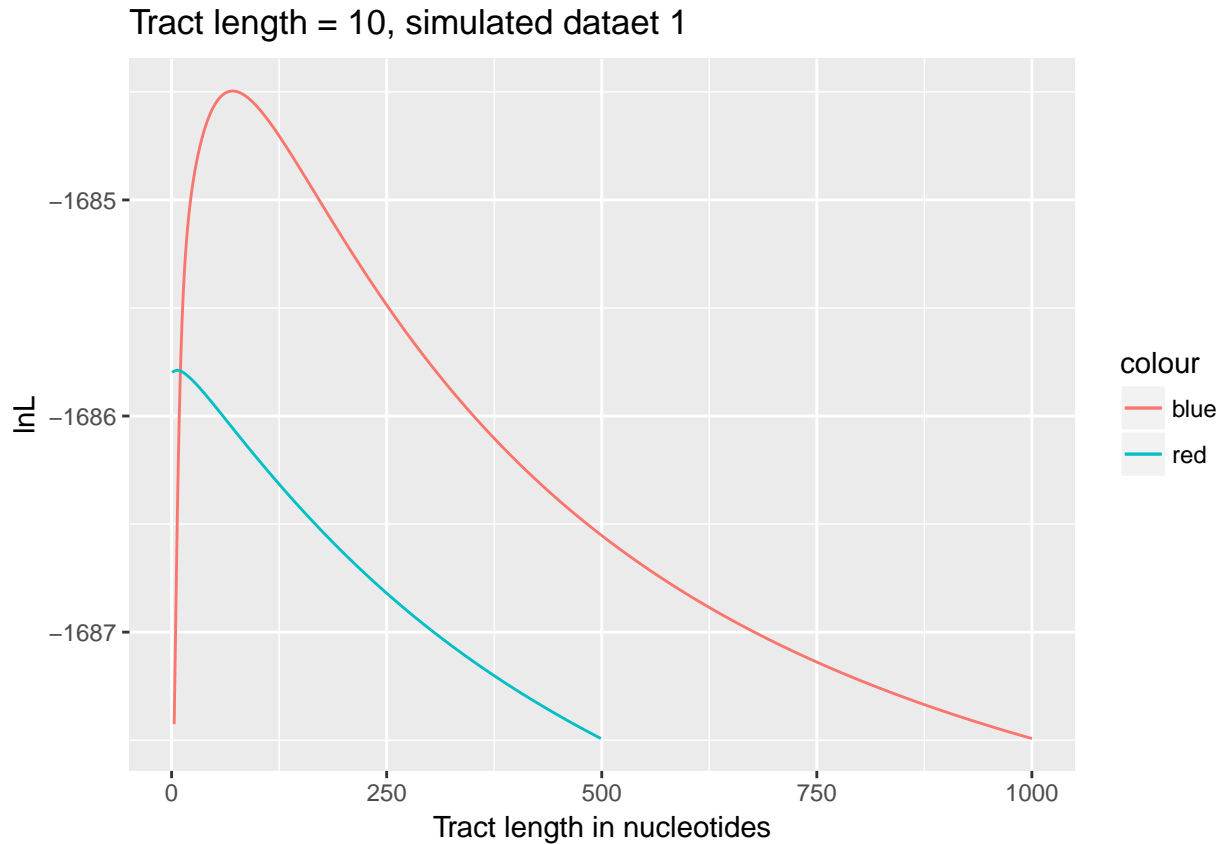




Now plot the two plots of lnL

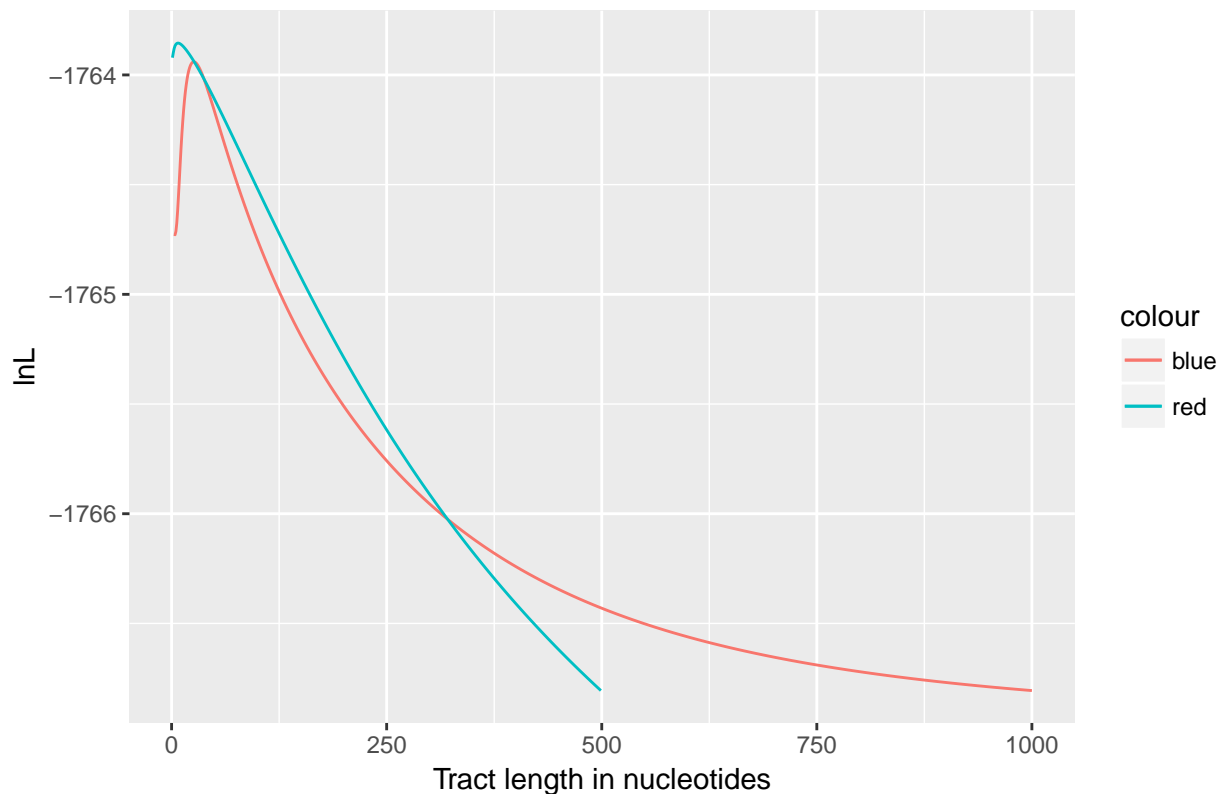
```
library(ggplot2)
# Tract length = 10
hmm.plot <- read.table("./plot/Tract_10.0/sim_1/HMM_YDR418W_YEL054C_lnL_sim_1_1D_surface.txt")
PSJS.plot <- read.table("./plot/Tract_10.0/sim_1/PSJS_HKY_rv_sim_1_Tract_10.0_lnL_1D_surface.txt")
```

```
ggplot(mapping = aes(x = 3.0*exp(-hmm.plot[,1]), y = hmm.plot[, 2], colour = "blue")) + geom_line() +
  geom_line(aes(x = exp(-PSJS.plot[,1]), y = PSJS.plot[, 2]/488 + min(hmm.plot[, 2]) - min(PSJS.plot[, 2],
    colour = "red")) +
  xlab("Tract length in nucleotides") +
  ylab("lnL") +
  ggtitle("Tract length = 10, simulated dataet 1")
```



```
hmm.plot <- read.table("./plot/Tract_10.0/sim_100/HMM_YDR418W_YELO54C_lnL_sim_100_1D_surface.txt")
PSJS.plot <- read.table("./plot/Tract_10.0/sim_100/PSJS_HKY_rv_sim_100_Tract_10.0_lnL_1D_surface.txt")
ggplot(mapping = aes(x = 3.0*exp(-hmm.plot[,1]), y = hmm.plot[, 2], colour = "blue")) + geom_line() +
  geom_line(aes(x = exp(-PSJS.plot[,1]), y = PSJS.plot[, 2]/488 + min(hmm.plot[, 2]) - min(PSJS.plot[, 2],
    colour = "red")) +
  xlab("Tract length in nucleotides") +
  ylab("lnL") +
  ggtitle("Tract length = 10, simulated dataet 100")
```

Tract length = 10, simulated dataet 100



Now see how estimates from the two approaches differ from the actual mean tract length in each simulated data set.

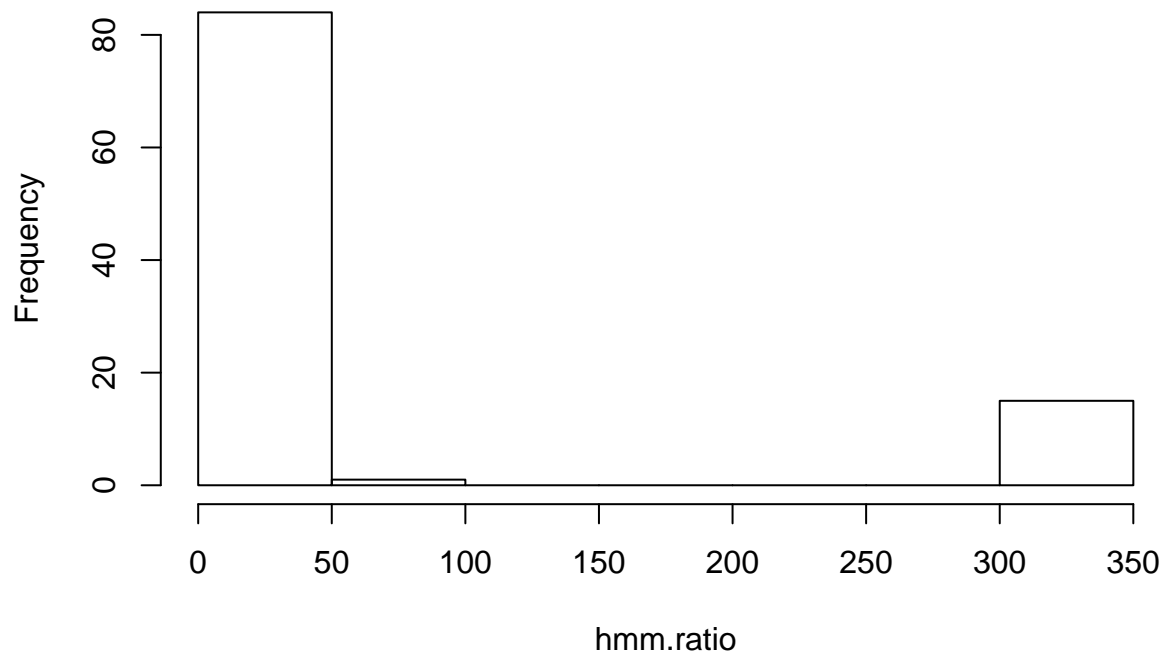
```
for(tract in Tract.list){
  sim.info <- get(paste("sim.tract.", toString(tract), sep = ""))
  # Show mean and sd
  print(c("empirical mean", mean(sim.info["mean tract length", ]),
    "geometric mean", tract,
    "empirical sd", mean(sim.info["sd tract length", ], na.rm = TRUE),
    "geometric sd", sqrt(tract^2-tract*3.0)))
  hmm.info <- get(paste("HMM_Tract_", toString(tract), "_plot", sep = ""))
  PSJS.info <- get(paste("PSJS_Tract_", toString(tract), "_summary", sep = ""))
  shared.col <- intersect(colnames(hmm.info), colnames(PSJS.info))

  # Now show the ratio of HMM estimated tract / actual mean tracts in simulation
  hmm.ratio <- hmm.info[1, shared.col]/sim.info[1, shared.col]
  hist(hmm.ratio, main = paste("HMM ratio Tract = ", toString(tract), sep = ""))
  print(c("HMM mean", mean(hmm.ratio), "HMM sd", sd(hmm.ratio)))

  # Now show the ratio of PSJS estimated tract / actual mean tracts in simulation
  PSJS.ratio <- PSJS.info["tract_length", shared.col]/sim.info[1, shared.col]
  hist(PSJS.ratio, main = paste("PSJS ratio Tract = ", toString(tract), sep = ""))
  print(c("PSJS mean", mean(PSJS.ratio), "PSJS sd", sd(PSJS.ratio)))
}
```

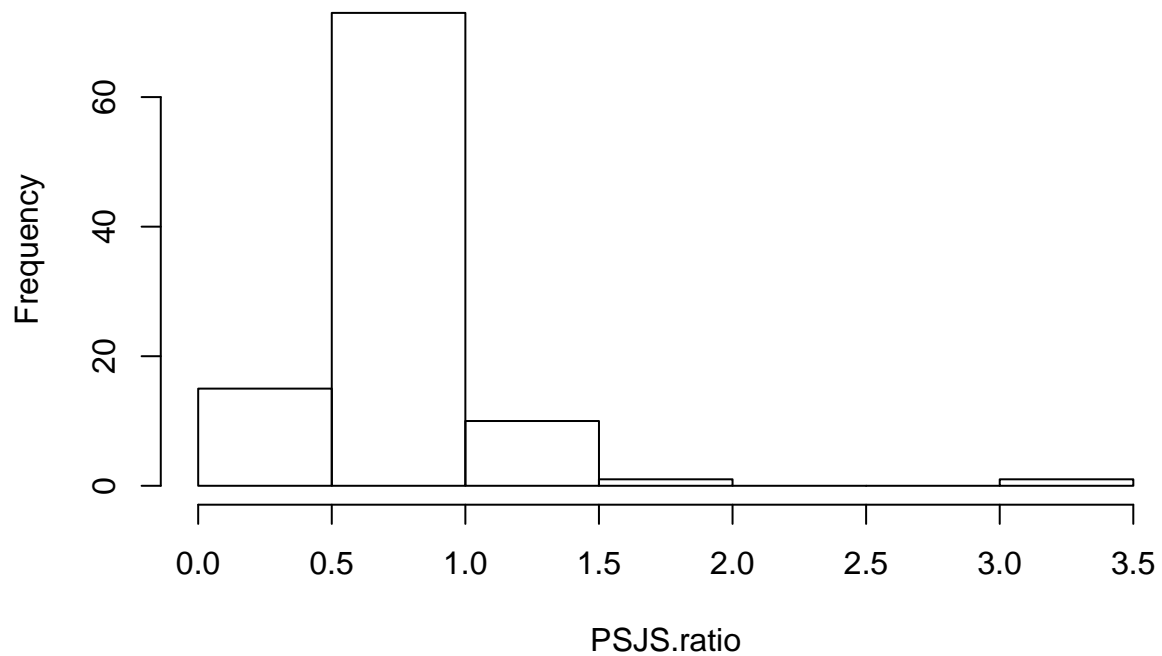
```
## [1] "empirical mean" "3"          "geometric mean" "3"
## [5] "empirical sd"   "0"          "geometric sd"   "0"
```

### HMM ratio Tract = 3



```
## [1] "HMM mean"          "53.3133333331312" "HMM sd"  
## [4] "118.442841347853"
```

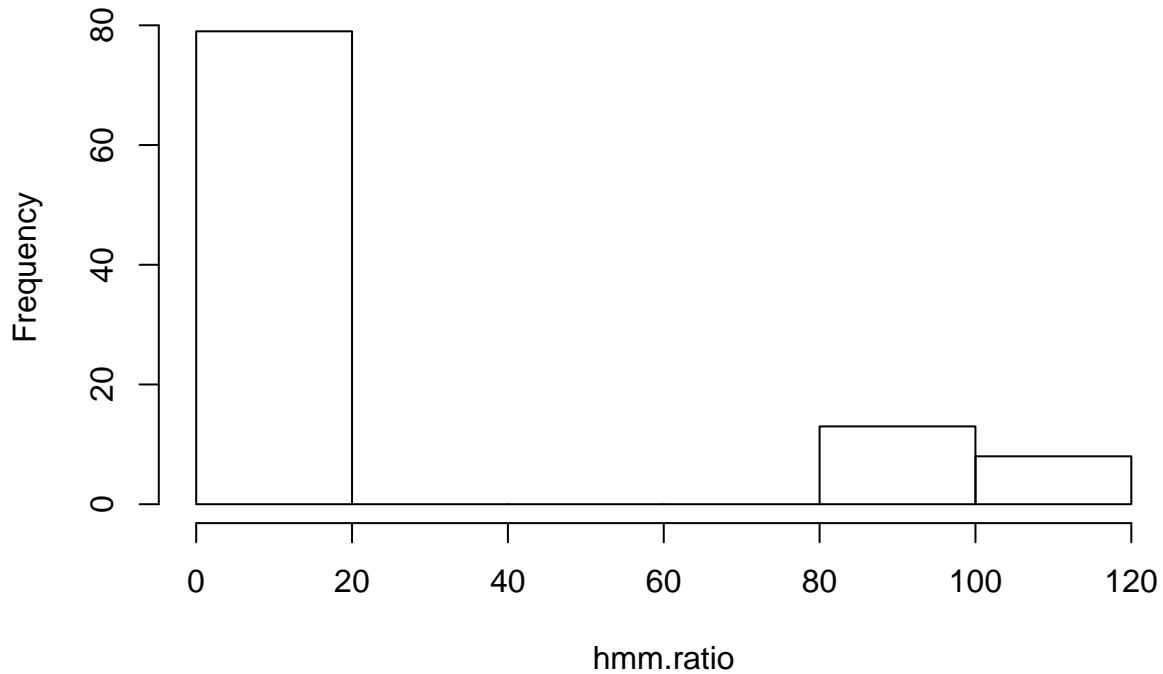
### PSJS ratio Tract = 3



```
## [1] "PSJS mean"          "0.765726712026039" "PSJS sd"  
## [4] "0.342915119293345"  
## [1] "empirical mean"     "10.0091602931023" "geometric mean"
```

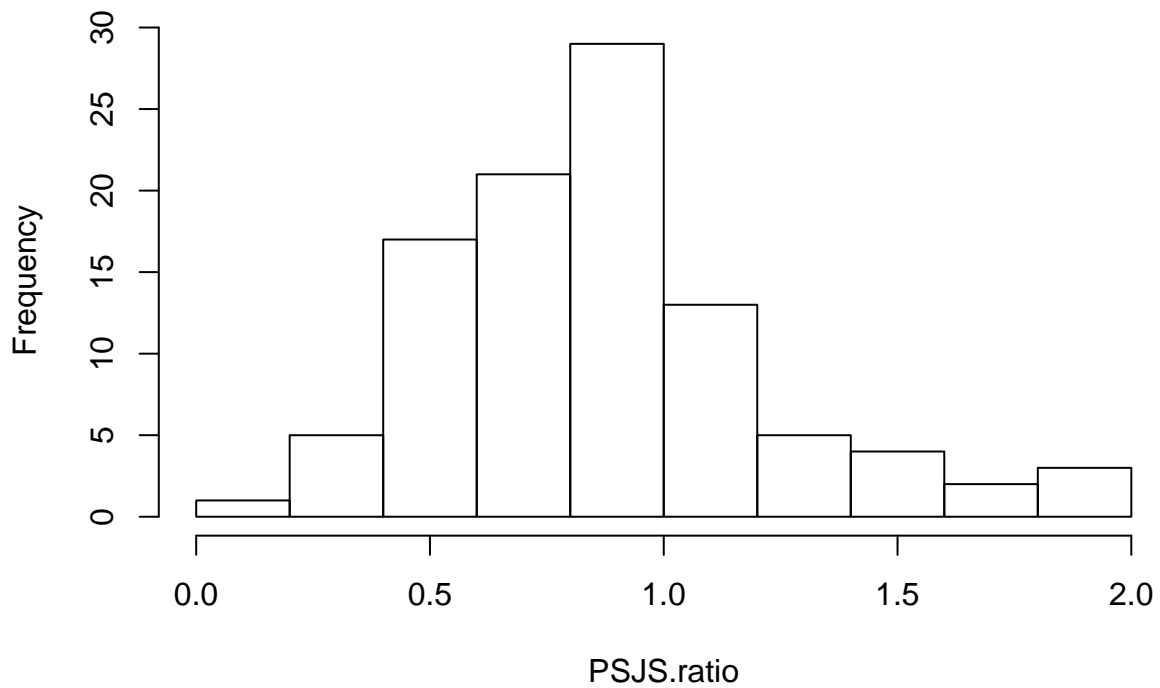
```
## [4] "10" "empirical sd" "8.37472754557179"  
## [7] "geometric sd" "8.36660026534076"
```

### HMM ratio Tract = 10



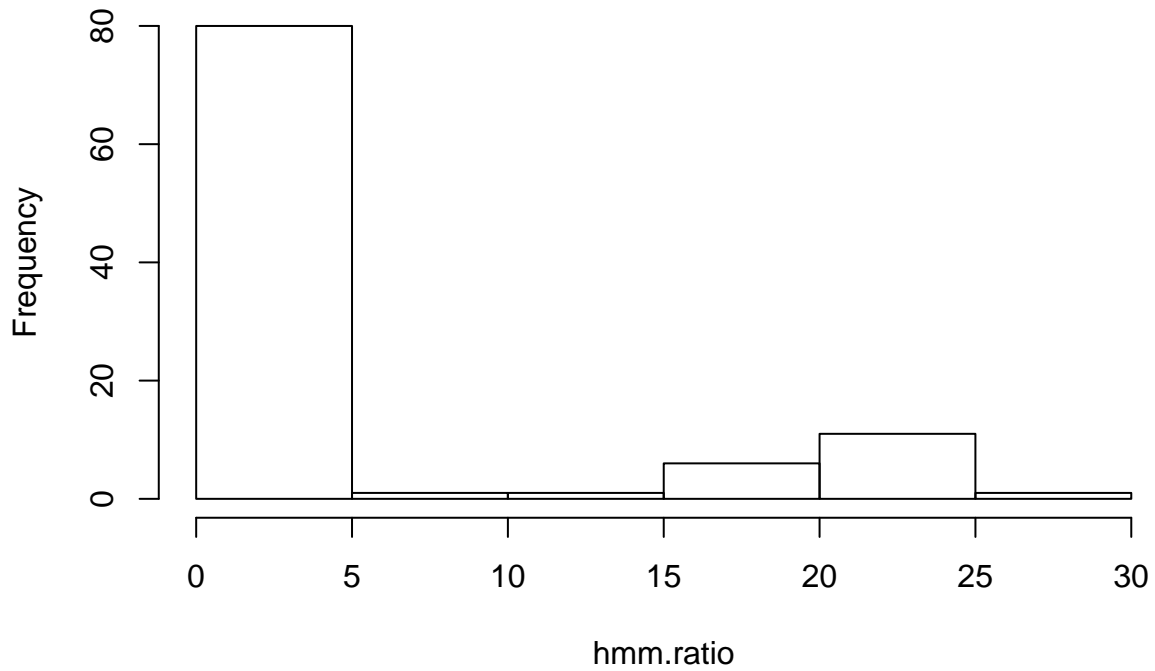
```
## [1] "HMM mean" "22.8463240019266" "HMM sd"  
## [4] "39.1183435578928"
```

### PSJS ratio Tract = 10



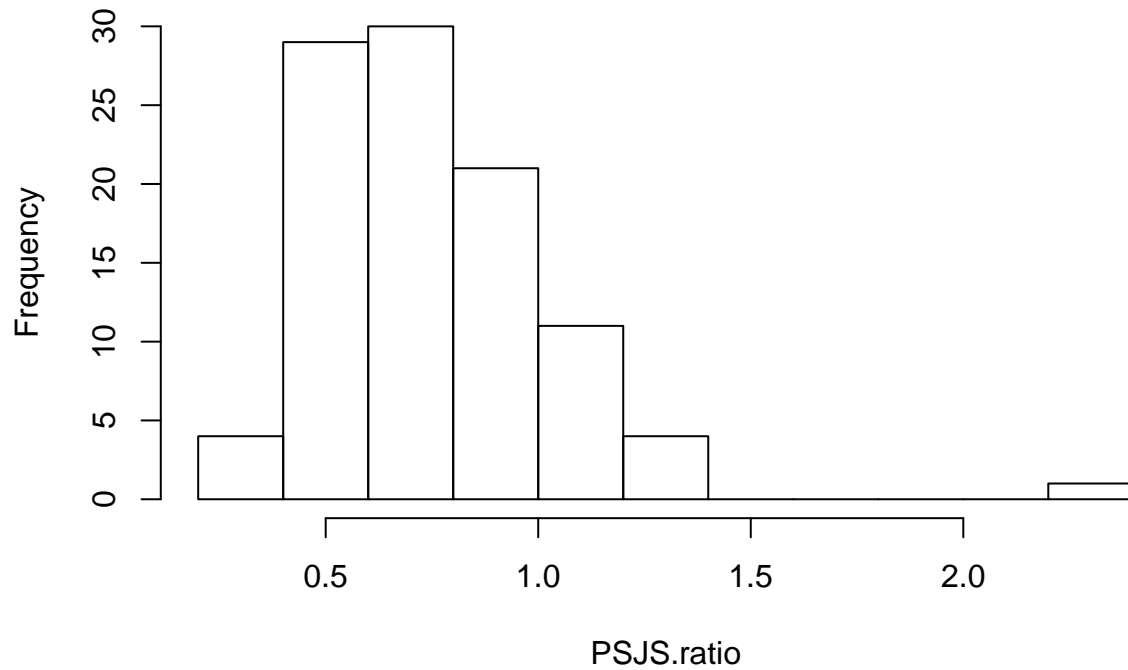
```
## [1] "PSJS mean"          "0.867271896676903" "PSJS sd"
## [4] "0.350196556681162"
## [1] "empirical mean"     "50.5630334377432"  "geometric mean"
## [4] "50"                 "empirical sd"       "48.5906932965629"
## [7] "geometric sd"       "48.4767985741633"
```

### HMM ratio Tract = 50



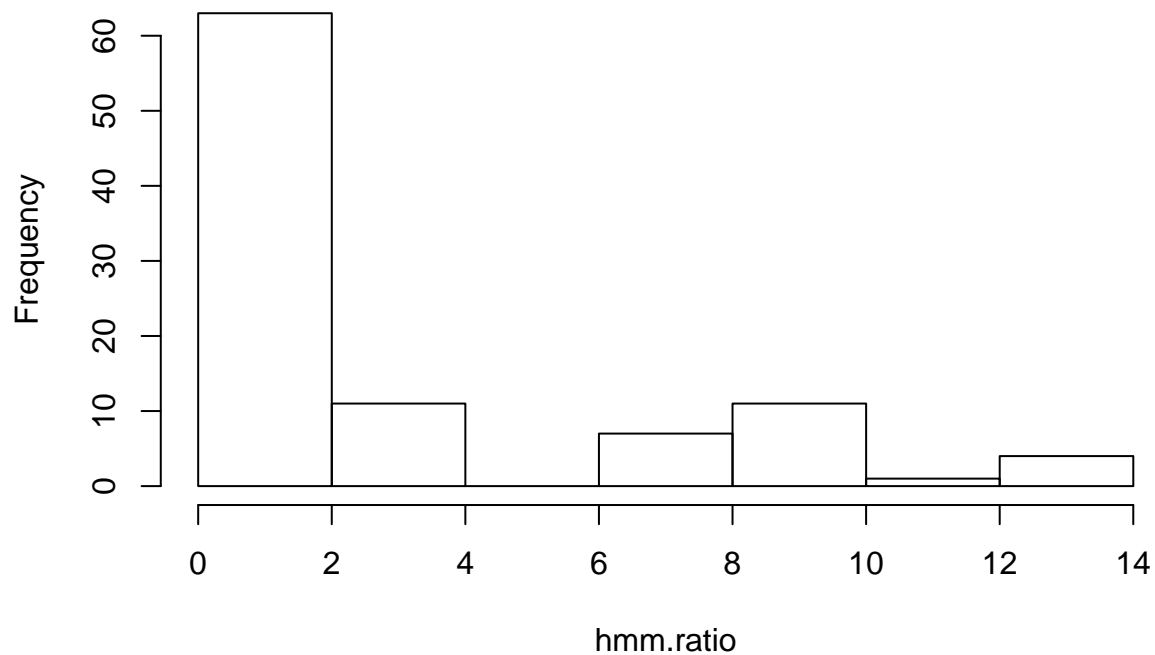
```
## [1] "HMM mean"          "5.02521520761096" "HMM sd"
## [4] "7.81334636408309"
```

### PSJS ratio Tract = 50



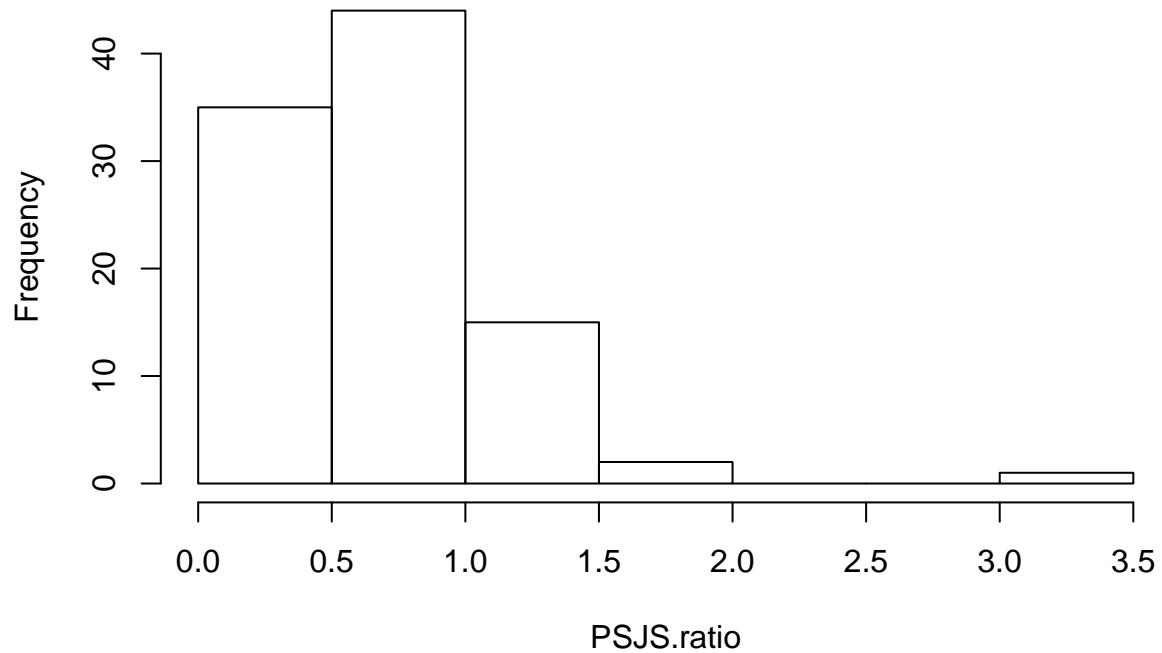
```
## [1] "PSJS mean"          "0.75123941605822"  "PSJS sd"
## [4] "0.283108335026374"
## [1] "empirical mean"     "98.5552302358531"  "geometric mean"
## [4] "100"                "empirical sd"       "90.6703458252144"
## [7] "geometric sd"       "98.488578017961"
```

### HMM ratio Tract = 100



```
## [1] "HMM mean"          "3.15364468735454" "HMM sd"
## [4] "3.65107531734578"
```

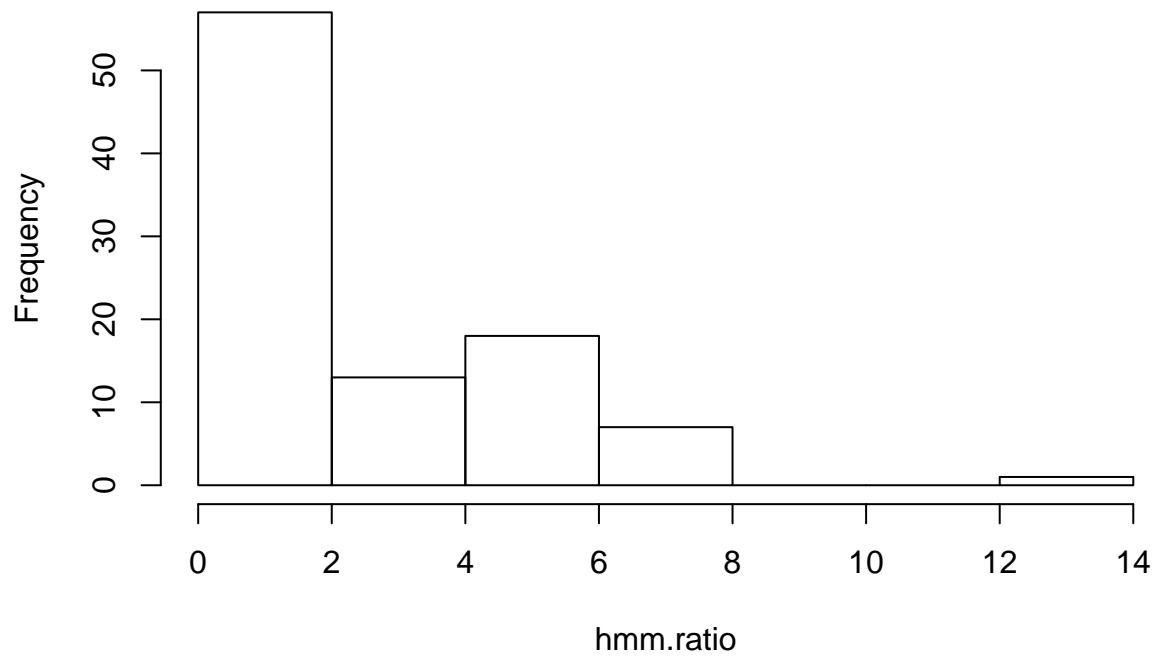
### PSJS ratio Tract = 100



```
## [1] "PSJS mean"          "0.705828559489391" "PSJS sd"
## [4] "0.451338788199344"
## [1] "empirical mean"     "198.259889997902" "geometric mean"
## [4] "200"                "empirical sd"      "181.444635451734"
## [7] "geometric sd"       "198.494332412792"
```

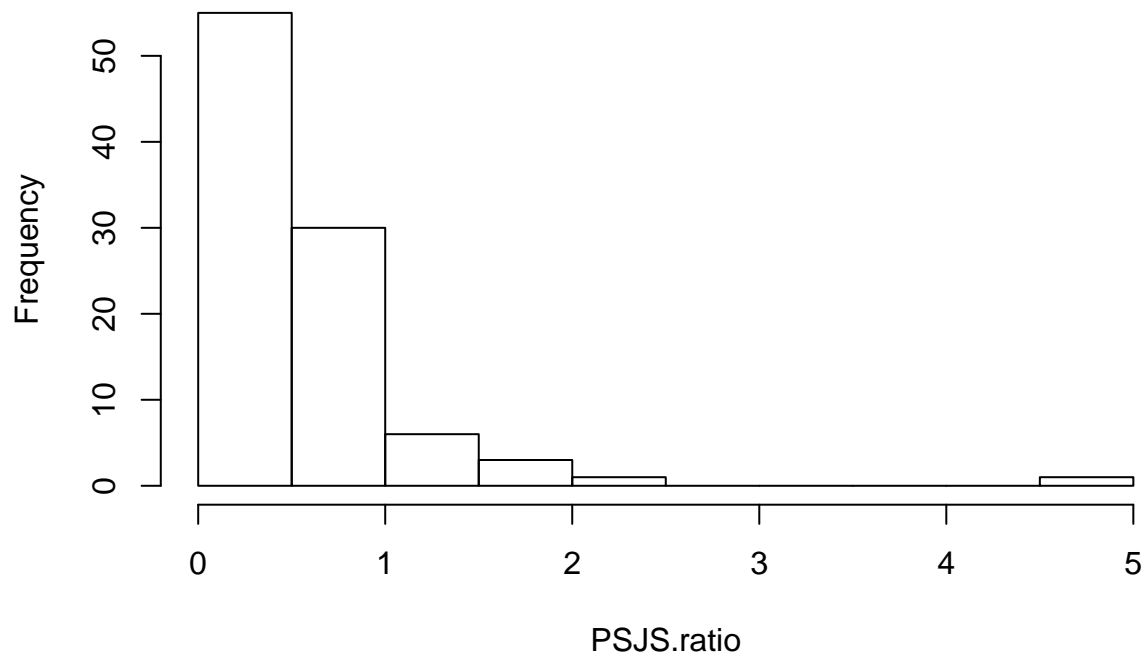


### HMM ratio Tract = 200



```
## [1] "HMM mean"          "2.25114401983335" "HMM sd"
## [4] "2.46929176842981"
```

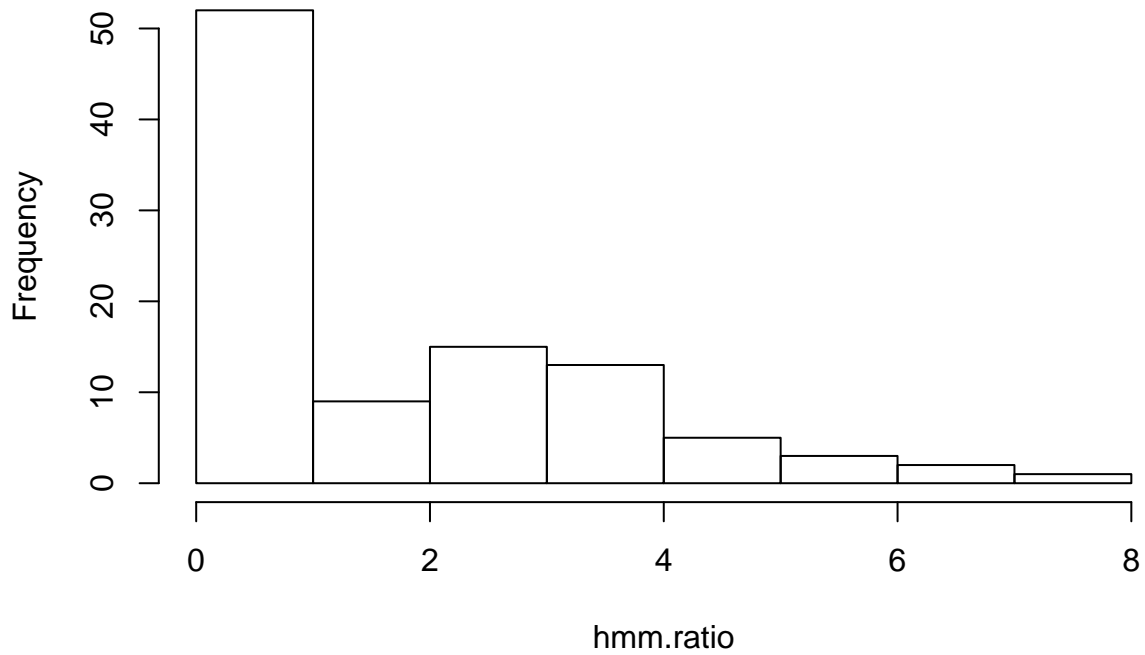
### PSJS ratio Tract = 200



```
## [1] "PSJS mean"          "0.594669364732503" "PSJS sd"
## [4] "0.594057156992222"
## [1] "empirical mean"     "300.306681568432" "geometric mean"
```

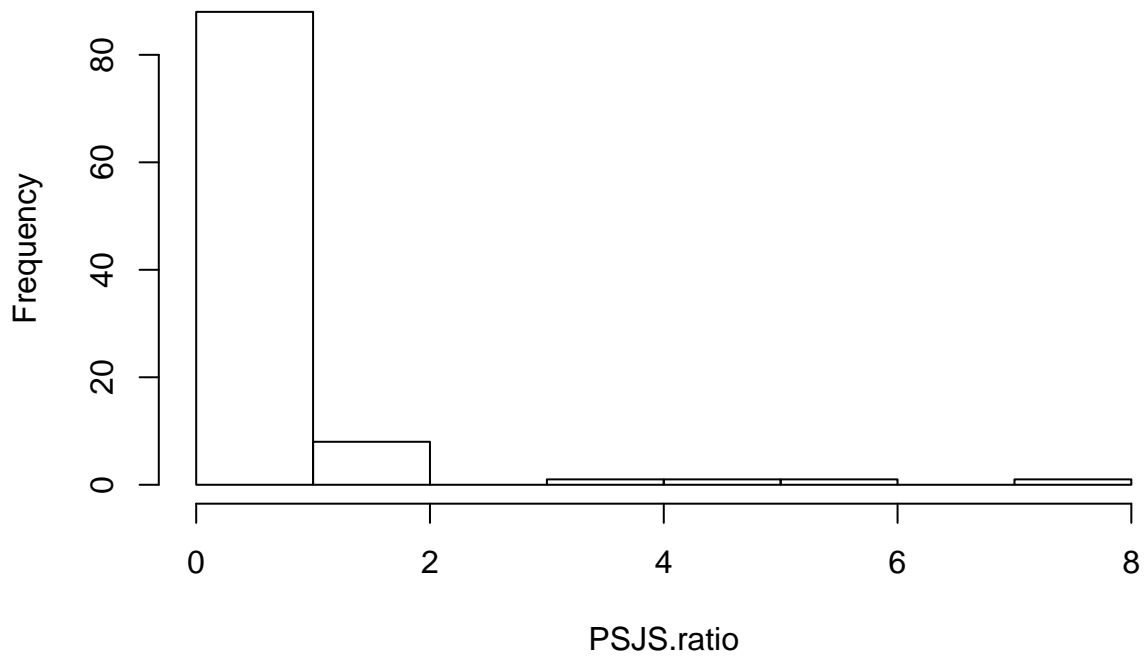
```
## [4] "300" "empirical sd" "270.918634615386"  
## [7] "geometric sd" "298.496231131986"
```

### HMM ratio Tract = 300



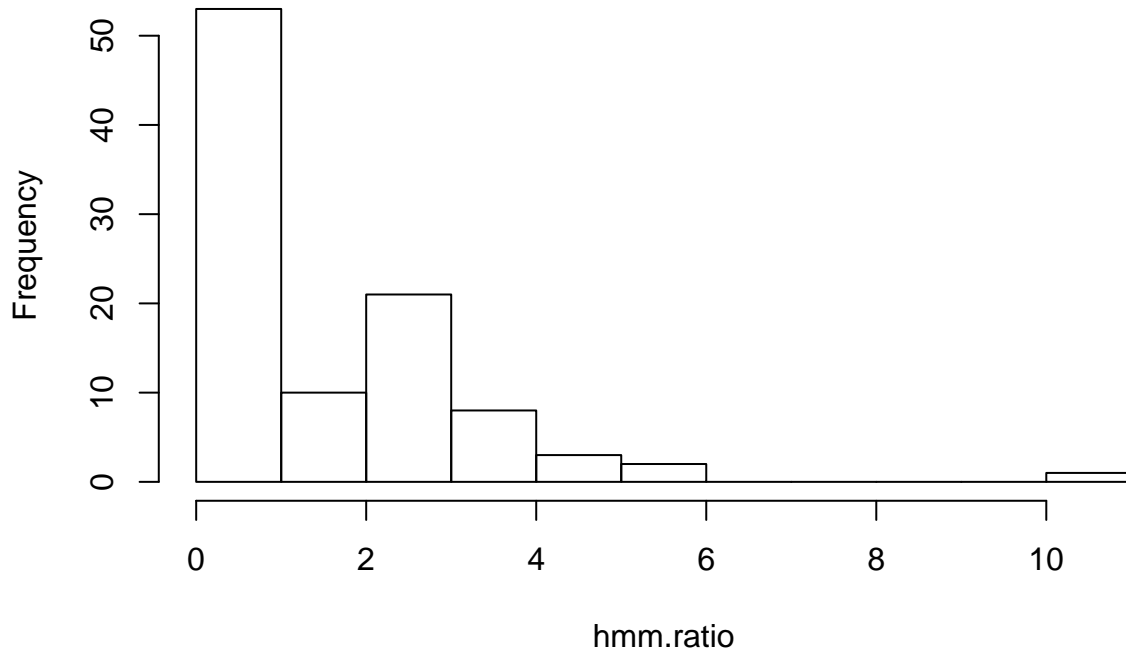
```
## [1] "HMM mean" "1.73254267738278" "HMM sd"  
## [4] "1.7636542177027"
```

### PSJS ratio Tract = 300



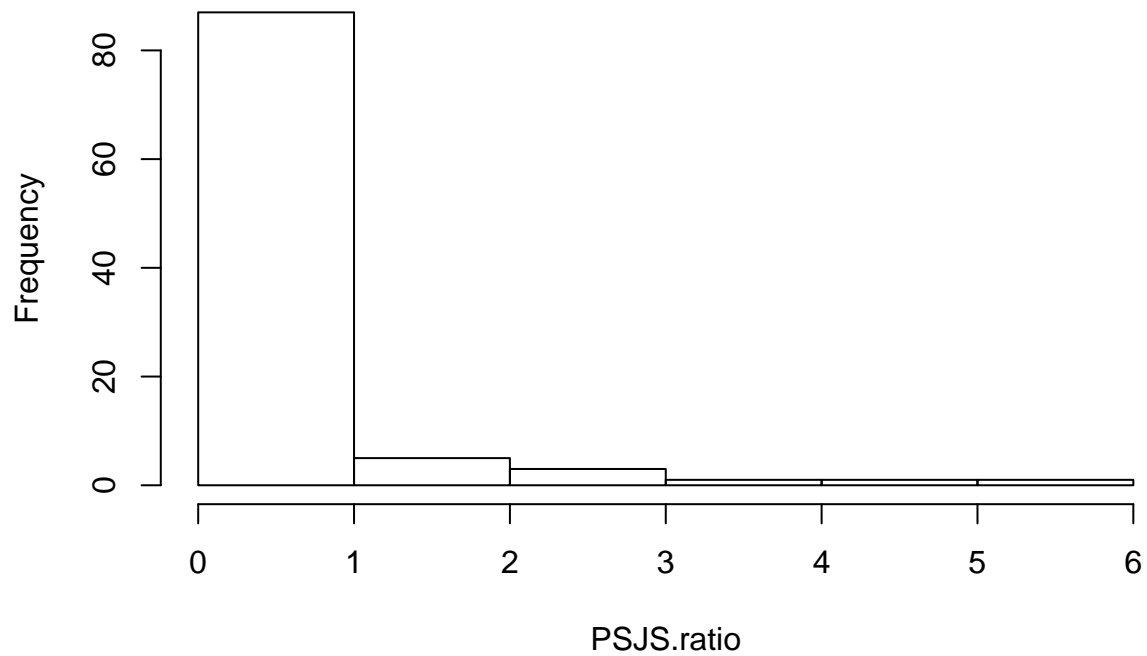
```
## [1] "PSJS mean"          "0.602489007078694" "PSJS sd"
## [4] "1.01256919318485"
## [1] "empirical mean"     "390.721396103896" "geometric mean"
## [4] "400"               "empirical sd"      "349.193048885713"
## [7] "geometric sd"       "398.497176903425"
```

### HMM ratio Tract = 400



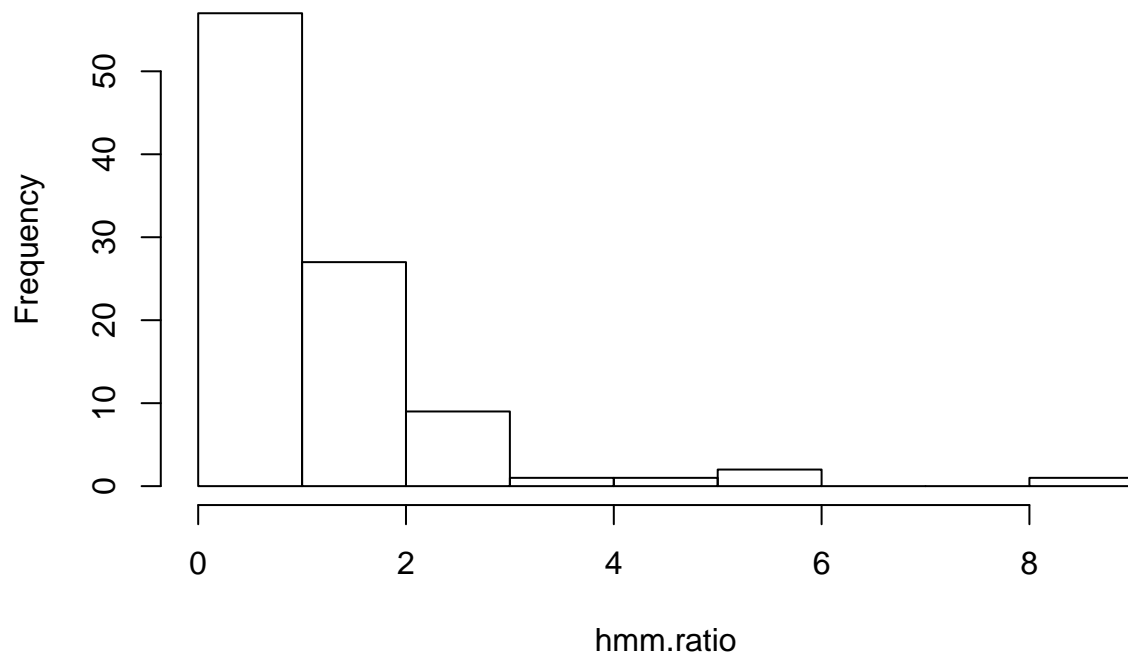
```
## [1] "HMM mean"          "1.45661938280292" "HMM sd"
## [4] "1.67860989325479"
```

### PSJS ratio Tract = 400



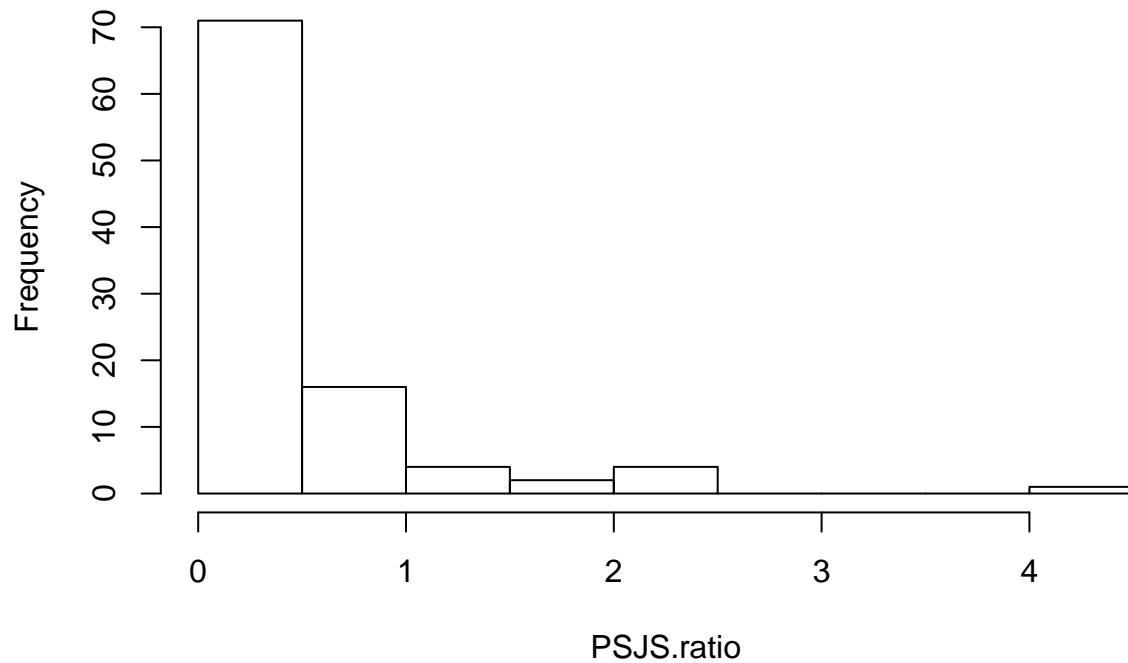
```
## [1] "PSJS mean"          "0.544780348828959" "PSJS sd"
## [4] "0.863921114745281"
## [1] "empirical mean"     "505.424964285714"  "geometric mean"
## [4] "500"                "empirical sd"       "435.966931089703"
## [7] "geometric sd"       "498.497743224581"
```

### HMM ratio Tract = 500



```
## [1] "HMM mean"          "1.11455605873148" "HMM sd"  
## [4] "1.38547997394801"
```

### PSJS ratio Tract = 500



```
## [1] "PSJS mean"          "0.477727688735254" "PSJS sd"  
## [4] "0.668483064587436"
```

save workspace now

```
save.image("./SimulationStudy.RData")
```