# age_distribution_of_cases
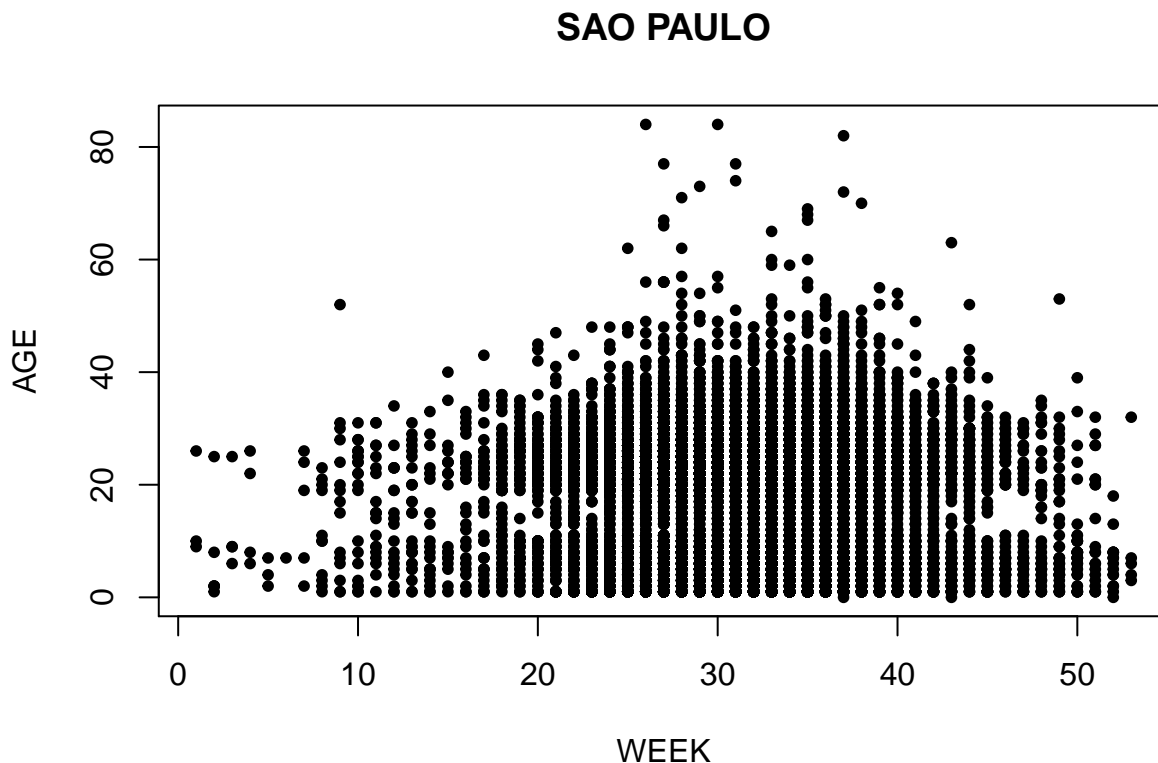
*Maggie Walters*

*June 22, 2017*

## Document Synopsis

This document aims to visually and quantitatively describe the different age distribution of cases in the counties represented in spm.data.long.csv. Densities of cases for age groups will be calculated and then compared between counties.

The difference between counties will be used in order to explore possible predictive relationships between population size or urbanness of a county and it's density of cases in different age classes.
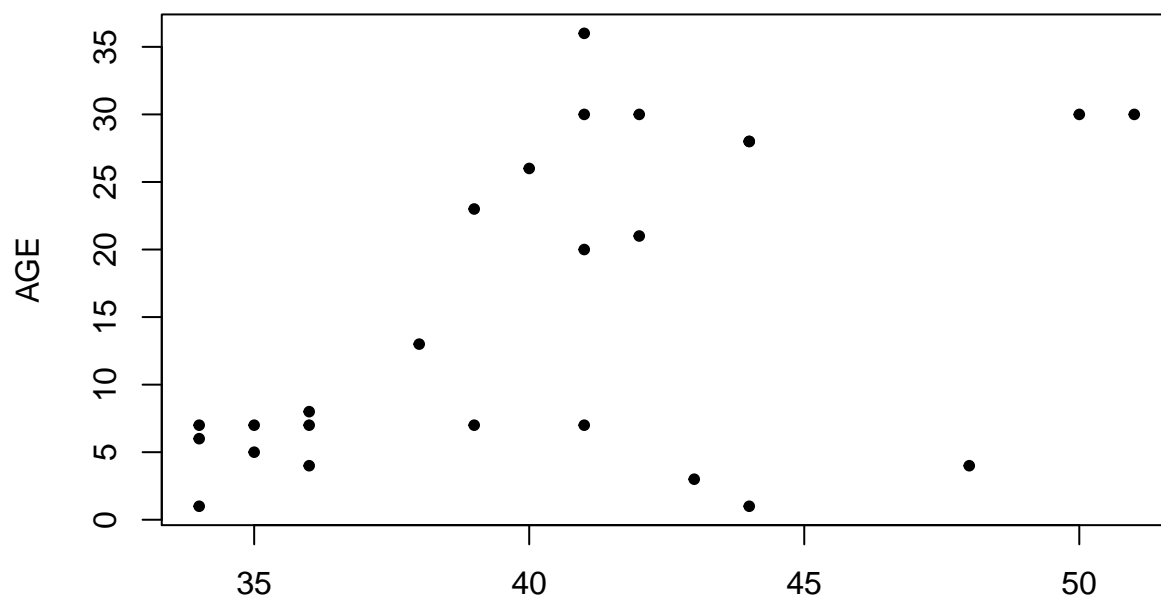
## Age distribution scatter plots by county

Scatter plots will be created in order to visually identify the age groups which frequently have a low density of cases. Depending on results, this may be limited to counties with number of cases above a certain threshold.

```
for(i in 1:length(county_vec)){
  subset_county <- subset(data, data$COUNTY == county_vec[i])
  if(length(subset_county$AGE) > 20){
    plot(subset_county$WEEK, subset_county$AGE, pch = 20,
         main = county_vec[i], xlab = "WEEK", ylab = "AGE")
  }
}
```
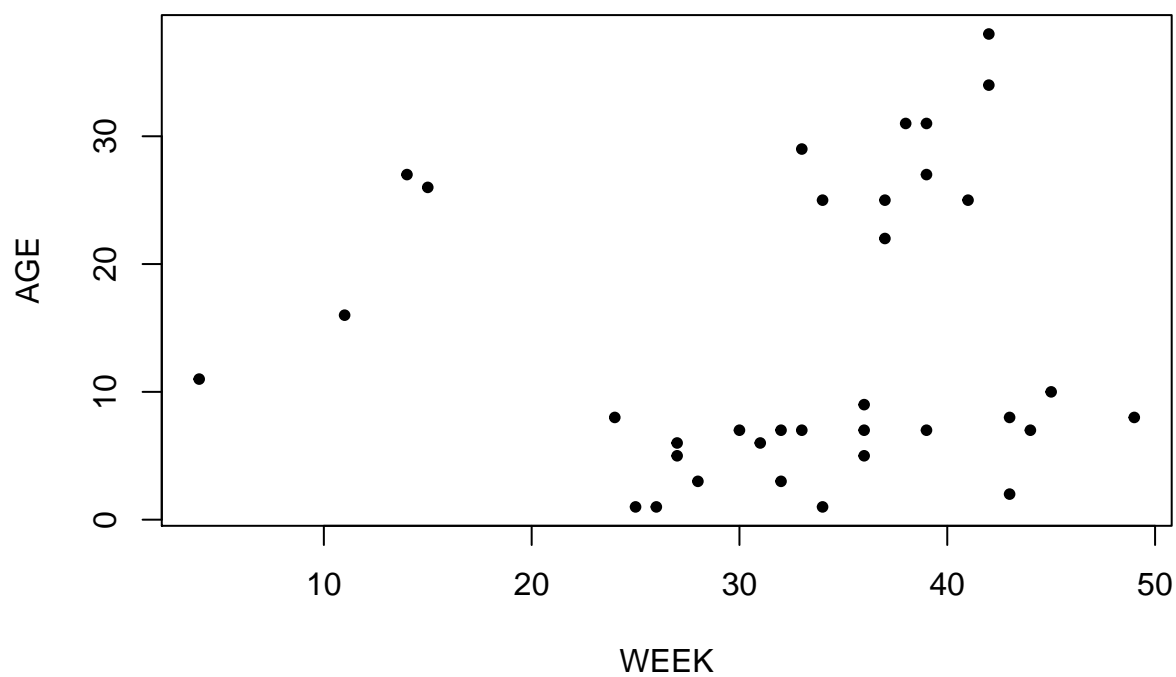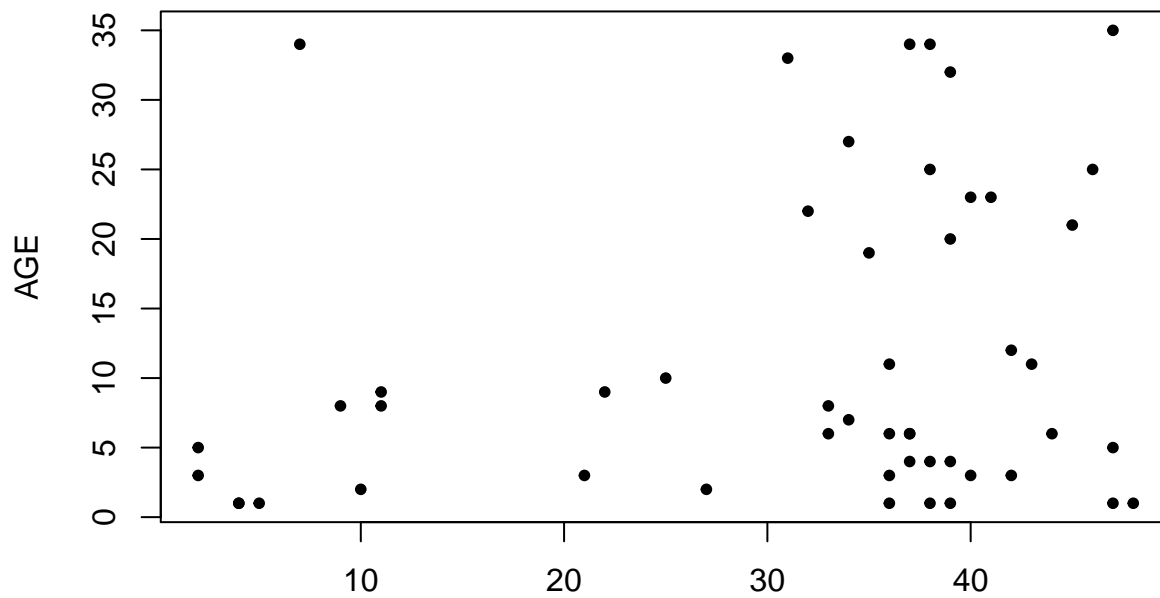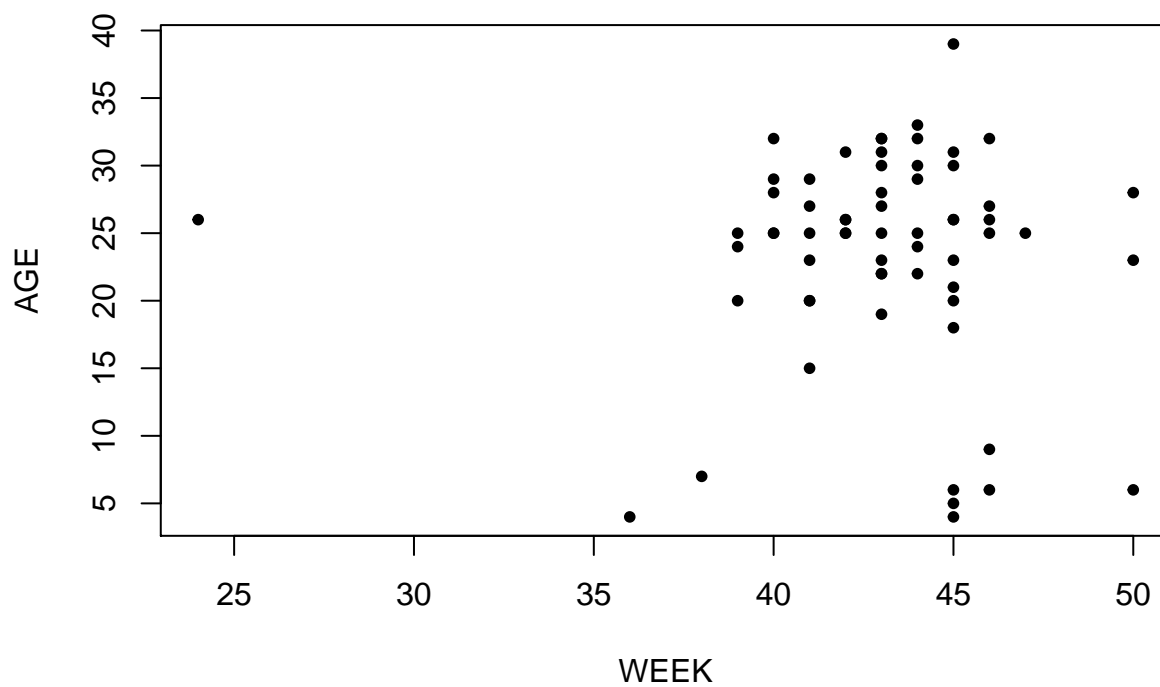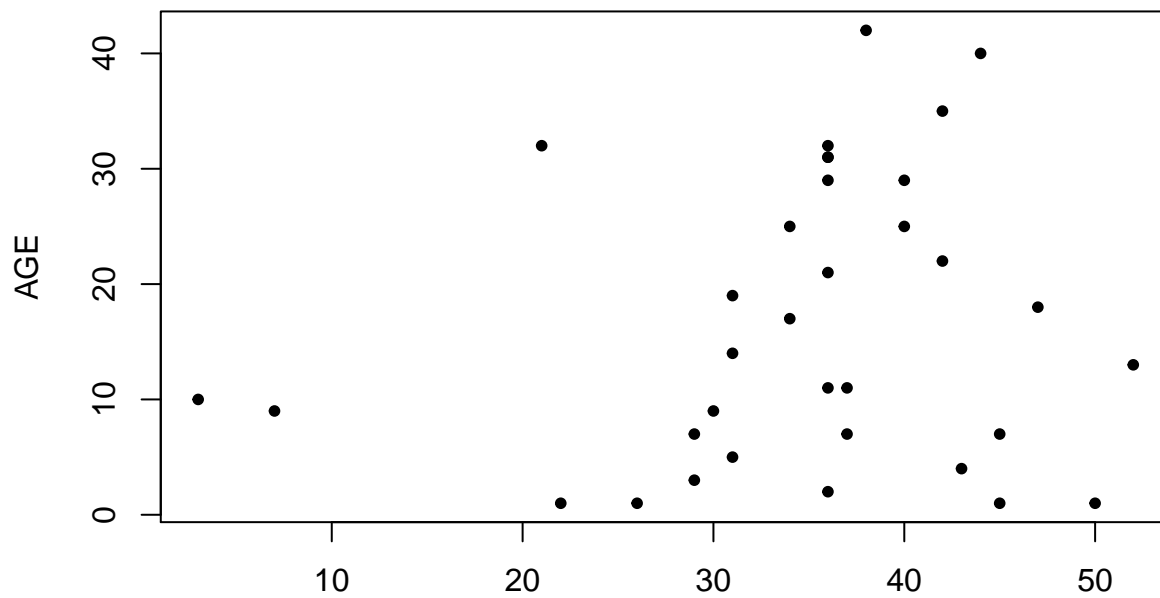
**SAO PAULO**

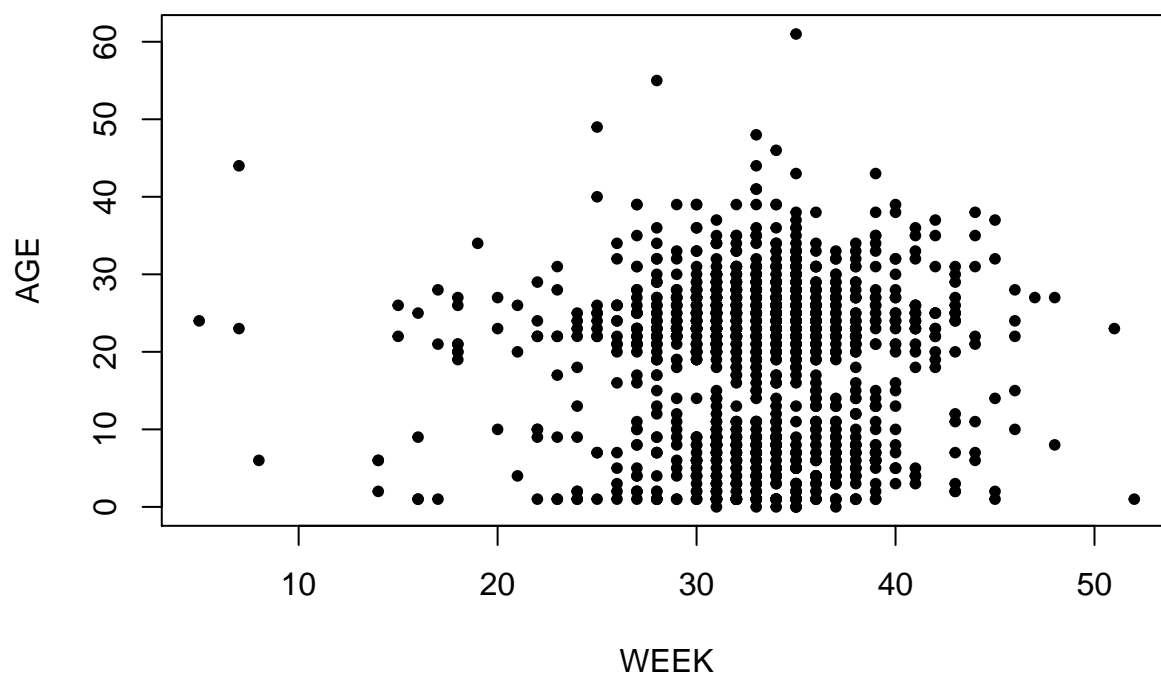# ARTUR NOGUEIRA



# BOTUCATU

**MOGI GUACU**



**ENGENHEIRO COELHO**

**BAURU**



**SAO CAETANO DO SUL**

## MORRO AGUDO



## OSASCO

**DIADEMA**

**POA**

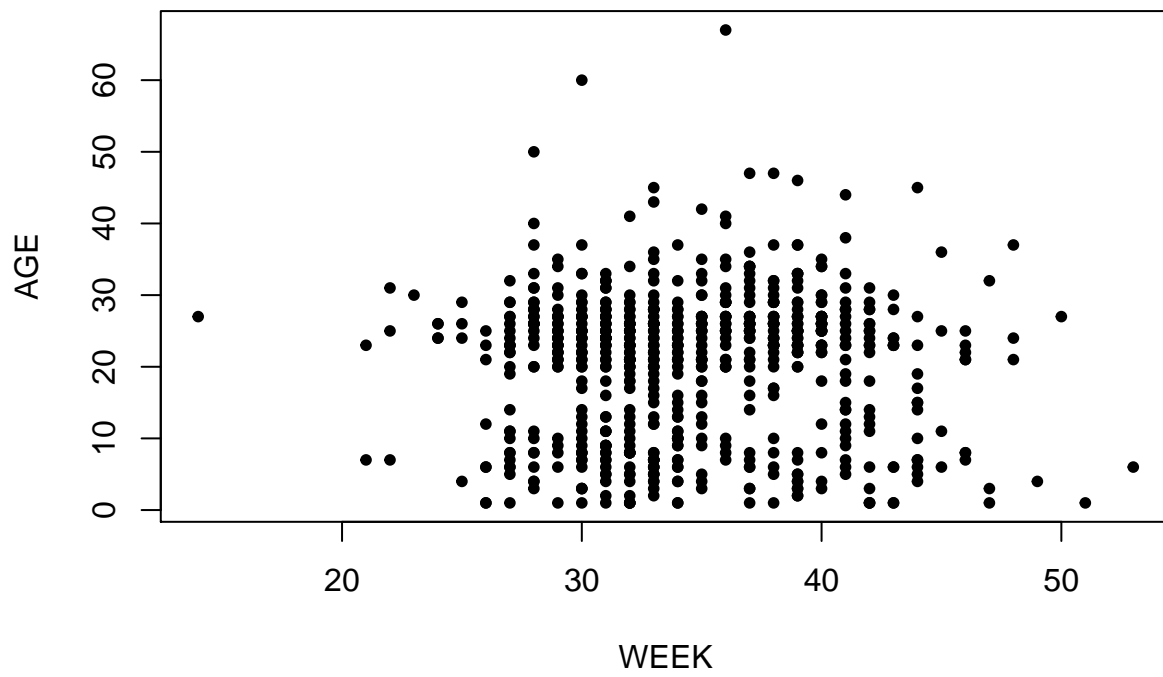## SAO BERNARDO DO CAMPO



WEEK

## GUARULHOS



WEEK

**CAJAMAR**



WEEK

**SANTO ANDRE**



WEEK

8

## PRESIDENTE PRUDENTE



## MARILIA

## FERRAZ DE VASCONCELOS



WEEK

## GUARUJA



WEEK

## CUBATAO



## SAO JOAO DA BOA VISTA

## ITAQUAQUECETUBA



WEEK

## GUARATINGUETA



WEEK

**SANTOS**



WEEK

**CAMPINAS**



WEEK

# ARARAQUARA



# MOGI MIRIM

**JANDIRA**



**COTIA**

**ARARAS**



WEEK

**PINDAMONHANGABA**



WEEK

16

# FRANCO DA ROCHA



# MAUA

**MOGI DAS CRUZES**



WEEK

**CARAPICUIBA**



WEEK

18

**SAO JOSE DO RIO PRETO**



**EMBU−GUACU**

# CATANDUVA



# SAO JOSE DOS CAMPOS

# SANTA CRUZ DAS PALMEIRAS



# EMBU DAS ARTES

## RIO GRANDE DA SERRA



## PONTAL

# RIBEIRAO PIRES



# ITAPECERICA DA SERRA

## SANTANA DE PARNAIBA



## BARUERI

**FRANCISCO MORATO**

AGE

WEEK

**CAIEIRAS**

AGE

WEEK

# JACAREI



WEEK

# ITU



WEEK

**ATIBAIA**



WEEK

**ITAPEVI**



WEEK

27

**MAIRIPORA**

**RIBEIRAO PRETO**

# LIMEIRA



WEEK

# JUNDIAI



WEEK

# RIO CLARO



# PIRACICABA

## ARUJA



WEEK

## BRAGANCA PAULISTA



WEEK

**CAMPO LIMPO PAULISTA**



WEEK

**TABOAO DA SERRA**



WEEK

**ITATIBA**



**AMERICANA**

# COSMOPOLIS



# SANTA BARBARA D'OESTE

# HORTOLANDIA



# PAULINIA

# INDAIATUBA

WEEK

AGE

# SALTO

WEEK

AGE

**SAO SEBASTIAO**



WEEK

**SAO VICENTE**



WEEK

**SERTAOZINHO**



WEEK

**SOROCABA**



WEEK

**SUMARE**



WEEK

**SUZANO**



WEEK

**TAUBATE**



WEEK

**TERRA ROXA**



WEEK

**UBATUBA**



**VARGEM GRANDE PAULISTA**

**VARZEA PAULISTA**



**VOTORANTIM**

**Preliminary Analysis**

## Confidence intervals for age group densities

Created two_matrix, three_matrix, four_matrix, and five_matrix which contain the density of cases in age groups of the corresponding span for each county. Density was just calculated as the number of cases in that specific class divided by the number of cases in that county total.

```r
#find appropriate age windows, look at age groups with ranges ranging from 2 to 5 years
density_age <- function(max.age, size){
  amount <- round(max.age / size)
  amount.vec <- rep(NA, amount)
  amount.vec <- as.data.frame(amount.vec)
  amount.vec[1,] <- length(which(subset_county$AGE <= size))
  for(i in 2:amount){
    before <- as.integer(sum(amount.vec[1: i - 1,]))
    amount.vec[i,] <- length(which(subset_county$AGE <= (size * i))) - before
  }
  density.vec <- amount.vec / length(subset_county$AGE)
  return(density.vec)
}


#matrix for age group of 2's
two_matrix <- matrix(rep(NA, round(84/2) * length(county_vec)), ncol = length(county_vec))
two_matrix <- as.data.frame(two_matrix)
colnames(two_matrix) <- county_vec
row.names(two_matrix) <- c("<2", "2-4", "4-6", "6-8", "8-10",
                           "10-12", "12-14", "14-16", "16-18", "18-20",
                           "20-22", "22-24", "24-26", "26-28", "28-30",
                           "30-32", "32-34", "34-36", "36-38", "38-40",
                           "40-42", "42-44", "44-46", "46-48", "48-50",
                           "50-52", "52-54", "54-56", "56-58", "58-60",
                           "60-62", "62-64", "64-66", "66-68", "68-70",
                           "70-72", "72-74", "74-76", "76-78", "78-80",
                           "80-82", "82-84")

for(j in 1:length(county_vec)){
  subset_county <- subset(data, data$COUNTY == county_vec[j])
  max.age <- max(subset_county$AGE)
  density <- unlist(density_age(max.age, 2))
  length_d <- length(density)
  density <-  c(density, rep(NA, (42 - length_d)))
  two_matrix[,j] <- density
}
#matrix for age groups of 3's
three_matrix <- matrix(rep(NA, round(84/3) * length(county_vec)), ncol = length(county_vec))
three_matrix <- as.data.frame(three_matrix)
colnames(three_matrix) <- county_vec
rownames(three_matrix) <- c("<3", "3-6", "6-9",
                            "9-12", "12-15", "15-18",
                            "18-21", "21-24", "24-27",
                            "27-30", "30-33", "33-36",
                            "36-39", "39-42", "42-45",
                            "45-48", "48-51", "51-54",
```
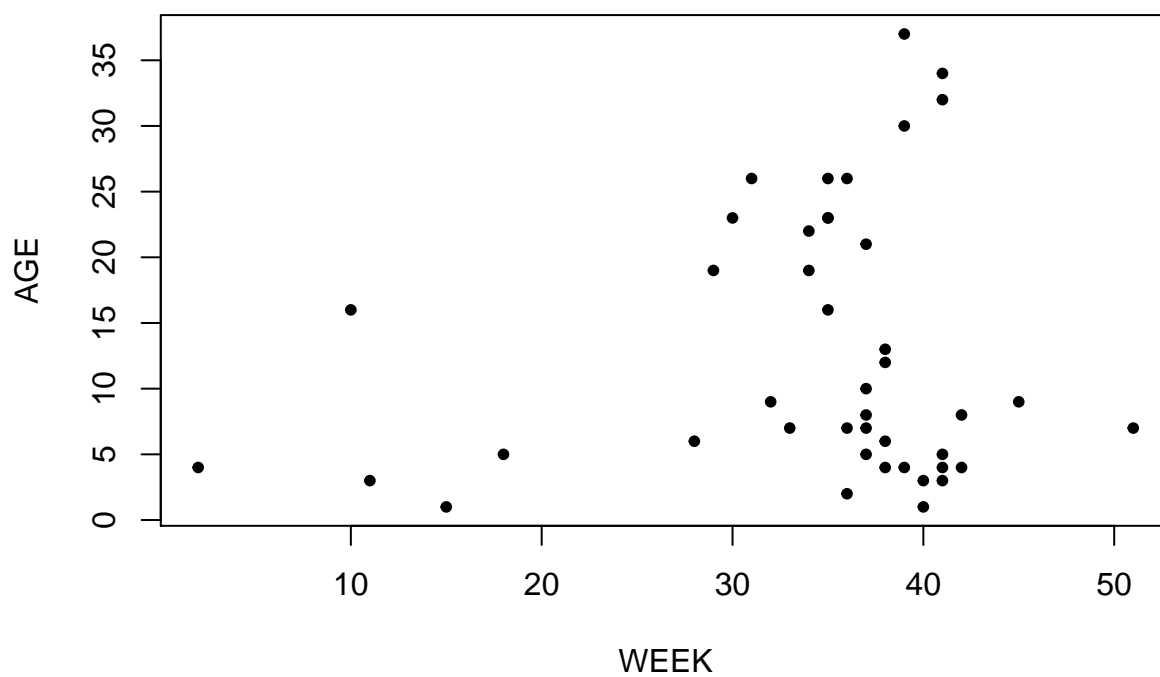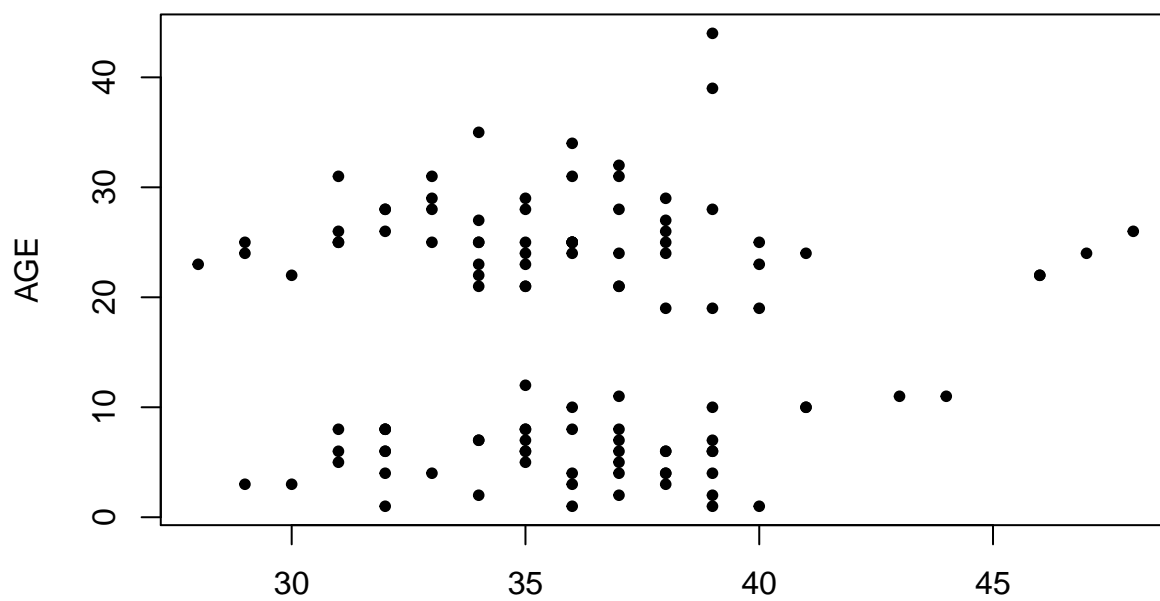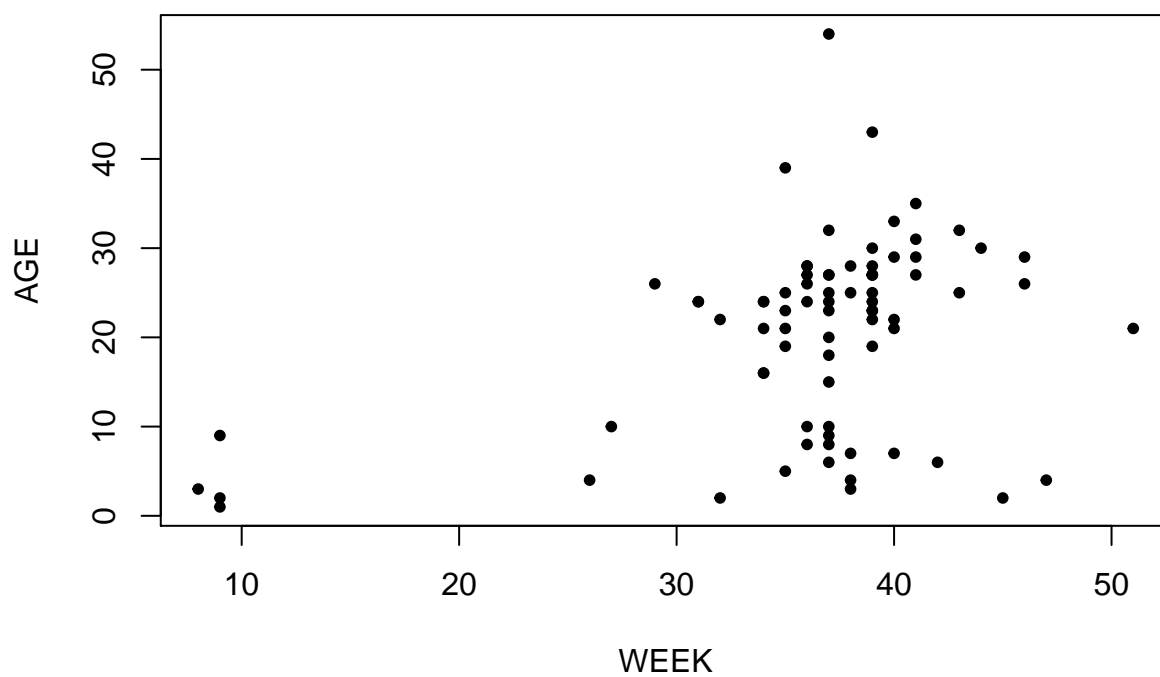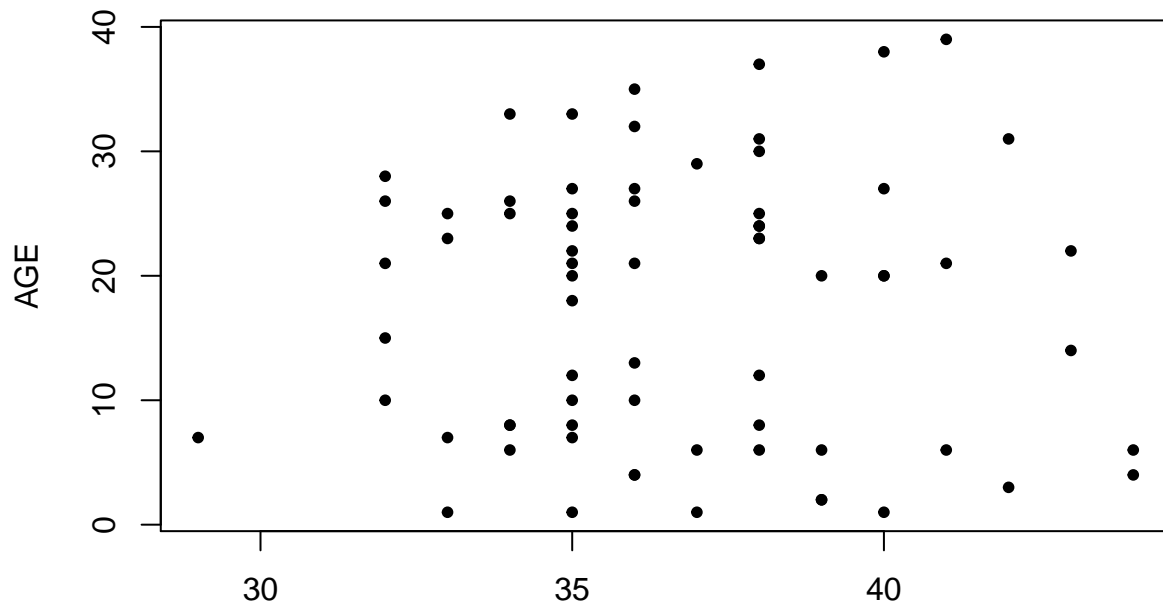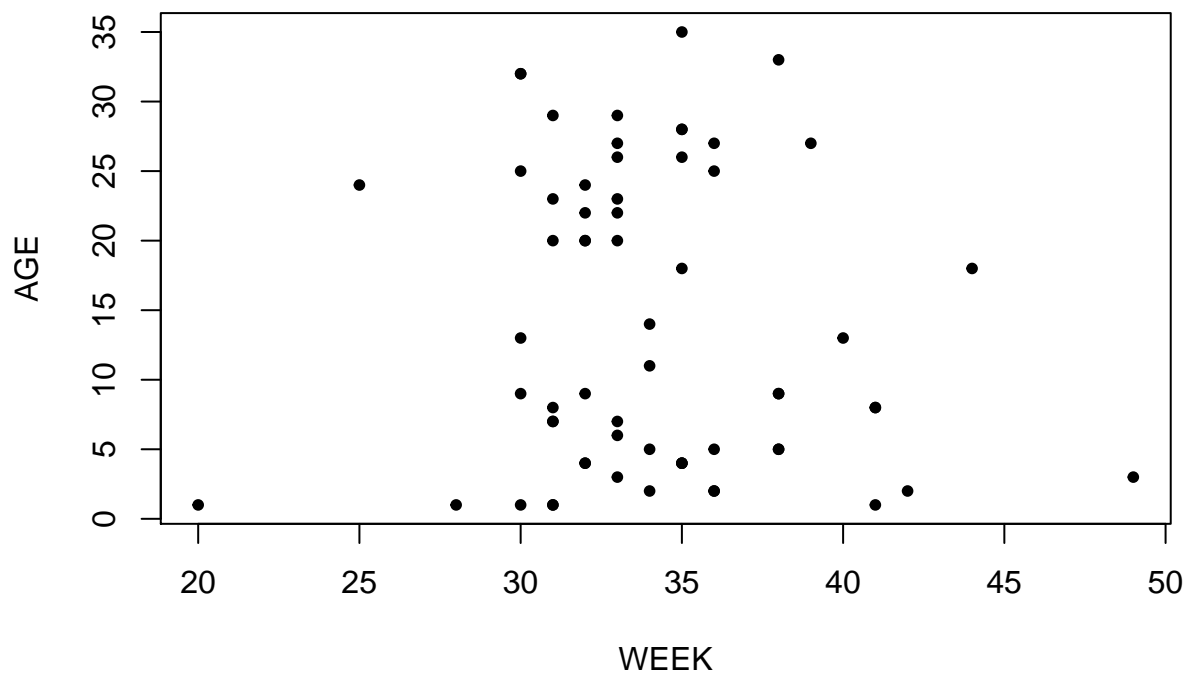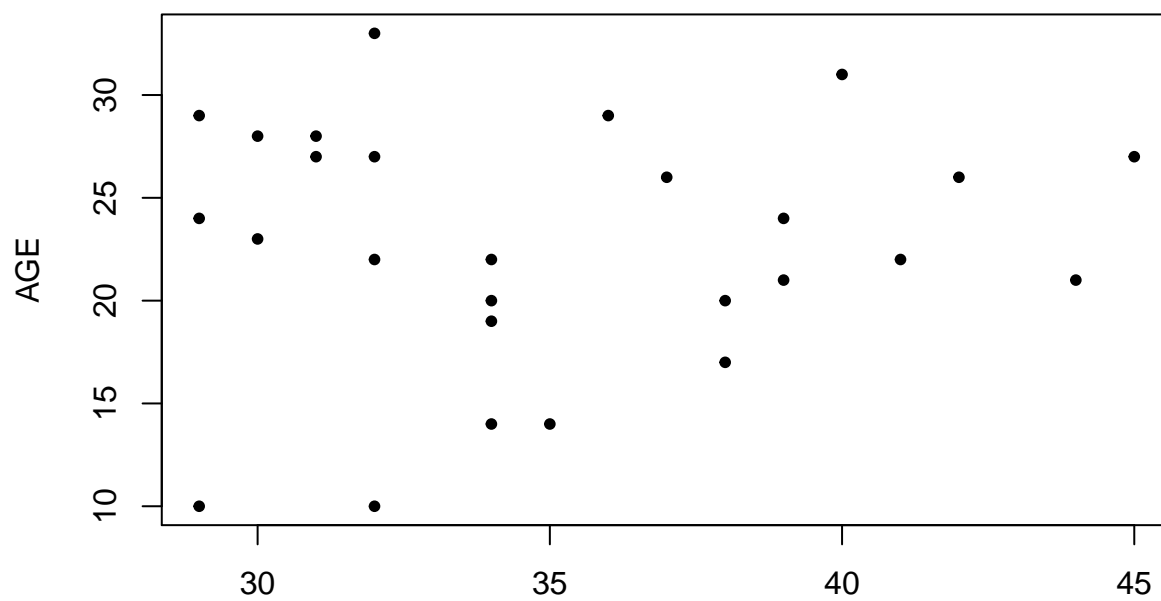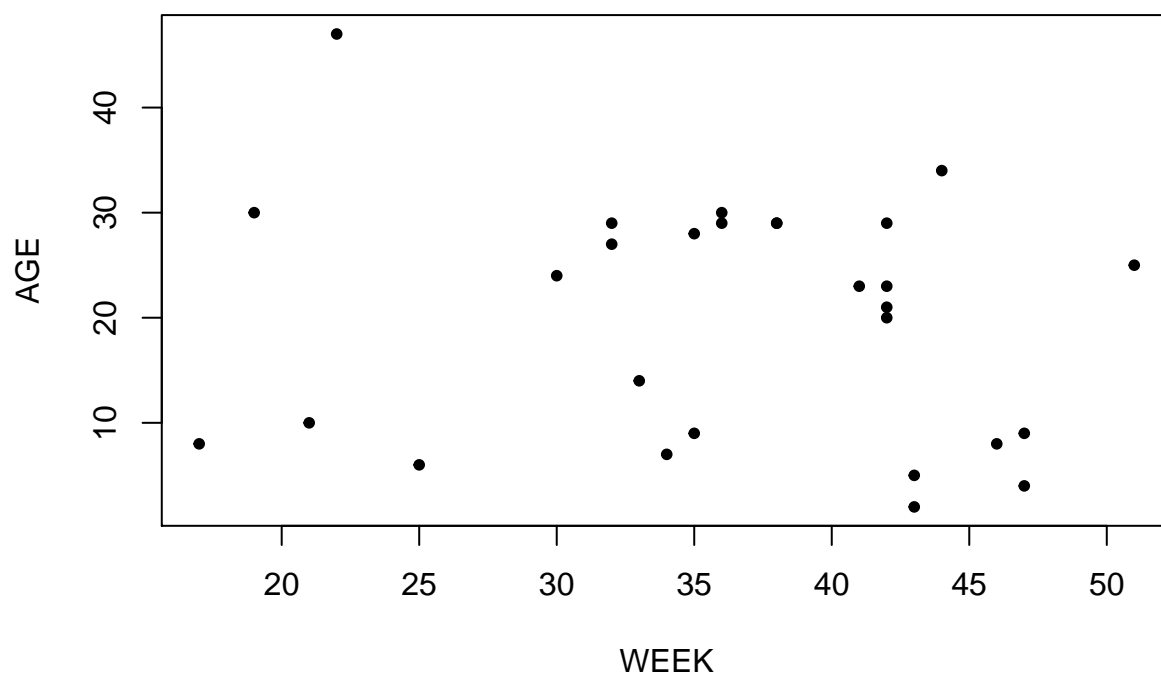
```r
                                  "54-57", "57-60", "60-63",
                                  "63-66", "66-69", "69-72",
                                  "72-75", "75-78", "78-81",
                                  "81-84")

for(j in 1:length(county_vec)){
  subset_county <- subset(data, data$COUNTY == county_vec[j])
  max.age <- max(subset_county$AGE)
  density <- unlist(density_age(max.age, 3))
  length_d <- length(density)
  density <-  c(density, rep(NA, (28 - length_d)))
  three_matrix[,j] <- density
}
#matrix for age groups of 4's
four_matrix <- matrix(rep(NA, 21 * length(county_vec)), ncol = length(county_vec))
four_matrix <- as.data.frame(four_matrix)
colnames(four_matrix) <- county_vec
row.names(four_matrix) <- c("<4", "4-8", "8-12", "12-16",
                                  "16-20", "20-24", "24-28", "28-32",
                                  "32-36", "36-40", "40-44", "44-48",
                                  "48-52", "52-56", "56-60", "60-64",
                                  "64-68", "68-72", "72-76", "76-80",
                                  "80-84")

for(j in 1:length(county_vec)){
  subset_county <- subset(data, data$COUNTY == county_vec[j])
  max.age <- max(subset_county$AGE)
  density <- unlist(density_age(max.age, 4))
  length_d <- length(density)
  density <-  c(density, rep(NA, (21 - length_d)))
  four_matrix[,j] <- density
}
#matrix for age groups of 5's
five_matrix <- matrix(rep(NA, round(84/5) * length(county_vec)), ncol = length(county_vec))
five_matrix <- as.data.frame(five_matrix)
colnames(five_matrix) <- county_vec
row.names(five_matrix) <- c("<5", "5-10", "10-15", "15-20",
                                  "20-25", "25-30", "30-35", "35-40",
                                  "40-45", "45-50", "50-55", "55-60",
                                  "60-65", "65-70", "70-75", "75-80",
                                  "80-85")

for(j in 1:length(county_vec)){
  subset_county <- subset(data, data$COUNTY == county_vec[j])
  max.age <- max(subset_county$AGE)
  density <- unlist(density_age(max.age, 5))
  length_d <- length(density)
  density <-  c(density, rep(NA, (17 - length_d)))
  five_matrix[,j] <- density
}
```
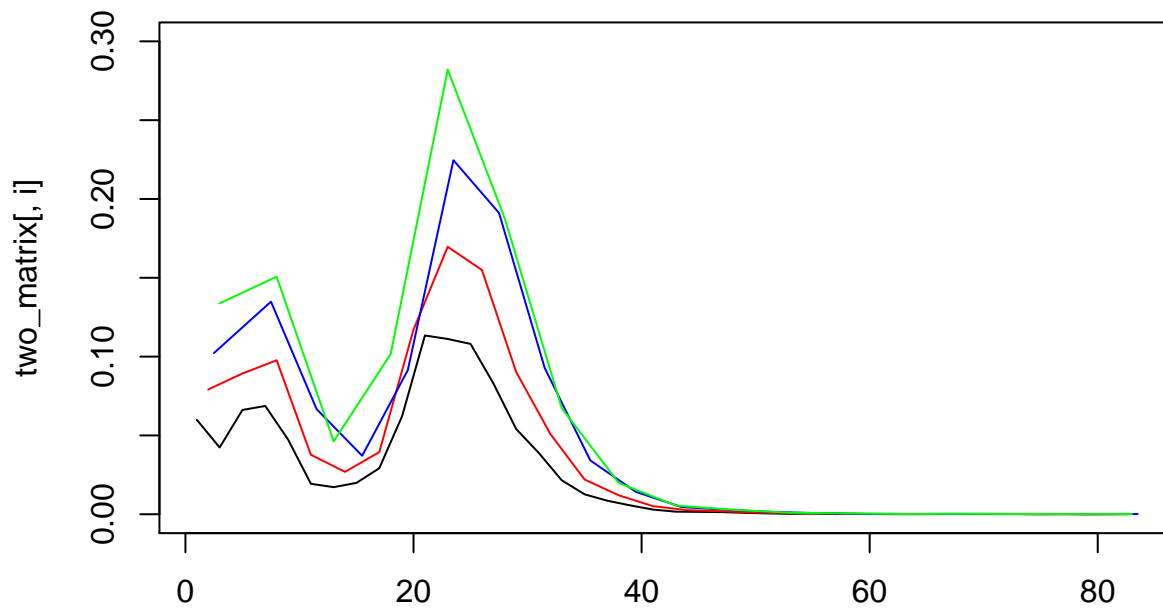
**Preliminary Analysis**

It appears that the size of the window does have a siginficant (colloquially) effect on the density in that close age classes have pretty different densities. Going to compare densities on a line graph.
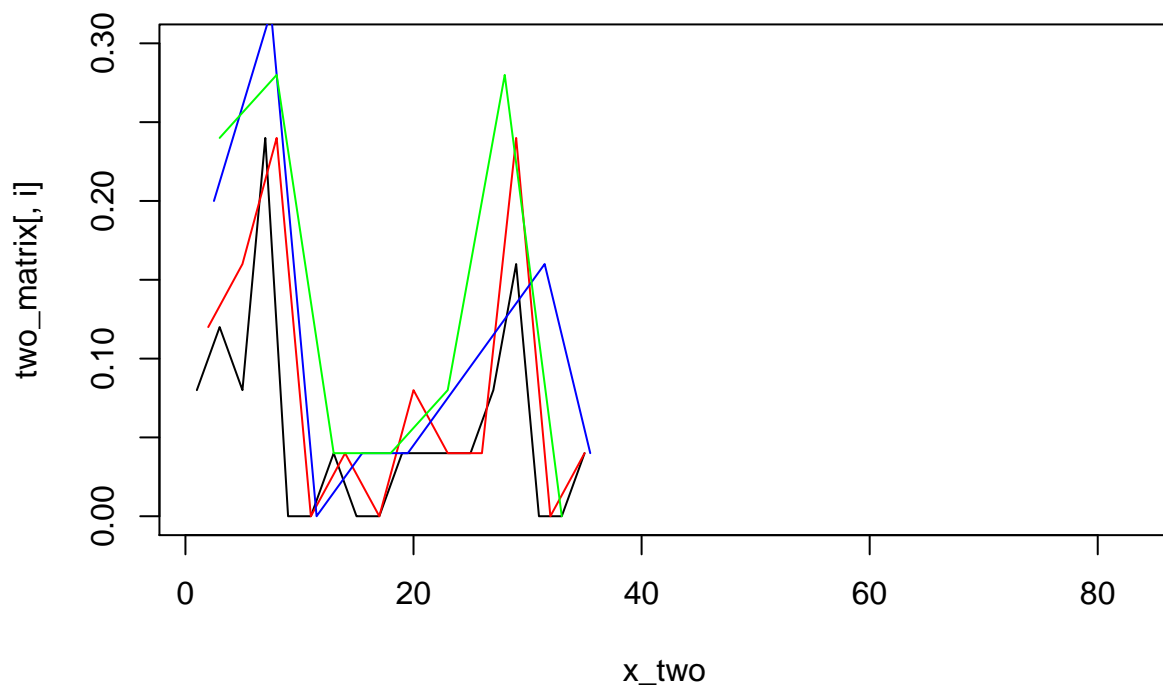
## Comparison of windows

```r
#x axis for twos
x_two <- rep(NA, 42)
x_two[1] <- 1
for(i in 2:42){
  x_add <- seq(1:41)
  x_two[i] <- i + x_add[i - 1]
}
#x axis for threes
x_three <- rep(NA, 28)
x_three[1] <- 2
for(i in 2:28){
  x_add <- seq(1:28)
  x_three[i] <- 2 * i + x_add[i] -1
}
#x axis for fours
x_four <- rep(NA, 21)
x_four[1] <- 2.5
for(i in 2:21){
  x_add <- seq(1:21)
  x_four[i] <- 3 * i + x_add[i] + 0.5 -1
}
#x axis for fives
x_five <- rep(NA, 17)
x_five[1] <- 3
for(i in 2:17){
  x_add <- seq(1:17)
  x_five[i] <- 4 * i + x_add[i] -2
}

for(i in 1:length(county_vec)){
  subset_county <- subset(data, data$COUNTY == county_vec[i])
  if(length(subset_county$AGE) > 20){
    {plot(x_two, two_matrix[,i], type = "l", ylim = c(0, 0.3), main = county_vec[i])
     lines(x_three, three_matrix[,i], col = "red")
     lines(x_four, four_matrix[,i], col = "blue")
     lines(x_five, five_matrix[,i], col = "green")}
  }
}
```
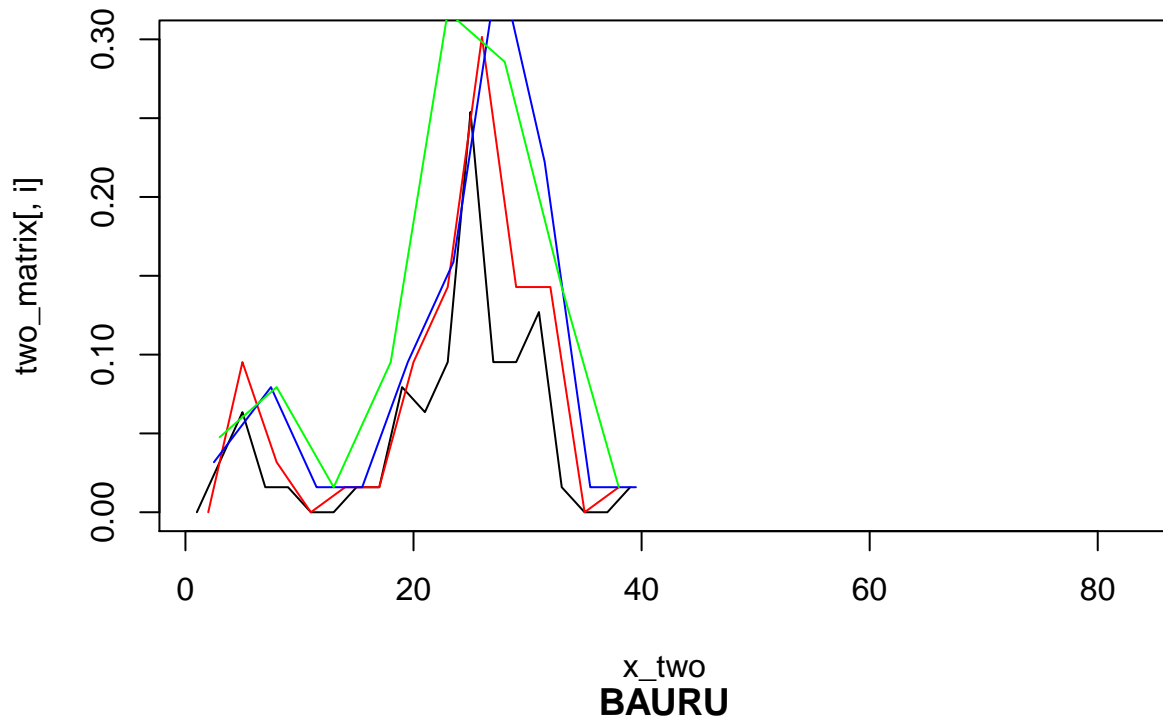
# SAO PAULO



x_two

# ARTUR NOGUEIRA



x_two

**BOTUCATU**



x_two

**MOGI GUACU**



x_two

# ENGENHEIRO COELHO



# BAURU

## SAO CAETANO DO SUL



x_two

## MORRO AGUDO



x_two

## OSASCO



x_two

## DIADEMA



x_two

**POA**


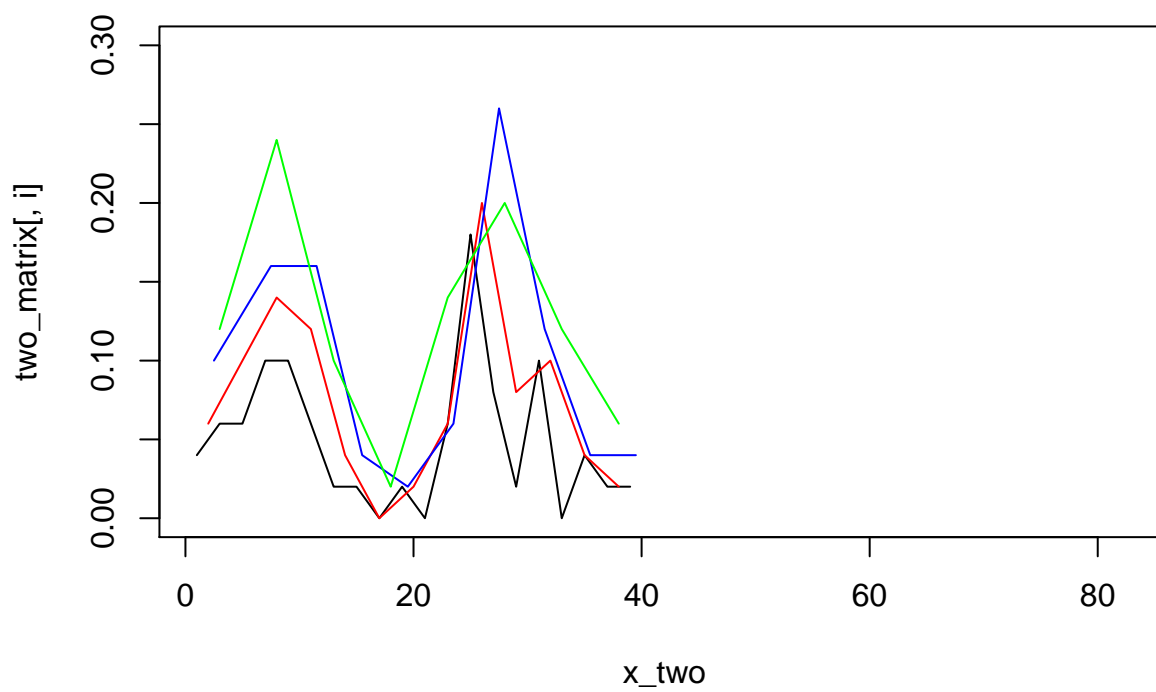
x_two

**SAO BERNARDO DO CAMPO**



x_two

**GUARULHOS**



x_two

**CAJAMAR**



x_two

# SANTO ANDRE



# PRESIDENTE PRUDENTE

# MARILIA



x_two

# FERRAZ DE VASCONCELOS



x_two

**GUARUJA**



x_two

**CUBATAO**



x_two

## SAO JOAO DA BOA VISTA



x_two

## ITAQUAQUECETUBA



x_two

56

**CAMPINAS**



x_two

**ARARAQUARA**



x_two
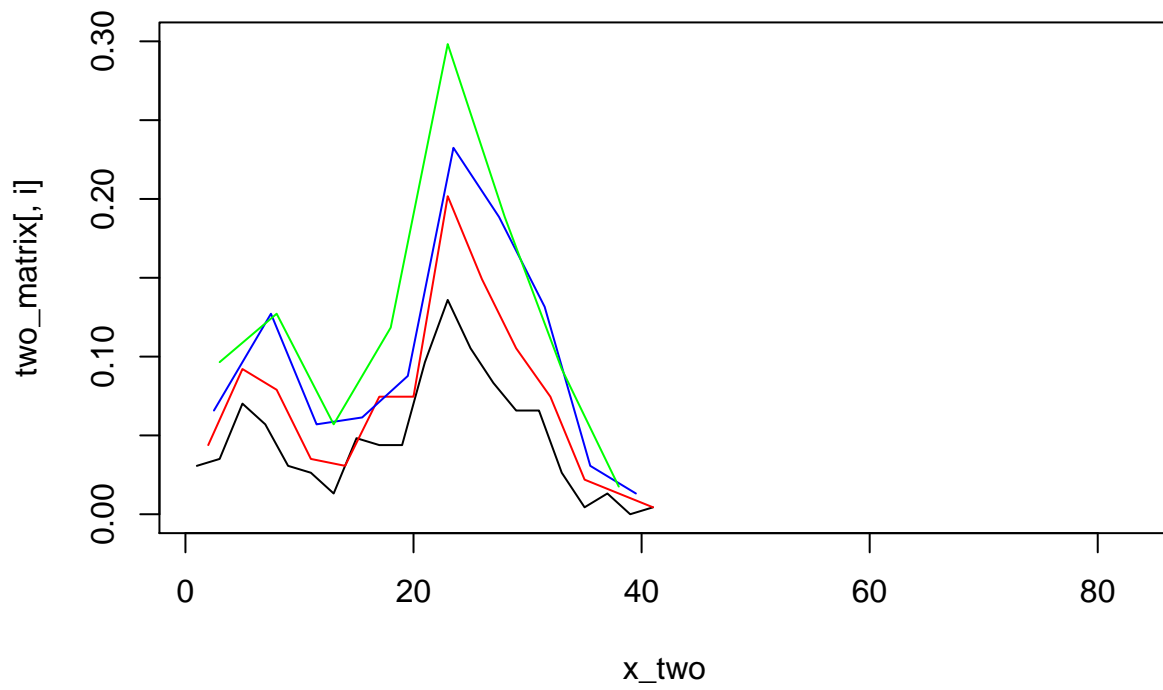
# MOGI MIRIM



# JANDIRA

# COTIA



x_two

# ARARAS



x_two

# PINDAMONHANGABA



# FRANCO DA ROCHA

**MAUA**



x_two

**MOGI DAS CRUZES**



x_two

**CARAPICUIBA**



x_two

**SAO JOSE DO RIO PRETO**



x_two

63

**EMBU−GUACU**



x_two

**CATANDUVA**



x_two

# SAO JOSE DOS CAMPOS



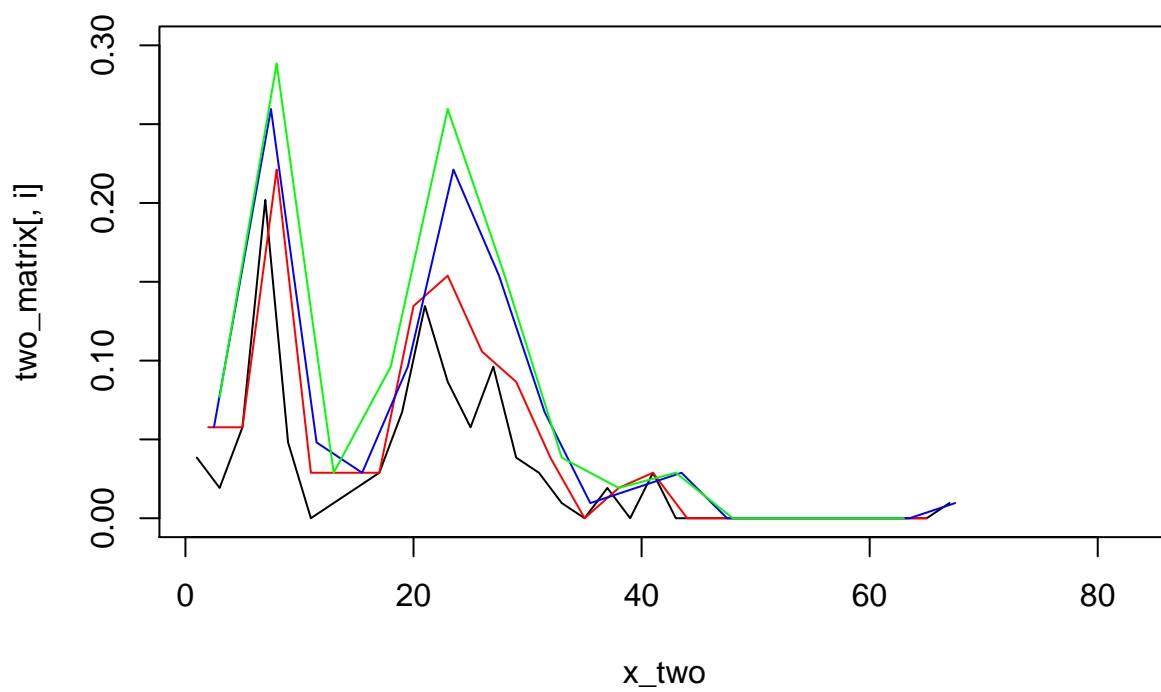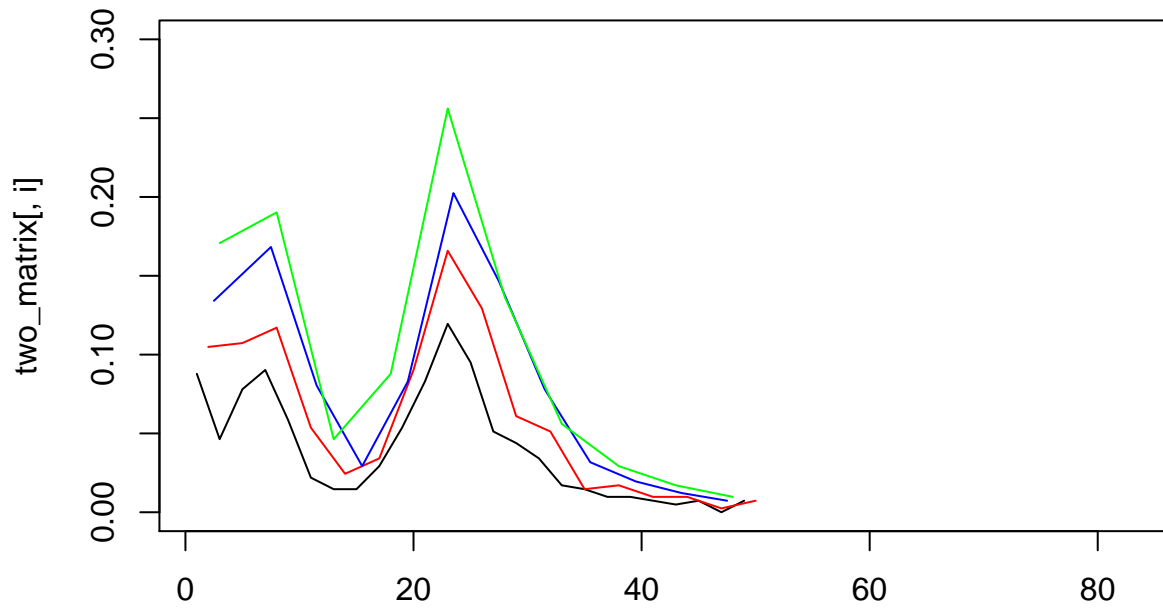# SANTA CRUZ DAS PALMEIRAS

# EMBU DAS ARTES


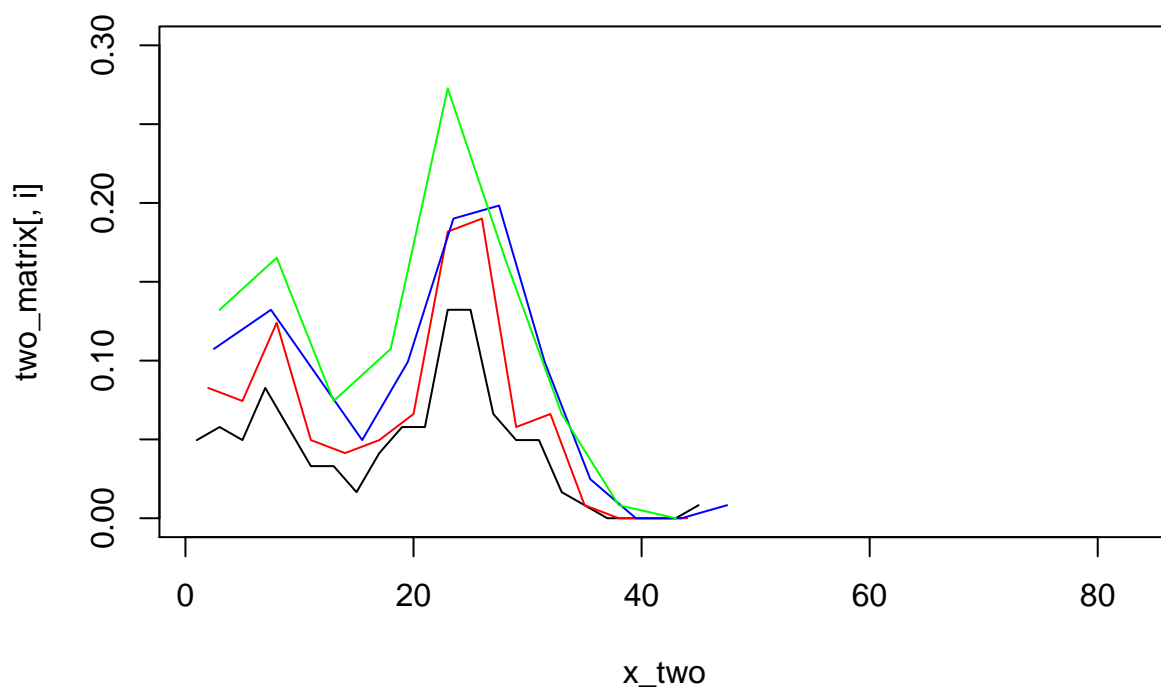
x_two

# RIO GRANDE DA SERRA



x_two

# PONTAL



x_two

# RIBEIRAO PIRES



x_two

# ITAPECERICA DA SERRA



x_two

# SANTANA DE PARNAIBA



x_two

**BARUERI**



x_two

**FRANCISCO MORATO**



x_two

69

## CAIEIRAS



x_two

## JACAREI



x_two

**ITU**



**ATIBAIA**

# ITAPEVI



x_two

# MAIRIPORA



x_two

# RIBEIRAO PRETO



# LIMEIRA

**JUNDIAI**

**RIO CLARO**

**PIRACICABA**

**ARUJA**

# BRANGANCA PAULISTA



x_two

# CAMPO LIMPO PAULISTA



x_two

## TABOAO DA SERRA



## ITATIBA

**AMERICANA**



**COSMOPOLIS**

# SANTA BARBARA D'OESTE



# HORTOLANDIA

**SALTO**



x_two

**SAO SEBASTIAO**



x_two

## SAO VICENTE



## SERTAOZINHO

**SOROCABA**



x_two

**SUMARE**



x_two

83

**SUZANO**



x_two

**TAUBATE**



x_two

84

# TERRA ROXA



# UBATUBA

## VARGEM GRANDE PAULISTA



## VARZEA PAULISTA

## VOTORANTIM



```r
#two
mean_two <- rep(NA, 42)
sd_two <- rep(NA, 42)
for(i in 1:42){
  mean_two[i] <- rowMeans(two_matrix[i,], na.rm = TRUE)
  sd_two[i] <- apply(two_matrix, 1, sd, na.rm = TRUE)[i]
}

#three
mean_three <- rep(NA, 28)
sd_three <- rep(NA, 28)
for(i in 1:28){
  mean_three[i] <- rowMeans(three_matrix[i,], na.rm = TRUE)
  sd_three[i] <- apply(three_matrix, 1, sd, na.rm = TRUE)[i]
}

#four
mean_four <- rep(NA, 21)
sd_four <- rep(NA, 21)
for(i in 1:21){
  mean_four[i] <- rowMeans(four_matrix[i,], na.rm = TRUE)
  sd_four[i] <- apply(four_matrix, 1, sd, na.rm = TRUE)[i]
}

#five
mean_five <- rep(NA, 17)
sd_five <- rep(NA, 17)
for(i in 1:17){
  mean_five[i] <- rowMeans(five_matrix[i,], na.rm = TRUE)
```

```r
  sd_five[i] <- apply(five_matrix, 1, sd, na.rm = TRUE)[i]
}

#RMSP = 1
urban <- subset(data, data$RMSP == "1")
urban_counties <- as.character(unique(urban$COUNTY))

which.urban <- rep(NA, length(urban_counties))
for(i in 1:length(urban_counties)){
 which.urban[i] <- which(county_vec[] == urban_counties[i])
}
two_urban_mat <- two_matrix[which.urban]
three_urban_mat <- three_matrix[which.urban]
four_urban_mat <- three_matrix[which.urban]
five_urban_mat <- five_matrix[which.urban]


#RMSP = 0
rural <- subset(data, data$RMSP == "0")
rural_counties <- as.character(unique(rural$COUNTY))

which.rural <- rep(NA, length(rural_counties))
for(i in 1:length(rural_counties)){
  which.rural[i] <- which(county_vec[] == rural_counties[i])
}
two_rural_mat <- two_matrix[which.rural]
two_rural <- matrix(rep(NA, 3 * 42 * length(rural_counties)), ncol = 3)
colnames(two_rural) <- c("COUNTY", "CLASS", "DENSITY")
two_rural <- as.data.frame(two_rural)
#fill counties
two_rural[,1] <- rep(rural_counties, 42)

#fill age class
two_rural[seq(1:length(rural_counties)),2] <- x_two[1]
for(i in 1:length(x_two)){
  x <- length(rural_counties) * i + seq(1:length(rural_counties))
  two_rural[x,2] <- x_two[i+1]
}

#fill density
for(i in 1:length(rural_counties)){
  two_rural[i,3] <- two_rural_mat[1,i]
}
for(i in 1:length(rural_counties)){
  for(j in 2:length(x_two)){
    two_rural[(i + 348 * (j - 1) ),3]  <- two_rural_mat[j,i]
  }
  }
two_rural$CLASS <- as.factor(two_rural$CLASS)
two_rural$COUNTY <- as.factor(two_rural$COUNTY)


three_rural_mat <- three_matrix[which.rural]
four_rural_mat <- three_matrix[which.rural]
```

```
five_rural_mat <- five_matrix[which.rural]
```

**Preliminary analysis**

Created a method (via for loops) to make a matrix which is compatible with ANOVAs.

Up next:

- Look for correlation or association between denstiy and class on both a rural and urban level.
  - This may be best to do with the five_rural_mat and five_urban_mat because they have fewer factor levels. Starting with urban may also be more promising because it will have fewer factors for county.
- Look into whether the moving windows/data smoothing can be used for the windows of age groups. Try to understand how this works.
- Plot row means with row sd error bars for the different windows.