# STA 108 Project 1

## Mahir Oza, Dylan M. Ang, In Seon Kim

### 1/31/2022

Part 1. Fitting Regression Models

```
# define active physicians columns
phys=CDI$V8
# define predictor variables total pop, hospital beds, and personal income
pop=CDI$V5
beds=CDI$V9
inc=CDI$V16
```

  a. Regression for 3 Predictors

Number of Active Physicians (Y) vs Total Population (Xp)

```
popFit=lm(phys~pop)
b0pop=popFit$coefficients[1]
b1pop=popFit$coefficients[2]
```

Y=b0pop+b1pop(Xp)

regression function: Y=-110.6348+0.002795425(Xp)

Number of Active Physicians (Y) vs Hospital Beds (Xh)

```
bedFit=lm(phys~beds)
b0bed=bedFit$coefficients[1]
b1bed=bedFit$coefficients[2]
```

Y=b0bed+b1bed(Xh)

regression function: Y=-95.93218+0.7431164(Xh)

Number of Active Physicians (Y) vs Total Personal Income (Xi)

```
incFit=lm(phys~inc)
b0inc=incFit$coefficients[1]
b1inc=incFit$coefficients[2]
```
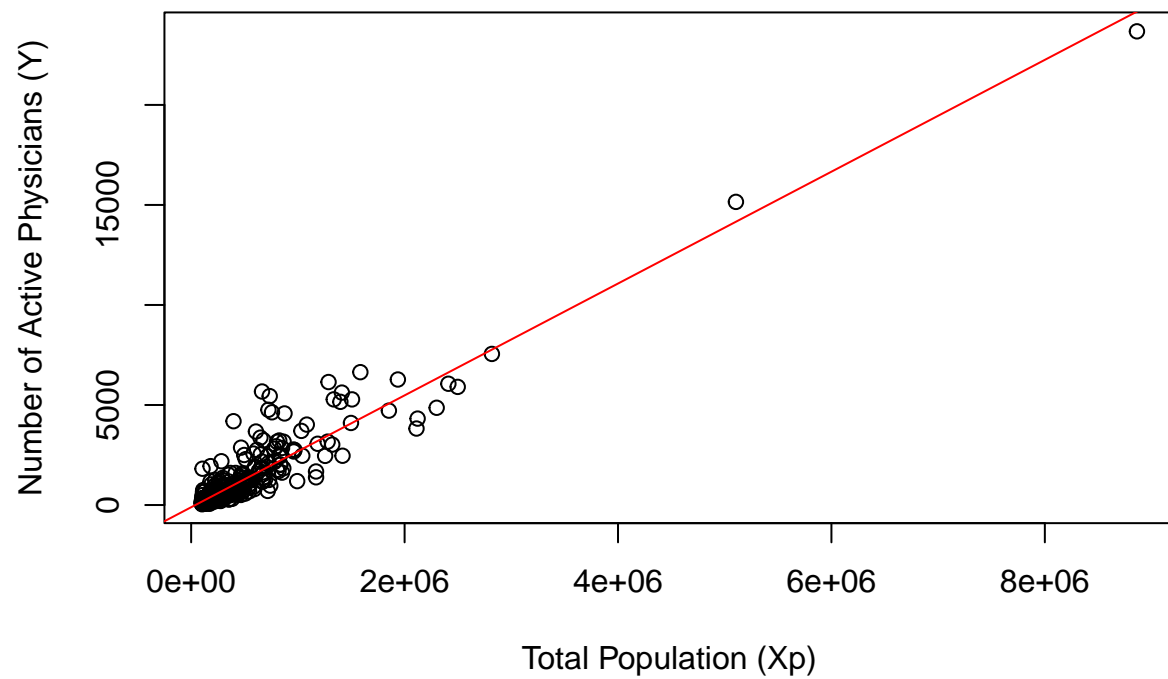
Y=b0inc+b1inc(Xi)

regression function: Y=-48.39485+0.1317012(Xi)
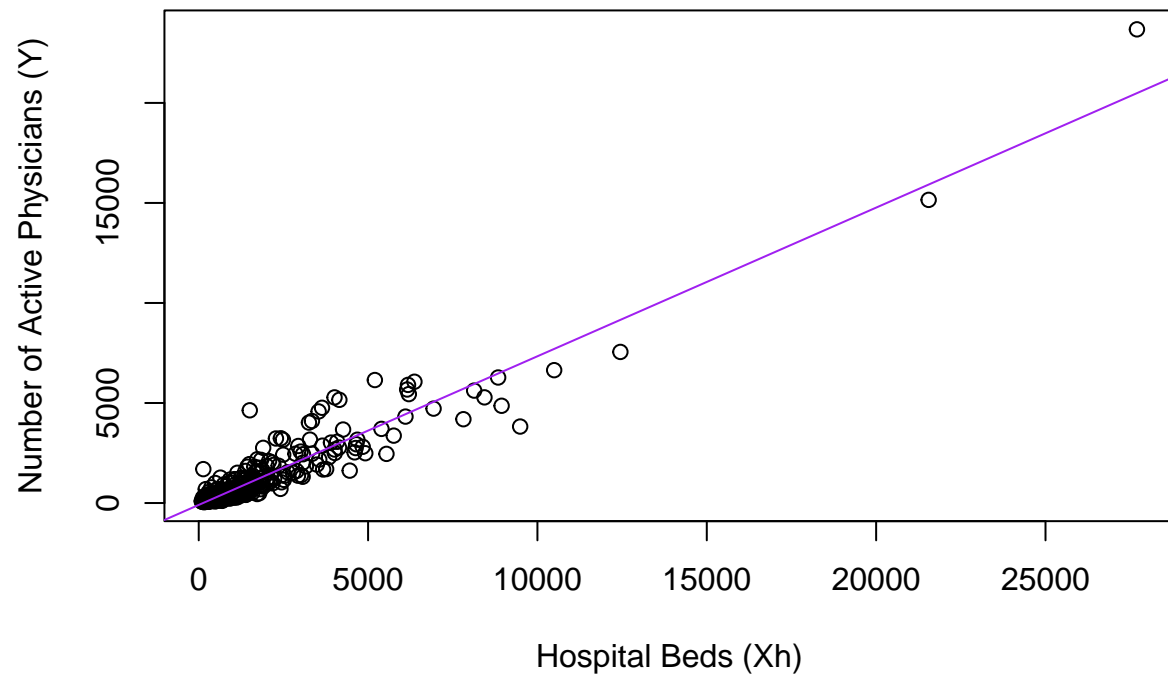
  b. Plots

Number of Active Physicians (Y) vs Total Population (Xp)

```
plot(x=pop,y=phys,xlab="Total Population (Xp)", ylab="Number of Active Physicians (Y)")
abline(popFit,col="red")
```
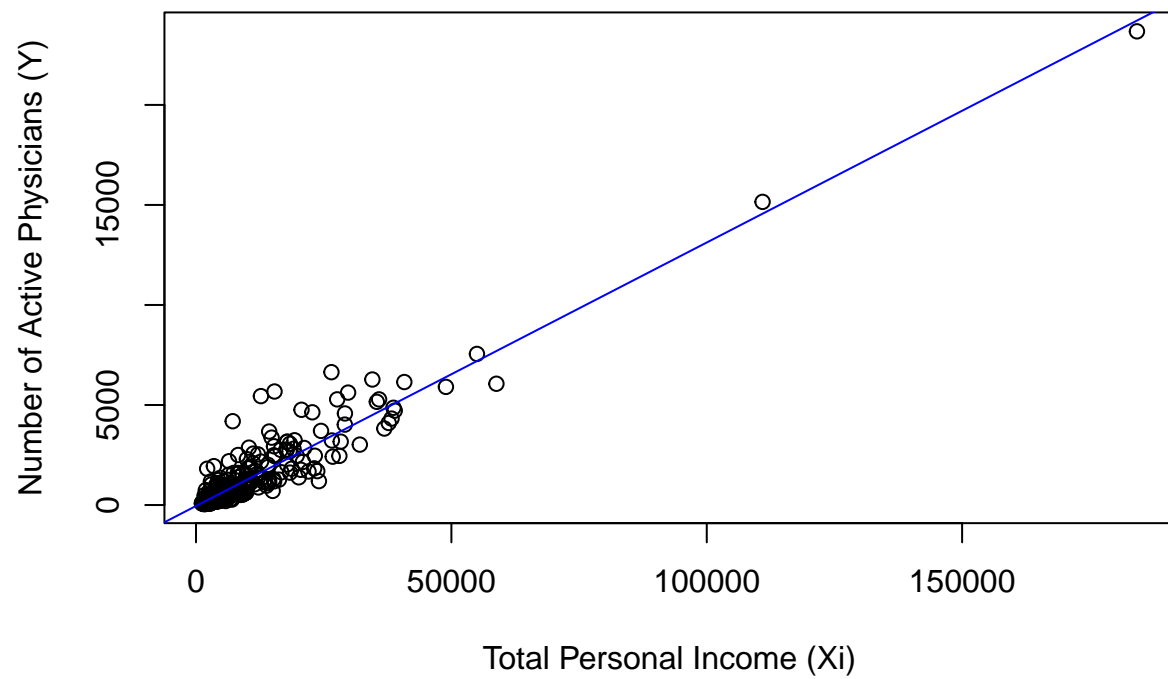
Number of Active Physicians (Y) vs Hospital Beds (Xh)

```
plot(x=beds,y=phys,xlab="Hospital Beds (Xh)", ylab="Number of Active Physicians (Y)")
abline(bedFit,col="purple")
```
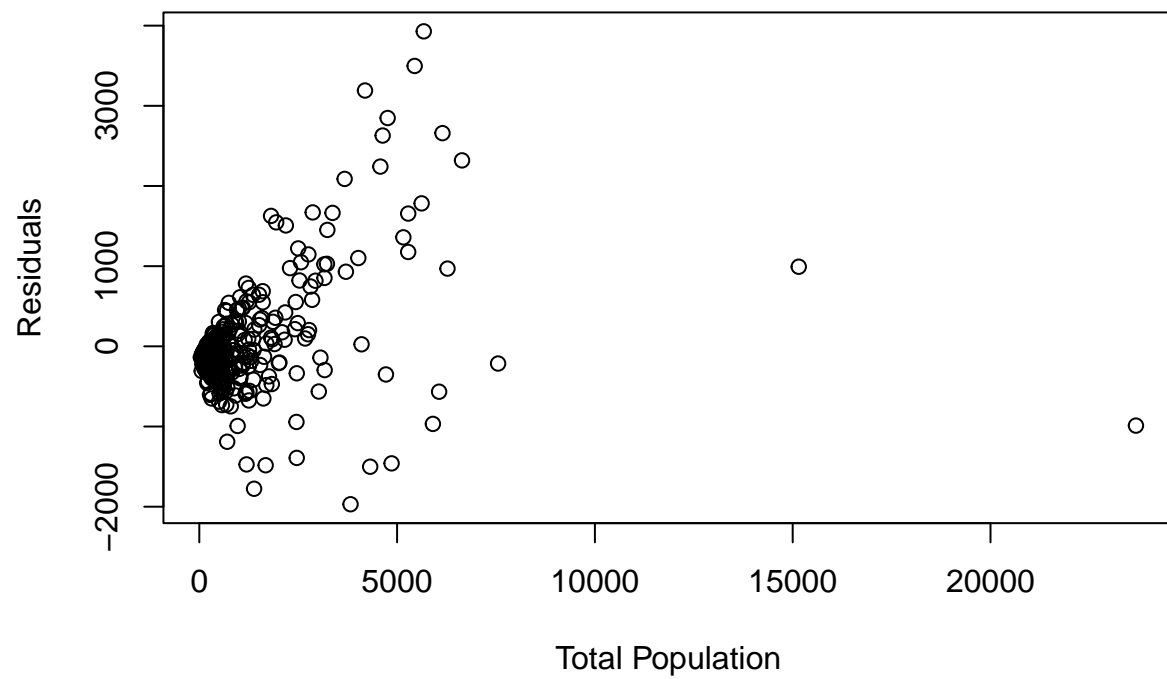
Number of Active Physicians (Y) vs Total Personal Income (Xi)

```
plot(x=inc,y=phys,xlab="Total Personal Income (Xi)", ylab="Number of Active Physicians (Y)")
abline(incFit,col="blue")
```
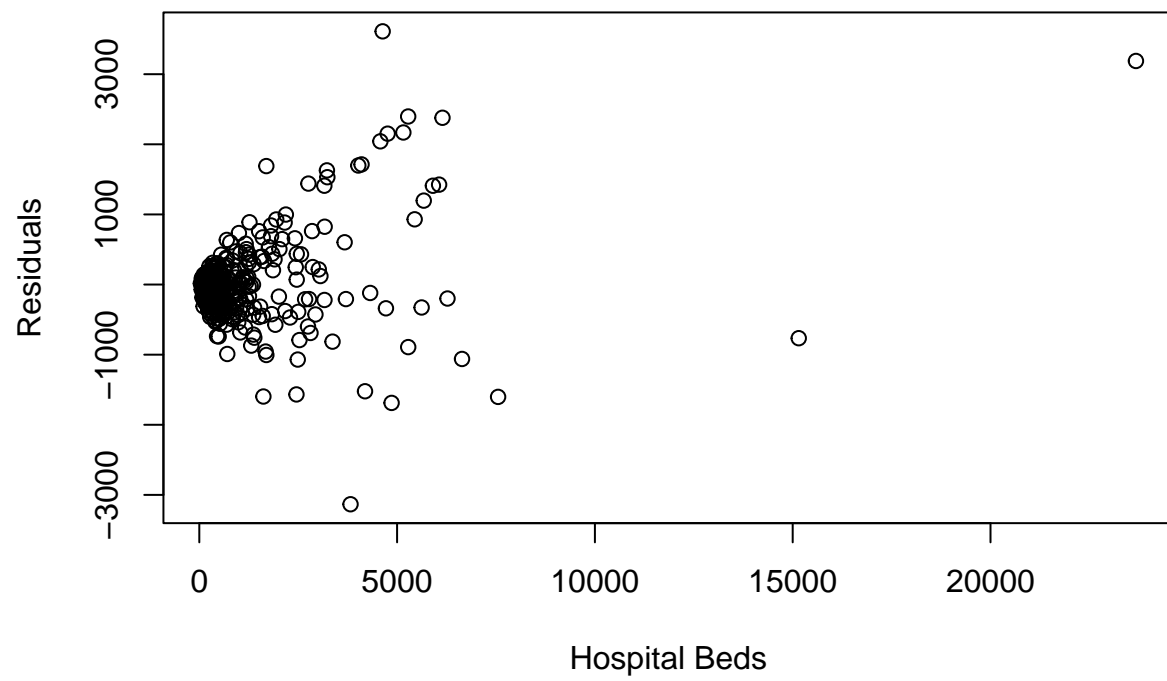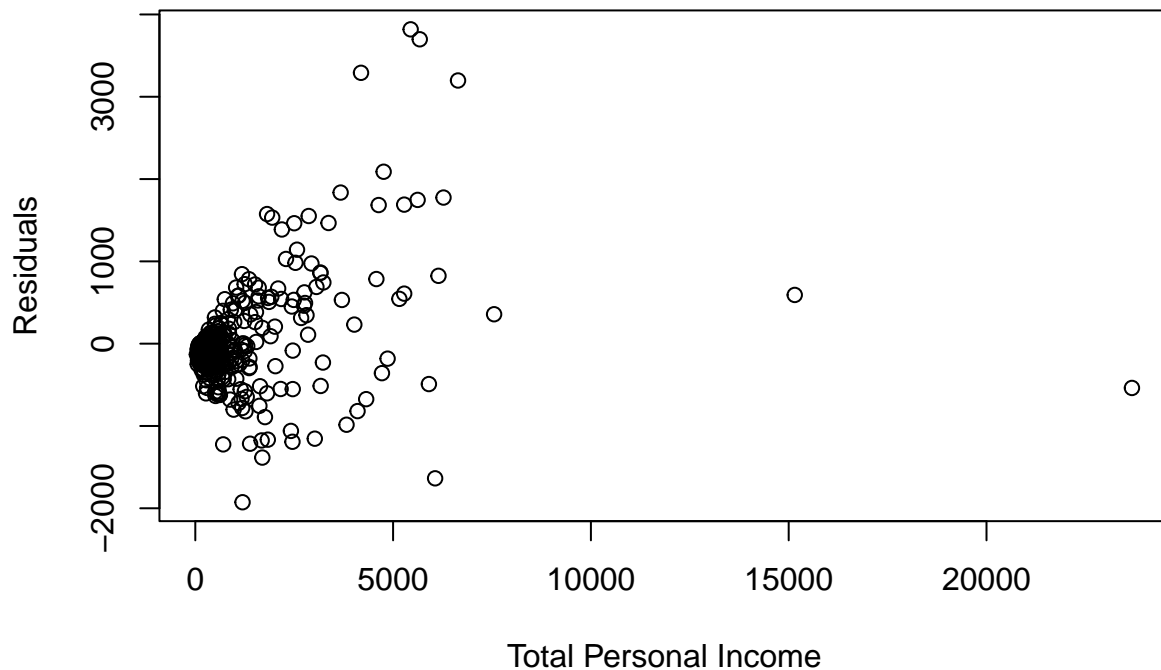
Residuals

```
popE=(phys-(b0pop+b1pop*pop))
bedE=(phys-(b0bed+b1bed*beds))
incE=(phys-(b0inc+b1inc*inc))
# Residuals Plot
plot(x=phys,y=popE,xlab="Total Population", ylab="Residuals")
```

```
plot(x=phys,y=bedE,xlab="Hospital Beds", ylab="Residuals")
```

```
plot(x=phys,y=incE,xlab="Total Personal Income", ylab="Residuals")
```

```
sum(popE)
```

```
## [1] 5.846914e-10
```

```
sum(bedE)
```

```
## [1] 9.003145e-10
```

```
sum(incE)
```

```
## [1] 1.212669e-09
```

The residual plot for all 3 predictor variables as a function of the number of active physicians shows no distinguishable mathematical trend nor any pattern. Also the sum of all the residuals for the 3 predictor variables are extremely small values, close to 0. These show that a linear regression model is appropriate in showing the relationship of the data. In addition to these arguments for a linear regression relationship being a good fit for the data, we can see that the data itself appears to demonstrate a linear relationship. Adding the regression function to the data further shows that the data of all 3 predictor variables follows the behavior of this regression function and further supports a linear regression relationship for all 3 predictor variables for the number of active physicians.

    c. MSE

```
npop=length(pop)
nbed=length(beds)
ninc=length(inc)
SSEp=sum((popE)^2)
SSEb=sum((bedE)^2)
SSEi=sum((incE)^2)
```

```
MSEp=SSEp/(npop-2)
MSEb=SSEb/(nbed-2)
MSEi=SSEi/(ninc-2)
```

The mean of squared errors of the number of hospital beds as the predictor for number of active physicians in a county, is 310,191.88. This is the lower MSE value compared to predictors total personal income, which had an MSE of 324,539.39, and total population of the county, which had the highest MSE of 372,203.50. This shows the mean squared errors in which the number of hospital beds in a county predicting the number of active physicians is lowest and therefore has the lowest variability around the fitted regression line.

Part 2. Measuring Linear Associations

R^2

```
popFit2=summary(popFit)
bedFit2=summary(bedFit)
incFit2=summary(incFit)
r2p=popFit2$r.squared
r2b=bedFit2$r.squared
r2i=incFit2$r.squared
```

For predicting the number of active physicians in a county, the r^2 value for total population as the predictor is 0.8841, for number of hospital beds is 0.9034, and for total personal income is 0.8989. Based on these values, the number of hospital beds as the predictor has the largest coefficient of determination meaning that 90.34% of the variation in the number of active physicians is due to introducing number of hospital beds into the regression model. This shows that the number of hospital beds accounts for the largest reduction in variability in the number of active physicians.