## CIS 320 Data Analysis Final Report:

Professor: Phillip D. Thomas

Billy La

Introduction:

This report is about the final step of my data, which I am going to analysis. My final report include all my data is being analysis. I use excel, SPSS statistic, SPSS Amos and R to analysis all my data. I analysis the population estimate. The data is taking from United States Census Bureau. There are one set of data and the data is about the population estimate on sumlev, region, devision, state, census2010 population estimate, estimatebase of year 2010, population estimate of year 2014-2015, death of year 2014 -2015, birth of years 2014-2015.

Hypothesis:

I am going to find the relationshiop between the number of population between year 2014-2015. The purpose is to monitor and the grow of us population between the year of 2014-2015. Also, I will do the monitor on the number of death and the number of birth on years 2014-2015. I will use calculation of statistic to see the different between those 2 unit.
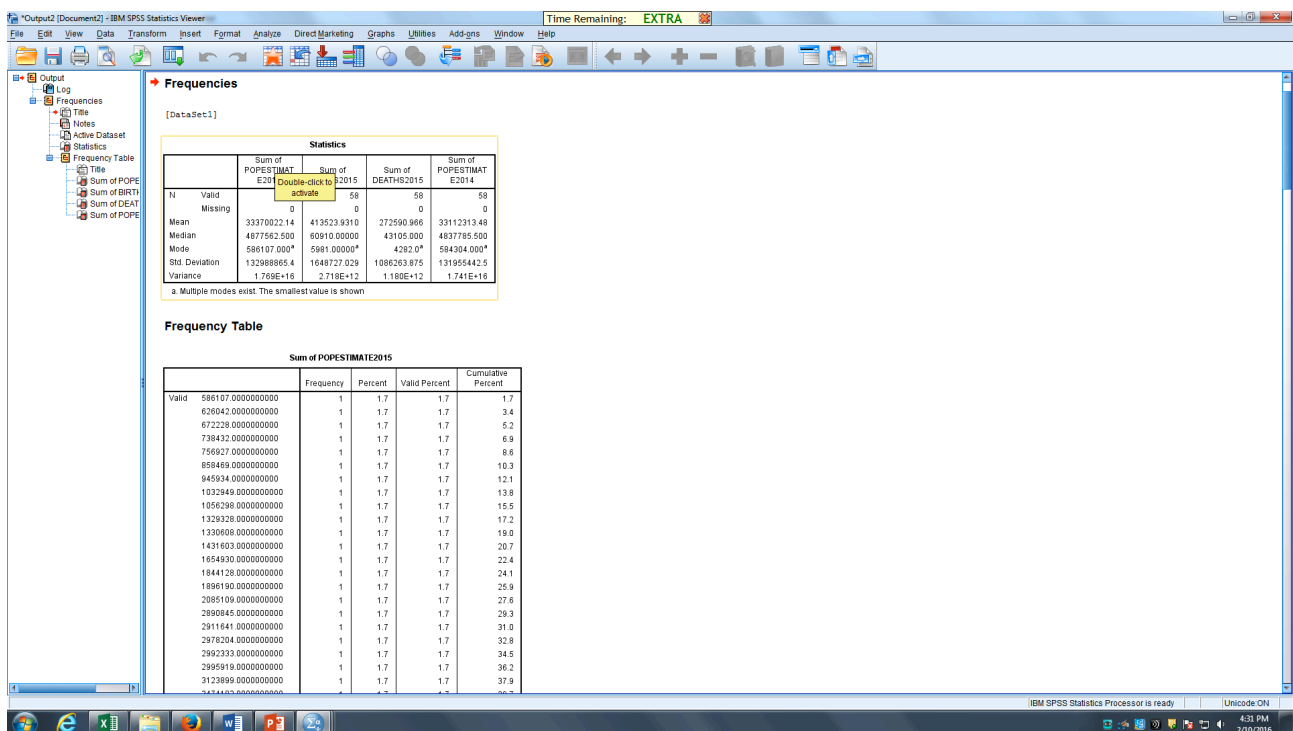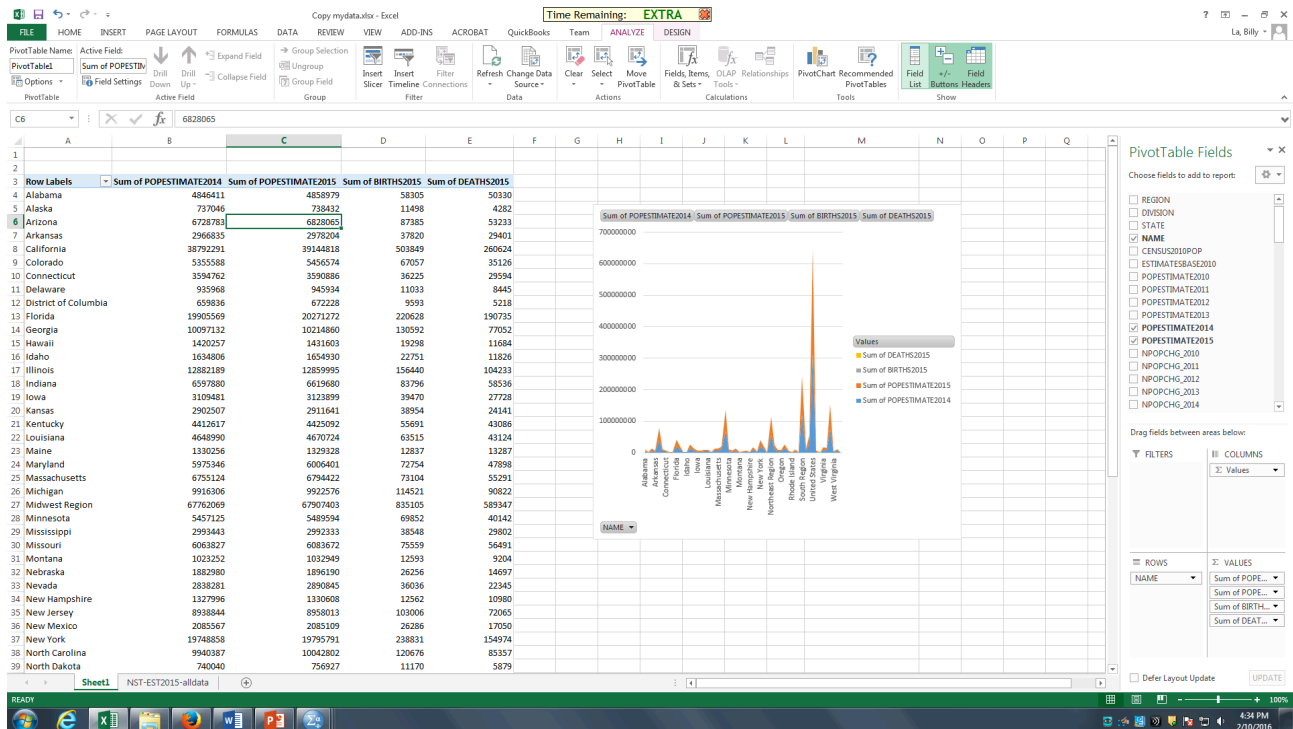
Goals:

I am going to user the pivot table to separate and summary my data include the sumofdeath and sumofbirth of year 2014-2015. Also, the data of populaiton on year 2014-2015 are include in pivot table. After finish the pivot table, I save my data as .csv file. After the save of my data, I insert to IBM spss to do the freqency and t-test for my data. The resourece is going to show accroding to the picture below. After done with the spss analysis, I will use IBM amos for the regression linear to predict the number of death, birth and the populartion of the future year. The final step, which is using R to graph the division of region and summary all the data.
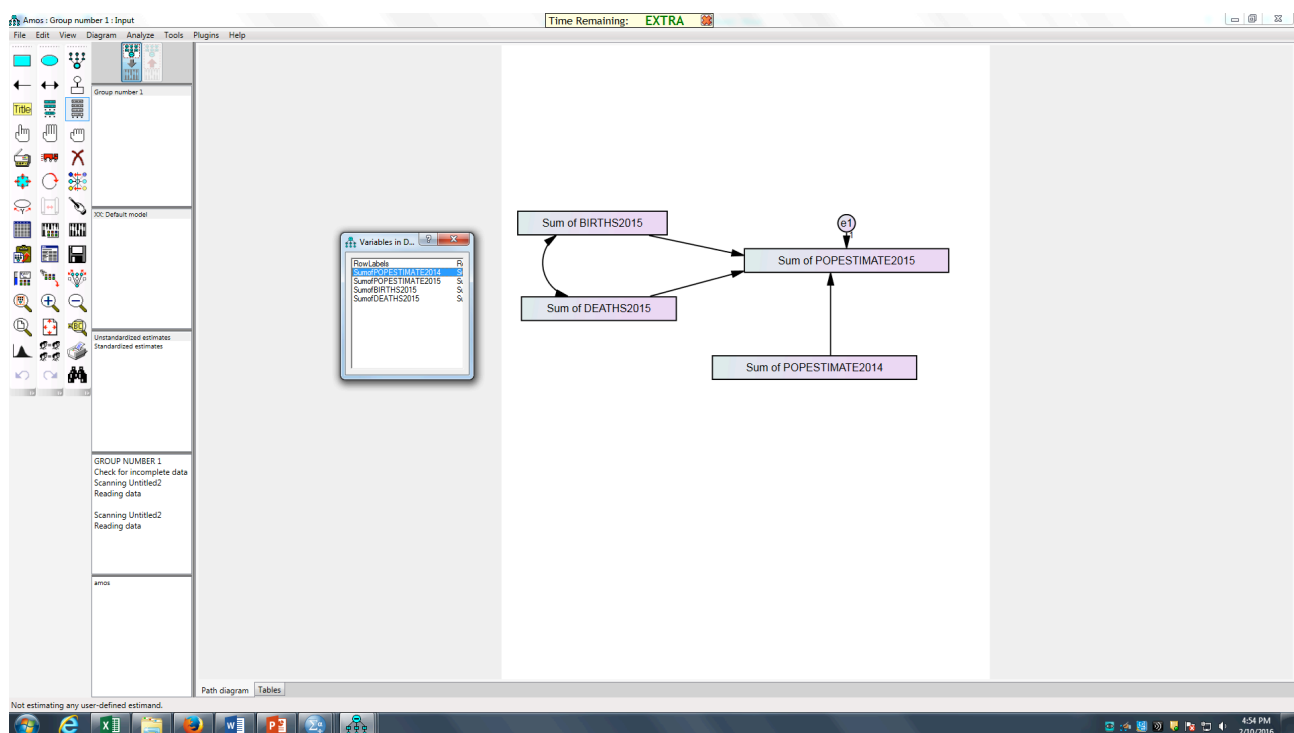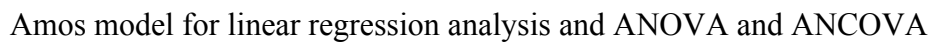
*Collect and manage data*

*The population estimate data is collected from* [http://www.census.gov/popest/data/index.html](http://www.census.gov/popest/data/index.html) *The time domain is Sep. 17, 2014*



Build the model

Excel, IBM SPSS, and R are the tools for my data analysis. The data that shows the Sum of popluation of years 2014-2015, Sum of BIRTHS2015, Sum of DEATHS2015 from all the state have been choose and processed to pivot table on Exel. Also, those data have been processed to SPSS to calculate the frequency include mean, median, mode, standard deviation, and variance. My data is a quantative data set.

Acroding to the statistics table, the population is increasing thought years 2014-2015 and the birth

of years 2015 far more than the death of 2015.

I use compare mean in one sample to t-test to compare the mean between the population between

2014-2015 and between the death and birth in years 2015.



Amos model for linear regression analysis and ANOVA and ANCOVA

The next step which is the final step. I am going to insert my data into R and do the summary and show a head of my data.

RStudio

Console ~/Desktop/

```
  unused argument (hearder = TRUE)
> cis<-read.csv("/Users/Bla/Desktop/Workbook2.csv", hearder = T)
Error in read.table(file = file, header = header, sep = sep, quote = qu
ote,  :
  unused argument (hearder = TRUE)
> cis<-read.csv("/Users/Bla/Desktop/Workbook2.csv", header  = T)
> summary(cis)
     SUMLEV       REGION    DIVISION      STATE
 Min.   :10.00   0: 1    5      : 9   Min.   : 0.00
 1st Qu.:40.00   1:10    8      : 8   1st Qu.:12.00
 Median :40.00   2:13    4      : 7   Median :27.00
 Mean   :38.07   3:18    1      : 6   Mean   :27.18
 3rd Qu.:40.00   4:14    0      : 5   3rd Qu.:41.00
 Max.   :40.00   X: 1    3      : 5   Max.   :72.00
                         (Other):17
        NAME      CENSUS2010POP       ESTIMATESBASE2010
 Alabama   : 1   Min.   :   563626   Min.   :   563767
 Alaska    : 1   1st Qu.:  1852994   1st Qu.:  1853011
 Arizona   : 1   Median :  4625364   Median :  4625401
 Arkansas  : 1   Mean   : 16315130   Mean   : 16315798
 California: 1   3rd Qu.:  9535483   3rd Qu.:  9535692
 Colorado  : 1   Max.   :308745538   Max.   :308758105
 (Other)   :51
  POPESTIMATE2014     POPESTIMATE2015      BIRTHS2014
 Min.   :   584304   Min.   :   586107   Min.   :   5965
 1st Qu.:  1882980   1st Qu.:  1896190   1st Qu.:  26147
 Median :  4829160   Median :  4858979   Median :  58334
 Mean   : 16846616   Mean   : 16977731   Mean   : 208948
 3rd Qu.:  9916306   3rd Qu.:  9922576   3rd Qu.: 114148
 Max.   :318907401   Max.   :321418820   Max.   :3958107

   BIRTHS2015       DEATHS2014        DEATHS2015
 Min.   :   5981   Min.   :   4172   Min.   :   4282
 1st Qu.:  26256   1st Qu.:  15841   1st Qu.:  15582
 Median :  58305   Median :  43063   Median :  43086
 Mean   : 210389   Mean   : 137946   Mean   : 138687
 3rd Qu.: 114521   3rd Qu.:  76120   3rd Qu.:  77052
 Max.   :3985924   Max.   :2611362   Max.   :2625033

>
```

Environment | History

Import Dataset | List

Global Environment

**Data**

cis                57 obs. of 13 variables

Files | Plots | Packages | Help | Viewer

R: Box Plots | Find in Topic

boxplot {graphics}                    R Documentation

# Box Plots

**Description**

Produce box-and-whisker plot(s) of the given (grouped) values.

**Usage**

```
boxplot(x, ...)

## S3 method for class 'formula'
boxplot(formula, data = NULL, ..., subset, na.action =

## Default S3 method:
boxplot(x, ..., range = 1.5, width = NULL, varwidth = 
        notch = FALSE, outline = TRUE, names, plot = TR
        border = par("fg"), col = NULL, log = "",
        pars = list(boxwex = 0.8, staplewex = 0.5, out
        horizontal = FALSE, add = FALSE, at = NULL)
```

| | SUMLEV | REGION | DIVISION | STATE | NAME | CENSUS2010POP |
|---|---|---|---|---|---|---|
| 1 | 10 | 0 | 0 | 0 | United States | 308745538 |
| 2 | 20 | 1 | 0 | 0 | Northeast Region | 55317240 |
| 3 | 20 | 2 | 0 | 0 | Midwest Region | 66927001 |
| 4 | 20 | 3 | 0 | 0 | South Region | 114555744 |
| 5 | 20 | 4 | 0 | 0 | West Region | 71945553 |
| 6 | 40 | 3 | 6 | 1 | Alabama | 4779736 |

| | ESTIMATESBASE2010 | POPESTIMATE2014 | POPESTIMATE2015 | BIRTHS2014 |
|---|---|---|---|---|
| 1 | 308758105 | 318907401 | 321418820 | 3958107 |
| 2 | 55318348 | 56171281 | 56283891 | 631620 |
| 3 | 66929897 | 67762069 | 67907403 | 833294 |
| 4 | 114562953 | 119795010 | 121182847 | 1521933 |
| 5 | 71946907 | 75179041 | 76044679 | 971260 |
| 6 | 4780127 | 4846411 | 4858979 | 58334 |

| | BIRTHS2015 | DEATHS2014 | DEATHS2015 |
|---|---|---|---|
| 1 | 3985924 | 2611362 | 2625033 |
| 2 | 635486 | 476784 | 479649 |
| 3 | 835105 | 591344 | 589347 |
| 4 | 1534496 | 1016353 | 1023601 |
| 5 | 980837 | 526881 | 532436 |
| 6 | 58305 | 50228 | 50330 |

The population estimate 2014 data. We will look for potential outliers in the data

```
reusing))
  undefined columns selected
> boxplot.stats(cis$POPESTIMATE2014)$out
[1] 318907401  56171281  67762069 119795010  75179041  38792291
[7]  26979078
>
```

❖I may change the coef argument to 3 (it is 1.5 by default) to identify suspected outliers.

```
> boxplot.stats(cis$POPESTIMATE2014, coef = 3)$out
[1] 318907401  56171281  67762069 119795010  75179041  38792291
>
```

❖I am going to compare the potential outliers of population estimate of year 2014 and 2015

```
> boxplot.stats(cis$POPESTIMATE2015)$out
[1] 321418820  56283891  67907403 121182847  76044679  39144818
[7]  27469114
> boxplot.stats(cis$POPESTIMATE2015, coef = 3)$out
[1] 321418820  56283891  67907403 121182847  76044679  39144818
```
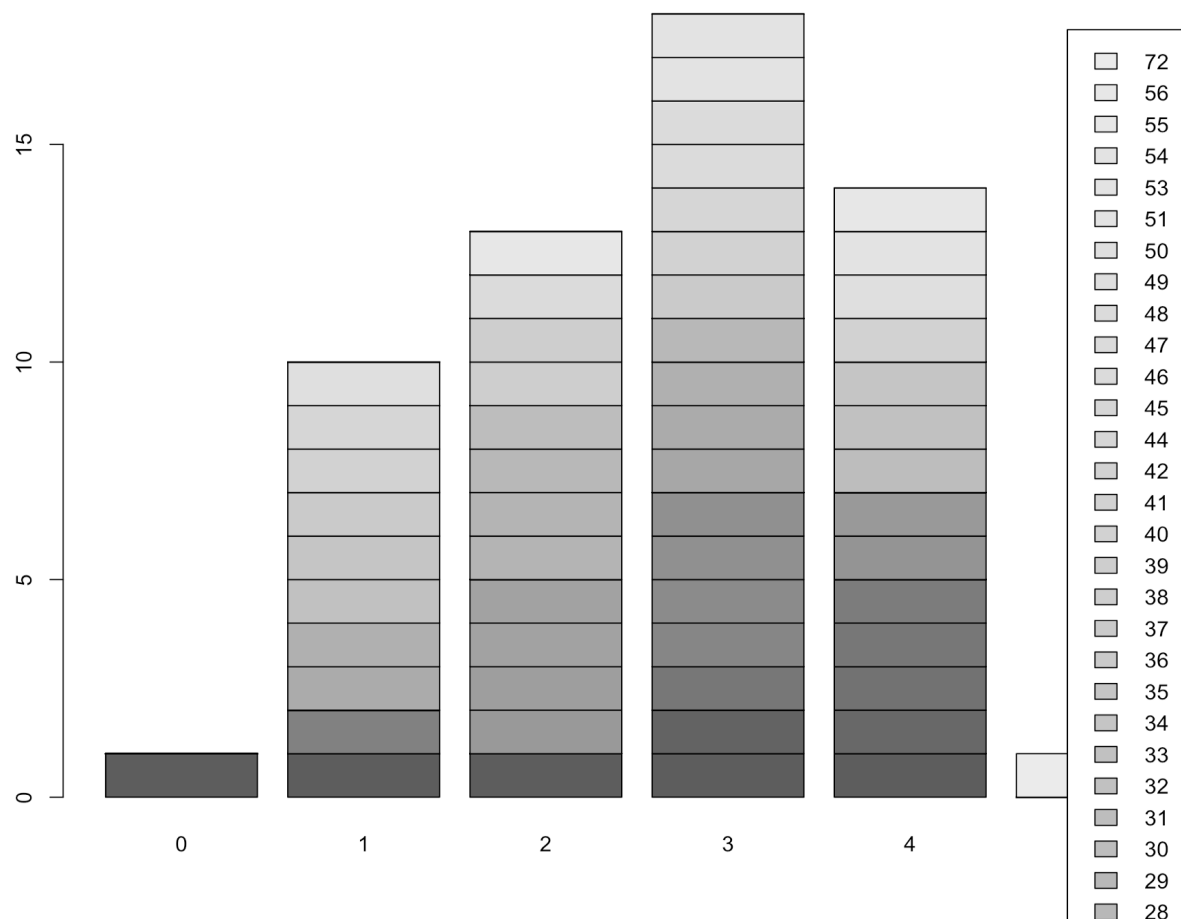
The plot of the different between the state of region and state of division

     barplot(table(cis$STATE,cis$REGION), legend.text=TRUE)
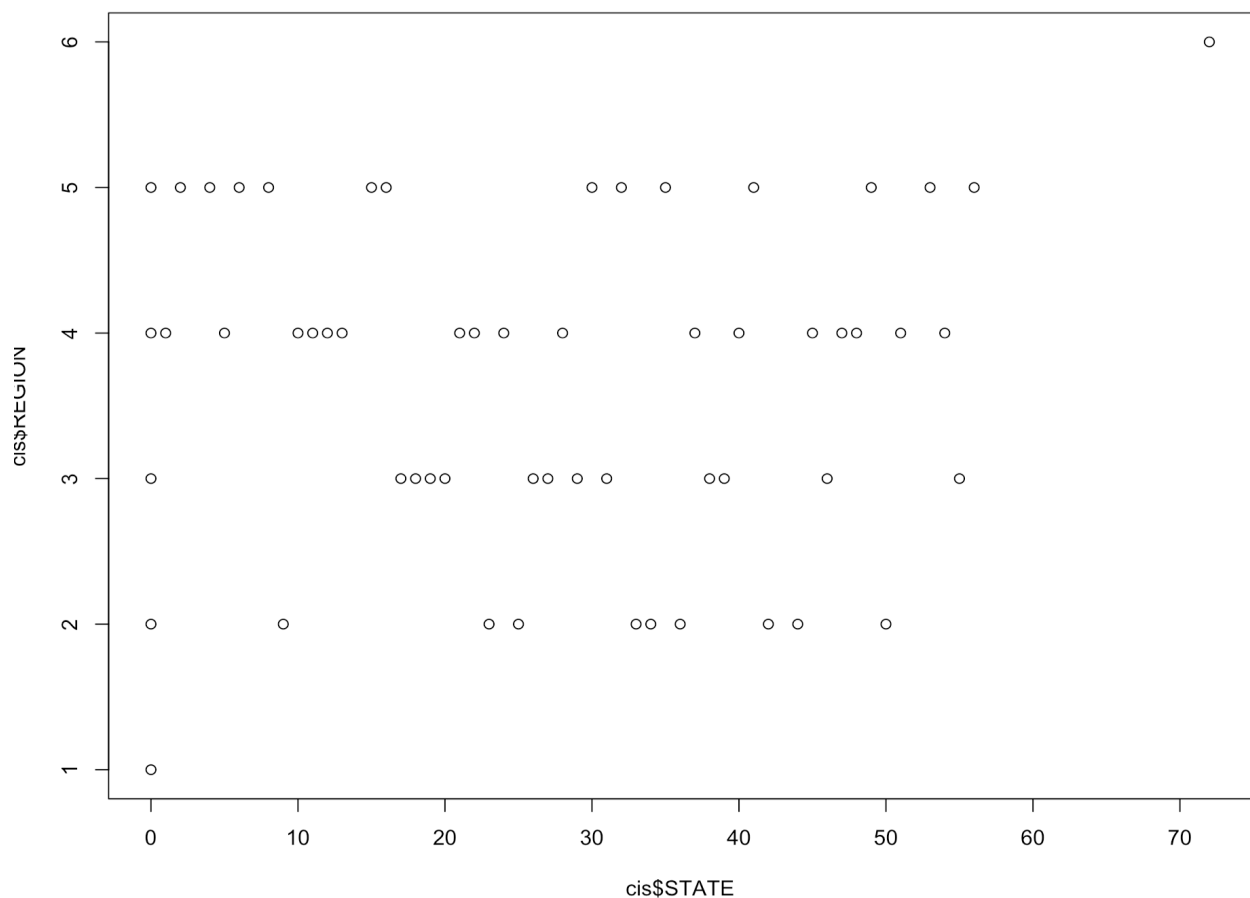     plot(cis$STATE, cis$REGION)
     The region devide by 0-1

0 is United States, 1 is Northeast Region, 2 is Midwest Region, 3 is South Region,4 is West Region

**Conclusion:** Back to the regression table and estimate, The births is going to increase 5.652 and the death will increase .610. Also, the population will increase as well to .932. Those numbers shows us the big different between the death and the births. The population of years 2014-2015 is increasing. Accroding to all the data that I was analysis, the population of US is increasing every year. On the other hand, the division of south is more than other. It divide in 3 colum. Across to the chart, we can see south is the higher bar. South is the place with most people living.