

DataVizA Tutorial: Plotting many variables: Solutions

Department of Econometrics and Business Statistics, Monash University

Tutorial 6

Swiss Data

1. Plot time series line plots of yearly Swiss exports to Germany (DE), the USA (US), China (CN) and India (IN). Facet by country.

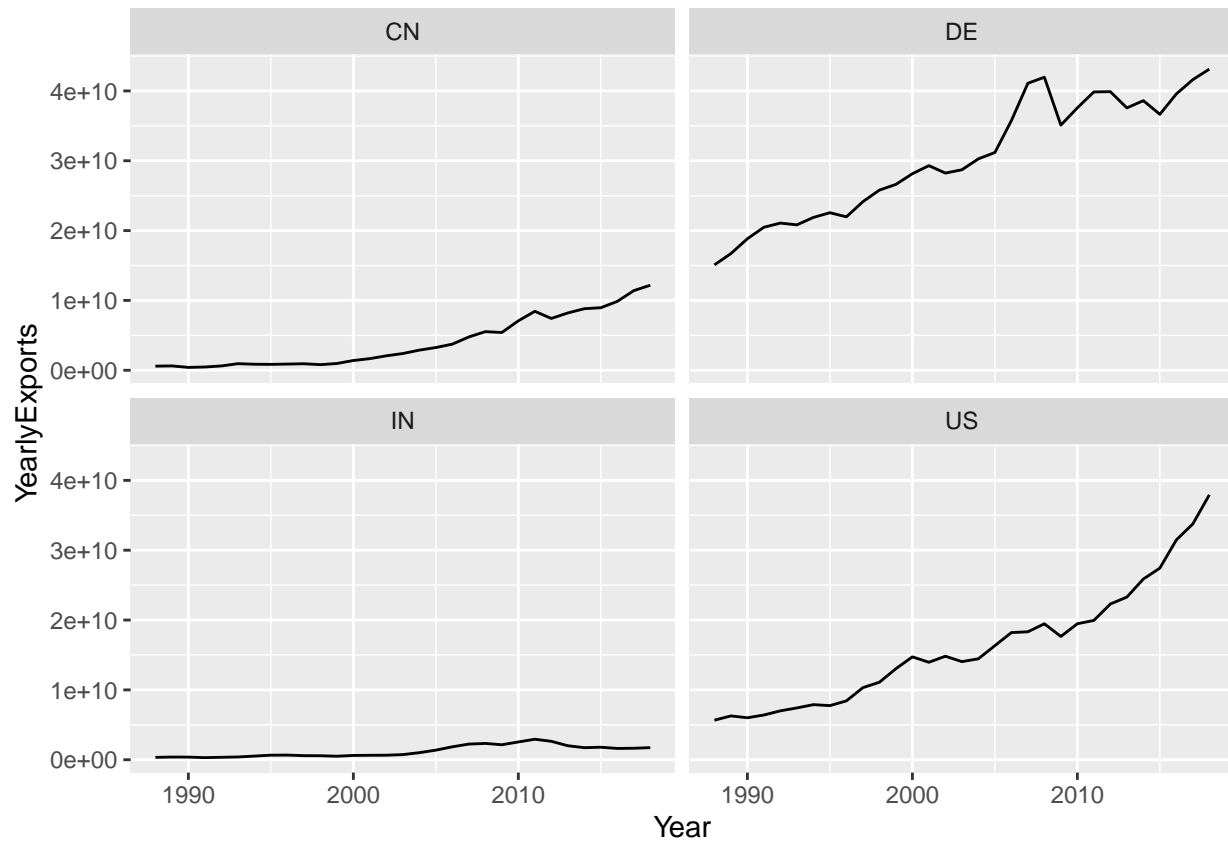
#The following block of code was all covered in the previous tutorial

```
library(tidyverse)
SwissWide<-read_csv('SwissExportsFull.csv')%>%
  rename(`NA`=X154)%>%
  pivot_longer(cols=c(-Date,-Year),
               names_to = 'Country',
               values_to = 'Exports')%>%
  group_by(Year,Country)%>%
  summarise(YearlyExports=sum(Exports))>SwissYearly
```

```
## Warning: Missing column names filled in: 'X154' [154]
```

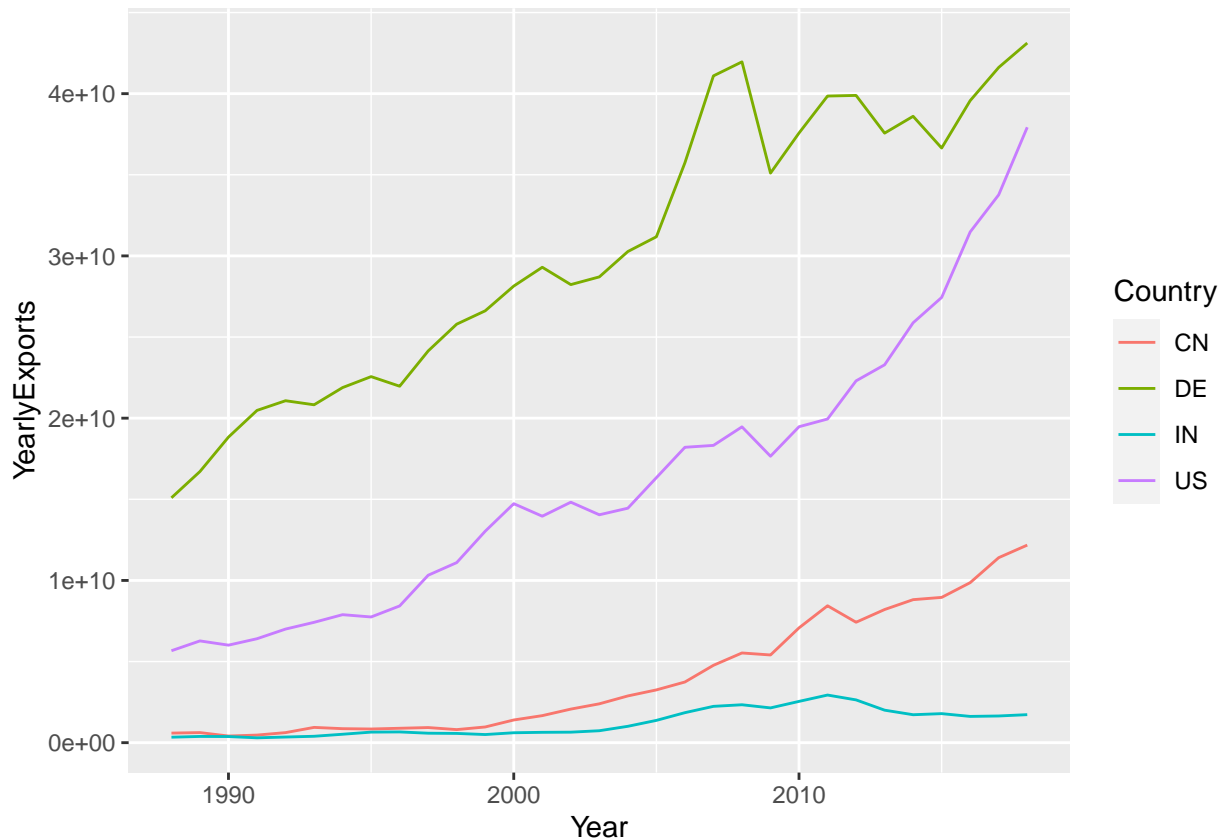
#New code below

```
SwissYearly%>%
  filter(Country %in% c('DE','US','CN','IN'))%>%
  ggplot(aes(x=Year,y=YearlyExports))+
  geom_line()+
  facet_wrap(~Country)
```



- Plot these four lines on a single plot with each country in a different colour. Hint: Use the aesthetic `col`.

```
SwissYearly%>%
  filter(Country %in% c('DE', 'US', 'CN', 'IN'))%>%
  ggplot(aes(x=Year, y=YearlyExports, col=Country))+
  geom_line()
```



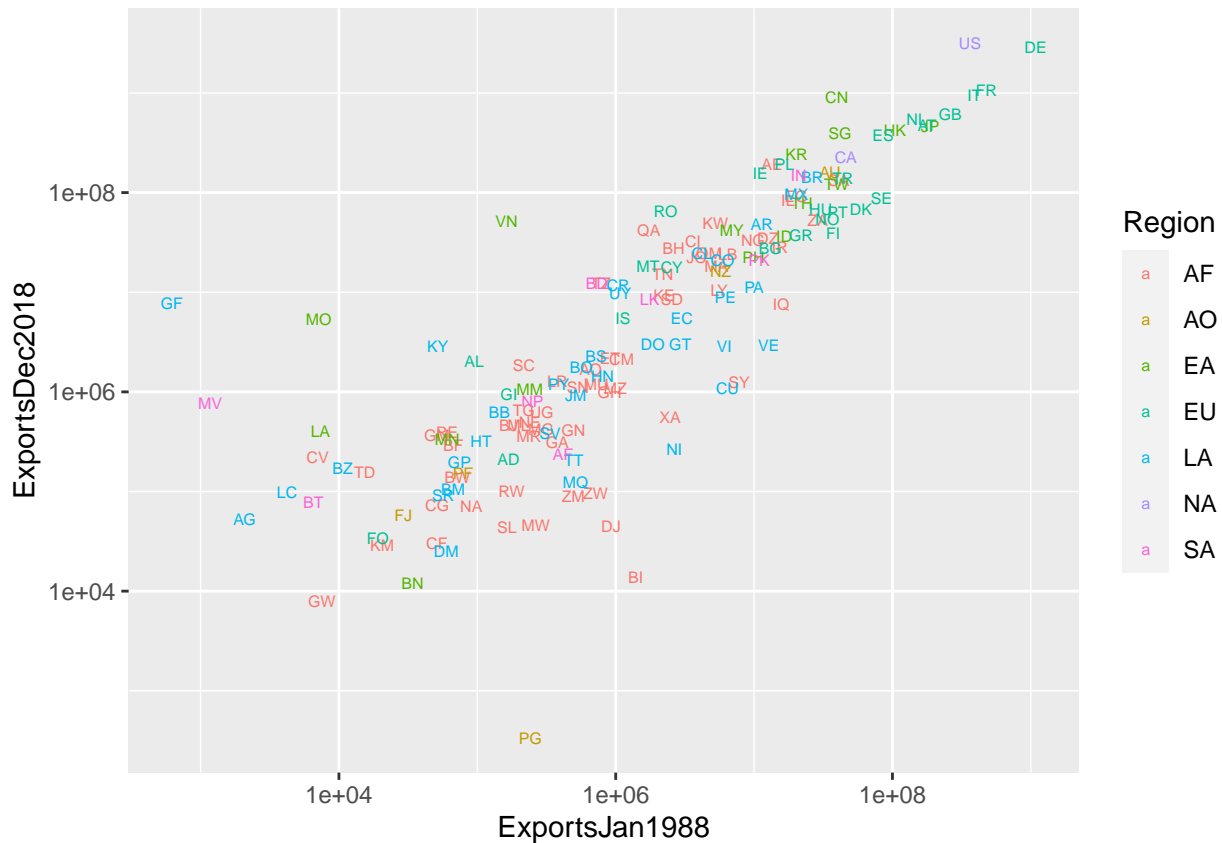
3. Comment on these plots

Both plots show the same information, Swiss exports are mostly trending upwards, Germany and the US are much bigger trading partners for Switzerland than China and India. Any seasonality has been lost by taking the yearly aggregate.

The second plot allows us to more easily line up the impact of the GFC. There was a drop in exports to all markets but this was most pronounced in Germany and hardly there at all in China and India. Post GFC the recovery in exports has been quickest and strongest for the US and China. Exports to Germany were volatile but have picked up in the last three years while growth in exports to India has stagnated

4. Construct the same scatterplot from Question 3 of tutorial 3. Use text to represent the country and color to represent the region. Note these regions are very loosely defined for instance Mexico could easily be in North America, and the Central Asian (CA) could be better named the former Soviet region, etc.

```
SwissExp<-readRDS('SwissExport.rds')
SwissExp%>%
  filter((ExportsJan1988!=0)&(ExportsDec2018!=0))%>%
  ggplot(aes(x=ExportsJan1988,
             y=ExportsDec2018,
             label=Country,
             col=Region))+
  geom_text(size=2)+
  scale_x_log10()+scale_y_log10()
```

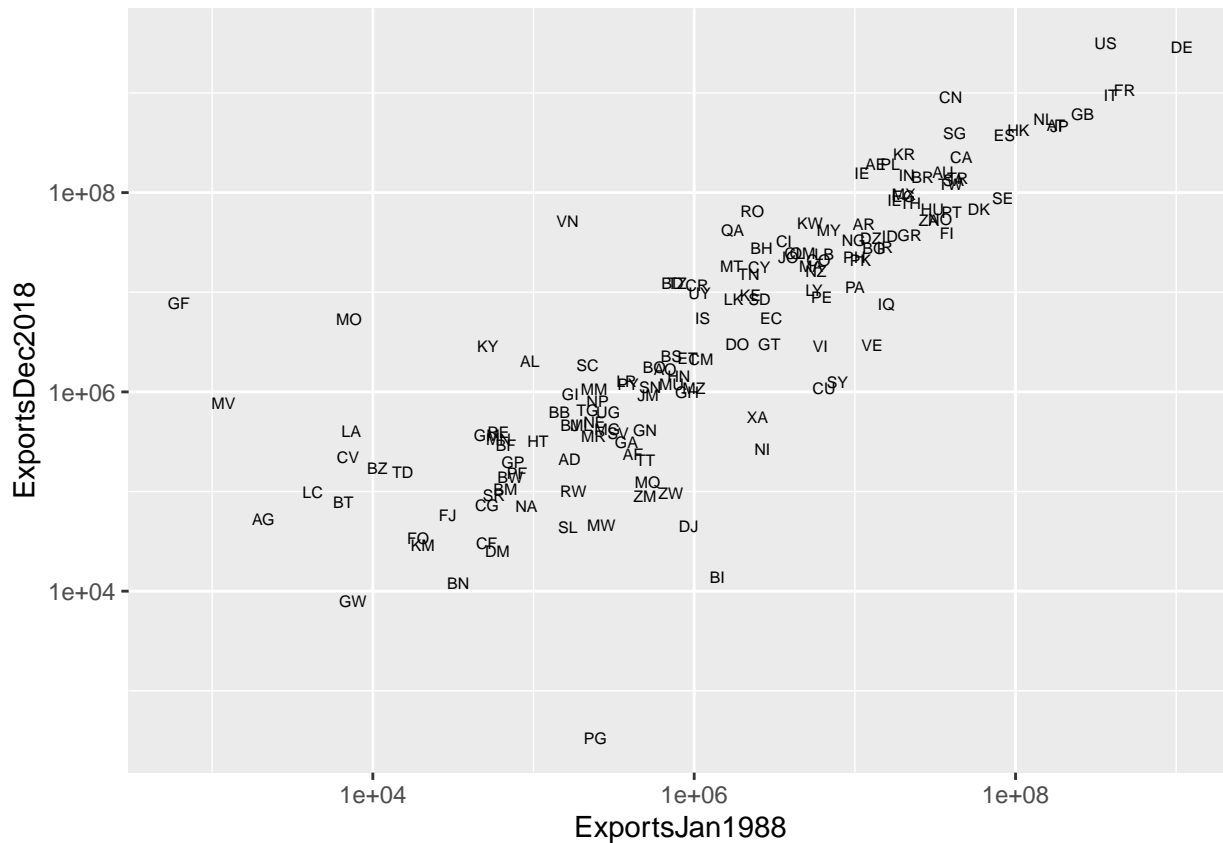


2. Discuss this plot.

Some interesting things to note are that the many EU nations occupy the top right hand side of the plot. This is to be expected since Switzerland is in the middle of Europe and exports mostly to its neighbours. Note that Frech Guyana (GF) went from having almost no exports in January 1988 to a moderately high level of exports in December 2018, while Papua New Guinea (PG) went in the opposite direction. It should be noted that monthly data are more volatile and a more accurate understanding could be obtained by looking at a yearly average.

3. Since the previous plots suffer from a little overplotting, repeat the analysis but facet by region.

```
library(tidyverse)
SwissExp %>%
  filter((ExportsJan1988!=0)&(ExportsDec2018!=0)) %>%
  ggplot(aes(x=ExportsJan1988,
             y=ExportsDec2018,
             label=Country))+
  geom_text(size=2)+
  scale_x_log10()+scale_y_log10()
```



```
facet_wrap(~Region,nrow = 4)
```

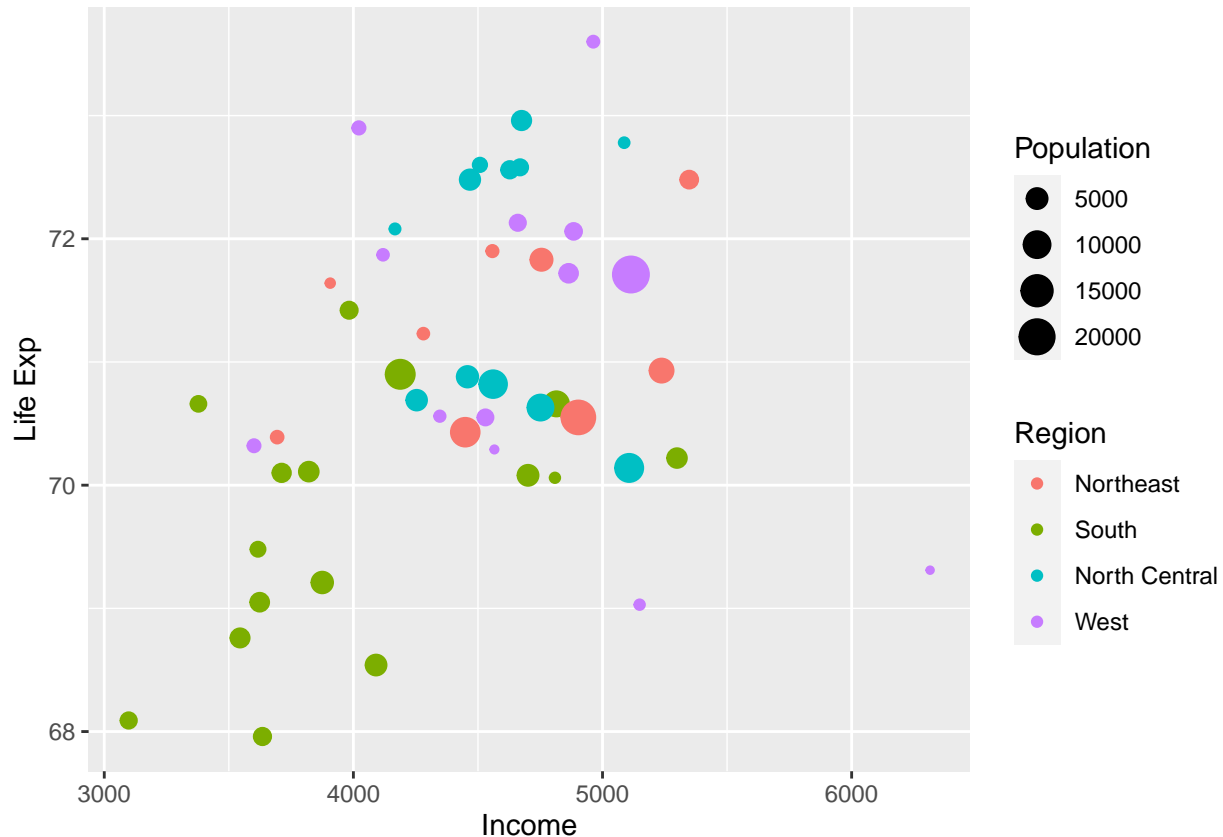
```
## <ggproto object: Class FacetWrap, Facet, gg>
##   compute_layout: function
##   draw_back: function
##   draw_front: function
##   draw_labels: function
##   draw_panels: function
##   finish_data: function
##   init_scales: function
##   map_data: function
##   params: list
##   setup_data: function
##   setup_params: function
##   shrink: TRUE
##   train_scales: function
##   vars: function
##   super: <ggproto object: Class FacetWrap, Facet, gg>
```

U.S. States Data

The data in the file *USStateData.rds* contains information on demography and geography for the states of the United States. These data were measured in 1977.

4. Construct a bubble chart of income per capita (x axis) against life expectancy (y axis). The size of the bubble should reflect the State population, and the color should reflect the region.

```
USState<-readRDS('USStateData.rds')
ggplot(USState,
  aes(x=Income,
    y='Life Exp', #Note use of ``
    col=Region, #Region mapped to color
    size=Population))+ #Country mapped to label
  geom_point()
```



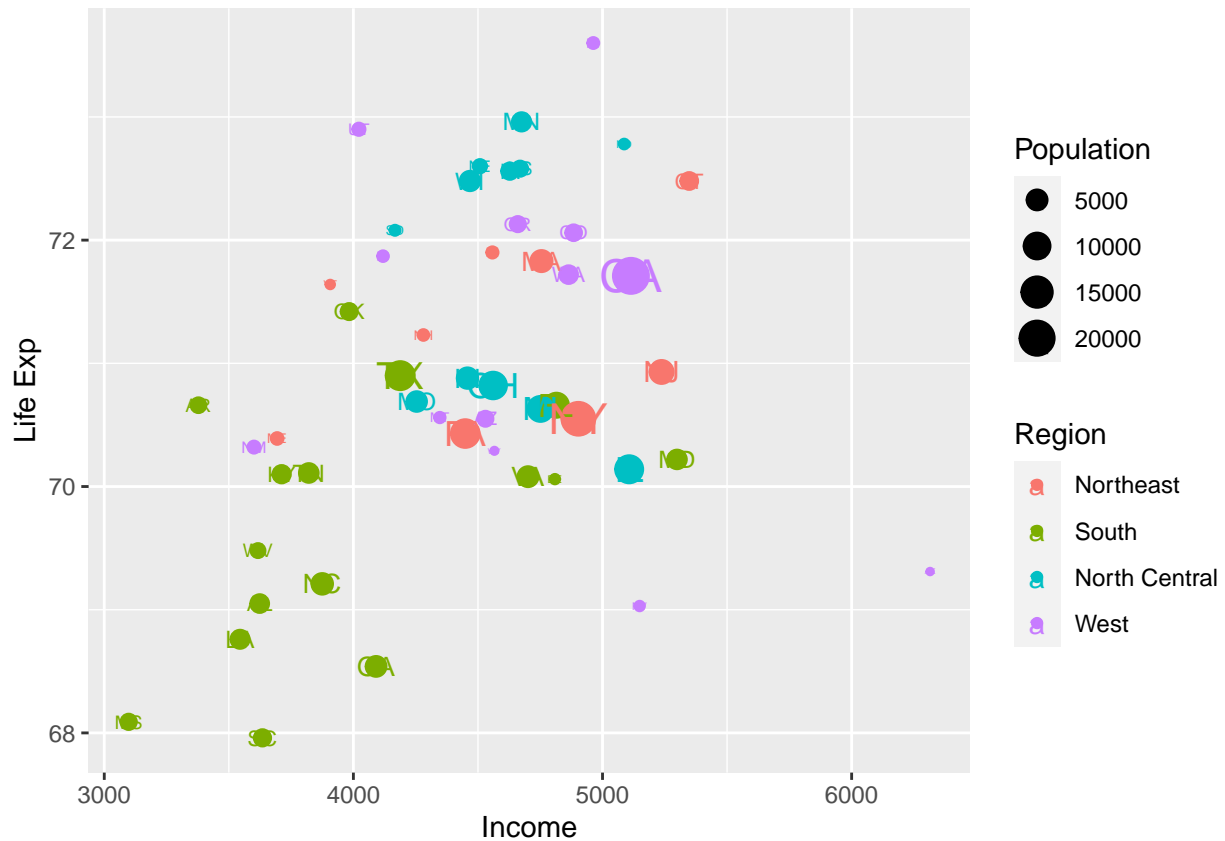
5. Discuss this plot.

Some things to notice are the following. The South region has mostly lower income states with lower life expectancy. Generally the larger states have higher income and moderately high life expectancy. There is a small state that is an outlier for having the highest per capita income.

6. Add the state abbreviations as text on the plot

*#The obvious thing to do would be to take the previous code
#and add geom_text and a label aesthetic*

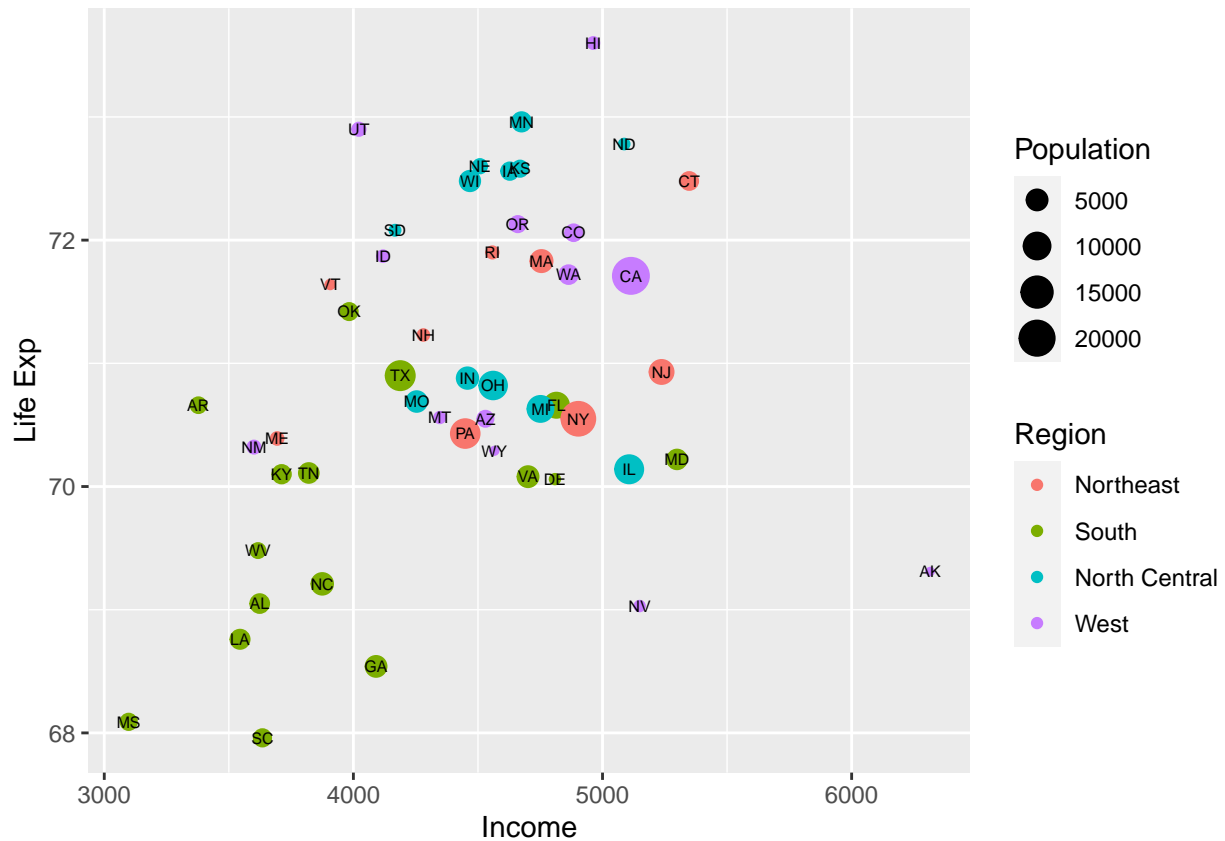
```
ggplot(USState,
  aes(x=Income,
    y='Life Exp', #Note use of ``
    col=Region, #Region mapped to color
    size=Population,
    label=Abbreviation))+ #Country mapped to label
  geom_point()+
  geom_text()
```



*#This looks bad since the text labels do not need to inherit color
#or size.*

- If you are not satisfied with your answer to question 6, redo the question. Hint: you may need to specify some aesthetics within a geom.

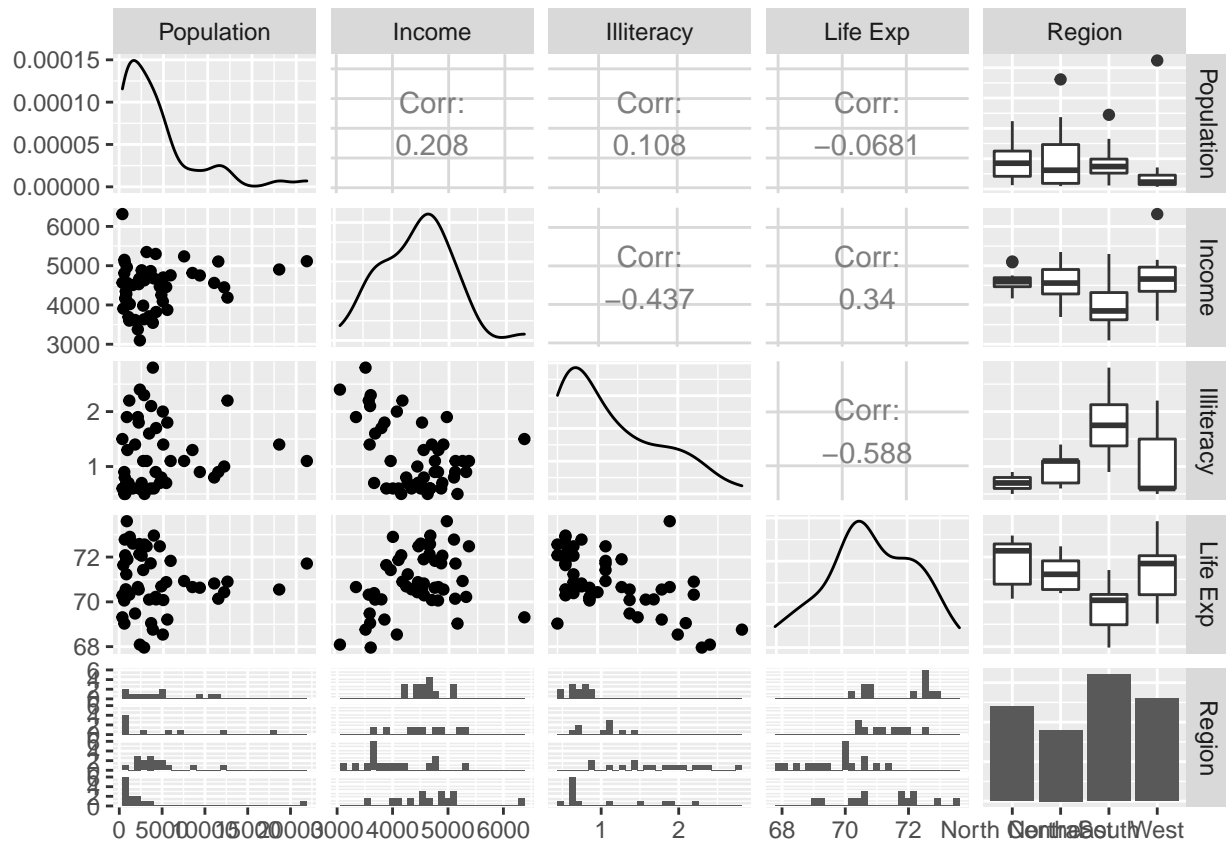
```
ggplot(USState,
  aes(x=Income,
    y=`Life Exp`))+ #Country mapped to label
  geom_point(aes(col=Region,
    size=Population))+ #color only in point
  geom_text(aes(label=Abbreviation),size=2)
```



*#Note we can now clearly see more things. For instance the
 #outlying state for income is Alaska (AK). It is not too suprising
 #since Alaska produces oil and has a small population. Similarly
 #it is not too suprising that life expectancy is lower since
 #Alaska is very cold.*

- The file `USStateRed.rds` contains the same data but with the variables `State`, `Abbreviation`, `Frost`, `Murder` and `HS Grad` stripped out. Load this data and plot a ggpairs plot.

```
library(GGally)
StateRed<-readRDS('USStateRed.rds')
ggpairs(StateRed)
```

9. Discuss some interesting features of the plot.

Just a few things are the following. The marginal distributions (densities) show that state population is right skewed, most states are small but there are a few big ones. Life expectancy is correlated highly with other socioeconomic variables such as income and illiteracy but not as much with either State size or area. The correlation between state area and state population is surprisingly close to zero. However the scatterplot suggests that an outlier (Alaska again) may be having an impact on this result. The Western states tend to have skewed distributions for illiteracy and life expectancy, most Western states do well (low illiteracy, high life expectancy) but there are a few western states with extreme values in the opposite direction.