

Exploring Relationships

- [Correlation in Scatterplots](#)
- [Changing ggplot themes](#)
- [What to do with outliers](#)
- [Smoother and linear regression with and without confidence interval bands](#)
- [Creating a linear model and understanding linear regression output](#)
- [Ggally for plotting pairs of variables](#)
- [Correlation plot to assess collinearity](#)
- [Performing multiple regression analysis and understanding regression plots](#)
- [Add interactivity with Plotly](#)
- [Line and dot-and-line and bar charts with Food Stamps Data](#)
- [Color Brewer](#)
- [Week 6 Homework Assignment](#)

Click here for the [rpubs link](#) to this document with the following notes.

Correlation and Scatterplots

Create a Scatterplot

In this example, look at US crime rates at the state level, in 2005, with rates per 100,000 population for crime types such as murder, robbery, and aggravated assault, as reported by the Census Bureau. There are 7 crime types in total. The dataset is clean to begin with.

```
library(tidyverse)
library(ggfortify)
library(htmltools)
library(plotly) ##

crime <- read_csv('http://datasets.flowingdata.com/crimeRatesByState2005.csv')

## Parsed with column specification:
## cols(
##   state = col_character(),
##   murder = col_double(),
##   forcible_rape = col_double(),
##   robbery = col_double(),
##   aggravated_assault = col_double(),
##   burglary = col_double(),
```

```
## larceny_theft = col_double(),
## motor_vehicle_theft = col_double(),
## population = col_double()
## )

# source: U.S. Census Bureau and Nathan Yau
```

Check out the first few lines

state	murder	forcible_rape	robbery	aggravated_assault	burglary	larceny_theft
<chr>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>
United States	5.6	31.7	140.7	291.1	726.7	2286.3
Alabama	8.2	34.3	141.4	247.8	953.8	2650.0
Alaska	4.8	81.1	80.9	465.1	622.5	2599.1
Arizona	7.5	33.8	144.4	327.4	948.4	2965.2
Arkansas	6.7	42.9	91.1	386.8	1084.6	2711.2
California	6.9	26.0	176.1	317.3	693.3	1916.5

6 rows | 1-7 of 9 columns

Notice

The data has a column for the state and then the rest are rates for various crimes. Now make a quick scatterplot.

Map variables in the data onto the X and Y axes and change the axes labels and theme

The default gray theme of ggplot2 has a rather academic look. See [here](#) and [here](#) for how to use the theme option to customize individual elements of a chart. Use one of the ggplot2 built-in themes, and then customize the fonts.

Theme minimal

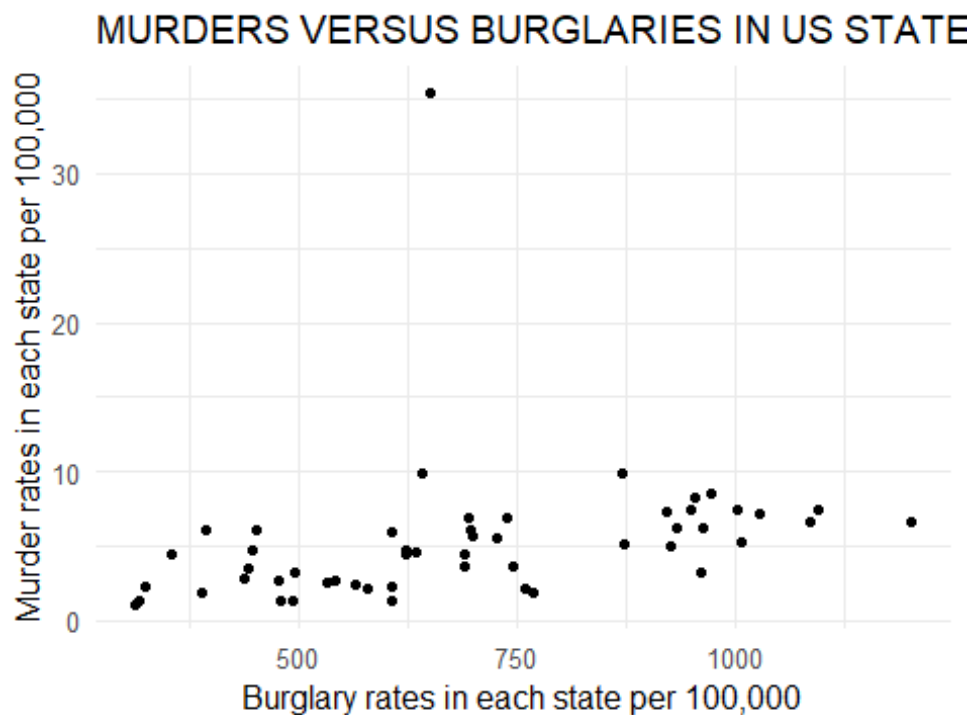
```
ggplot(crime, aes(x = burglary, y = murder)) +
  xlab("Burglary rates in each state per 100,000") +
```

```
ylab("Murder rates in each state per 100,000") +  
theme_minimal(base_size = 12) # Change the theme
```



Add a caption

```
p1 <- ggplot(crime, aes(x = burglary, y = murder)) +  
  labs(title = "MURDERS VERSUS BURGLARIES IN US STATES PER 100,000",  
        caption = "Source: U.S. Census Bureau and Nathan Yau") +  
  xlab("Burglary rates in each state per 100,000") +  
  ylab("Murder rates in each state per 100,000") +  
  theme_minimal(base_size = 12)  
p1 + geom_point()
```



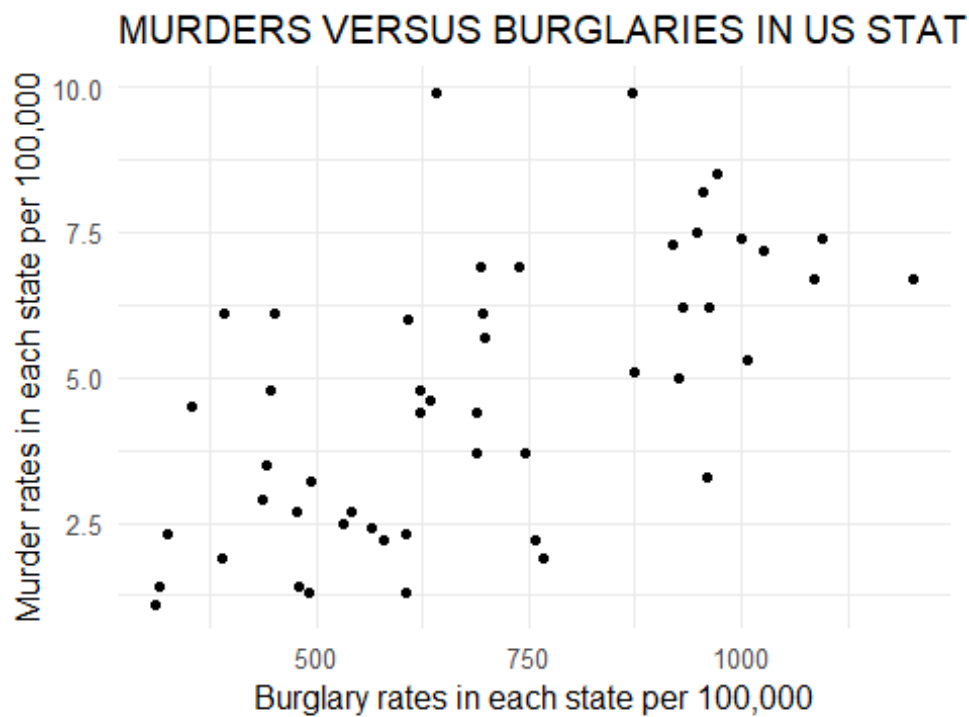
Source: U.S. Census Bureau and Nathan Yau

What is going on with the outlier?

The one point far higher than the rest represents Washington, D.C., which had a much higher murder rate of 35.4. The states with the next highest murder rate at that time were Louisiana and Maryland at 9.9 per 100,000.

Remove D.C. and US averages and replot:

```
crime2 <- crime[crime$state != "District of Columbia",]
crime2 <- crime2[crime2$state != "United States",]
p2 <- ggplot(crime2, aes(x = burglary, y = murder)) +
  labs(title = "MURDERS VERSUS BURGLARIES IN US STATES PER 100,000",
    caption = "Source: U.S. Census Bureau and Nathan Yau") +
  xlab("Burglary rates in each state per 100,000") +
  ylab("Murder rates in each state per 100,000") +
  theme_minimal(base_size = 12)
p2 + geom_point()
```

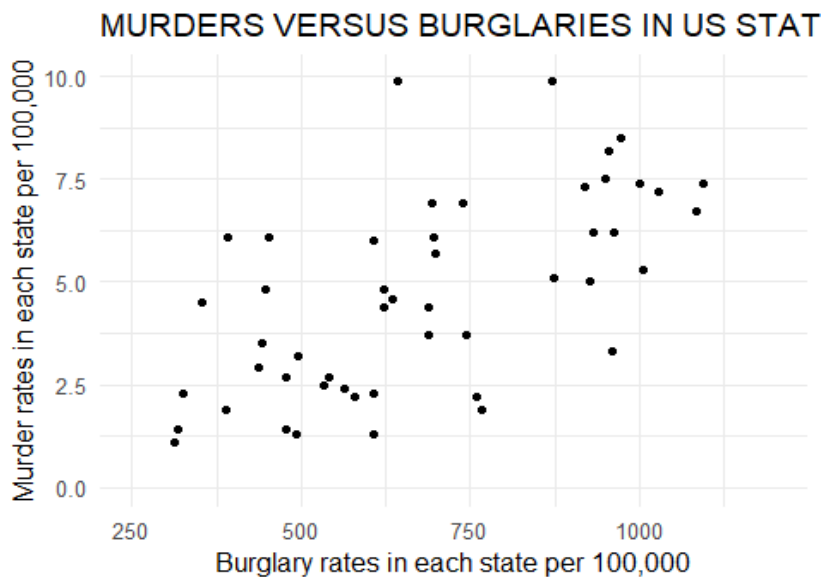


Source: U.S. Census Bureau and Nathan Yau

Now the scatterplot appears to show a correlation

Fix the axes to start at 0.

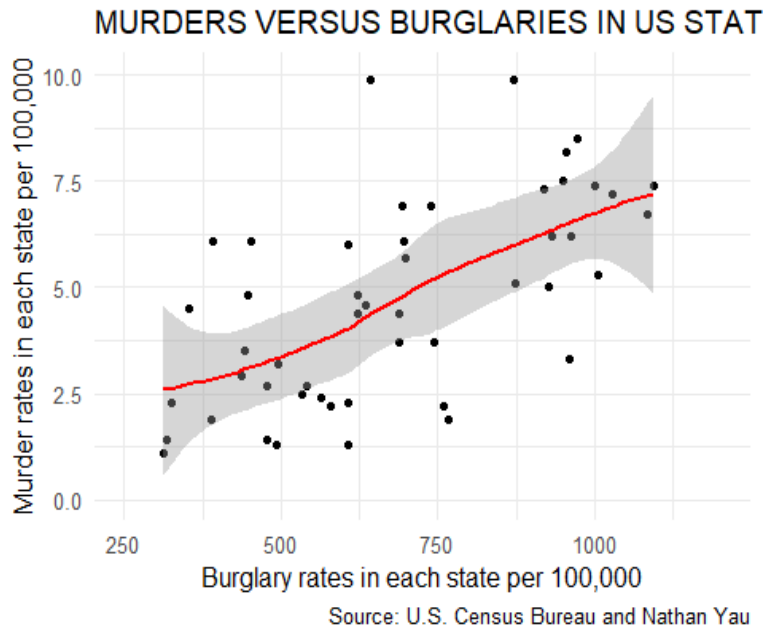
```
p3 <- p2 + xlim(250,1200)+ ylim(0,10)
p3 + geom_point()
```



Source: U.S. Census Bureau and Nathan Yau

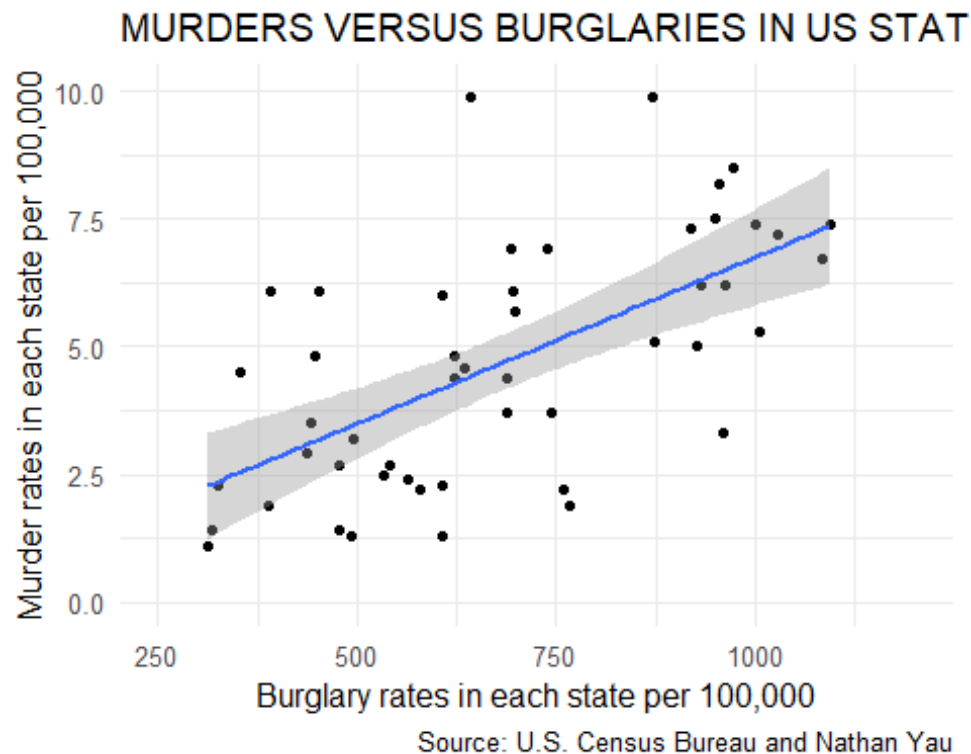
Add a smoother in red with a confidence interval

```
p4 <- p3 + geom_point() + geom_smooth(color = "red")
p4
```



Add a linear regression with confidence interval

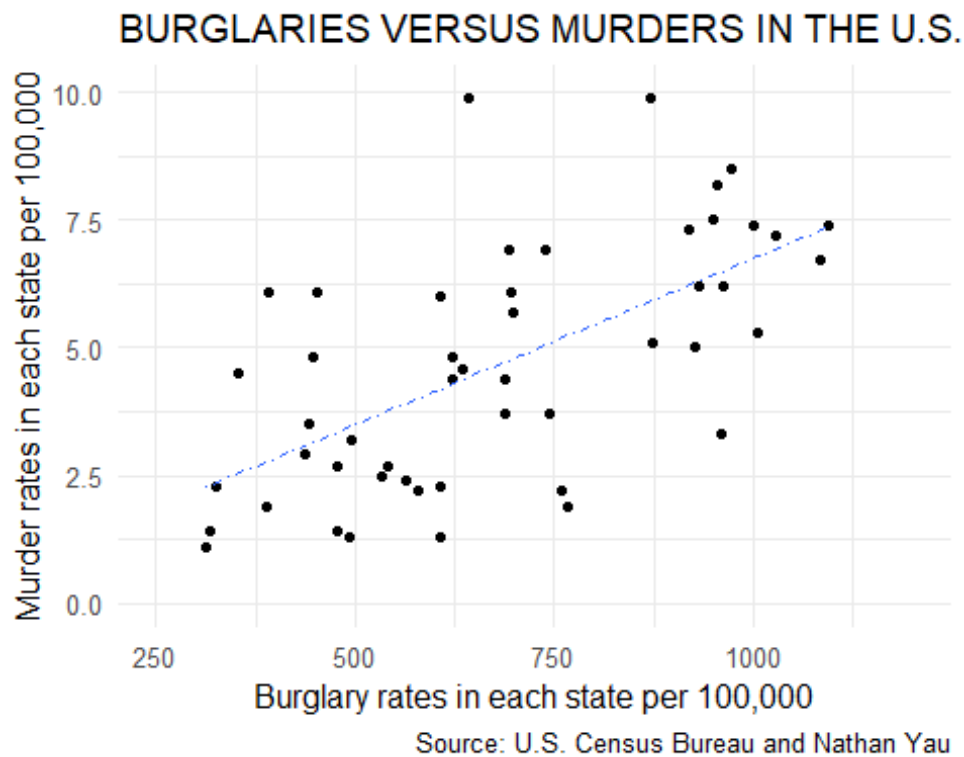
```
p5 <- p3 + geom_point() + geom_smooth(method='lm',formula=y~x)
p5
```



Add a title, make the line dashed, and remove the confidence interval band

The command `se = FALSE` takes away the CI band

```
p6 <- p3 + geom_point() + geom_smooth(method='lm', formula=y~x, se = FALSE, linetype="dotdash", size = 0.3) +  
  ggtitle("BURGLARIES VERSUS MURDERS IN THE U.S.")  
p6
```



What is the linear equation of that linear regression model?

In the form, $y = mx + b$, we use the command, `lm(y~x)`, meaning, fit the predictor variable x into the model to predict y . Look at the values of (Intercept) and murder. The column, Estimate gives the value you need in your linear model. The column for Pr(>|t|) describes whether the predictor is useful to the model. The more asterisks, the more the variable contributes to the model.

Calculate the correlation coefficient and model summary

```
cor(crime2$burglary, crime2$murder)  
[1] 0.6231757
```

```
fit1 <- lm(murder ~ burglary, data = crime2)
summary(fit1)
```

Call:
lm(formula = murder ~ burglary, data = crime2)

Residuals:

Min	1Q	Median	3Q	Max
-3.2924	-1.2156	-0.2142	1.1749	5.4978

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	0.395519	0.825748	0.479	0.634
burglary	0.006247	0.001132	5.521	1.34e-06 ***

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 1.87 on 48 degrees of freedom
Multiple R-squared: 0.3883, Adjusted R-squared: 0.3756
F-statistic: 30.48 on 1 and 48 DF, p-value: 1.342e-06

What does the output mean?

Cor stands for “correlation”. This is a value between (inclusively) -1 and 1. The correlation coefficient tells how strong or weak the correlation is. Values closer to +/- 1 are strong correlation (the sign is determined by the linear slope), values close to +/- 0.5 are weak correlation, and values close to zero have no correlation.

The model has the equation:

$$\text{murder} = 0.0062(\text{burglary}) + 0.396$$

The slope may be interpreted in the following: For each additional burglary per 100,000, there is a predicted increase of 0.006 murders.

The p-value on the right of burglary has 3 asterisks which suggests it is a meaningful variable to explain the linear increase in murders.

But we also need to look at the Adjusted R-Squared value. It states that about 38% of the variation in the observations may be explained by the model. In other words, 62% of the variation in the data is likely not explained by this model.

What about more variables?

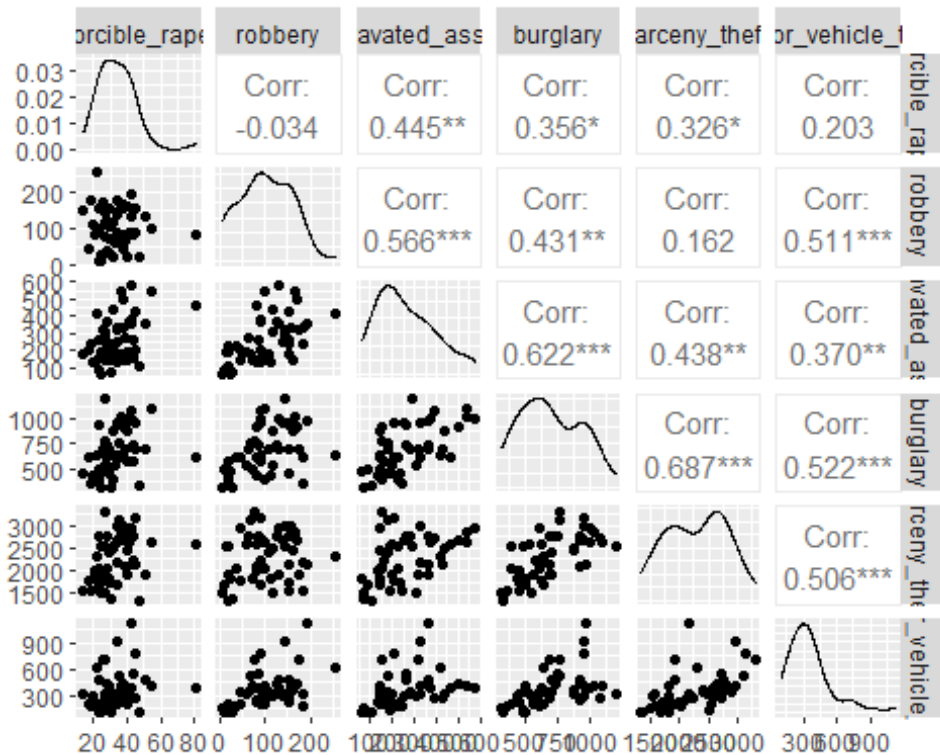
Can a model with more predictors also be used? What would we be trying to predict?

Is there an easier way to compare multiple variables using a scatterplot matrix?

Check out the pairwise comparisons with density curves and correlation output

Library GGally

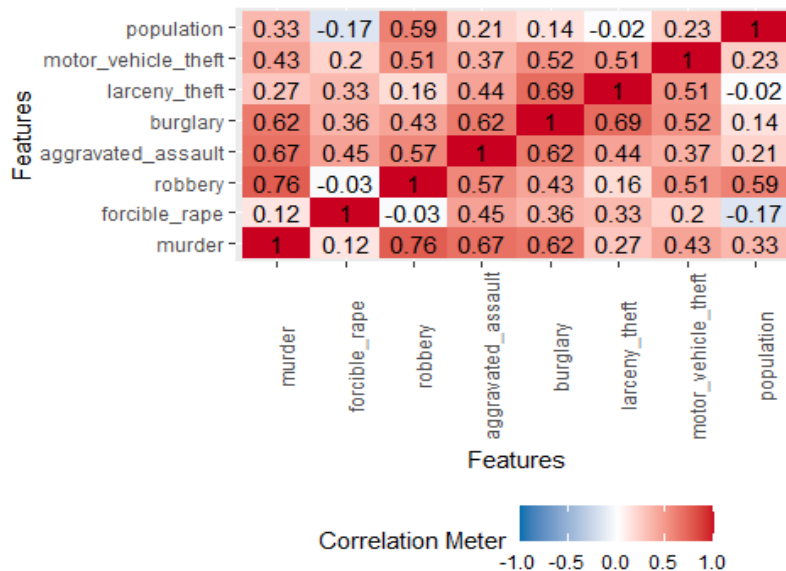
```
library(GGally)
ggpairs(crime2, columns = 3:8) # only include predictor variables in the matrix
```



Another method: Use a correlation plot to explore the correlation among all variables

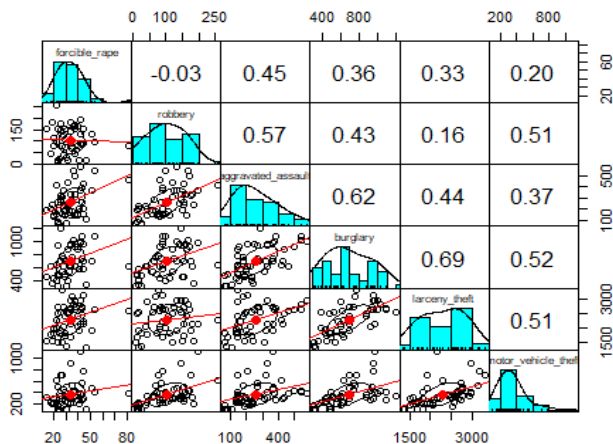
This correlation plot shows similar pairwise results as above, but in a heatmap of correlation values.

```
#install.packages("DataExplorer")
library(DataExplorer)
plot_correlation(crime2)
```



A third option to explore correlations using library(psych)

```
library(psych)
pairs.panels(crime2[3:8], # plot distributions and correlations for all the data
a
gap = 0,
pch = 21,
lm = TRUE)
```



Collinearity

The key goal of multiple regression analysis is to isolate the relationship between EACH INDEPENDENT VARIABLE and the DEPENDENT VARIABLE.

COLLINEARITY means explanatory variables are correlated and thus NOT INDEPENDENT. The more correlated the variables, the more difficult it is to change one variable without changing the other. This is important to

keep in mind. The two different matrices gave slightly different correlation information. We are concerned with dependence of 2 or more variables.

The two variables with the highest correlation of 0.68 or 0.69 are burglary and larceny_theft.

Now try to make a multiple regression model.

With multiple regression, there are several strategies for comparing variable inputs into a model. I will show you backward elimination. In backward elimination, start with all possible predictor variables with your response variable. In this case, we will use: burglary forcible_rape aggravated_assault larceny_theft motor_vehicle_theft Perform a model fit with all predictors.

1. Look at the p-value for each variable - if it is relatively small (< 0.10), then it is likely contributing to the model.
2. Check out the residual plots. A good model will have a relatively straight horizontal red line across the scatterplot between residuals plotted with fitted values (see below for a good residuals plot). You can also look at the other plots (Normal QQ, Scale-Location, and Residuals vs Leverage), but for now we will focus on the residual vs. fitted plot. The more curved the red line, the more likely that a better model exists.
3. Look at the output for the Adjusted R-Squared value at the bottom of the output. The interpretation is:

__% (from the adjusted r-squared value) of the variation in the observations may be explained by this model. The higher the adjusted R-squared value, the better the model. We use the adjusted R-squared value because it compensates for more predictors mathematically increasing the normal R-squared value.

Model fit2

```
fit2 <- lm(murder ~ robbery + burglary + forcible_rape + aggravated_assault + larceny_theft + motor_vehicle_theft, data = crime2)
summary(fit2)
```

Call:

```
lm(formula = murder ~ robbery + burglary + forcible_rape + aggravated_assault + larceny_theft + motor_vehicle_theft, data = crime2)
```

Residuals:

Min	1Q	Median	3Q	Max
-2.7088	-0.7961	-0.0508	0.6757	3.4723

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)	
(Intercept)	0.9940985	1.0835208	0.917	0.364014	
robbery	0.0194331	0.0052193	3.723	0.000567	***
burglary	0.0041431	0.0013339	3.106	0.003352	**
forcible_rape	-0.0126884	0.0210395	-0.603	0.549627	
aggravated_assault	0.0045161	0.0023433	1.927	0.060576	.
larceny_theft	-0.0007841	0.0005622	-1.395	0.170246	
motor_vehicle_theft	-0.0002426	0.0012751	-0.190	0.849982	

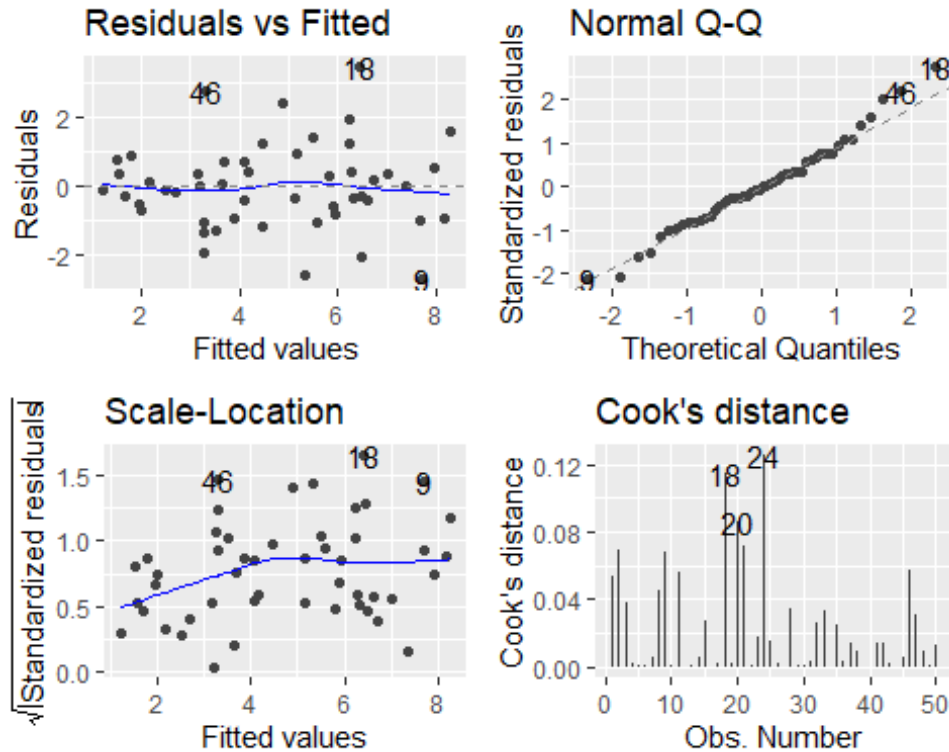
Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 1.338 on 43 degrees of freedom

Multiple R-squared: 0.7193, Adjusted R-squared: 0.6801

F-statistic: 18.36 on 6 and 43 DF, p-value: 1.949e-10

```
autoplot(fit2, 1:4, nrow=2, ncol=2)
```



What does these diagnostic plots mean?

1. Residual plot essentially indicates whether a linear model is appropriate - you can see this by the blue line showing relatively horizontal. If it is not relatively horizontal, a linear plot may not be appropriate.
2. QQPlot indicates whether the distribution is relatively normal. Observations that might be outliers are indicated by their row number.
3. Scale-Location indicates homogeneous variance (homoscedasticity). Influential observations that are skewing the variance distribution are indicated.
4. Cook's Distance indicates which outliers have high leverage, meaning that some outliers may not cause the model to violate basic assumptions required for the regression analysis (see #1-3). If outliers have high leverage, then they may be causing problems for your model. You can try to remove those observations, especially if they appear in any of the other 3 plots above.

What are we really trying to predict?

If we are trying to predict murder rates, then we can see if any of the predictor variables contribute to this model. Note the adjusted R-squared value is 68.01% The only variable that does not appear to be as significant as the others is motor_vehicle_theft. So drop that and re-run the model.

Model fit3

```
fit3 <- lm(murder ~ robbery + burglary + forcible_rape + aggravated_assault + lar  
ceny_theft, data = crime2)  
summary(fit3)
```

Call:

```
lm(formula = murder ~ robbery + burglary + forcible_rape + aggravated_assault +  
    larceny_theft, data = crime2)
```

Residuals:

Min	1Q	Median	3Q	Max
-2.6923	-0.7545	-0.0751	0.6404	3.4836

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)	
(Intercept)	1.0611101	1.0134089	1.047	0.300785	
robbery	0.0189486	0.0045060	4.205	0.000126	***
burglary	0.0041189	0.0013131	3.137	0.003044	**
forcible_rape	-0.0134321	0.0204456	-0.657	0.514623	
aggravated_assault	0.0046152	0.0022596	2.042	0.047124	*
larceny_theft	-0.0008229	0.0005181	-1.588	0.119349	

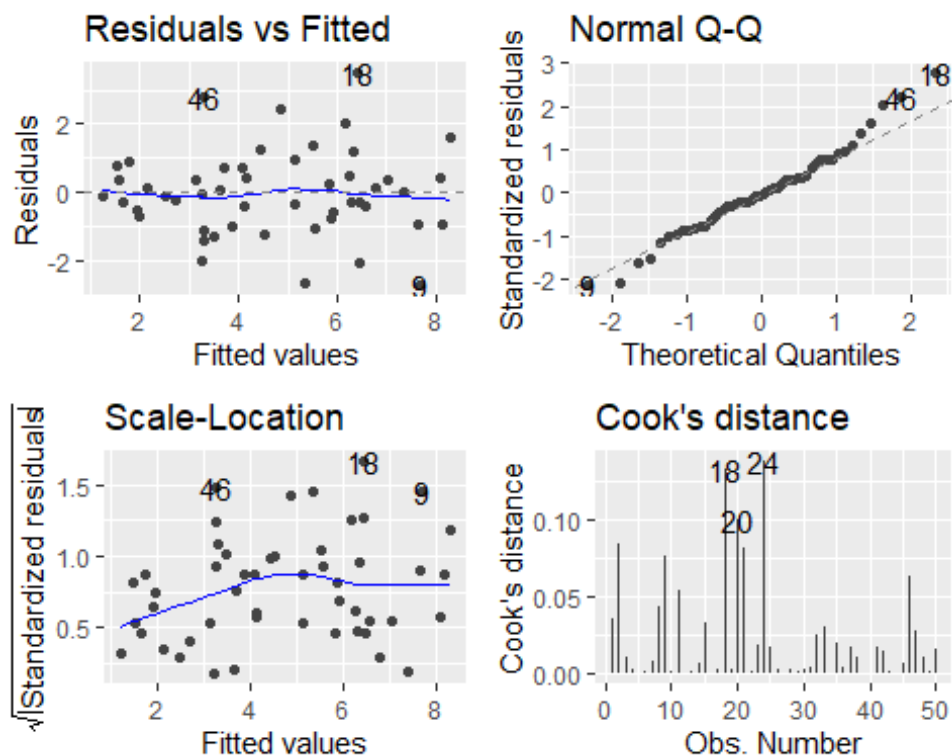
Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 1.324 on 44 degrees of freedom

Multiple R-squared: 0.719, Adjusted R-squared: 0.6871

F-statistic: 22.52 on 5 and 44 DF, p-value: 3.917e-11

```
autoplot(fit3, 1:4, nrow=2, ncol=2)
```



Drop motor_vehicle_theft - the adjusted R-squared value improved slightly to 68.7%.

Maybe try removing forcible rape since it had a large p-value of 0.51. Don't forget to check the diagnostic plots.

Model fit4

```
fit4 <- lm(murder ~ robbery + burglary + aggravated_assault + larceny_theft, data
= crime2)
summary(fit4)
```

Call:

```
lm(formula = murder ~ robbery + burglary + aggravated_assault +
    larceny_theft, data = crime2)
```

Residuals:

Min	1Q	Median	3Q	Max
-2.6290	-0.7670	-0.0601	0.4779	3.6348

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)	
(Intercept)	0.7555163	0.8946439	0.844	0.40286	
robbery	0.0201084	0.0041195	4.881	1.36e-05	***
burglary	0.0040134	0.0012950	3.099	0.00334	**
aggravated_assault	0.0039521	0.0020089	1.967	0.05533	.
larceny_theft	-0.0008325	0.0005146	-1.618	0.11268	

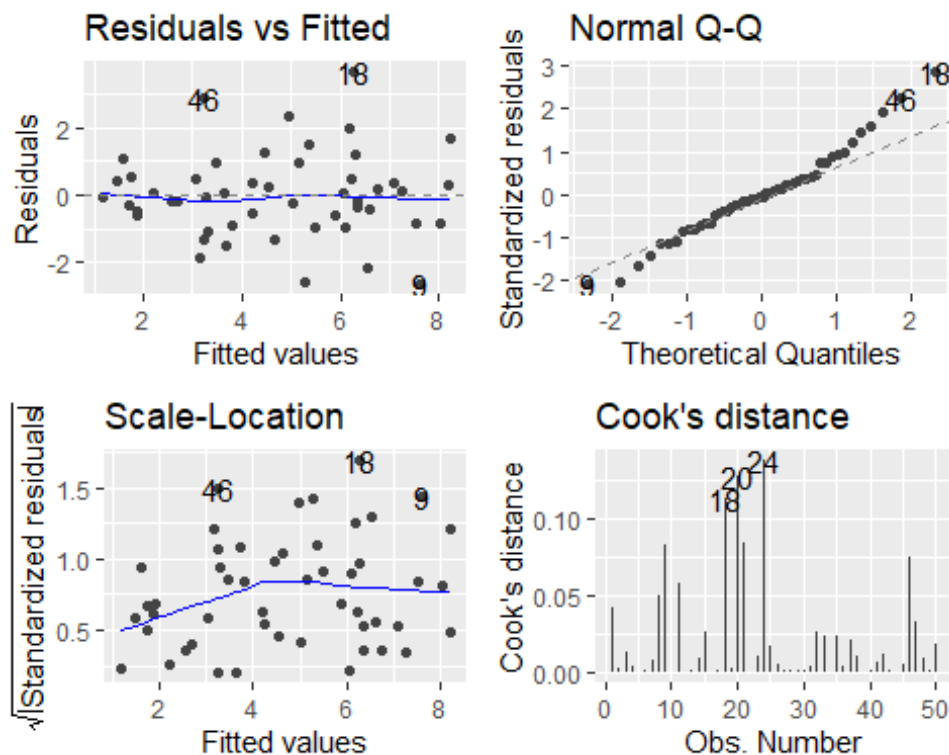
Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 1.315 on 45 degrees of freedom

Multiple R-squared: 0.7163, Adjusted R-squared: 0.691

F-statistic: 28.4 on 4 and 45 DF, p-value: 8.396e-12

```
autoplot(fit4, 1:4, nrow=2, ncol=2)
```



Interesting!!

The adjusted R-squared went up to 69.1%. The residuals plot looks about the same.

One final model - the simplest (parsimonious) by removing larceny_theft.

Model fit5

```
fit5 <- lm(murder ~ robbery + burglary + aggravated_assault, data = crime2)
summary(fit5)
```

Call:

```
lm(formula = murder ~ robbery + burglary + aggravated_assault,
    data = crime2)
```

Residuals:

Min	1Q	Median	3Q	Max
-2.6434	-0.7535	-0.0107	0.7229	3.7420

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	-0.330470	0.601764	-0.549	0.5855
robbery	0.021669	0.004075	5.318	3e-06 ***
burglary	0.002732	0.001042	2.621	0.0118 *
aggravated_assault	0.003570	0.002030	1.759	0.0853 .

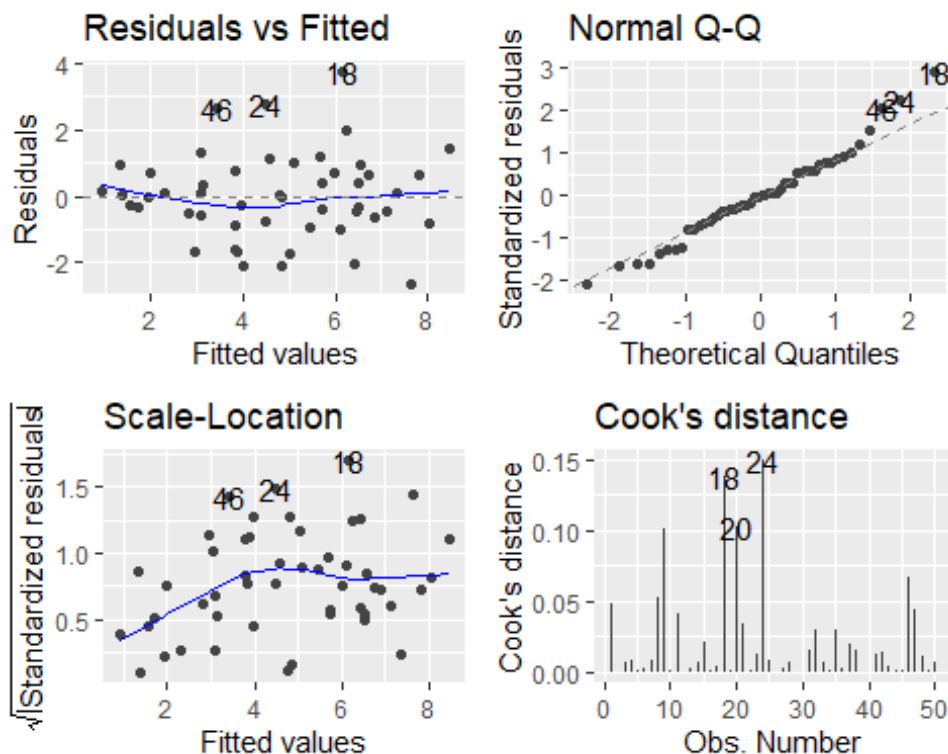
Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 1.338 on 46 degrees of freedom

Multiple R-squared: 0.6998, Adjusted R-squared: 0.6802

F-statistic: 35.74 on 3 and 46 DF, p-value: 4.451e-12

```
autoplot(fit5, 1:4, nrow=2, ncol=2)
```



The residuals plot shows observations 24 and 18 have an effect on the residuals plot as well having high scale-location values.

Mississippi is 24 Louisiana is 18

Try the last model (fit6), but remove those 2 observations:

```
crime3 <- crime2[-c(18,24),]  
fit6 <- lm(murder ~ robbery + burglary + aggravated_assault, data = crime3)  
summary(fit6)
```

Call:


```
lm(formula = murder ~ robbery + burglary + aggravated_assault,
   data = crime3)
```

Residuals:

	Min	1Q	Median	3Q	Max
	-2.43348	-0.56500	-0.01404	0.88995	2.61186

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	-0.1464862	0.5216248	-0.281	0.7802
robbery	0.0227896	0.0035316	6.453	7.27e-08 ***
burglary	0.0020752	0.0009364	2.216	0.0319 *
aggravated_assault	0.0036188	0.0018328	1.974	0.0546 .

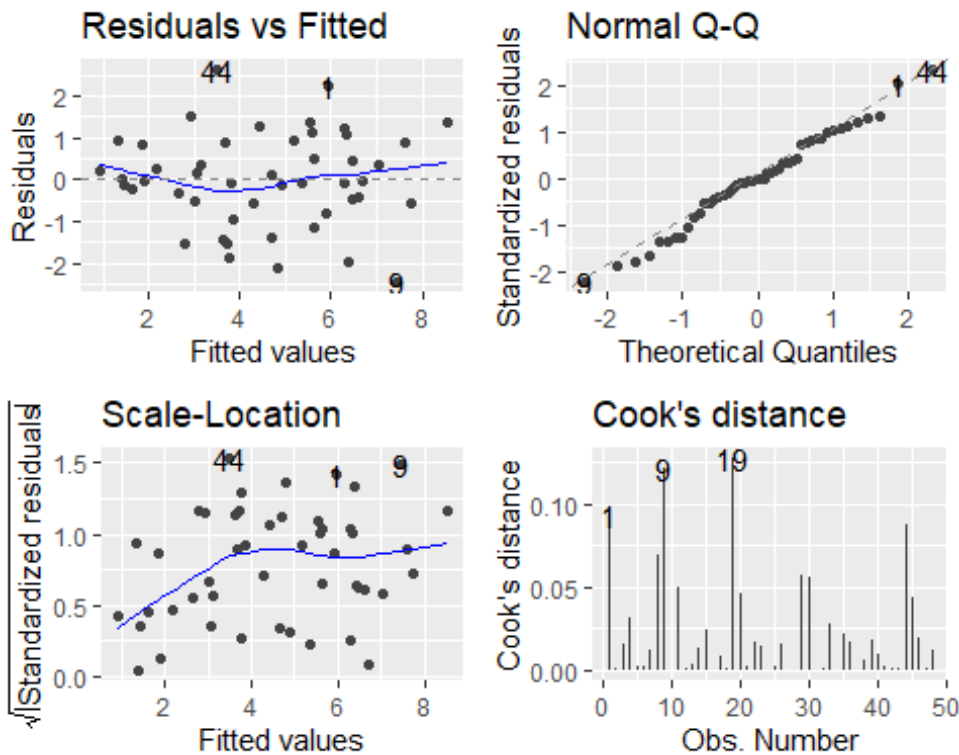
Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 1.156 on 44 degrees of freedom

Multiple R-squared: 0.7547, Adjusted R-squared: 0.738

F-statistic: 45.13 on 3 and 44 DF, p-value: 1.762e-13

```
autoplot(fit6, 1:4, nrow=2, ncol=2)
```



The adjusted R^2 went up to about 73%, which is an improvement. The residuals plot does not seem to have improved.

One last attempt - we can compare the last models to see if removing `larceny_theft` is an improvement on the model using ANOVA

ANOVA (analysis of variance) compares 2 models, one simpler than the other. If the result is a small p-value, then the larger model is better than the smaller model

ANOVA

```
anova(fit5, fit4)
```

Analysis of Variance Table

Model 1: murder ~ robbery + burglary + aggravated_assault

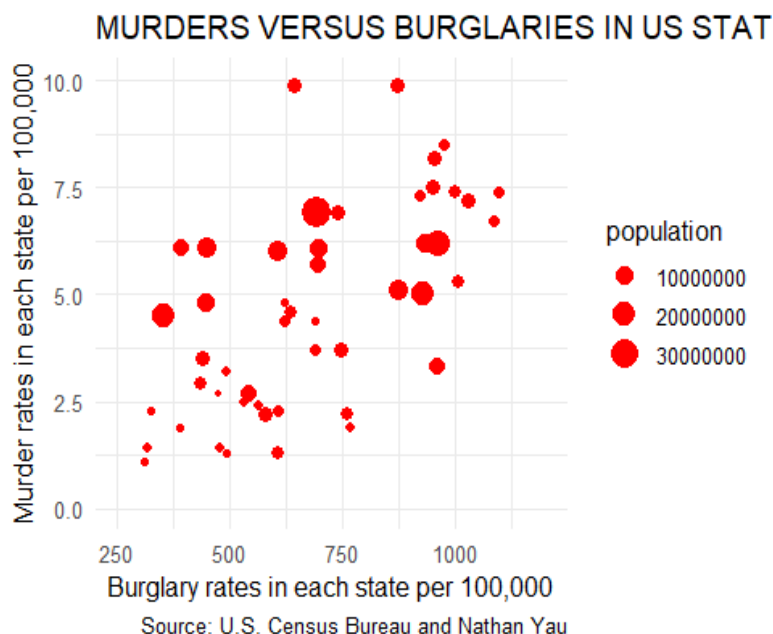
Model 2: murder ~ robbery + burglary + aggravated_assault + larceny_theft

	Res.Df	RSS	Df	Sum of Sq	F	Pr(>F)
1	46	82.381				
2	45	77.852	1	4.5284	2.6175	0.1127

We can see that the p-value is large, so there is no compelling evidence that larceny_theft contributes significantly to the model.

Back to simply murders and burglaries - bring in the state's population as a size of the circle

```
options(scipen = 999)
p2 +
  geom_point(aes(size = population), color = "red") + xlim(250,1200) + ylim(0,10)
+
  labs(title = "MURDERS VERSUS BURGLARIES IN US STATES PER 100,000",
caption = "Source: U.S. Census Bureau and Nathan Yau") +
  xlab("Burglary rates in each state per 100,000") +
  ylab ("Murder rates in each state per 100,000") +
  theme_minimal(base_size = 12)
```



Finally, add some interactivity to the plot with plotly

Warning: the ggplotly function does not render in a quarto document. It only renders in a markdown document.

```
p <- ggplot(crime2, aes(x = burglary, y = murder, size = population, text = paste("state:", state))) +  
  geom_point(alpha = 0.5, color = "red") + xlim(250,1200) + ylim(0,10) +  
  ggtitle("BURGLARIES VERSUS MURDERS IN THE U.S.", subtitle = "Sizes of circles are proportional to state populations") +  
  xlab("Burglary rates in each state per 100,000") +  
  ylab("Murder rates in each state per 100,000") +  
  theme_minimal(base_size = 12)  
p <- ggplotly(p)  
p
```

Make a series of charts from food stamps data

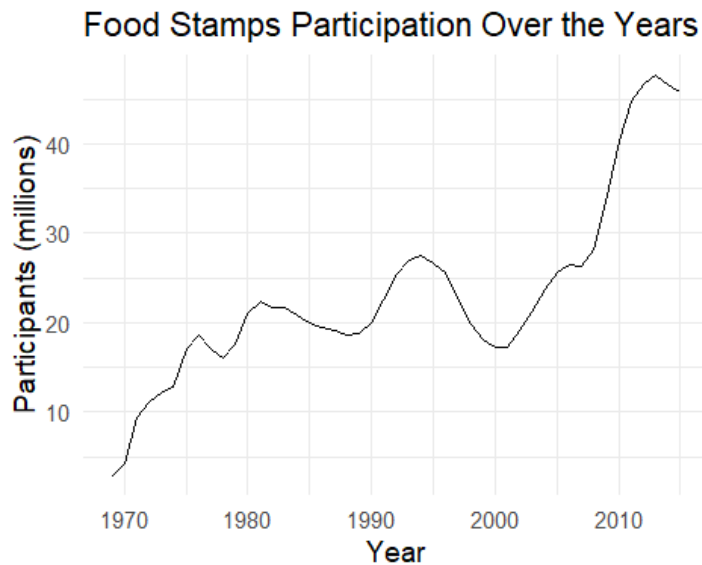
Now we will explore a series of other geom functions using the food stamps data.

Load the foodstamps data, map variables onto the X and Y axes, and save chart template

```
# load data  
food_stamps <- read_csv("food_stamps.csv")
```

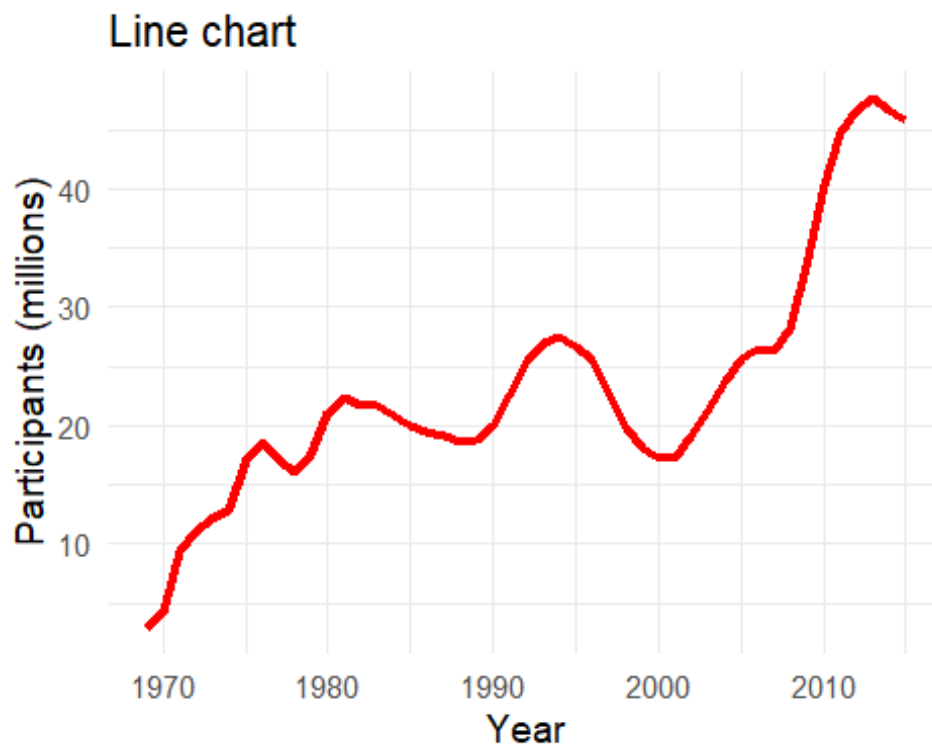
Use foodstamps dataset to make a line chart

```
food_stamps_chart <- ggplot(food_stamps, aes(x = year, y = participants)) +  
  labs(title = "Food Stamps Participation Over the Years") +  
  geom_line() +  
  xlab("Year") +  
  ylab("Participants (millions)") +  
  theme_minimal(base_size = 14)  
food_stamps_chart
```



Customize the line, and add a title

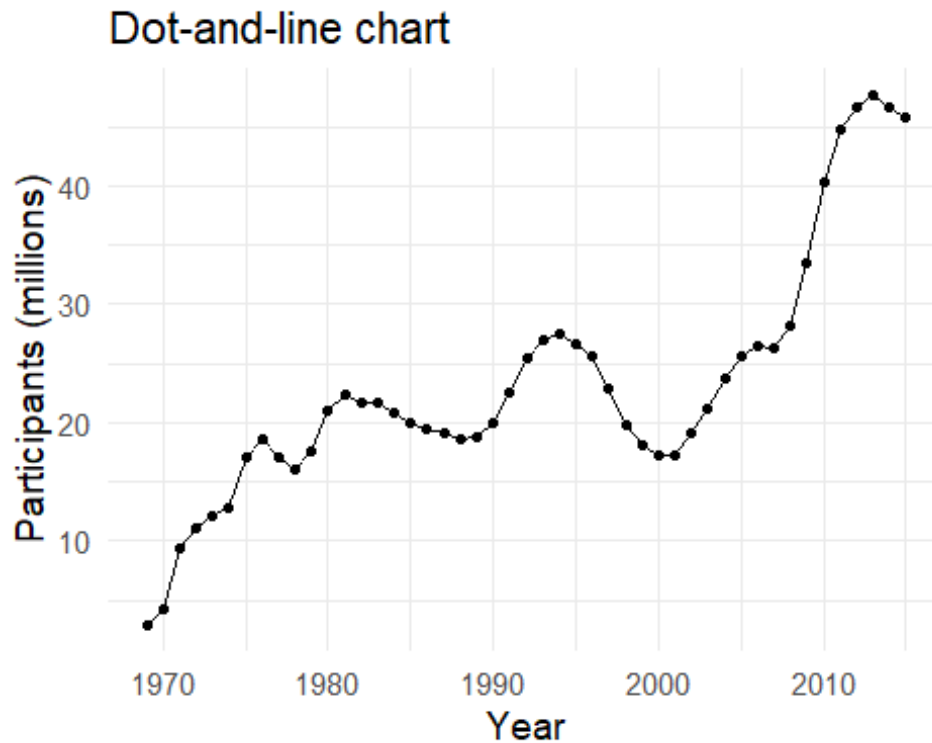
```
food_stamps_chart +  
  geom_line(size = 1.5, color = "red") +  
  ggtitle("Line chart")
```



Add a second layer to make a dot-and-line chart

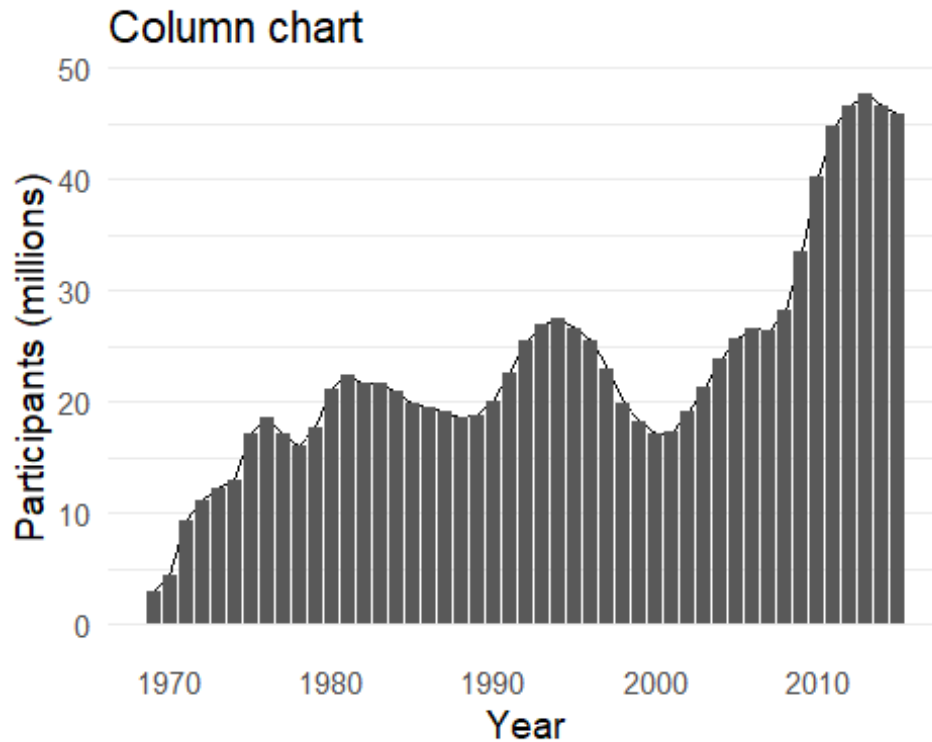
```
food_stamps_chart +  
  geom_line() +
```

```
geom_point() +  
ggtitle("Dot-and-line chart")
```



Make a column chart, then flip its coordinates to make a bar chart

```
# Make a column chart  
food_stamps_chart +  
  geom_bar(stat = "identity") +  
  ggtitle("Column chart") +  
  theme(panel.grid.major.x = element_blank(),  
        panel.grid.minor.x = element_blank())
```

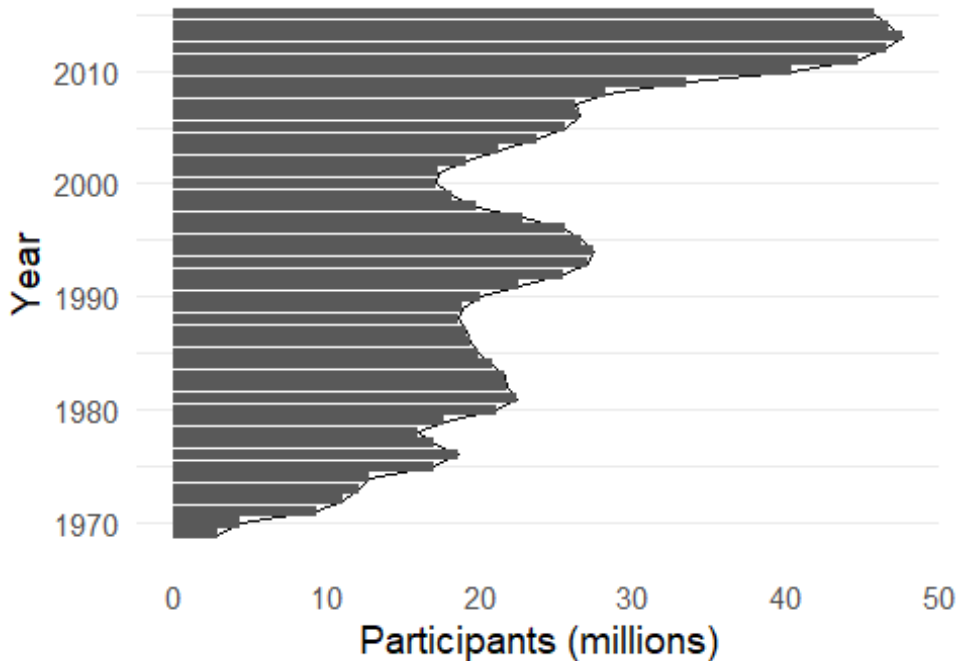


`geom_bar` works a little differently to the geoms we have considered previously. If you have not mapped data values to the Y axis with `aes`, its default behavior is to set the heights of the bars by counting the number of records for values along the X axis. If you have mapped a variable to the Y axis, and want the heights of the bars to represent values in the data, use you must use `stat="identity"`.

`coord_flip` switches the X and Y axes.

```
# Make a bar chart
food_stamps_chart +
  geom_bar(stat = "identity") +
  ggtitle("Bar chart") +
  theme(panel.grid.major.x = element_blank(),
        panel.grid.minor.x = element_blank()) +
  coord_flip()
```

Bar chart



The difference between color and fill

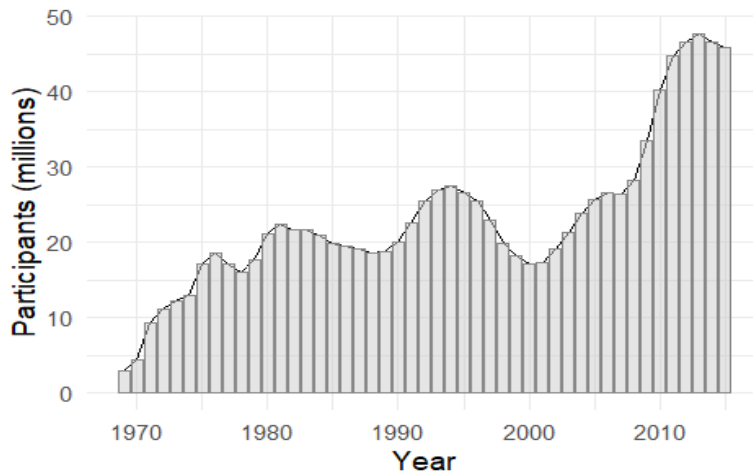
For some geoms, notably `geom_bar`, you can set color for their outline as well as the interior of the shape.

When setting colors, `color` refers to the outline, `fill` to the interior of the shape.

use both color and fill

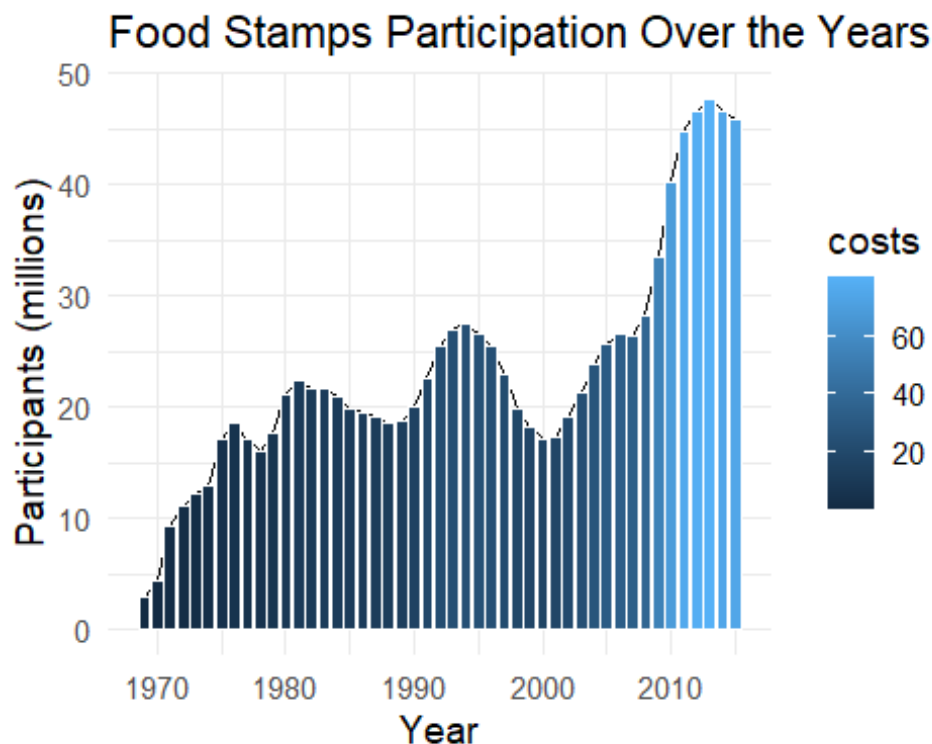
```
# set color and fill
food_stamps_chart +
  geom_bar(stat = "identity", color = "#888888", fill = "#CCCCCC", alpha = 0.5) +
  ggtitle("Column chart")
```

Column chart



Map color to the values of a continuous variable

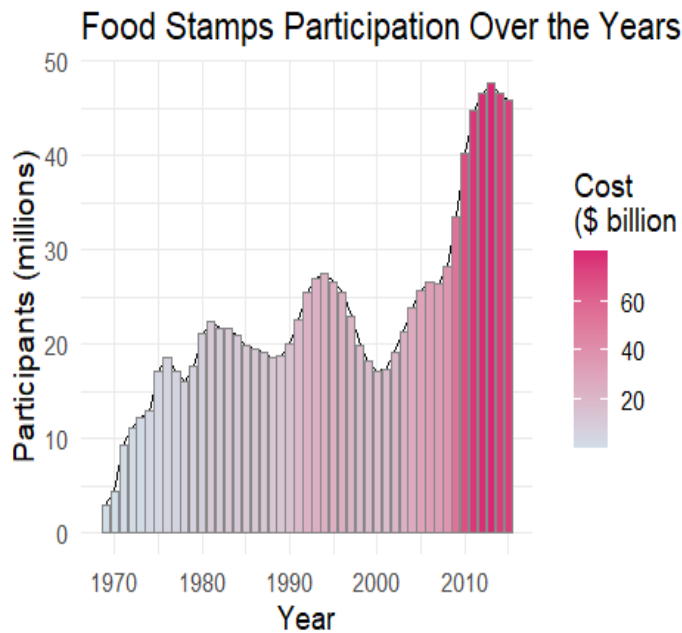
```
# fill the bars according to values for the cost of the program
food_stamps_chart +
  geom_bar(stat = "identity", color = "white", aes(fill = costs))
```



This code uses an aes mapping to color the bars according values for the costs of the program, in billions of dollars. ggplot2 recognizes that costs is a continuous variable, but its default sequential scheme applies more intense blues to lower values, which is counterintuitive.

Use a ColorBrewer sequential color palette

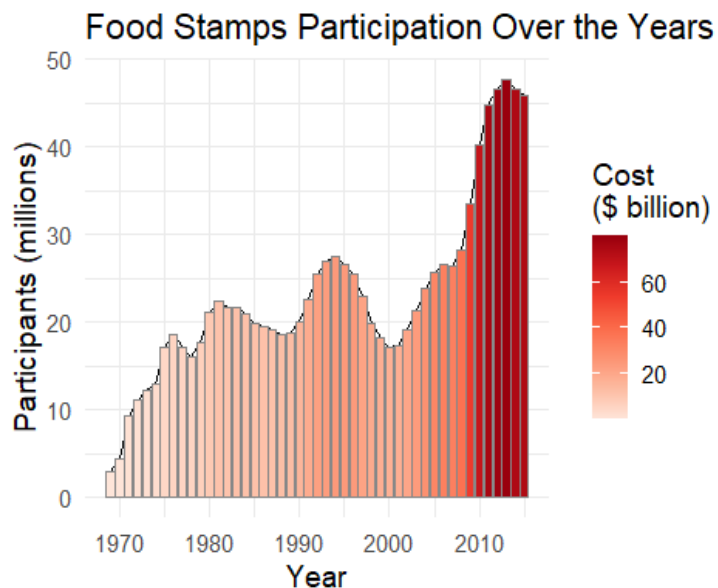
```
# use a colorbrewer gradient levels for intensity
food_stamps_chart +
  geom_bar(stat = "identity", color = "#888888", aes(fill = costs)) +
  scale_fill_gradient(name = "Cost\n($ billion", low = "#d1dee8", high = "#d92774")
  )
```

`scale_fill_distiller` (and `scale_color_distiller`) work like `scale_color_brewer`, but set color gradients for ColorBrewer's sequential and diverging color palettes; `direction = 1` ensures that larger numbers are mapped to more intense colors (`direction = -1` reverses the color mapping). Try changing the code I have: `scale_fill_gradient()` to `scale_fill_distiller` with different directions (1 or -1).

Scale_fill_distiller

```
food_stamps_chart +
  geom_bar(stat = "identity", color = "#888888", aes(fill = costs)) +
  scale_fill_distiller(name = "Cost\n($ billion)", palette = "Reds", direction =
1)
```



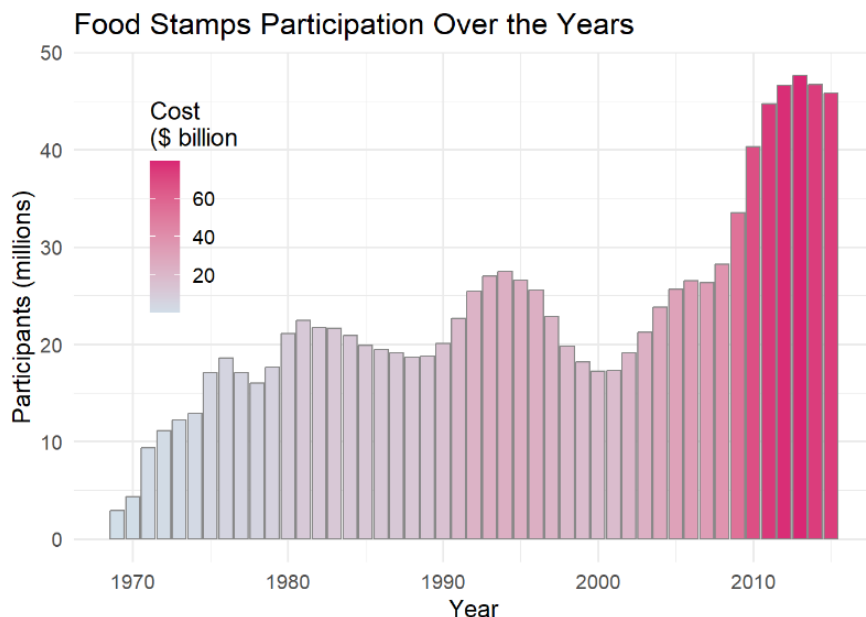
Notice also the in the title for the legend. This introduces a new line.

Control the position of the legend

This code uses the theme function to moves the legend from its default position to the right of the chart to use some empty space on the chart itself.

Move the legend

```
food_stamps_chart +  
  geom_bar(stat="identity", color = "#888888", aes(fill=costs)) +  
  scale_fill_gradient(name = "Cost\n($ billion", low = "#d1dee8", high = "#d92774") +  
  theme(legend.position=c(0.15,0.7))
```

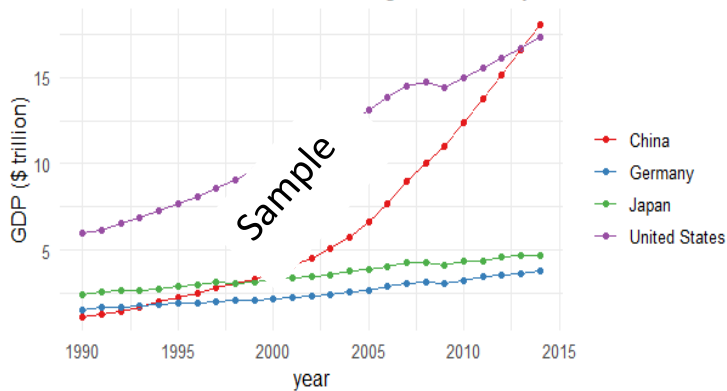


The coordinates for the legend are given as a list: The first number sets the horizontal position, from left to right, on a scale from 0 to 1; the second number sets the vertical position, from bottom to top, again on a scale from 0 to 1.

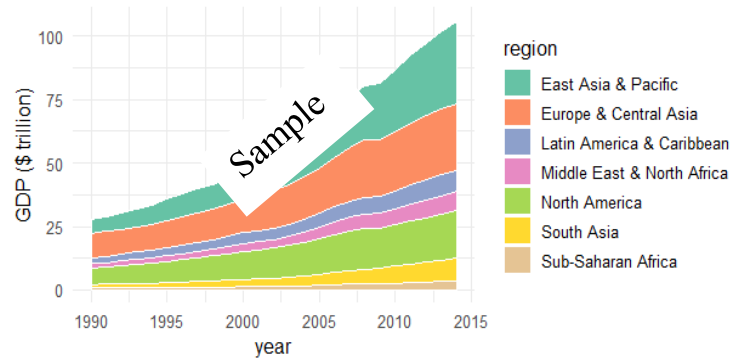
Week 6 Homework Assignment

1. **(Ungraded)** Complete copying notes on scatterplotting, correlation, and regression analysis.
2. **(Worth up to 10 points for each chart)** Use **dplyr** and **ggplot2** to process data and draw these two charts (shown below) from the Nations dataset. You do NOT need to incorporate interactivity, but you can, if you want to challenge yourself.

China's Rise to Become the Largest Economy



GDP by World Bank Region



Details for Nations Dataset Charts Assignment

- For both charts, you will first need to create a new variable in the data, using `mutate` from **dplyr**, giving the GDP of each country in trillions of dollars, by multiplying `gdp_percap` by `population` and dividing by a trillion.
- Draw both charts with **ggplot2**.
- For the first chart, you will need to `filter` the data with **dplyr** for the four desired countries. When making the chart with **ggplot2** you will need to add both `geom_point` and `geom_line` layers, and use the `Set1` `ColorBrewer` palette using: `scale_color_brewer(palette = "Set1")`.
- For the second chart, using **dplyr** you will need to `group_by` `region` and `year`, and then `summarize` on your mutated value for `gdp` using `summarise(GDP = sum(gdp, na.rm = TRUE))`. (There will be null values, or NAs, in this data, so you will need to use `na.rm = TRUE`).
- Each region's area will be generated by the command `geom_area()`
- When drawing the chart with **ggplot2**, you will need to use the `Set2` `ColorBrewer` palette using `scale_fill_brewer(palette = "Set2")`
- Think about the difference between `fill` and `color` when making the chart, and where the above `fill` command needs to go in order for the regions to fill with the different colors when making the chart, and put a very thin white line around each area.

Knit your code for each chart and save your work in `rpubs`. Submit the link on the assignment dropbox by **11:59 pm on Tuesday, ____**.