# Measurement and visualization, 2

Frank Edwards

10/19/21

- Survey methods with randomization
- Administrative data and agency surveys
- Unit, item non-response
- Desirability bias
- Latent variables, latent groups
- More visualization

- A census records information about a population, with measurement for each individual or unit in the population
- A survey samples from a population to make an inference about population characteristics

```
n<-1e6
marbles<-data.frame(
  color = c(rep("blue", n),
            rep("green", n),
            rep("red", n)))
table(marbles$color)


##
##    blue   green     red
## 1000000 1000000 1000000
```
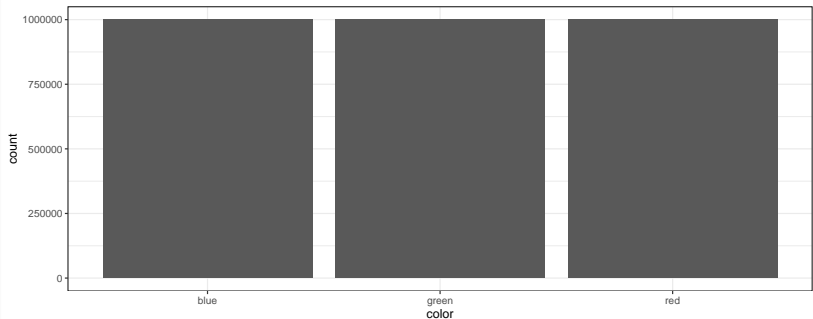
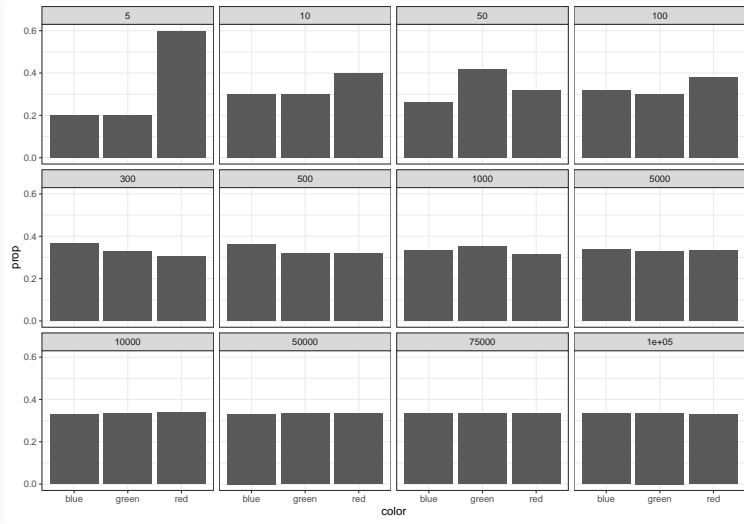## How could we know how many of each color are in the enormous bag of marbles?

- Count them all (tedious!)
- Sample

```
ggplot(marbles, aes(x = color)) + geom_bar()
```

# How many random draws is enough to accurately measure the characteristics of 3 million marbles?

- With a sufficiently large sample and equal probability of sampling for all units in the population, a simple random sample allows for unbiased measurement of population characteristics.

- With a sufficiently large sample and equal probability of sampling for all units in the population, a simple random sample allows for unbiased measurement of population characteristics.
- Identical motivation for randomization in experiments

- With a sufficiently large sample and equal probability of sampling for all units in the population, a simple random sample allows for unbiased measurement of population characteristics.
- Identical motivation for randomization in experiments
- Such a sample is representative of the population across both measured and unmeasured characteristics

## Stratified random sampling

If we wish to learn about particular sub-populations (i.e. geographies), we can use multi-stage or stratified sampling

## Stratified random sampling

If we wish to learn about particular sub-populations (i.e. geographies), we can use multi-stage or stratified sampling

1. Randomly sample larger units (geographic) or select larger units of interest purposively

## Stratified random sampling

If we wish to learn about particular sub-populations (i.e. geographies), we can use multi-stage or stratified sampling

1. Randomly sample larger units (geographic) or select larger units of interest purposively
2. Randomly sample individuals within these larger units

If we wish to learn about particular sub-populations (i.e. geographies), we can use multi-stage or stratified sampling

1. Randomly sample larger units (geographic) or select larger units of interest purposively
2. Randomly sample individuals within these larger units

EXAMPLE: The American Community Survey (simplified)

If we wish to learn about particular sub-populations (i.e. geographies), we can use multi-stage or stratified sampling

1. Randomly sample larger units (geographic) or select larger units of interest purposively
2. Randomly sample individuals within these larger units

EXAMPLE: The American Community Survey (simplified)

1. Take a list of all US Census tracts

If we wish to learn about particular sub-populations (i.e. geographies), we can use multi-stage or stratified sampling

1. Randomly sample larger units (geographic) or select larger units of interest purposively
2. Randomly sample individuals within these larger units

EXAMPLE: The American Community Survey (simplified)

1. Take a list of all US Census tracts
2. Randomly sample households within tract based on complete list of addresses (sampling frame)

## Stratified random sampling

If we wish to learn about particular sub-populations (i.e. geographies), we can use multi-stage or stratified sampling

1. Randomly sample larger units (geographic) or select larger units of interest purposively
2. Randomly sample individuals within these larger units

EXAMPLE: The American Community Survey (simplified)

1. Take a list of all US Census tracts
2. Randomly sample households within tract based on complete list of addresses (sampling frame)
3. Randomly sample adults within household, conduct survey

1. Unit non-response

1.  Unit non-response
2.  Item non-response

1. Unit non-response
2. Item non-response
3. Lying (of various sorts)

Individual (or organization) doesn't respond to the survey

Individual (or organization) doesn't respond to the survey

- How are surveys actually administered?
- Response rates are generally low (and decreasing!)

Completely random non-response does not bias inference

Individual (or organization) doesn't respond to the survey

- How are surveys actually administered?
- Response rates are generally low (and decreasing!)

Completely random non-response does not bias inference

When would non-response be an issue?

Individual takes the survey, but refuses to answer (skips) a particular question

- Why might this occur?

Individual takes the survey, but refuses to answer (skips) a particular question

- Why might this occur?
- When would this be a problem?

- Social desirability bias
    - Did you vote? Remember HW 1?
    - Are you a racist?
    - What kinds of crimes do you like to do?

Examining non-response in a survey of exposure to violence in Afghanistan

```
### survey of
afghan<-read_csv("https://raw.githubusercontent.com/kosukeimai/qss/master/MEASUREMENT/afghan.csv")
afghan.village<-read_csv("https://raw.githubusercontent.com/kosukeimai/qss/master/MEASUREMENT/afghan-village.cs
```

# Explore the variables in afghan

```
table(afghan$province)
```

```
##
## Helmand   Khost   Kunar   Logar Uruzgan
##     855     630     396     486     387
```

```
table(afghan$district)
```

```
##
##        Asadabad           Bak   Baraki Barak    Chapa Dara        Dangam
##              54            54           180           108            63
##        Dihrawud       Garmser     Ghaziabad  Khas Uruzgan        Khoshi
##             117           225            63           117            54
##           Khost   Lashkar Gah     Musa Qala      Naw Zad     Puli Alam
##             243           108           225           216           252
##        Qalandar Shahidi Hassas         Spira          Tani       Washer
##              63           153            99           171            81
##        Wata Pur
##             108
```

# Explore the variables in afghan

```
ggplot(afghan, aes(x=age)) + geom_histogram()
```

# Explore the variables in afghan

```
ggplot(afghan, aes(x = educ.years)) + geom_bar()
```

```
table(afghan$employed)
```

```
##
##    0    1
## 1149 1605
```

```
table(afghan$income)
```
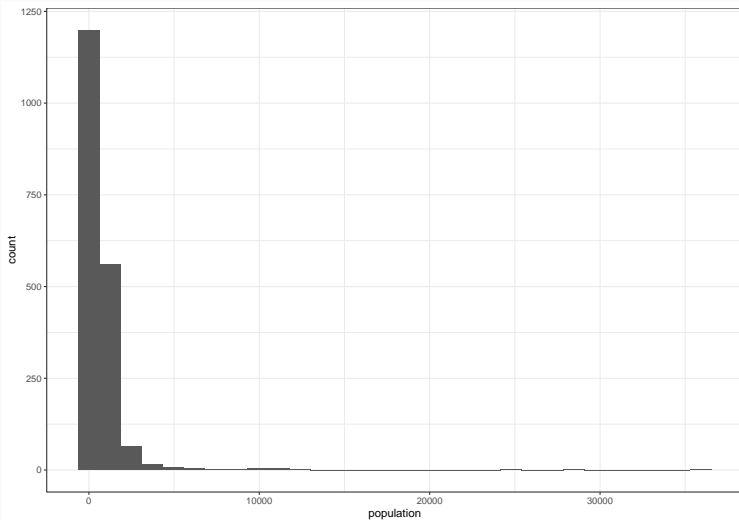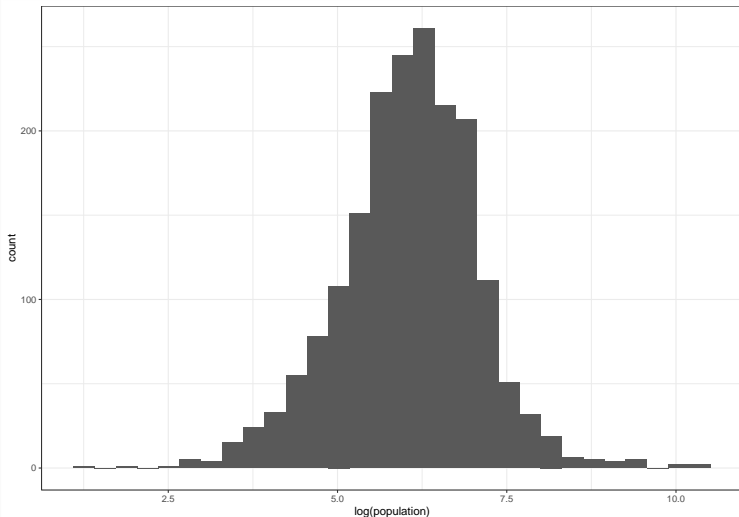
```
##
##   10,001-20,000    2,001-10,000   20,001-30,000 less than 2,000     over 30,000
##             616            1420              93             457              14
```

```
## for ordered categorical
afghan<-afghan %>%
  mutate(income =
           factor(income,
                  levels = c(
                    "less than 2,000", "2,001-10,000",
                    "10,001-20,000", "20,001-30,000",
                    "over 30,000")))
```

```
table(afghan$income)
```

```
##
## less than 2,000    2,001-10,000   10,001-20,000   20,001-30,000     over 30,000
##             457            1420             616              93              14
```

# Explore the variables in afghan

```
afghan %>%
  select(employed, violent.exp.ISAF, violent.exp.taliban) %>%
  summary()
```

```
##     employed        violent.exp.ISAF violent.exp.taliban
##  Min.   :0.0000   Min.   :0.0000   Min.   :0.0000
##  1st Qu.:0.0000   1st Qu.:0.0000   1st Qu.:0.0000
##  Median :1.0000   Median :0.0000   Median :0.0000
##  Mean   :0.5828   Mean   :0.3749   Mean   :0.3289
##  3rd Qu.:1.0000   3rd Qu.:1.0000   3rd Qu.:1.0000
##  Max.   :1.0000   Max.   :1.0000   Max.   :1.0000
##                   NA's   :25       NA's   :54
```

## Explore the variables in afghan.village

```
head(afghan.village)
```

```
## # A tibble: 6 x 3
##    altitude population village.surveyed
##       <dbl>      <dbl>            <dbl>
## 1    1959.        197                1
## 2    2426.        744                0
## 3    2237.        179                1
## 4    1692.        225                0
## 5    1928.        379                0
## 6    1195.        617                0
```

# Explore the variables in afghan.village

```
ggplot(afghan.village, aes(x=altitude)) +
  geom_histogram()
```

# Explore the variables in afghan.village

```
ggplot(afghan.village, aes(x=population)) +
  geom_histogram()
```

# Explore the variables in afghan.village: logs help!

```
ggplot(afghan.village, aes(x=log(population))) +
  geom_histogram()
```

```
mean(afghan.village$village.surveyed)
```

```
## [1] 0.1094421
```

# Is the sampling representative of villages?

```
ggplot(afghan.village,
        aes(x=village.surveyed==1,
            y = altitude)) +
  geom_boxplot()
```
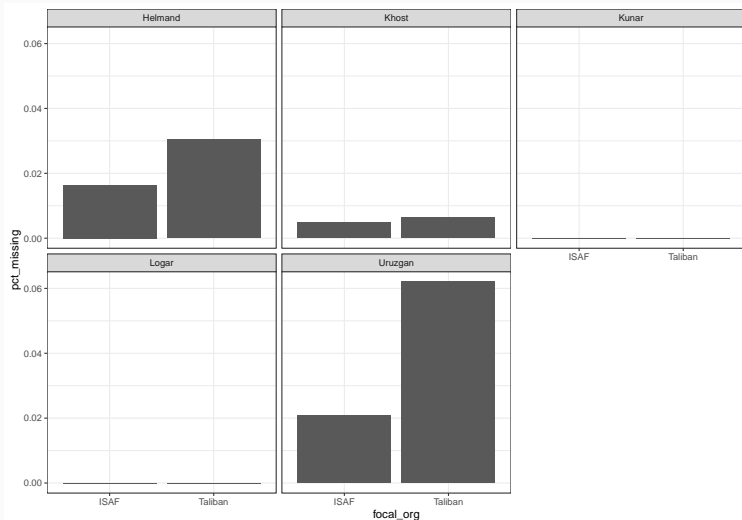
## Is the sampling representative of villages?

```
ggplot(afghan.village,
       aes(x=village.surveyed==1,
           y = log(population))) +
  geom_boxplot()
```

# Is the sampling representative of villages? Alternative plot

```
ggplot(afghan.village,
       aes(x=log(population))) +
  geom_histogram() +
  facet_wrap(~ village.surveyed, scales = "free")
```

# Does item non-response bias estimates of violence by region?

- Unit non-responses can bias survey estimates

- Unit non-responses can bias survey estimates
- Item non-response can bias survey estimates

- Unit non-responses can bias survey estimates
- Item non-response can bias survey estimates
- Social desirability can bias survey estimates

- Unit non-responses can bias survey estimates
- Item non-response can bias survey estimates
- Social desirability can bias survey estimates
- Errors induced by these biases can lead to incorrect conclusions (see polling consensus on 2016 election)

Returning to the IPV example

```
ipv<-read_csv("./data/dhs_ipv.csv")
## on your machine, path is /slides/data/
head(ipv)
```

```
## # A tibble: 6 x 7
##   beat_burnfood beat_goesout sec_school no_media country     year region
##           <dbl>        <dbl>      <dbl>    <dbl> <chr>      <dbl> <chr>
## 1           4.4         18.6       25.2      1.5 Albania     2008 Middle East a~
## 2           4.9         19.9       67.7      8.7 Armenia     2000 Middle East a~
## 3           2.1         10.3       67.6      2.2 Armenia     2005 Middle East a~
## 4           0.3          3.1       46        6.4 Armenia     2010 Middle East a~
## 5          12.1         42.5       74.6      7.4 Azerbaijan  2006 Middle East a~
## 6          NA           NA         24       41.9 Bangladesh  2004 Asia
```

# Look at bivariate relationships

```
ggplot(ipv, aes(x = sec_school, y =no_media, color = region)) + geom_point()
```
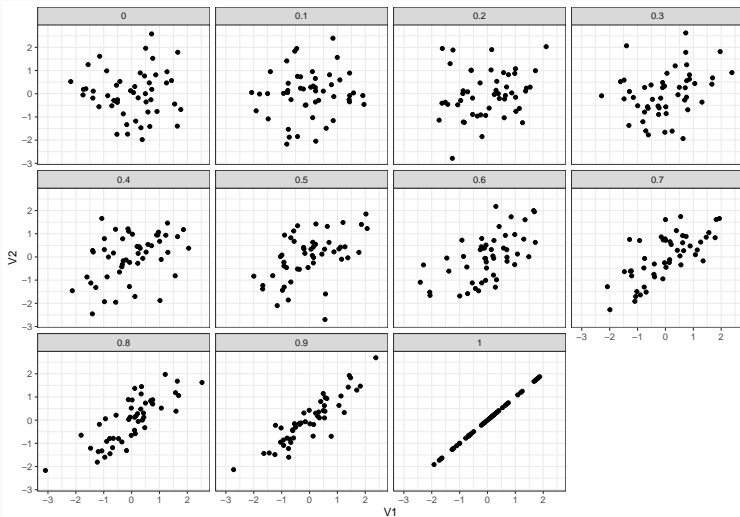
# Is there a change in sec_school by region over time across this sample? Does time matter here?
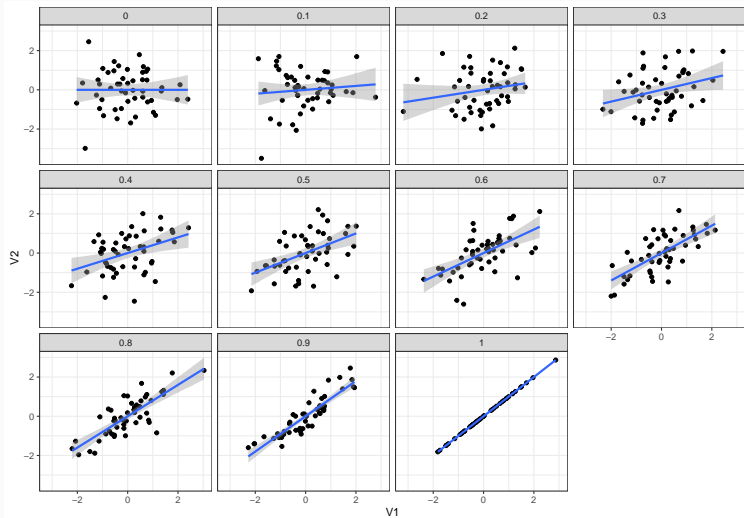
```r
ipv_ts<-ipv %>%
  group_by(region, year) %>%
  summarise(sec_school=mean(sec_school))

ggplot(ipv_ts, aes(x=year, y = sec_school, color = region)) +
  geom_line()
```

## Correlation (math time): Z-scores

First, we need the variables to be comparable, so we transform them to be on a standard deviation scale.

A *z*-score scales a variable measures the number of standard deviations an observation is away from it's mean.

$$z \text{ score of } x_i = \frac{x_i - \bar{x}}{S_x}$$

Where $\bar{x}$ is the mean, and $S_x$ is the standard deviation of variable *x*. Z scores have a mean zero, and a range defined by the range of the data on a standard deviation scale.

For a normally (Gaussian) distributed variable, this will typically range between $[-3, 3]$

In R, we can transform a numeric into a *z*-score using `scale()`

# Z-scores in R

```
ipv_scale<-ipv %>%
  mutate(sec_scale = scale(sec_school)) %>%
  select(sec_school, sec_scale)
summary(ipv_scale)
```

```
##    sec_school      sec_scale.V1
## Min.   : 3.10   Min.   :-1.345006
## 1st Qu.:10.18   1st Qu.:-0.898292
## Median :22.40   Median :-0.126408
## Mean   :24.40   Mean   : 0.000000
## 3rd Qu.:34.90   3rd Qu.: 0.662840
## Max.   :74.60   Max.   : 3.169492
## NA's   :3       NA's   :3
```

# Z-scores in R

```
ggplot(ipv_scale, aes(x=sec_school, y=sec_scale)) + geom_point()
```

Correlation measures the degree to which two variables are associated
with each other. We often use the letter *r* to denote a correlation.

$$r(x, y) = \frac{1}{n} \sum_{i=1}^{n} \frac{x_i - \bar{x}}{S_x} \times \frac{y_i - \bar{y}}{S_y}$$

Note that this is equal to the average of the product of the *z − scores* of *x*
and *y*

In R, you can use `cor()`

# Returning to our example: Are sec_school and no_media correlated?

```r
cor(ipv$sec_school, ipv$no_media, use="complete")
```

```
## [1] -0.6077951
```

```r
## z score method
mean(scale(ipv$sec_school) * scale(ipv$no_media), na.rm=TRUE)
```

```
## [1] -0.6084724
```

# Clustering

Data often *cluster* based on unobserved or unobservable characteristics. We can use *classification methods* to try to uncover these latent structures in data.

*k*-means is a straightforward method we can use to identify *k* latent groupings in our data, based on proximity of observations for specified variables.

An algorithm is a sequential set of steps used to solve a problem.

A *centroid* is the mean value of a cluster within a group.

1. Choose the initial centroids for each of the *k* clusters
2. Assign each observation to the cluster with the nearest centroid
3. Assign a new centroid based on the within-cluster mean for assigned observations
4. Repeat steps 2 and 3 until the cluster assignments no longer change

We arbitrarily choose the number of clusters *k*, and R randomly selects starting centroid values for step 1.

```
ipv_scale<-ipv %>%
  select(sec_school, no_media) %>%
  mutate(sec_school = scale(sec_school),
         no_media = scale(no_media)) %>%
  filter(!(is.na(sec_school)), !(is.na(no_media)))

ipv_kmeans<-kmeans(ipv_scale,
                   centers = 3,
                   nstart=10)
```

## Working with the k-means object

```
ipv_kmeans


## K-means clustering with 3 clusters of sizes 72, 17, 46
##
## Cluster means:
##   sec_school   no_media
## 1 -0.6135490  0.7910351
## 2  1.8803248 -1.1669709
## 3  0.2071354 -0.7678187
##
## Clustering vector:
##   [1] 3 2 2 2 2 1 1 1 1 1 3 3 1 1 1 1 3 1 1 1 1 1 3 3 1 1 1 1 3 3 3 3 3 3 3
##  [38] 2 1 1 1 1 2 1 3 3 1 1 2 2 1 3 3 3 3 1 3 3 3 2 2 3 3 2 1 3 1 1 1 1 1 1 1 3
##  [75] 1 1 1 2 3 1 1 3 3 3 1 1 3 1 1 1 1 1 3 3 3 2 3 3 3 3 3 3 1 1 1 3 3 3 3 3
## [112] 1 1 2 1 1 1 1 1 1 1 1 2 1 1 1 1 3 3 2 2 1 1 1 1 1 1
##
## Within cluster sum of squares by cluster:
## [1] 52.858371  8.657364 24.241269
##  (between_SS / total_SS =  68.3 %)
##
## Available components:
##
## [1] "cluster"      "centers"      "totss"         "withinss"      "tot.withinss"
## [6] "betweenss"    "size"         "iter"          "ifault"
```

## Pull out what we need from the list

```r
ipv_clusters<-ipv %>%
  filter(!(is.na(sec_school)), !(is.na(no_media))) %>%
  mutate(cluster = factor(ipv_kmeans$cluster))

centers<-data.frame(ipv_kmeans$centers)

centers

##   sec_school   no_media
## 1 -0.6135490  0.7910351
## 2  1.8803248 -1.1669709
## 3  0.2071354 -0.7678187
```
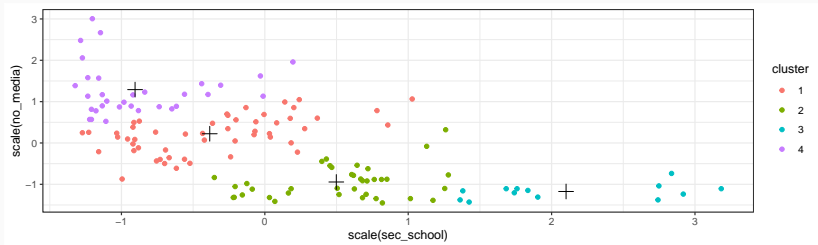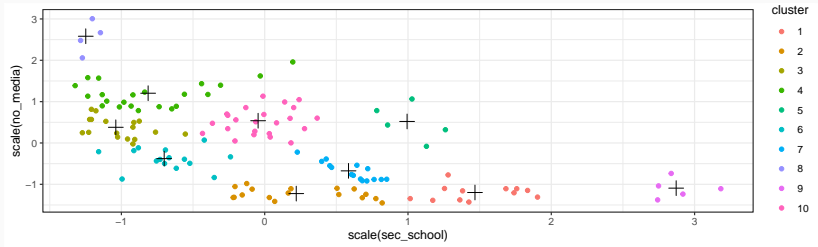
# Plot it!

```
ggplot(ipv_clusters, aes(x = scale(sec_school),
                         y =scale(no_media),
                         color = cluster)) +
  geom_point() +
  geom_point(data = centers,
             aes(x=sec_school,
                 y=no_media),
             color = "black", size = 4, shape = 3)
```

# What if we thought there were 4 clusters?

- Measurement and design matter!
- Always check your data, and think about how unit and item non-response may inform your conclusions
- Think about desirability and other forms of response bias as you interpret your results
- Design visuals and exploratory analyses to check hypotheses about what's going on in the data
- Think about the structure of your data, use descriptive statistics like correlations to describe relationships
- Think about latent structures in your data to capture clustering

Lab: more data visualization with ggplot

```
head(mpg)
```

```
## # A tibble: 6 x 11
##   manufacturer model displ  year   cyl trans      drv     cty   hwy fl    class
##   <chr>        <chr> <dbl> <int> <int> <chr>      <chr> <int> <int> <chr> <chr>
## 1 audi         a4      1.8  1999     4 auto(l5)   f        18    29 p     compa~
## 2 audi         a4      1.8  1999     4 manual(m5) f        21    29 p     compa~
## 3 audi         a4      2    2008     4 manual(m6) f        20    31 p     compa~
## 4 audi         a4      2    2008     4 auto(av)   f        21    30 p     compa~
## 5 audi         a4      2.8  1999     6 auto(l5)   f        16    26 p     compa~
## 6 audi         a4      2.8  1999     6 manual(m5) f        18    26 p     compa~
```
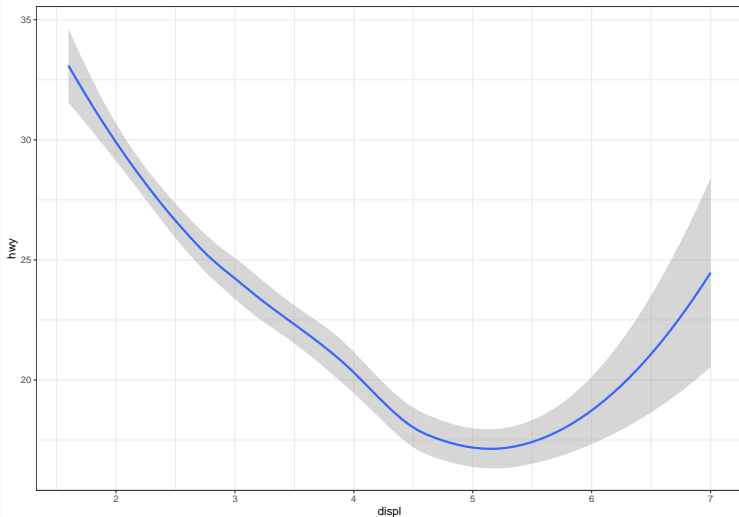
# How are these plots similar?

```
ggplot(mpg,
       aes(x = displ, y = hwy)) +
  geom_point()
```

# How are these plots similar?

```
ggplot(mpg,
       aes(x = displ, y = hwy)) +
  geom_smooth()
```
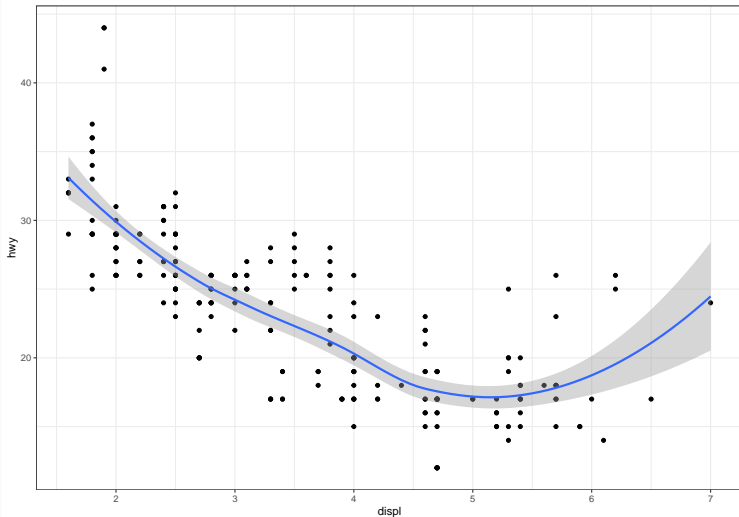
# Geometric objects (geoms) map data to visual objects

```
ggplot(mpg,
       aes(x = displ, y = hwy)) +
  geom_point()
```
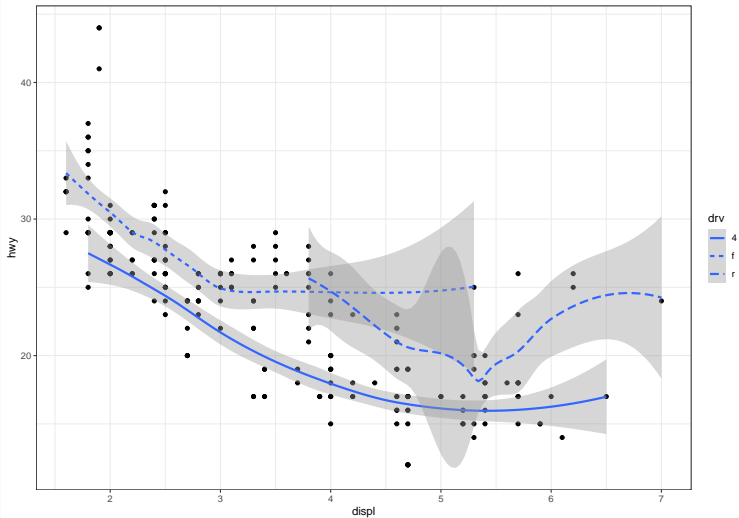
# We can layer geoms

```
ggplot(mpg,
       aes(x = displ, y = hwy)) +
  geom_point() +
  geom_smooth()
```
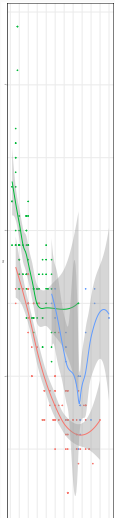
# Add aesthetics to map variables to visual objects

```
ggplot(mpg,
       aes(x = displ, y = hwy, lty = drv)) +
  geom_point() +
  geom_smooth()
```
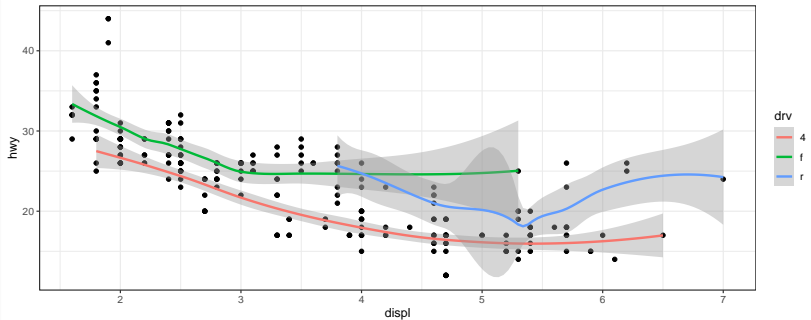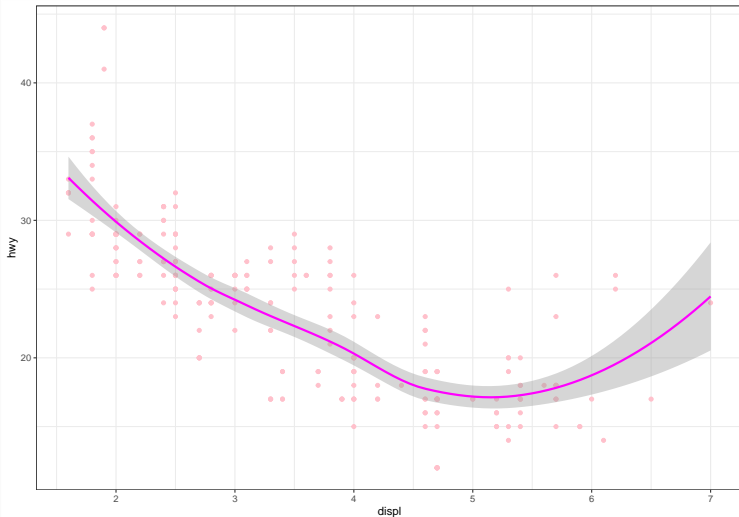
```
ggplot(mpg,
       aes(x = displ, y = hwy, color = drv)) +
  geom_point() +
  geom_smooth()
```

```
## What's the difference from the prior plot?
## could i make this more compact?
ggplot(mpg) +
  geom_point(aes(x = displ, y = hwy)) +
  geom_smooth(aes(x = displ, y = hwy, color = drv))
```
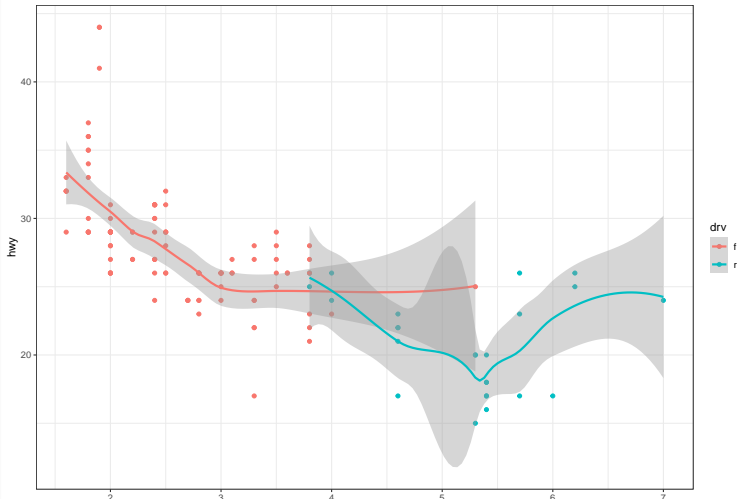
# Modifying visual objects without variable mapping

```
ggplot(mpg,
       aes(x = displ, y = hwy)) +
  geom_point(color = "pink") +
  geom_smooth(color = "magenta")
```

# Modifying data prior to plotting

```
ggplot(mpg %>%
       filter(drv !="4"),
       aes(x = displ, y = hwy, color = drv)) +
  geom_point() +
  geom_smooth()
```
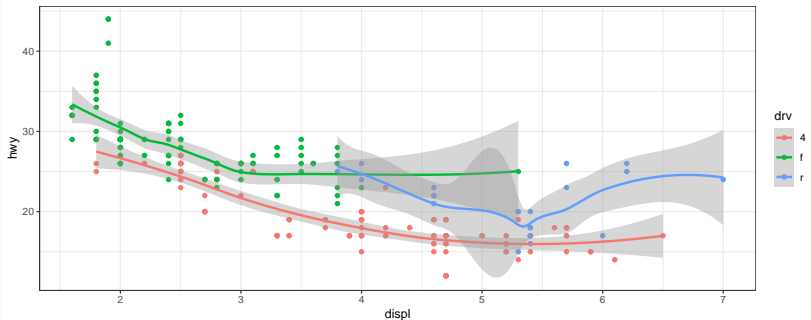
- Using the `afghan` data, visualize the relationship between age and educ.years. What is the best geom for examining this relationship?
- Add an additional aesthetic for province
- Add an additional aesthetic for age
- Remove unemployed individuals prior to plotting
- Do you think there is a correlation here? Add a smooth to examine (hint, use method = "lm")

## This is not pretty
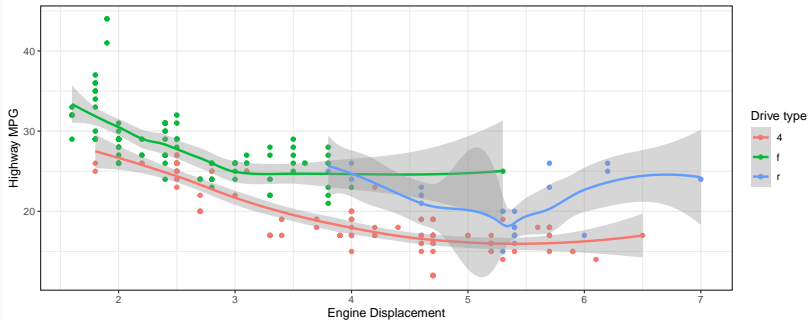
```
ggplot(mpg,
       aes(x = displ, y = hwy, color = drv)) +
  geom_point() +
  geom_smooth()
```
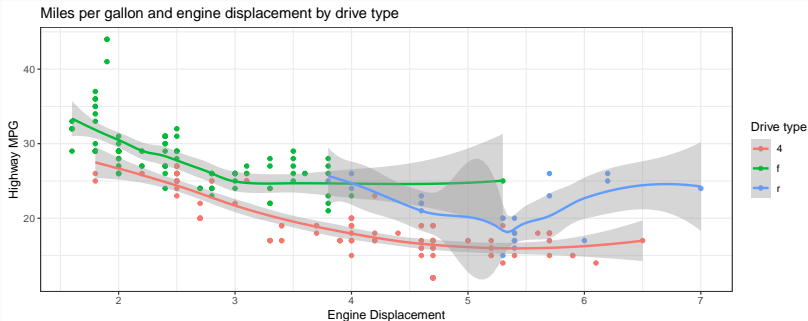
## This is better!

```
ggplot(mpg,
       aes(x = displ, y = hwy, color = drv)) +
  geom_point() +
  geom_smooth() +
  labs(x = "Engine Displacement",
       y = "Highway MPG",
       color = "Drive type")
```
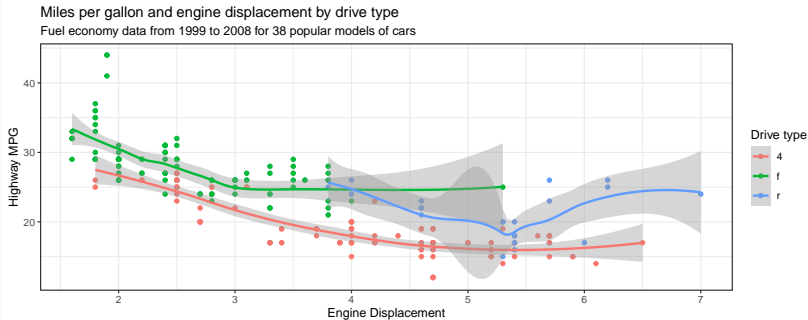
# Titles

```
ggplot(mpg,
       aes(x = displ, y = hwy, color = drv)) +
  geom_point() +
  geom_smooth() +
  labs(x = "Engine Displacement",
       y = "Highway MPG",
       color = "Drive type",
       title = "Miles per gallon and engine displacement by drive type")
```
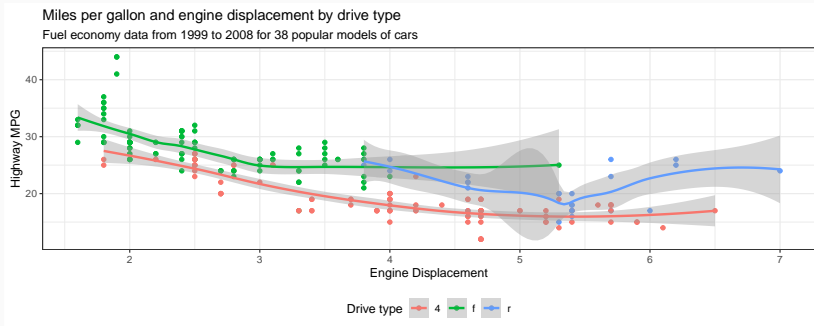
# Subtitles

```
ggplot(mpg,
       aes(x = displ, y = hwy, color = drv)) +
  geom_point() +
  geom_smooth() +
  labs(x = "Engine Displacement",
       y = "Highway MPG",
       color = "Drive type",
       title = "Miles per gallon and engine displacement by drive type",
       subtitle = "Fuel economy data from 1999 to 2008 for 38 popular models of cars")
```
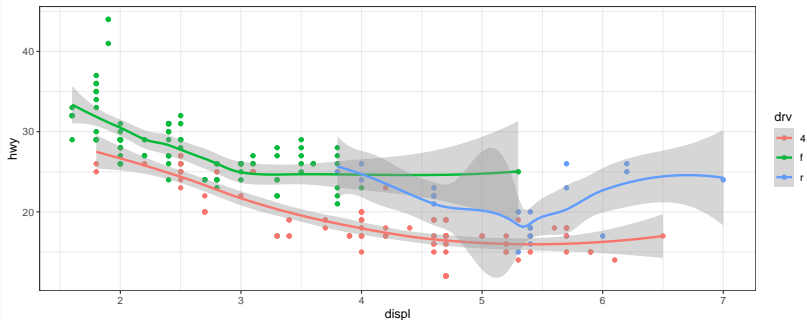
```r
ggplot(mpg,
       aes(x = displ, y = hwy, color = drv)) +
  geom_point() +
  geom_smooth() +
  labs(x = "Engine Displacement",
       y = "Highway MPG",
       color = "Drive type",
       title = "Miles per gallon and engine displacement by drive type",
       subtitle = "Fuel economy data from 1999 to 2008 for 38 popular models of cars") +
  theme(legend.position = "bottom")
```



Miles per gallon and engine displacement by drive type
Fuel economy data from 1999 to 2008 for 38 popular models of cars

If two vehicles have identical MPG and displ, they will overlap, and we can't actually see them. A jitter adds a small amount of noise to help us see all the data.

```
ggplot(mpg,
       aes(x = displ, y = hwy, color = drv)) +
  geom_point() +
  geom_smooth()
```

# Improving scatterplots

```
ggplot(mpg,
       aes(x = displ, y = hwy, color = drv)) +
  geom_jitter() +
  geom_smooth()
```