# Quantitative Sociological Analysis

# Inferential Statistics
# Hypothesis Testing and Bivariate Statistics

Part 7

April 3, 2025

# Assignment 2: Task A

(A) Firm X's sample has a mean age of 24.23 with a standard deviation (sd) of 10.07. Thus, on average, a given respondent's age is plus or minus 10.07 years from the mean age of 24.23 years.

Run Example 1 in RScript "Assignment_2" to see how the above results were obtained.

Run the line of code that begins with "age_OurFirm" to read our firm's sampled ages into R, and modify the RScript to answer the following questions:

1. What is the mean age of our firm's sample?

**Mean age = 19.82**

2. What is the standard deviation (sd) of our firm's sample age?

**sd age = 4.09**

3. In a sentence or two, interpret the standard deviation (sd) of our firm's sample age.

**The ages in the sample vary, on average, by 4.09 years from the mean of 19.82.**

- Some suggested sd had to do with range

- While range is also a measure of dispersion, it does not directly determine the sd

- Consider the following example:

---

FirmX: min 7, max 42, range = 42-7 = 35

```
age_FirmX<-c(7,10,11,16,17,18,19,19,20,20,21,23,22,26,28,30,33,35,37,39,40,42)
```

HAge: min 6, max 43, range = 43-6 = 37

```
HAge<-c(6,24,24,24,24,24,24,24,24,24,24,24,24,24,24,24,24,24,24,24,24,43)
```

FirmX: mean age = 24.23, sd = 10.07

```
> mean(age_FirmX)        > sd(age_FirmX)
[1] 24.22727             [1] 10.07085
```

HAge: mean age = 24.05, sd =    5.71

```
> mean(HAge)             > sd(HAge)
[1] 24.04545             [1] 5.711119
```
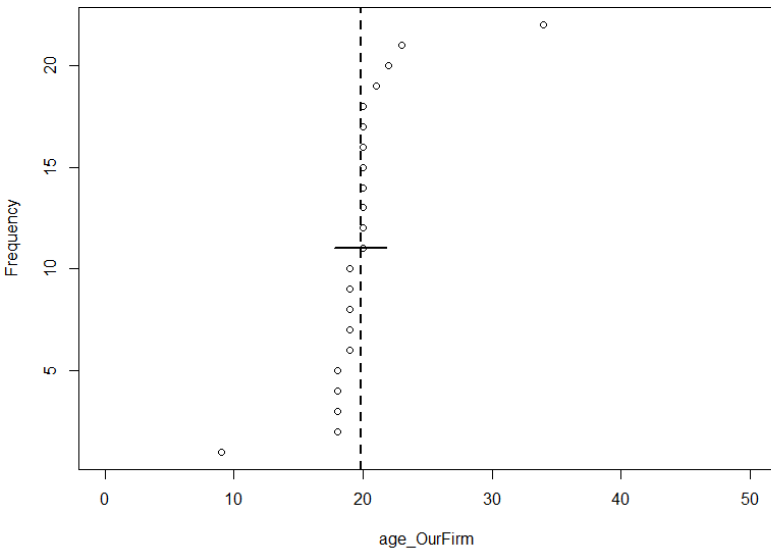
Let's see how sd measures avg. spread around a mean…

# Assignment 2: Task A cont.

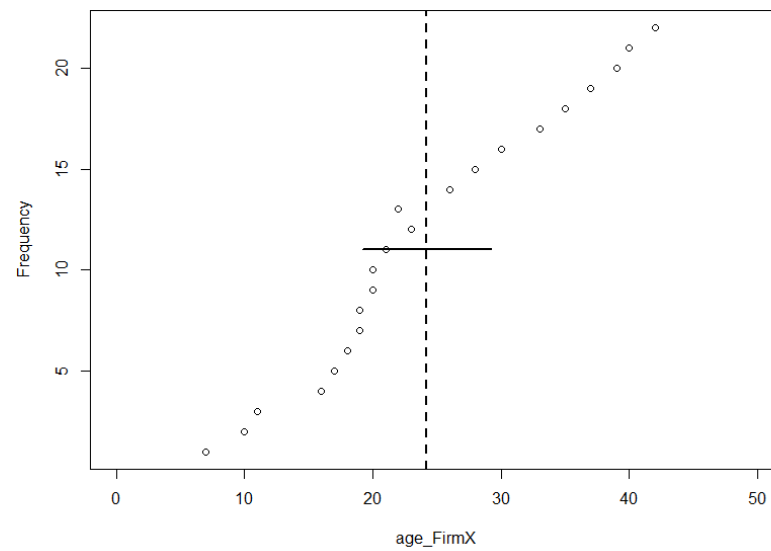$$sd = \sqrt{\frac{\sum(X_i - \bar{X})^2}{n-1}}$$
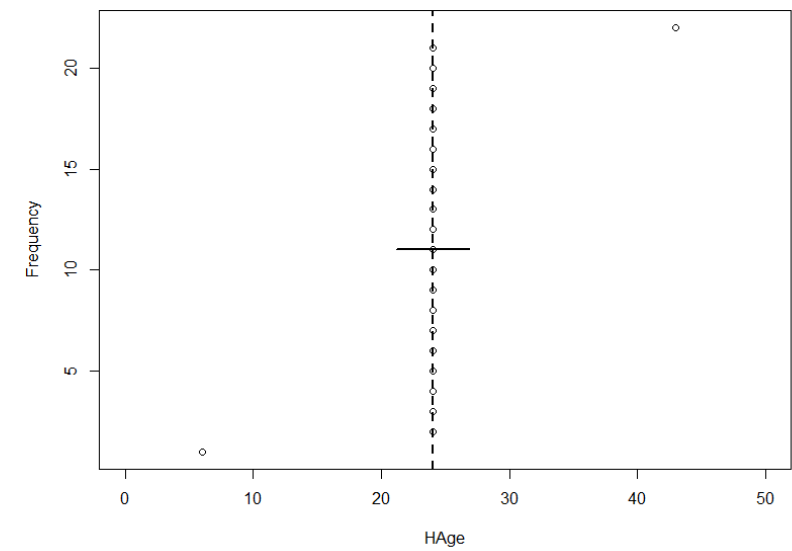


$\bar{X} = 19.82$     $sd = 4.09$

$\bar{X} = 24.23$     $sd = 10.07$

$\bar{X} = 24.05$     $sd = 5.71$

# Assignment 2: Task B

(B) Firm X's standard error (SE) for their sample age is 2.15.

Run Example 2 in RScript "Assignment_2" to see how the above results were obtained.

See how the standard error (SE) is a function of the standard deviation (sd) and the sample size (N):

$$SE = \frac{sd}{\sqrt{N}}$$

A smaller standard error (SE) means the sample estimate is more precise.

4. What is the standard error (SE) for our firm's sample age?

**SE age = 0.87**

5. Briefly explain why our firm's standard error (SE) for age is relatively smaller compared to Firm X.

**The SE for our sample's age is smaller than Firm X's because our sd is smaller. We have the same sample size as FirmX (N=22) but less spread, variation, around our sample mean.**

6. Briefly explain which firm's age estimate you expect is more precise, and why?

**I expect our firm's age estimate to be more precise because our SE is smaller.**

- Some suggested SE had to do with accuracy, implying how well a sample statistic estimates the population

- precision ≠ accuracy

When population parameters are unknown, we can never be certain how well a sample statistic estimates the population (accuracy).

When population parameters are known then features like age, sex, and race in survey data can be used to assess how well a sample estimates the population.

In general, relatively larger random samples are theoretically more generalizable, lead to more accurate estimates of population parameters.

This type of "accuracy" in terms of representativeness, generalizability, is distinct from the reliability and validity of a variable in terms of measurement.

# Assignment 2: Task C

(C) Firm X's margin of error (MoE) for their sample age is 6.08.

Run Example 3 in RScript "Assignment_2" to see how the above results were obtained.

See how the margin of error (MoE) is the function of a critical value from a probability distribution, like a Z or t distribution, and the standard error of a sample statistic, like a mean.

$$MoE = t \times SE$$

7. What is the margin of error (MoE) for our firm's sample mean age?

**MoE age = 2.47**

8. Briefly explain why our firm's MoE for mean age is relatively smaller compared to Firm X.

**The MoE is smaller for our firm's sample age because we had a smaller SE compared to Firm X, and the critical value, t score, was the same because the sample sizes were the same (N=22).**

- MoE expands on SE by incorporating a critical value from a probability distribution, which accounts for random sampling variability at a specific confidence, alpha, level.

- The MoE quantifies the expected range within which the true population parameter is likely to fall, given the sampling variability.

Again, when the population variability is unknown, we can never be certain how well a sample statistic estimates the population (accuracy).

Consider that a sample mean ($\bar{X}$) that matches the population ($\mu$) can often be obtained in ways by which the sample standard deviation does not match the population standard deviation ($\sigma$)

For example:

FirmX's mean age = 24.23, sd = 10.07

HAge mean age = 24.05, sd = 5.71

# Assignment 2: Task D

(D) Firm X's 95% confidence interval (CI) for their sample mean age is (18.15,30.31). Thus, based on their sample, it is expected that 95% of all possible random samples (N=22) from the Netflix consumer population would contain a mean age between 18.15 and 30.31.

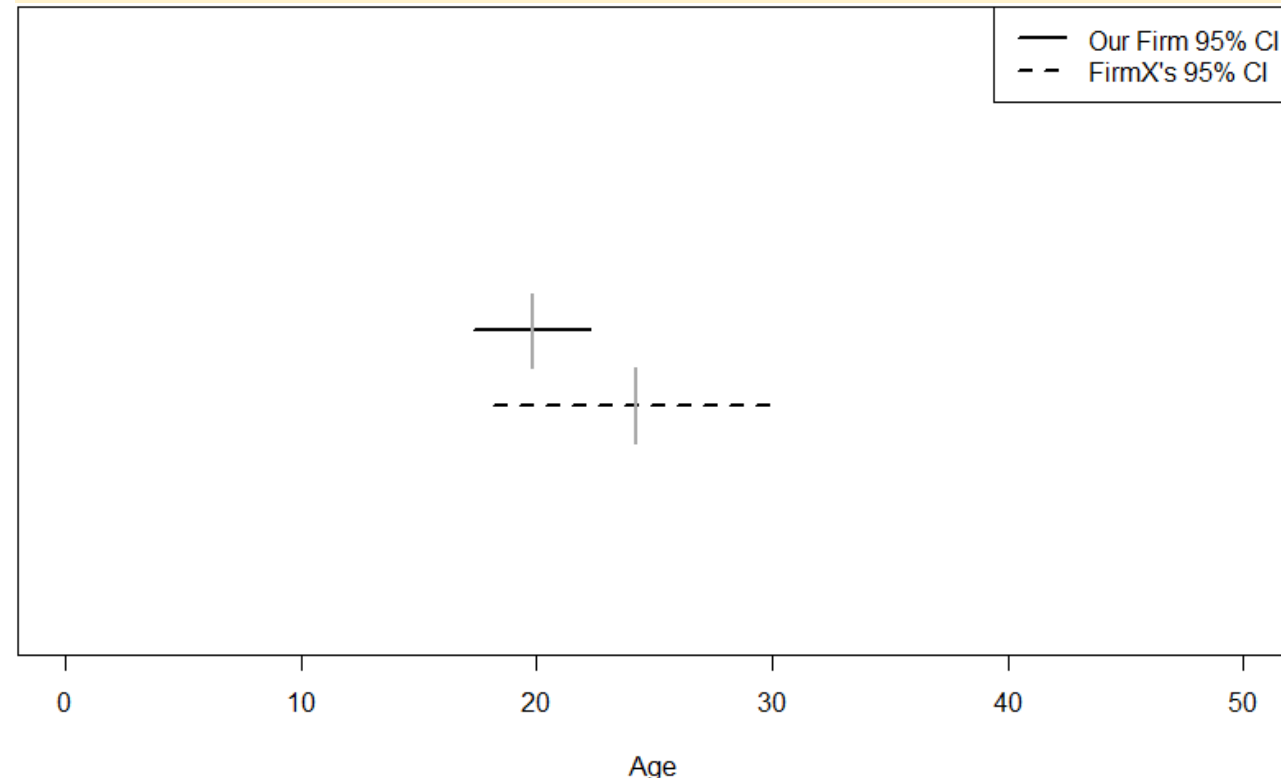Run Example 4 in RScript "Assignment_2" to see how the above results were obtained.

See how a confidence interval (CI) for a sample statistic is a function of its margin of error (MoE).

$$CI = \hat{\theta} \pm MoE$$

9. What is the 95% confidence interval (CI) for our firm's sample mean age?

**95% CI age = (17.35,22.29)**

10. In a sentence or two, interpret the 95% confidence interval (CI) for our firm's sample mean age.

**Theoretically, 95% of all possible random samples (N=22) from the Netflix consumer population would contain a mean age between 17.35 and 22.29, given the sample variability.**

11. Briefly explain why our firm's 95% confidence interval (CI) for age is relatively narrower, more precise, compared to Firm X.

**Our 95% CI for age is smaller because we had a smaller SE compared to Firm X.**

- In view of the previous tasks, recognize how a confidence interval is a function of sampling variability and sample size.

- This tells us about precision, but not accuracy

Recall that very different age distributions were obtained from random samples of the same size (N=22) from the Netflix consumer population.

However, the 95% CIs for each sample's statistic overlap.



Legend: —— Our Firm 95% CI, – – FirmX's 95% CI

Age axis: 0, 10, 20, 30, 40, 50

# Assignment 2: Task E

(E) Our firm's sample age appears to differ in some ways compared to Firm X's sample age. However, is our firm's sample mean age statistically different from Firm X's?

Two-sample t-test: two-tailed

$H_0$: Our firm's sample mean age = Firm X's sample mean age

$H_a$: Our firm's sample mean age ≠ Firm X's sample mean age

Run the two-sample t-test in RScript "Assignment_2" to obtain results, where alpha is set at 0.05.

If the p-value is < than 0.05 then reject $H_0$

If the p-value is > than 0.05 then fail to reject $H_0$

12. In a sentence or two, interpret the results from this two-sample t-test.

**We fail to reject the null hypothesis because the obtained p-value of 0.07 is greater than 0.05.**

- Our sample's mean age is not statistically different from FirmX's at the alpha=0.05 level

- Our sample statistic is more precise, but no way to tell which sample statistic is more accurate

If the hypothesis test obtained a p-value ≤ 0.05, then all we would know is that these two sample means are statistically different from one another at the 0.05 level.

No way to determine accuracy unless population parameters are known, or informed hypotheticals.

Theoretically, as sample size increases so does a sample statistic's precision and accuracy.

# Univariate vs bivariate statistics

- Thus far, we've only worked with one single variable at a time
  - univariate statistics

- This is useful for summarizing sample data in terms of central tendency and dispersion, or
  - accounting for uncertainty by estimating a sample statistic's precision and
    - its accuracy if population parameters are known

- Bivariate statistics: analysis of two variables simultaneously to determine their association
  - see how quantitative analytic techniques and statistical tests differ by level of measurement

See RScript_3 in R + RStudio instructions module on Canvas

```
284 ▾ ######################Bivariate Statistics###################################
285
286   load("E:/1_UK/2_Teaching/SOC303/Data/GSS.RData")
287   # First, attaching GSS like with Netflix data
288   attach(GSS)
289
290   ###
291   # Crosstabulation
292   ###
293   # examine two variables simultaneously, to begin to explore their association
294   # see how approaches may differ in terms of level of measurement(s)
```

# Continuous-continuous

- Is mean years of education greater among relatively more recent birth cohorts?

- GSS includes 34 random samples with surveys administered from 1972-2004 (N=64,555)
    - birth cohort ranges from 1883-2004

The distribution of education ($\bar{X} = 12.98\ (3.19)$) looks like this,   `313   table(educ_yrs)`

| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 |
|---|---|---|---|---|---|---|---|---|---|----|----|----|----|----|----|----|----|----|----|----|
| 147 | 42 | 138 | 252 | 303 | 388 | 794 | 847 | 2631 | 1995 | 2778 | 3610 | 19307 | 5268 | 7177 | 2919 | 8738 | 2102 | 2576 | 978 | 1565 |

but we are interested in education by cohort....

# Continuous-continuous

```
335   table(cohort,educ_yrs)
```

This contingency table, or crosstab, for years of education by cohort is massive.

Summarizing these data in terms of mean years of education would make understanding these data more manageable.

However, this would also include 122 cohorts, rows in the table

Let's consider mean years of education by cohort depicted in a plot…

# Continuous-continuous

- Is mean years of education greater among relatively more recent birth cohorts?
    - It looks so, but what's up with…



Mean Years of Education by Birth Cohort: GSS 1972-2022 N=64,555

$H_0: mean\ years\ of\ edu \leq among\ more\ recent\ cohorts$
$H_a: mean\ years\ of\ edu > among\ more\ recent\ cohorts$

Last week we determined whether one mean was statistically different from another mean.

The present hypothesis test involves multiple means, which requires a different type of hypothesis test – more complex than we are ready for now, but we will build up to this.

Let's first explore some different bivariate combinations by level of measurement …

# Continuous-nominal

- Is mean years of education greater among males compared females?

$H_0: male \ \bar{X}_{edu} \leq female \ \bar{X}_{edu}$
$H_a: male \ \bar{X}_{edu} > female \ \bar{X}_{edu}$

The distribution of education by sex looks like this    `359  table(female,educ_yrs)`

| female | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 |
|--------|----|----|----|-----|-----|-----|-----|-----|------|------|------|------|-------|------|------|------|------|------|------|-----|-----|
| 0 | 74 | 28 | 75 | 133 | 156 | 182 | 385 | 400 | 1121 | 838 | 1180 | 1546 | 7938 | 2158 | 3249 | 1341 | 4092 | 1029 | 1201 | 546 | 910 |
| 1 | 73 | 14 | 63 | 119 | 147 | 206 | 409 | 447 | 1510 | 1157 | 1598 | 2064 | 11369 | 3110 | 3928 | 1578 | 4646 | 1073 | 1375 | 432 | 655 |



Males: Years of Education (mean=13.11, sd=3.34)



Females: Years of Education (mean=12.87, sd=3.06)

See how considering within-group proportions can be helpful in bivariate statistics…

# Continuous-nominal

- Is mean years of education greater among males compared females?
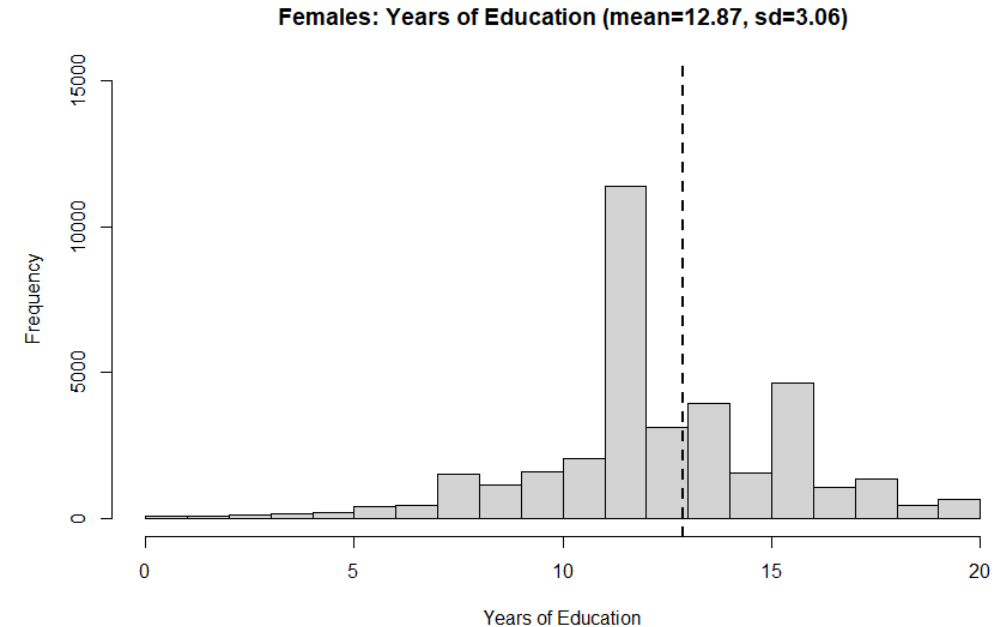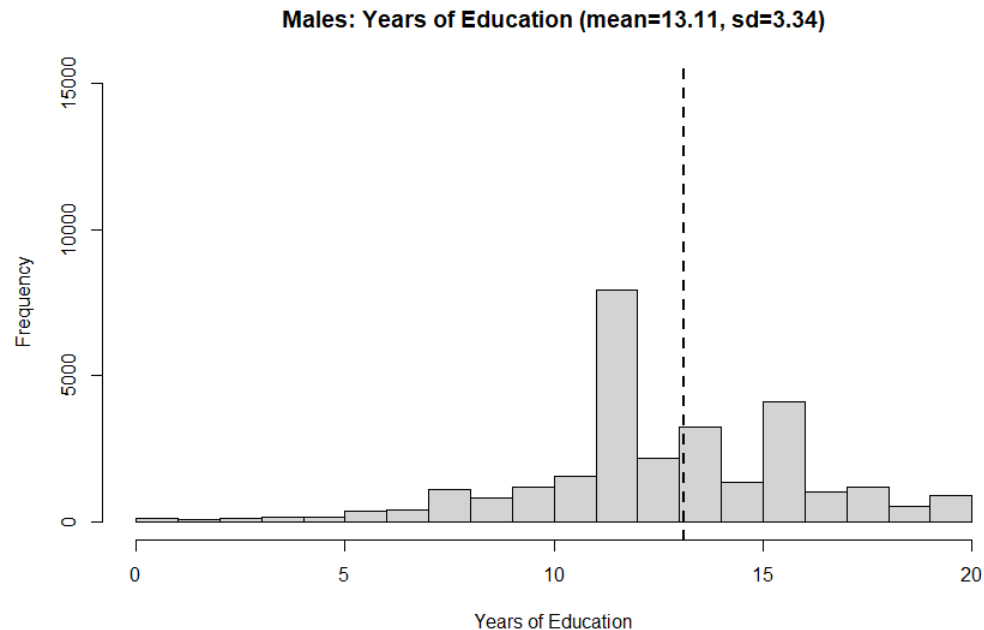
$H_0: male\ \bar{X}_{edu} \leq female\ \bar{X}_{edu}$

$H_a: male\ \bar{X}_{edu} > female\ \bar{X}_{edu}$

The distribution of education by sex looks like this
`381  prop.table(table(female,educ_yrs))`

| female | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 |
|--------|---|---|---|---|---|---|---|---|---|---|----|----|----|----|----|----|----|----|----|----|----|
| 0 | 0.00115 | 0.00043 | 0.00116 | 0.00206 | 0.00242 | 0.00282 | 0.00596 | 0.00620 | 0.01737 | 0.01298 | 0.01828 | 0.02395 | 0.12296 | 0.03343 | 0.05033 | 0.02077 | 0.06339 | 0.01594 | 0.01860 | 0.00846 | 0.01410 |
| 1 | 0.00113 | 0.00022 | 0.00098 | 0.00184 | 0.00228 | 0.00319 | 0.00634 | 0.00692 | 0.02339 | 0.01792 | 0.02475 | 0.03197 | 0.17611 | 0.04818 | 0.06085 | 0.02444 | 0.07197 | 0.01662 | 0.02130 | 0.00669 | 0.01015 |



Males: Years of Education (mean=13.11, sd=3.34)

Females: Years of Education (mean=12.87, sd=3.06)

See how considering within-group proportions can be helpful in bivariate statistics…

# Continuous-nominal

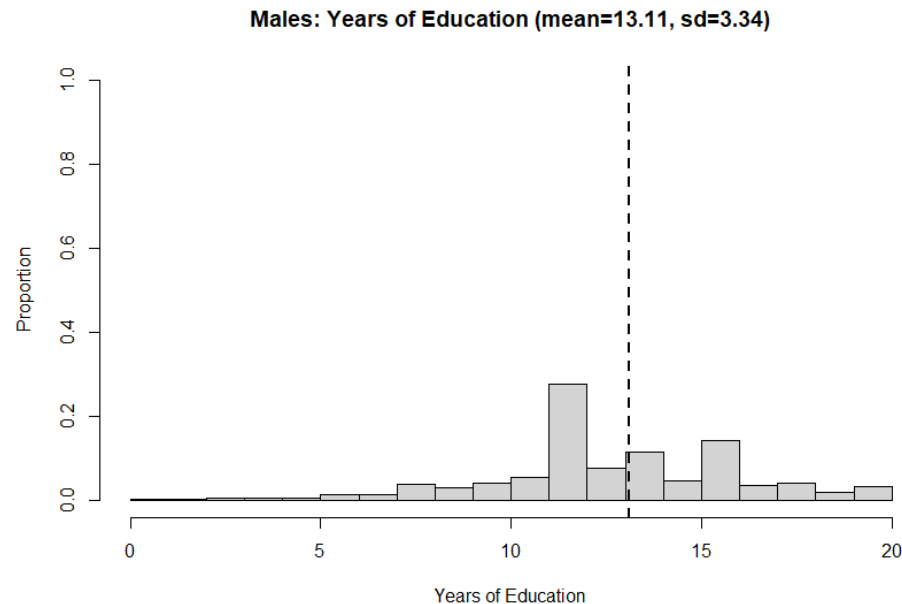- Is mean years of education greater among males compared females?

$H_0: male\ \bar{X}_{edu} \leq female\ \bar{X}_{edu}$
$H_a: male\ \bar{X}_{edu} > female\ \bar{X}_{edu}$

two-sample t-test

$$t = \frac{\bar{X}_1 - \bar{X}_2}{sd_p\sqrt{\frac{1}{n_1} + \frac{1}{n_2}}}$$

$sd_p$: pooled standard deviation

$$sd_p = \sqrt{\frac{(n_1 - 1)sd_1^2 + (n_2 - 1)sd_2^2}{n_1 + n_2 - 2}}$$

```
406   t.test(educ_yrs~female,alternative="greater")

data:  educ_yrs by female
t = 9.2, df = 58657, p-value <2e-16
alternative hypothesis: true difference in means between group 0 and group 1 is greater than 0
95 percent confidence interval:
 0.1928    Inf
sample estimates:
mean in group 0 mean in group 1
        13.11           12.87
```

**Mean Years of Education by Sex**



The difference in means is small, but statistically significant.

Why might that be? Hint: consider n, SE, and MoE

Male      95% CI (13.07,13.15)   ` 416   t.test(educ_yrs[female==0],conf.level=(0.95))`

Female   95% CI (12.84,12.90)   ` 418   t.test(educ_yrs[female==1],conf.level=(0.95))`

Let's further explore different bivariate combinations by level of measurement ...

# Ordinal-nominal

- Unlike continuous variables, the mean is not very meaningful for categorical variables

- Let's consider whether happiness differs by sex

  - 1=not very happy, 2=pretty happy, 3=very happy

The distribution of happiness by sex looks like this

```
454  table(happy,female)
```

```
          female
happy       0       1
    1    3842    4909
    2   16169   19990
    3    8571   11074
```



**Happiness: Males**

**Happiness: Females**

See how considering within-group proportions can be helpful in bivariate statistics...

# Ordinal-nominal

- Does happiness differ by sex?

$H_0$: happiness is independent of sex

$H_a$: happiness is dependent on sex

The distribution of happiness by sex looks like this

```
468  prop.table(table(happy,female))
```

```
           female
happy        0        1
    1  0.060  0.076
    2  0.250  0.310
    3  0.133  0.172
```

Margins can be an incredibly helpful tool

```
472  addmargins(prop.table(table(happy,female),2),1)
```

```
           female
happy      0      1
    1    0.13   0.14
    2    0.57   0.56
    3    0.30   0.31
  Sum   1.00   1.00
```



Happiness: Males



Happiness: Females

We need a different type of hypothesis test to determine if happiness differs by sex…

# Hypothesis testing: chi-squared

Same steps as with any hypothesis test,

1. Make assumptions and meet test requirements
   - random sample; level of measurement; sample size; are parameters un/known
2. State the null ($H_0$) and alternative ($H_a$) hypothesis
   - $H_0$ no difference;  $H_a$ there is a difference
3. Choose a significance level (critical value)
   - e.g., $\alpha = 0.05$
4. Compute the test statistic
5. Draw a conclusion and interpret the test results
   - If $x^2 > critical\ value$ then reject $H_0$

but different procedures for computing the test statistic…

# Chi-squared test of independence

4a. Create a bivariate contingency table

4b. Find the discrepancies between observed ($f_o$) and expected ($f_e$) counts

4c. Sum the discrepancies between $f_o$ and $f_e$

$$\chi^2 = \sum \frac{(f_o - f_e)^2}{f_e}$$

4d. Determine statistical significance

$$df = (rows - 1) \times (\text{columns} - 1)$$

but different procedures for computing the test statistic…

# Chi-squared test of independence: step 4a

4a. Create a bivariate contingency table, sometimes called a crosstab

| | Columns (IV) | | |
|---|---|---|---|
| Rows (DV) | Column 2 | Column 1 | |
| Row 1 | cell a | cell b | Row Marginal 3 |
| Row 2 | cell c | cell d | Row Marginal 2 |
| | Column Marginal 1 | Column Marginal 2 | N |

Let's check out our example for happiness by sex…

# Chi-squared test of independence: step 4a

4a. Create a bivariate contingency table, sometimes called a crosstab

|  | Columns (IV) | | |
|---|---|---|---|
|  | male | female |  |
| not too happy | 3842 | 4909 | 8751 |
| pretty happy | 16169 | 19990 | 36159 |
| very happy | 8571 | 11074 | 19645 |
|  | 28582 | 35973 | 64555 |

Probability theory: if two events are independent then their joint probability equals the product of their respective marginal probabilities. So, ….

# Chi-squared test of independence: step 4b

4b. Find the discrepancies between frequency observed ($f_o$) and frequency expected ($f_e$)

$$\chi^2 = \sum \frac{(f_o - f_e)^2}{f_e}$$

$f_o$ = the cell frequencies observed in the bivariate table
$f_e$ = the cell frequencies that would be expected if the
variables were independent

$$f_e = \frac{row\ marginal\ \times column\ marginal}{N}$$

Let's look at this in another way…

# Chi-squared test of independence: step 4b

|                          | $Y = 0$ | $Y = 1$ | Total (Marginals for Y) |
| ------------------------ | ------- | ------- | ----------------------- |
| $X = 0$                  | $a$     | $b$     | $a + b$                 |
| $X = 1$                  | $c$     | $d$     | $c + d$                 |
| Total (Marginals for X)  | $a + c$ | $b + d$ | $a + b + c + d = n$     |

- If X and Y are independent, then $p(x, y) = p(x) \times p(y)$

- So, $p(x = 0, y = 0) = p(x = 0) \times p(y = 0) = \dfrac{a+c}{n} \times \dfrac{a+b}{n}$

- Expected cell count is: $\dfrac{(O_{x=0})(O_{y=0})}{n}$

Let's check out our example for happiness by sex…

# Chi-squared test of independence: step 4b-c

4b. Find the discrepancies between frequency observed ($f_o$) and frequency expected ($f_e$)

| | Columns (IV) | | |
|---|---|---|---|
| | male | female | Total |
| not too happy | 3842 | 4909 | 8751 |
| Expected | $\frac{(8751)(28582)}{64555}$ = 3875 | $\frac{(8751)(35973)}{64555}$ = 4876 | |
| $\frac{(O-E)^2}{E}$ | 0.27 | 0.22 | 0.49 |
| pretty happy | 16169 | 19990 | 36159 |
| Expected | $\frac{(36159)(28582)}{64555}$ = 16009 | $\frac{(36159)(35973)}{64555}$ = 20149 | |
| $\frac{(O-E)^2}{E}$ | 1.59 | 1.26 | 2.85 |
| very happy | 8571 | 11074 | 19645 |
| Expected | $\frac{(19645)(28582)}{64555}$ = 8698 | $\frac{(19645)(35973)}{64555}$ = 10947 | |
| $\frac{(O-E)^2}{E}$ | 1.85 | 1.47 | 3.32 |
| Total | 28582 | 35973 | 64555 |

4c.

$\chi^2 = 0.27 + 0.22 + 1.59 + 1.26 + 1.85 + 1.47 = 6.66$

# Chi-squared test of independence: step 4d

- Use the $x^2$ probability distribution table to identify critical values
  - like the Z and t tables for a sample mean

- Find the critical value that corresponds with the degrees of freedom ($df$), and
  - your selected alpha ($\alpha$), significance level

$$df = (rows - 1) \times (\text{columns} - 1)$$

Like with the Z- and t-tests, I won't have you do this, but let's check it out…

# Chi-squared test of independence: step 5

5. Draw a conclusion and interpret the test results

- If $x^2 > critical\ value$ then reject $H_0$
  - $df = (3-1) \times (2-1) = 2 \times 1 = 2$

$df = (rows - 1) \times (columns - 1)$

  - $\alpha = 0.05$ , or 95% confidence level

Significance level (α)

| Degrees of freedom (df) | .99 | .975 | .95 | .9 | .1 | .05 | .025 | .01 |
|---|---|---|---|---|---|---|---|---|
| 1 | -------- | 0.001 | 0.004 | 0.016 | 2.706 | 3.841 | 5.024 | 6.635 |
| 2 | 0.020 | 0.051 | 0.103 | 0.211 | 4.605 | 5.991 | 7.378 | 9.210 |
| 3 | 0.115 | 0.216 | 0.352 | 0.584 | 6.251 | 7.815 | 9.348 | 11.345 |
| 4 | 0.297 | 0.484 | 0.711 | 1.064 | 7.779 | 9.488 | 11.143 | 13.277 |
| 5 | 0.554 | 0.831 | 1.145 | 1.610 | 9.236 | 11.070 | 12.833 | 15.086 |

$P(\chi^2)$

0.5

0.4

df = 1
df = 2
df = 3
df = 4
df = 5

0.3

0.2

0.1

0  1  2  3  4  5  6  7  8  9  10  11  12  $\chi^2$

As the df increases the chi square distribution tends to normal distribution.
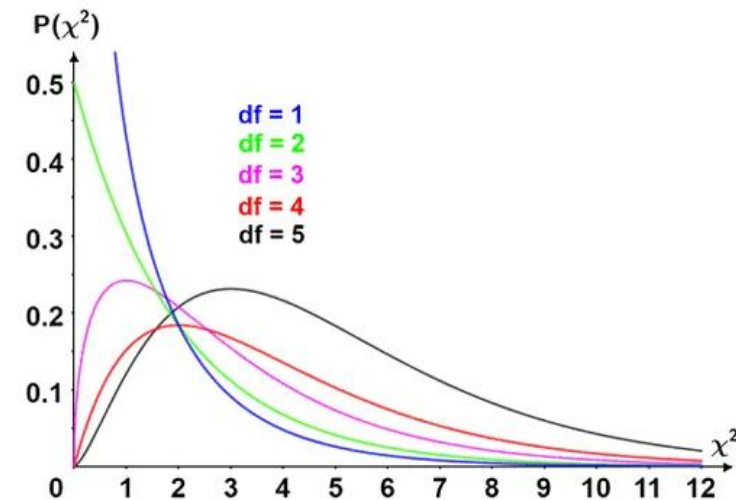
- 6.66 > 5.99, so reject $H_0$
  - there is a statistically significant association between happiness and biologically assigned sex
- Or, just use the chisq.test command in RStudio

```
493    chisq.test(happy_by_sex)
```

```
        Pearson's Chi-squared test

data:  happy_by_sex
X-squared = 7, df = 2, p-value = 0.04
```

Again, this is a large sample, so hypothesis tests are likely to be statistically significant

# Chi-squared test of independence: limitations

- $x^2$ test statistics are unstable with small cell counts
  - Rule of thumb: don't use when expected cell counts <5
    - Instead, could use Fisher's exact test

- $x^2$ tests are more sensitive to sample size than other test statistics
  - e.g., in $\frac{(O-E)^2}{E}$ numerator grows faster than the denominator as $n$ increase