

gretl Notes

Week 4

Module 2 Week 4

Intro

- Dataset used: `National Election Study 2008.gdt`
 - Open with `File -> Open data -> User file...`
- Open up a new script editor: `File -> Script files -> New script -> gretl script`
- Checkout summary statistics:
 - ```
1 | summary --simple
```
  - Most Variables are "thermometer" values: higher => more positive feelings of a respondent
- How to get summary statistics across categories:
  - ```
1 | summary obama_therm --by=race_ethnicity
```
- Give `race_ethnicity` variable a label (from data documentation):
 - ```
1 | setinfo race_ethnicity --description="1=White, 2=Black, 3=Hispanic, 4=Other"
```

### Two sample hypothesis test for means, with independent sample

- Does the mean for the big business thermometer differ across whether or not someone has the Internet?
  - $H_0: \mu_{bigbus\_withinternet} - \mu_{bigbus\_withoutinternet} = 0$   
 $H_a: \mu_{bigbus\_withinternet} - \mu_{bigbus\_withoutinternet} \neq 0$
  - 5 - no internet, 1 - internet:
    - ```
1 | freq internet
```

- ```

1 Frequency distribution for internet, obs 1-2323
2
3 frequency rel. cum.
4
5 1 1615 69.58% 69.58% *****
6 5 706 30.42% 100.00% *****
7
8 Missing observations = 2 (0.09%)

```

- Fix coding:

- ```
1 series netaccess = (internet==1)
```

- Let `gret1` know it's discrete:

- ```
1 discrete netaccess
```

- Fixed version:

- ```
1 freq netaccess
```

- ```

1 Frequency distribution for netaccess, obs 1-2323
2
3 frequency rel. cum.
4
5 0 706 30.42% 30.42% *****
6 1 1615 69.58% 100.00% *****
7
8 Missing observations = 2 (0.09%)

```

- Summary stats on `bigbusiness` by `netaccess` to see the means we'll compare:

- ```
1 summary bigbus_therm --by=netaccess
```

- ```

1 Summary statistics for bigbus_therm, by value of netaccess
2
3 netaccess = 0 (n = 706):
4 Mean 59.994
5 Median 60.000
6 Minimum 0.0000
7 Maximum 100.00
8 Standard deviation 22.078
9 C.V. 0.36800
10 Skewness -0.28223
11 Ex. kurtosis 0.42036
12 5% percentile 15.000
13 95% percentile 100.00

```

```

14 Interquartile range 20.000
15 Missing obs. 0
16
17 netaccess = 1 (n = 1615):
18 Mean 53.182
19 Median 50.000
20 Minimum 0.0000
21 Maximum 100.00
22 Standard deviation 20.595
23 C.V. 0.38725
24 Skewness -0.15362
25 Ex. kurtosis 0.41778
26 5% percentile 15.000
27 95% percentile 85.000
28 Interquartile range 30.000
29 Missing obs. 0

```

- Is this a real difference in means or just by chance?

- ```
1 anova bigbus_therm netaccess
```

- ```

1 Analysis of Variance, response = bigbus_therm, treatment =
 netaccess:
2
3 Sum of squares df Mean square
4
5 Treatment 22797.6 1 22797.6
6 Residual 1.02821e+06 2319 443.384
7 Total 1.05101e+06 2320 453.02
8
9 F(1, 2319) = 22797.6 / 443.384 = 51.4172 [p-value 1e-12]
10
11 Level n mean std. dev
12
13 0 706 59.9943 22.078
14 1 1615 53.182 20.595
15
16 Grand mean = 55.2542

```

- Conclusion:

- `p-value 1e-12`

- p-value is approximately 0. Therefore, we reject  $H_0$ , we believe there is a statistically sign difference in the mean of the big business therm between those with and without Internet access

- $T\text{-Stat} = \sqrt{F\text{-Stat}}$ :

- ```
1 ? eval sqrt(51.4172)
2 7.1705788
```

- Use the T-value and the T-Distribution instead of the F-Distribution:

- ```
1 ? pvalue t 2319 7.1705788
2 t(2319): area to the right of 7.17058 = 4.99898e-13
3 (two-tailed value = 9.99795e-13; complement = 1)
```

- Conclusion:

- p-value =  $9.99795e - 13 \Rightarrow$  Almost the same as the earlier F-Stat-based p-value =  $1e - 12$ , as expected!
  - Can use `anova` output to perform t-tests!

## Chi-squared test for independence between two nominal variables

- Are the variables internet access and gun ownership related or independent?

- $H_0$ : `netaccess` and `gunown` are independent
  - $H_a$ : `netaccess` and `gunown` are not independent

- ```
1 ? freq gunown
2
3 Frequency distribution for gunown, obs 1-2323
4
5           frequency    rel.    cum.
6
7      1           651    28.63%   28.63% *****
8      5          1623    71.37%  100.00% *****
9
10 Missing observations = 49 ( 2.11%)
```

- The coding is messed up again, need to re-code (5 = no, 1 = yes):

- ```
1 ? series gunown = (gunown==1)
2 Replaced series gunown (ID 56)
3 ? freq gunown
4
5 Frequency distribution for gunown, obs 1-2323
6
7 frequency rel. cum.
8
9 0 1623 71.37% 71.37% *****
10 1 651 28.63% 100.00% *****
11
12 Missing observations = 49 (2.11%)
```

- Perform the Chi-Squared Test for Independence:

```

1 ? xtab netaccess gunown
2
3 Cross-tabulation of netaccess (rows) against gunown (columns)
4
5 [0][1] TOT.
6
7 [0] 547 149 696
8 [1] 1074 502 1576
9
10 TOTAL 1621 651 2272
11
12 51 missing values
13
14 Pearson chi-square test = 25.7635 (1 df, p-value = 3.8591e-07)

```

- $df = (cols - 1) * (rows - 1)$

- $p\text{-value} = 3.8591e - 07 \Rightarrow$  Reject the null hypothesis, the variables `netaccess` and `gunown` are related to one another

- Look at the marginals by `gunown` to investigate the relationship:

```

1 ? xtab netaccess gunown --column
2
3 Cross-tabulation of netaccess (rows) against gunown (columns)
4
5 [0][1] TOT.
6
7 [0] 33.7% 22.9% 30.6%
8 [1] 66.3% 77.1% 69.4%
9
10 TOTAL 1621 651 2272
11
12 51 missing values
13
14 Pearson chi-square test = 25.7635 (1 df, p-value = 3.8591e-07)

```

- Conclusion:

- Those who own a gun seem to be more likely to have internet access.

- GUI: `View` -> `Cross Tabulation` -> Move `netaccess` and `gunown` to the right -> `OK`

## One-way ANOVA

- Is there a difference in the mean of the federal government term across race/ethnicity?

```
1 ? anova fedgov_term race_ethnicity
2
3 Analysis of Variance, response = fedgov_term, treatment =
 race_ethnicity:
4
5 Sum of squares df Mean square
6
7 Treatment 81808.3 3 27269.4
8 Residual 1.03246e+06 2311 446.759
9 Total 1.11427e+06 2314 481.533
10
11 F(3, 2311) = 27269.4 / 446.759 = 61.0384 [p-value 5.66e-38]
12
13 Level n mean std. dev
14
15 1 1159 46.3658 19.925
16 2 569 58.8067 22.351
17 3 509 57.8625 22.260
18 4 78 48.141 21.987
19
20 Grand mean = 52.0112
```

- Conclusion:

- $p\text{-value} = 5.66e - 38 \Rightarrow$
- Reject  $H_0$ , there is a statistically significant difference in mean of `fedgov_term` across race/ethnicity categories
- GUI Anova: `Model` -> `Other linear models` -> `ANOVA` -> Choose Response ( `fedgov_term` ) and Treatment ( `race_ethnicity` ) Variables (No block variable, that's for 2-way ANOVA) -> `OK`