

Introduction to Statistics in STATA

Basic univariate descriptive statistics. Basics of bivariate inferential statistics.

Tomáš Oleš

Department of Economic Policy
Faculty of Economics and Finance

February 1, 2025

Agenda

- Obtain basic univariate descriptive statistics using STATA.
- Obtain basic bivariate inferential statistics using STATA.
- Perform simple bivariate analyses using STATA.

Why Do We Need Statistics?



Why Do We Need Statistics?

- "Facts are stubborn, but statistics are more pliable." - Mark Twain
- "Before the curse of statistics fell upon mankind we lived a happy, innocent life, full of merriment and go and informed by fairly good judgment." - Hilaire Belloc

What is Univariate Descriptive Statistics?

Univariate statistics describe and summarize a single variable in a dataset.

These include:

- Measures of Central Tendency: mean, median, mode.
- Measures of Dispersion: variance, standard deviation, range, interquartile range

BRACE YOURSELF



... math is coming!

Univariate Descriptive Statistics

- Measures of Central Tendency:
 - Mean: $\bar{x} = \frac{\sum x_i}{n}$
 - Median: Middle value of ordered data
 - Mode: Most frequently occurring value
- Measures of Dispersion:
 - Variance: $s^2 = \frac{\sum (x_i - \bar{x})^2}{n-1}$
 - Standard Deviation: $s = \sqrt{s^2}$
 - Range: $\max(x) - \min(x)$
 - Interquartile Range: $Q_3 - Q_1$

Descriptive Statistics in STATA

```
*Load built-in dataset
```

```
sysuse nlsw88, clear
```

```
*Basic summary statistics
```

```
summarize wage hours
```

```
*Detailed statistics with percentiles
```

```
summarize wage, detail
```

Mean and standard deviation provide insights into central tendency and dispersion. The 'detail' option includes percentiles, skewness, and kurtosis.

Frequency and Percentiles

*Frequency table for categorical variables

```
tabulate race
```

```
tabulate industry, sort
```

*Percentiles and quartiles

```
centile wage, centile(25 50 75)
```

- 'tabulate' provides counts and percentages.
- 'centile' helps understand data distribution.

Visualizing Data: Histograms

```
*Histogram for wage  
histogram wage, bin(20) normal
```

- Helps visualize distribution.
- Can compare against a normal curve.
- Adjust 'bin()' for granularity.

Skewness and Kurtosis

Skewness: Measures symmetry of distribution.

$$\text{Skewness} = \frac{\sum (x_i - \bar{x})^3}{(n - 1)s^3} \quad (1)$$

Kurtosis: Measures tail heaviness.

$$\text{Kurtosis} = \frac{\sum (x_i - \bar{x})^4}{(n - 1)s^4} \quad (2)$$

Skewness and Kurtosis in STATA

- Positive skew: Right-tailed distribution.
- Negative skew: Left-tailed distribution.
- Kurtosis > 3 : Heavy tails (leptokurtic).
- Kurtosis < 3 : Light tails (platykurtic).

```
*Test for skewness and kurtosis  
sktest wage
```

Boxplots for Outlier Detection

```
*Boxplot for wage  
graph box wage
```

- Median line represents the central tendency.
- Whiskers show variability.
- Outliers appear as individual points.

Descriptive Statistics and Graphs: workout1.dta dataset

```
*Back to using the workout1 dataset
```

```
use datasets/workout1,clear
```

```
encode v07, gen(v07_num) *turns v07 into numeric
```

```
*Shows frequency distributions
```

```
tab v07_num
```

```
fre v07_num
```

```
hist v07_num, discrete percent addlabel xlabel(1/2, valuelabel noticks)
```

```
graph pie, over(v07_num) plabel(_all percent)
```

```
*Open a STATA installed dataset
```

```
sysuse auto,clear
```

```
*Summary statistics
```

```
sum price
```

```
sum price, d
```

```
mean price
```

```
tabstat price weight length, stats(mean sd range count) by(foreign)
```

```
tabstat price weight length, stats(mean sd range count) by(foreign) col(  
    stats) nototal
```

```
tab foreign rep78, sum(mpg)
```

Plotting

```
*Open a stata-installed dataset
sysuse nlsw88,clear
hist wage, frequency
replace race=. if race==3 *category 3 set to missing
graph box wage, by(race)
```

Bivariate Inferential Statistics

These include:

- Correlation
- t-test
- ANOVA
- Chi-squared test

BRACE YOURSELF



... math is coming!

Correlation Analysis

Definition: Examines the relationship between two continuous variables.

Formula: Pearson Correlation Coefficient

$$r = \frac{\sum(x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum(x_i - \bar{x})^2} \cdot \sqrt{\sum(y_i - \bar{y})^2}} \quad (3)$$

```
*Open STATA-installed data
```

```
sysuse nlsw88, clear
```

```
*Correlation analysis
```

```
pwcorr wage ttl_exp, star(0.05) obs
```

```
corr wage ttl_exp
```

Interpretation: A moderate positive/negative correlation between wage and experience ($r = 0.27$, $p < 0.05$).

Independent t-test

Definition: Tests if the means of a variable differ between two independent groups.

$$t = \frac{\bar{x}_1 - \bar{x}_2}{\sqrt{s_p^2 \left(\frac{1}{n_1} + \frac{1}{n_2} \right)}} \quad (4)$$

where

$$s_p^2 = \frac{(n_1 - 1)s_1^2 + (n_2 - 1)s_2^2}{n_1 + n_2 - 2} \quad (5)$$

***Independent t-test**

```
ttest wage, by(collgrad)
ttest wage, by(collgrad) unequal
sdtest wage, by(collgrad)
```

Interpretation: The mean hourly wage of those with a college degree differs non-significantly/significantly from those without ($t(2244) = -13, p < 0.001$).

Analysis of Variance (ANOVA)

Definition: Tests for differences between more than two independent means.

Formula: F-ratio

$$F = \frac{\text{Between-group variance}}{\text{Within-group variance}} = \frac{MS_{\text{between}}}{MS_{\text{within}}} \quad (6)$$

*ANOVA

```
tab race, sum(wage)
anova wage race
pwcompare race, pveffects
```

Interpretation: There is a non-significant/significant difference in mean hourly wages across racial groups.

Chi-Squared Test

Definition: Tests the relationship between two categorical variables.

Formula: Chi-Squared Statistic

$$\chi^2 = \sum \frac{(O_i - E_i)^2}{E_i} \quad (7)$$

where:

- O_i is the observed frequency,
- E_i is the expected frequency.

```
*Chi-square test
```

```
tab union collgrad, col chi2
```

Interpretation: There is a non-significant/significant relationship between union membership and having a college degree.

- Laffers, L. (2021). *Draft poznámok k predmetu Moderná Aplikovaná regresia 1*. UMB Banská Bystrica.
- Mehmetoglu, M. and Jakobsen, T. G. (2022). *Applied Statistics using Stata: a Guide for the Social Sciences*. Sage.