

# Data Analysis with Stata 14.1 Cheat Sheet

For more info see Stata's reference manual (stata.com)

Results are stored as either **r** -class or **e** -class. See [Programming Cheat Sheet](#)

## Summarize Data

Examples use `auto.dta` (`sysuse auto, clear`) unless otherwise noted

**univar** price mpg, **boxplot**

calculate univariate summary, with box-and-whiskers plot

**stem** mpg

return stem-and-leaf display of mpg

**summarize** price mpg, **detail**

calculate a variety of univariate summary statistics

**ci** mpg price, **level**(99)

compute standard errors and confidence intervals

**correlate** mpg price

return correlation or covariance matrix

**pwcorr** price mpg weight, **star**(0.05)

return all pairwise correlation coefficients with sig. levels

**mean** price mpg

estimates of means, including standard errors

**proportion** rep78 foreign

estimates of proportions, including standard errors for categories identified in varlist

**ratio**

estimates of ratio, including standard errors

**total** price

estimates of totals, including standard errors

## Statistical Tests

**tabulate** foreign rep78, **chi2 exact expected**

tabulate foreign and repair record and return chi<sup>2</sup> and Fisher's exact statistic alongside the expected values

**ttest** mpg, **by**(foreign)

estimate t test on equality of means for mpg by foreign

**prtest** foreign == 0.5

one-sample test of proportions

**ksmirnov** mpg, **by**(foreign) **exact**

Kolmogorov-Smirnov equality-of-distributions test

**ranksom** mpg, **by**(foreign) **exact**

equality tests on unmatched data (independent samples)

**anova** systolic drug

analysis of variance and covariance

**pwmean** mpg, **over**(rep78) **pveffects mcompare**(tukey)

estimate pairwise comparisons of means with equal variances include multiple comparison adjustment

## Estimation with Categorical & Factor Variables

more details at <http://www.stata.com/manuals14/u25.pdf>

### CONTINUOUS VARIABLES

measure something

### CATEGORICAL VARIABLES

identify a group to which an observations belongs

### INDICATOR VARIABLES

**T F** denote whether something is true or false

### OPERATOR

|     |                                |
|-----|--------------------------------|
| i.  | specify indicators             |
| ib. | specify base indicator         |
| c.  | treat variable as continuous   |
| o.  | omit a variable or indicator   |
| #   | specify interactions           |
| ##  | specify factorial interactions |

### EXAMPLE

|  |   |
|--|---|
| <code>regress price i.rep78</code>                   | specify rep78 variable to be an indicator variable                                    |
| <code>regress price ib(3).rep78</code>               | set the third category of rep78 to be the base category                               |
| <code>regress price i.foreign#c.mpg i.foreign</code> | treat mpg as a continuous variable and specify an interaction between foreign and mpg |
| <code>regress price io(2).rep78</code>               | set rep78 as an indicator; omit observations with rep78 == 2                          |
| <code>regress price mpg c.mpg#c.mpg</code>           | create a squared mpg term to be used in regression                                    |
| <code>regress price c.mpg##c.mpg</code>              | create all possible interactions with mpg (mpg and mpg <sup>2</sup> )                 |

## Declare Data

By declaring data type, you enable Stata to apply data munging and analysis functions specific to certain data types

### TIME SERIES

**webuse sunspot, clear**

**tsset** time, **yearly**

declare sunspot data to be yearly time series

**tsreport**

report time series aspects of a dataset

**generate** lag\_spot = L1.spot

create a new variable of annual lags of sun spots

**tsline** spot

plot time series of sunspots

**arima** spot, **ar**(1/2)

estimate an auto-regressive model with 2 lags

### TIME SERIES OPERATORS

|    |                                      |     |  |
|----|--------------------------------------|-----|--|
| L. | lag $x_{t-1}$                        | L2. | 2-period lag $x_{t-2}$   |
| F. | lead $x_{t+1}$                       | F2. | 2-period lead $x_{t+2}$  |
| D. | difference $x_t - x_{t-1}$           | D2. | difference of difference $x_t - x_{t-1} - (x_{t-1} - x_{t-2})$ |
| S. | seasonal difference $x_t - x_{t-12}$ | S2. | lag-2 (seasonal difference) $x_t - x_{t-2}$                    |

### USEFUL ADD-INS

|                     |   |
|---------------------|---|
| <b>tscollap</b>     | compact time series into means, sums and end-of-period values |
| <b>carryforward</b> | carry non-missing values forward from one obs. to the next    |
| <b>tsspell</b>      | identify spells or runs in time series                        |

### SURVIVAL ANALYSIS

**webuse drugtr, clear**

**stset** studytime, **failure**(died)

declare survey design for a dataset

**stsum**

summarize survival-time data

**stcox** drug age

estimate a cox proportional hazard model

### PANEL / LONGITUDINAL

**webuse nlswork, clear**

**xtset** id year

declare national longitudinal data to be a panel

**xtdescribe**

report panel aspects of a dataset

**xtsum** hours

summarize hours worked, decomposing standard deviation into between and within components

**xtline** ln\_wage if id <= 22, **labeled**(#3)

plot panel data as a line plot

**xtreg** ln\_w c.age##c.age ttl\_exp, **fe vce**(robust)

estimate a fixed-effects model with robust standard errors

### SURVEY DATA

**webuse rnhanes2b, clear**

**svyset** psuid [pweight = finalwgt], **strata**(stratid)

declare survey design for a dataset

**svydescribe**

report survey data details

**svy:** mean age, **over**(sex)

estimate a population mean for each subpopulation

**svy, subpop**(rural): mean age

estimate a population mean for rural areas

**svy:** tabulate sex heartatk

report two-way table with tests of independence

**svy:** reg zinc c.age##c.age female weight rural

estimate a regression using survey weights

## 1 Estimate Models

stores results as **e** -class

**regress** price mpg weight, **robust**

estimate ordinary least squares (OLS) model on mpg weight and foreign, apply robust standard errors

**regress** price mpg weight if foreign == 0, **cluster**(rep78)

regress price only on domestic cars, cluster standard errors

**rreg** price mpg weight, **genwt**(reg\_wt)

estimate robust regression to eliminate outliers

**probit** foreign turn price, **vce**(robust)

estimate probit regression with robust standard errors

**logit** foreign headroom mpg, **or**

estimate logistic regression and report odds ratios

**bootstrap, reps**(100): **regress** mpg /\*

\*/ weight gear foreign

**jackknife** r(mean), **double**: **sum** mpg

estimate regression with bootstrapping  
jackknife standard error of sample mean

### ADDITIONAL MODELS

|                  |                        |                               |
|------------------|------------------------|-------------------------------|
| <b>pca</b>       | built-in Stata command | principal components analysis |
| <b>factor</b>    |                        | factor analysis               |
| <b>poisson</b>   | <b>nbreg</b>           | count outcomes                |
| <b>tobit</b>     |                        | censored data                 |
| <b>ivregress</b> | <b>ivreg2</b>          | instrumental variables        |
| <b>diff</b>      | user-written           | difference-in-difference      |
| <b>rd</b>        | ssc install ivreg2     | regression discontinuity      |
| <b>xtabond</b>   | <b>xtabond2</b>        | dynamic panel estimator       |
| <b>psmatch2</b>  |                        | propensity score matching     |
| <b>synth</b>     |                        | synthetic control analysis    |
| <b>oaxaca</b>    |                        | Blinder-Oaxaca decomposition  |

## 2 Diagnostics

not appropriate after `robust cluster()`

**estat** **hettest**

test for heteroskedasticity

**ovtest**

test for omitted variable bias

**vif**

report variance inflation factor

**dfbeta**(length)

calculate measure of influence

**rvfplot**, **yline**(0)

plot residuals against fitted values

**avplots**

plot all partial-regression leverage plots in one graph

**display** \_b[length]

return coefficient estimate or standard error for mpg from most recent regression model

**margins, dydx**(length) returns e-class information when `post` option is used

return the estimated marginal effect for mpg

**margins, eyex**(length)

return the estimated elasticity for price

**predict** yhat if **e**(sample)

create predictions for sample on which model was fit

**predict** double resid, **residuals**

calculate residuals based on last fit model

**test** mpg = 0

test linear hypotheses that mpg estimate equals zero

**lincom** headroom - length

test linear combination of estimates (headroom = length)