# Data Analysis
## with Stata          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**          *ssc install univar*
calculate univariate summary with box-and-whiskers plot

**stem** mpg
return stem-and-leaf display of mpg

**summarize** price mpg, **detail** —— *frequently used commands are highlighted in yellow*
calculate a variety of univariate summary statistics

**ci mean** mpg price, **level**(99) —— *for Stata 13: 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** price/mpg
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$^2$ 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

**ranksum** mpg, **by**(foreign)
equality tests on unmatched data (independent samples)

**anova** systolic drug          *webuse* systolic, clear
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 https://www.stata.com/manuals/u26.pdf*

CONTINUOUS VARIABLES
measure something

CATEGORICAL VARIABLES
identify a group to which an observation belongs

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

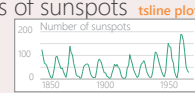| OPERATOR | DESCRIPTION | EXAMPLE | |
|---|---|---|---|
| i. | specify indicators | regress price i.rep78 | specify rep78 variable to be an indicator variable |
| ib. | specify base indicator | regress price ib(3).rep78 | set the third category of rep78 to be the base category |
| fvset | command to change base | fvset base frequent rep78 | set the base to most frequently occurring category for rep78 |
| c. | treat variable as continuous | regress price i.foreign#c.mpg i.foreign | treat mpg as a continuous variable and specify an interaction between foreign and mpg |
| o. | omit a variable or indicator | regress price io(2).rep78 | set rep78 as an indicator; omit observations with rep78 == 2 |
| # | specify interactions | regress price mpg c.mpg#c.mpg | create a squared mpg term to be used in regression |
| ## | specify factorial interactions | regress price c.mpg##c.mpg | create all possible interactions with mpg (mpg and mpg$^2$) |

## 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 sunspots    *tsline plot*

**tsline** spot
plot time series of sunspots


Number of sunspots

**arima** spot, **ar**(1/2)
fit an autoregressive 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$_{t-1}$ | S2. | lag-2 (seasonal difference) x$_t$-x$_{t-2}$ |

USEFUL ADD-INS
**tscollap**    compact time series into means, sums, and end-of-period values
**carryforward**    carry nonmissing values forward from one obs. to the next
**tsspell**    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

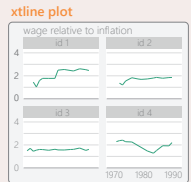**stcox** drug age
fit a Cox proportional hazards model

## 1 Fit models          *stores results as* **e**-class

**regress** price mpg weight, **vce**(robust)
fit ordinary least-squares (OLS) model on mpg, weight, and foreign, apply robust standard errors

**regress** price mpg weight **if** foreign == 0, **vce**(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
estimate regression with bootstrapping

**jackknife r**(mean): **sum** mpg
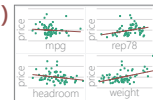jackknife standard error of sample mean

| ADDITIONAL MODELS | | |
|---|---|---|
| pca ← built-in Stata command | | principal components analysis |
| factor | | factor analysis |
| poisson • nbreg | | count outcomes |
| tobit | user-written | censored data |
| ivregress  ivreg2 | | instrumental variables |
| didregress | | difference-in-difference |
| rd  ssc install ivreg2 | | regression discontinuity |
| xtabond  xtdpdsys | | dynamic panel estimator |
| teffects psmatch | | propensity score matching |
| synth | | synthetic control analysis |
| oaxaca | | Blinder–Oaxaca decomposition |

## 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, **tlabel**(#3)
plot panel data as a line plot


xtline plot
wage relative to inflation

**xtreg** ln_w c.age##c.age ttl_exp, **fe vce**(robust)
fit a fixed-effects model with robust standard errors

## SURVEY DATA          *webuse* nhanes2b, 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

## 2 Diagnostics          *some are inappropriate with robust SEs*

**estat  hettest**    test for heteroskedasticity
**ovtest**    test for omitted-variable bias
**vif**    report variance inflation factor

**dfbeta**(length)          *Type* help regress postestimation plots *for additional diagnostic plots*
calculate measure of influence


**rvfplot, yline**(0)    plot residuals against fitted values


**avplots**    plot all partial-regression leverage plots in one graph

## 3 Postestimation          *commands that use a fitted model*

**regress** price headroom length          *Used in all postestimation examples*

**display** _b[length]          **display** _se[length]
return coefficient estimate or standard error for length from most recent regression model

**margins, dydx**(length)          *returns e-class information when post option is used*
return the estimated marginal effect for length

**margins, eyex**(length)
return the estimated elasticity for length

**predict** yhat if **e**(sample)
create predictions for sample on which model was fit

**predict** double resid, **residuals**
calculate residuals based on last fitted model

**test** headroom = 0
test linear hypotheses that headroom estimate equals zero

**lincom** headroom - length
estimate linear combination (headroom - length)