# 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**     *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)*
  **r** 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

**e** **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 $chi^2$ and Fisher's exact statistic alongside the expected values

**ttest** mpg, **by**(foreign)
  estimate t test on equality of means for mpg by foreign

**r** **prtest** foreign == 0.5
  one-sample test of proportions

**ksmirnov** mpg, **by**(foreign) **exact**
  Kolmogorov-Smirnov equality-of-distributions test

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

**anova** systolic drug     *webuse systolic, clear*
  analysis of variance and covariance

**e** **pwmean** mpg, **over**(rep78) **pveffects** **mcompare**(tukey)
  estimate pairwise comparisons of means with equal variances include multiple comparison adjustment

## Estimation with Categorical & Factor Variables

### 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 | 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 mpg²) |

---

## 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

**r** **tsreport**
  report time series aspects of a dataset

**generate** lag_spot = L1.spot
  create a new variable of annual lags of sun spots *tsline plot*

**tsline** spot
  plot time series of sunspots

**e** **arima** spot, **ar**(1/2)
  estimate an auto-regressive model with 2 lags

Number of sunspots
200
100
0
1850 1900 1950

#### 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 non-missing 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

**r** **stsum**
  summarize survival-time data

**e** **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**
  **r** report panel aspects of a dataset

**xtsum** hours
  summarize hours worked, decomposing standard deviation into between and within components     *xtline plot*

**xtline** ln_wage if id <= 22, **tlabel**(#3)
  plot panel data as a line plot

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

wage relative to inflation

### Survey Data     *webuse nhanes2b, clear*

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

**e** **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
  estimate regression with bootstrapping

**jackknife r**(mean), **double: 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** | | censored data |
| **ivregress** **ivreg2** | | instrumental variables |
| **diff** *user-written* | | difference-in-difference |
| **rd** *ssc install ivreg2* | | regression discontinuity |
| **xtabond** **xtabond2** | | dynamic panel estimator |
| **psmatch2** | | propensity score matching |
| **synth** | | synthetic control analysis |
| **oaxaca** | | Blinder-Oaxaca decomposition |

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

---

## 2   Diagnostics     *not appropriate after robust cluster( )*

**estat** **hettest**   test for heteroskedasticity
  **r** **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 mpg from most recent regression model

**margins, dydx**(length)     *returns e-class information when post option is used*
  **r** 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
  **r** test linear hypotheses that mpg estimate equals zero

**lincom** headroom - length
  test linear combination of estimates (headroom = length)

---

Tim Essam (tessam@usaid.gov) • Laura Hughes (lhughes@usaid.gov)
follow us @StataRGIS and @flaneuseks

inspired by RStudio's awesome Cheat Sheets (rstudio.com/resources/cheatsheets)

geocenter.github.io/StataTraining

updated June 2016