# Descriptives
*LPO 9951 | Fall 2015*

**PURPOSE**   Describing the data in your sample is one of the most important steps in reporting on your research. A reader that has a clear understanding of the patterns in your data will be able to readily understand your more complex analyses.

The key to descriptive statistics turns out to be the humble conditional mean: the mean of the dependent variable at various levels of the independent variable. Master the conditional mean and how to display it, and everyone will always remember your papers and presentations.

**HEADER**   Incidental to the lesson today, but important to set up correctly is the header. Notice that the plot and table files types are saved in global macros. With a quick switch at the top of the file, you can change the file format of the plots and tables that Stata saves. Very handy.

```
. global datadir "../data/"

. global plotdir "../plots/"

. global tabsdir "../tables/"

. // set plot and table types
. global gtype eps

. global ttype rtf

. // open up modified plans data
. use ${datadir}plans2, clear

. // use svyset to account for survey design
. svyset psu [pw = f1pnlwt], strat(strat_id) singleunit(scaled)

      pweight: f1pnlwt
          VCE: linearized
  Single unit: scaled
     Strata 1: strat_id
         SU 1: psu
        FPC 1: <zero>
```

## Tables

Every manuscript should include a table of descriptive statistics, listing the mean and standard error or standard deviation of every variable to be used in the dataset. In addition, tables should be used to convey crosstabs of two categorical variables. Most of your papers will also eventually include tables for regression results. Tables should be used sparingly for describing data: your best bet is almost always graphics.

For many categorical variables, however, tables may be your only option. In that case you need to think hard about two things:
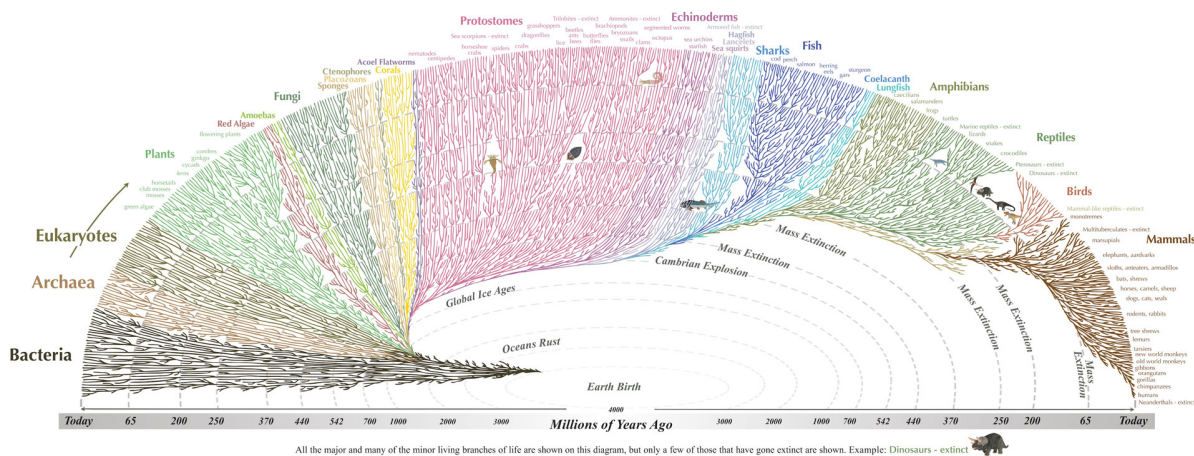
1. How can I best show patterns in the conditional mean of my dependent variable at different levels of my independent variables?
2. How can I best show relationships among key independent variables?

## Principles for displaying data

Tufte (2001) lists the following principles for describing data using graphics. He says they should:

- Show the data
- Induce the viewer to think about the substance rather than about the methodology, graphic design, the technology production, or something else.
- Avoid distorting what the data have to say.
- Present many numbers in a small space.
- Make large datasets coherent.
- Encourage the eye to compare different pieces of data.
- Reveal the data at several levels of detail, from a broad overview to fine structure.
- Serve a reasonably clear purpose: description, exploration, tabulation, or decoration.
- Be closely integrated with the statistical and verbal descriptions of a dataset.

### Tree of Life

*(credit: Leonard Eisenberg)*

## Describing variation and central tendency in continuous variables

### Plots

The two key tools for describing variation and central tendency in a continuous variable are the kernel density plot and the histogram. A histogram should be your first choice for most variables: the key decisions will be the number of bins or the frequency of the plot. Histograms can also be combined across levels using the onewayplot command.

The basic histogram is shown here.

```
. histogram bynels2m, name(hist_bynels2m) ///
>     xtitle("NELS-1992 Scale-Equated Math Score") ///
>     bin(100) /// we can try different bin widths
>     fraction
(bin=100, start=.1471, width=.006456)


. graph export ${plotdir}hist_bynels2m.$gtype, name(hist_bynels2m) replace
```

```
(file ../plots/hist_bynels2m.eps written in EPS format)
```

.

At the extreme end of the histogram is the "spike" plot, which has a single line for every level of the underlying variable.

```
. spikeplot bynels2m, name(spike_bynels2m) ///
>       xtitle("NELS-1992 Scale-Equated Math Score")


.
. graph export ${plotdir}spike_bynels2m.$gtype, name(spike_bynels2m) replace
(file ../plots/spike_bynels2m.eps written in EPS format)
```

Kernel density plots are a key tool for describing a continuous variable. The density of the variable can be compared to standard densities for visual comparison, like in the first plot below. Kernel density plots can be particularly illuminating when displayed across multiple levels of a categorical variable, as in the second plot below.

```
. kdensity bynels2m, name(kd_bynels2m) ///
>       xtitle("NELS Math Scores") ///
>       n(100) ///
>       bwidth(.025) ///
>       normal ///
>       normopts(lpattern(dash))

. graph export ${plotdir}kd_bynels2m.$gtype, name(kd_bynels2m) replace
(file ../plots/kd_bynels2m.eps written in EPS format)

. // kernel density plots of base year math score across gender
. kdensity bynels2m if bysex == 1, name(kd_bynels2m_cond) ///
>       xtitle("NELS Math Scores") ///
>       n(100) ///
>       bwidth(.025) ///
>       addplot(kdensity bynels2m if bysex == 2, ///
>               n(100) ///
>               bwidth(.025) ///
>               ) ///
>       legend(label(1 "Males") label(2 "Females")) ///
>       note("") ///
>       title("")

. graph export ${plotdir}kd_bynels2m_cond.$gtype, name(kd_bynels2m_cond) replace
(file ../plots/kd_bynels2m_cond.eps written in EPS format)
```

**QUICK EXERCISE**

> Display a histogram for byses1. Next show a kernel density plot with `byses1` for students whose father went to postsecondary education and overlay the density of students whose father did not go to postsecondary education.

**Tables**

For tables describing continuous variables, the industry standard is a table of means and standard errors or standard deviations. Below is a table of means and standard errors, nicely formatted.

```
. // get mean estimates using svy
. svy: mean bynels2m bynels2r byses1 byses2
(running mean on estimation sample)

Survey: Mean estimation

Number of strata =     361        Number of obs   =      15,512
Number of PSUs   =     751        Population size = 3,210,779
                                  Design df       =         390

--------------------------------------------------------------
             |             Linearized
             |      Mean   Std. Err.     [95% Conf. Interval]
-------------+------------------------------------------------
    bynels2m |  .4485513   .0026715      .443299     .4538036
    bynels2r |  .2947283   .0017403     .2913067     .2981498
      byses1 |  .0050858   .0145406    -.0235019     .0336736
      byses2 |  .0049615   .0143462    -.0232441     .0331671
--------------------------------------------------------------

. // store the estimates in a nice table using esttab
. esttab . using ${tabsdir}means_se.$ttype, ///    // . means all in current memory
>     not ///                              // do not include t-tests
>     replace ///                          // replace if it exists
>     nostar ///                           // no significance tests
>     label ///                            // use variable labels
>     main(b) ///                          // main = means
>     aux(se) ///                          // aux = standard errors
>     nonotes ///                          // no standard table notes
>     nonumbers ///                        // no column/model numbers
>     addnotes("Linearized estimates of standard errors in parentheses")
(output written to ../tables/means_se.rtf)
```

</td><td>        Mean</td></tr>

10th Grade Math Scores

   0.449</td></tr>

            </td><td>   (0.00267)</td></tr>

10th Grade Reading Scores

   0.295</td></tr>

            </td><td>   (0.00174)</td></tr>

4

SES v1

 0.00509</td></tr>

                </td><td>     (0.0145)</td></tr>

SES v2

 0.00496</td></tr>

                </td><td>     (0.0143)</td></tr>

Observations

    15512</td></tr>

Linearized estimates of standard errors in parentheses

You can also use standard deviations, which even in survey data like yours are based on the simple sample standard deviation calculation. Table 2 shows means and standard deviations.

```
. tabstat bynels2m bynels2r byses1 byses2, stat(sd) save

    stats |  bynels2m  bynels2r    byses1    byses2
---------+--------------------------------------
       sd |  .1353664  .0939987  .7429628  .7502604
--------------------------------------------------

. // grab matrix of sds and store in matrix
. mat mysd = r(StatTotal)

. // add to earlier results using estadd
. estadd matrix mysd

added matrix:
              e(mysd) :  1 x 4

. // save new table, this time with sds instead of ses
. esttab . using ${tabsdir}means_sd.$ttype, ///
>     not ///
>     replace ///
>     nostar ///
>     label ///
>     main(b) ///
>     aux(mysd) ///                       // NOTE: aux = standard deviations
>     nonumbers ///
>     nonotes ///
>     addnotes("Standard deviations in parentheses")
(output written to ../tables/means_sd.rtf)
```

</td><td>          Mean</td></tr>

10th Grade Math Scores

    0.449</td></tr>

                                </td><td>      (0.135)</td></tr>

10th Grade Reading Scores

    0.295</td></tr>

                                </td><td>      (0.0940)</td></tr>

SES v1

 0.00509</td></tr>

                                </td><td>      (0.743)</td></tr>

SES v2

 0.00496</td></tr>

                                </td><td>      (0.750)</td></tr>

Observations

    15512</td></tr>

Standard deviations in parentheses

## QUICK EXERCISE

Create a table that includes all of the above variables but also sex and race. Choose whether to include standard errors or standard deviations.

## Describing categorical or binary variables

### Plots

The first question is: do you need to? What can you present in a graphic that wouldn't be in a table or couldn't be described in the text? That said, some options include a histogram for categorical variables (bar chart) and a table of proportions.

Below is code for a barchart.

```
. histogram bystexp, name(bar_bystexp) ///
>     percent ///
>     addlabels ///
>     xlabel(-1 1 2 3 4 5 6 7, ///
>           value ///
>           angle(45) ///
>           labsize(vsmall) ///
>           ) ///
>     addlabopts(yvarformat(%4.1f)) ///
>     xtitle("")
(bin=41, start=-1, width=.19512195)

. graph export ${plotdir}bar_bystexp.$gtype, name(bar_bystexp) replace
(file ../plots/bar_bystexp.eps written in EPS format)
```

### Tables

A table of proportions is a way to display information regarding a single categorical variable. Here is a simple table of proportions.

```
. estpost svy: tabulate bystexp
(running tabulate on estimation sample)

Number of strata   =        361        Number of obs      =     15,236
Number of PSUs     =        751        Population size    =  3,210,779
                                       Design df          =        390


----------------------
how far    |
in school  |
student    |
thinks     |
will       |
get-compo  |
site       | proportion
-----------+-----------
 Don't Kn  |      .0967
 Less tha  |      .0089
       HS  |      .0688
     2 yr  |      .0657
  4 yr No  |      .0369
 Bachelor  |       .361
  Masters  |      .2002
```

7

```
  Advanced |        .1619
           |
     Total |           1
----------------------
  Key:  proportion  =  cell proportion

saved vectors:
            e(b) =  cell proportions
           e(se) =  standard errors of cell proportions
           e(lb) =  lower 95% confidence bounds for cell proportions
           e(ub) =  upper 95% confidence bounds for cell proportions
         e(deff) =  deff for variances of cell proportions
         e(deft) =  deft for variances of cell proportions
         e(cell) =  cell proportions
        e(count) =  weighted counts
          e(obs) =  number of observations

. esttab . using ${tabsdir}proportions.$ttype, ///
>     not ///
>     replace ///
>     nostar ///
>     label ///
>     main(b) ///
>     aux(se) ///
>     nonotes ///
>     nonumbers ///
>     addnotes("Linearized standard errors in parentheses")
(output written to ../tables/proportions.rtf)
```

                 </td><td>        Mean</td></tr>

Don't Know

  0.0967</td></tr>

                 </td><td>    (0.00311)</td></tr>


Less than HS

 0.00887</td></tr>

                 </td><td>    (0.00104)</td></tr>


HS

  0.0688</td></tr>

                 </td><td>    (0.00309)</td></tr>

2 yr

  0.0657</td></tr>

                </td><td>    (0.00332)</td></tr>


4 yr No Deg

  0.0369</td></tr>

                </td><td>    (0.00195)</td></tr>


Bachelors

   0.361</td></tr>

                </td><td>    (0.00507)</td></tr>


Masters

   0.200</td></tr>

                </td><td>    (0.00448)</td></tr>


Advanced

   0.162</td></tr>

                </td><td>    (0.00412)</td></tr>


Total

      1</td></tr>

                </td><td>          (0)</td></tr>

Observations

   15236</td></tr>

Linearized standard errors in parentheses

**QUICK EXERCISE**

> Create a nicely formatted bar chart and table showing the proportions of students with different levels of parental education.

## Describing relationships between two continuous variables

**Plots**

The scatterplot is the ultimate tool in describing relationships between two continuous variables. All scatterplots should always have the dependent variable on the $y$ axis and the independent variable on the $x$ axis. Both axes must be labeled clearly. Limits should be set based on the data: there's no need to include a 0 in every graph. This can in fact be counterproductive.

The code below produces a simple scatterplot of math and reading scores.

```
. graph twoway scatter bynels2m bynels2r, name(sc_math_read)

. graph export ${plotdir}sc_math_read.$gtype, name(sc_math_read) replace
(file ../plots/sc_math_read.eps written in EPS format)
```

Notice how there are too many dots; just blob reall. You may find it useful to plot only a random (representative) subsample of your data. Here is a better version of the above scatterplot that only uses 10% of the data.

```
. preserve                              // preserve data

. sample 10                             // sample random 10%
(14,544 observations deleted)

. graph twoway scatter bynels2m bynels2r, name(sc_math_read_10) ///
>   ytitle("NELS Math Scores") ///
>   xtitle("NELS Reading Scores") ///
>   msize(tiny)

. graph export ${plotdir}sc_math_read_10.$gtype, name(sc_math_read_10) replace
(file ../plots/sc_math_read_10.eps written in EPS format)

. restore                               // restore data
```

Finally, you can condition on another variable in order to add another level to your scatter plot. This can be done both with use of `if` statements, as in the first plot below, and `by` statements, as in the second plot.

```
. preserve                              // preserve data

. sample 25                             // sample random 25%
(12,120 observations deleted)

. graph twoway (scatter bynels2m byses1 if urm == 0, ///
>               mcolor("orange") ///
>               msize(vtiny) ///
>          ) ///
```

```
>                || scatter bynels2m byses1 if urm == 1, ///
>                   mcolor("green") ///
>                   msize(vtiny) ///
>                   msymbol(triangle) ///
>                   ytitle("NELS Math Scores") ///
>                   xtitle("SES") ///
>                   legend(order(1 "Non-Minority" 2 "Underrep Minority")) ///
>                   name(sc_complex)

. graph export ${plotdir}sc_complex.$gtype, name(sc_complex) replace
(file ../plots/sc_complex.eps written in EPS format)


. restore                                    // restore data


. graph twoway scatter bynels2m byses1, by(bystexp) ///
>     msize(*.01) ///
>     ytitle("NELS Math Scores") ///
>     xtitle("SES") ///
>     note("") ///
>     name(sc_cond)

. graph export ${plotdir}sc_cond.$gtype, name(sc_cond) replace
(file ../plots/sc_cond.eps written in EPS format)
```

You can also run a scatter plot across levels of a categorical variable if you suspect the underlying relationship may not be the same in each level of the categorical variable. The `matrix` plot helpfully with plot each combination of included varibles against each other to produe a sort of "small multiples" correlation plot.

```
. graph matrix bynels2m bynels2r byses1 byses2, name(matrix_plot) msize(vtiny)

. graph export ${plotdir}matrix_plot.$gtype, name(matrix_plot) replace
(file ../plots/matrix_plot.eps written in EPS format)
```

## Describing relationships between a categorical and a continuous variable

### Plots

There are multiple options for plotting the relationship between a categorical and a continuous variable. A particularly useful option is to plot the continuous variable as a series of boxplots, one for each level of the categorical variable.

**Boxplots**   For boxplots to be effective, they should be sorted by the median of the dependent variable. This contrast is shown in the two figures below.

```
. graph box bynels2m, over(byrace2, ///
>                     label(alternate ///
>                           labsize(tiny) ///
>                          ) ///
>                        ) ///
>                 name(box1)
```

```
. graph export ${plotdir}box1.$gtype, name(box1) replace
(file ../plots/box1.eps written in EPS format)

. graph box bynels2m, over(byrace2, ///
>                        label(alternate ///
>                                labsize(tiny) ///
>                                ) ///
>                        sort(1) ///
>                        ) ///
>                name(box2)

. graph export ${plotdir}box2.$gtype, name(box2) replace
(file ../plots/box2.eps written in EPS format)
```

**Dot plots**   Dot plots can also be useful for plotting the measure of central tendency across groups. In this case, we'll produce two plots, one each for reading and math scores, and then combine them into a single graphic.

```
. graph dot bynels2m, over(bypared, ///
>                        label(alternate ///
>                                labsize(tiny) ///
>                                ) ///
>                        ) ///
>                ytick(0(.10).80) ///
>                ylabel(0(.1).8) ///
>                ytitle("Math Scores") ///
>                marker(1, msymbol(O)) ///
>                name(dot_math)

. graph save ${plotdir}dot_math.gph, replace
(file ../plots/dot_math.gph saved)

. graph dot bynels2r, over(bypared, ///
>                        label(alternate ///
>                                labsize(tiny) ///
>                                ) ///
>                        ) ///
>                ytick(0(.10).80) ///
>                ylabel(0(.1).8) ///
>                ytitle("Reading Scores") ///
>                marker(1, msymbol(O)) ///
>                name(dot_read)

. graph save ${plotdir}dot_read.gph, replace
(file ../plots/dot_read.gph saved)

. // combine graphs into one plot
. graph combine ${plotdir}dot_math.gph ${plotdir}dot_read.gph, ///
>     name(dot_both) ///
>     colfirst

. // export combined graphics
. graph export ${plotdir}dot_both.$gtype, name(dot_both) replace
(file ../plots/dot_both.eps written in EPS format)
```

## Describing relationships between two categorical variables

**Plots**

The basic tool for comparing two categorical variables is the crosstabulation. In a crosstabulation we take a look at counts of the sample that are identified by their presence in cells created by the two categorical variables. There are several tools for plotting categorical variables, including tabplots, jittered plots, and heatmaps.

**Tabplots**   Below are examples of a tabplot, with both two and three dimensions.

```
. // first new recoded parental education level
. recode bypared (1/2 = 1) (3/5 = 2) (6 = 3) (7/8 = 4) (. = .), gen(newpared)
(14362 differences between bypared and newpared)

. label var newpared "Parental Education"

. label define newpared 1 "HS or Less" 2 "Less than 4yr" 3 "4 yr" 4 "Advanced"

. label values newpared newpared

. // next new recoded student expectations
. recode f1psepln (1/2 = 1) (3/4 = 2) (5 = 3) (6 = .) (. = .), gen(newpln)
(13995 differences between f1psepln and newpln)

. label var newpln "PS Plans"

. label define newpln 1 "No plans" 2 "VoTech/CC" 3 "4 yr"

. label values newpln newpln

. // tabplot of parental education against student plans
. tabplot newpared newpln, name(tabplot1) ///
>     percent(newpared) ///
>     showval ///
>     subtitle("")

. graph export ${plotdir}tabplot1.$gtype, name(tabplot1) replace
(file ../plots/tabplot1.eps written in EPS format)

. // tabplot again, but with new dimension
. tabplot newpared newpln, by(bysex) ///
>     percent(newpared) ///
>     showval ///
>     subtitle("") ///
>     name(tabplot2)

. graph export ${plotdir}tabplot2.$gtype, name(tabplot2) replace
(file ../plots/tabplot2.eps written in EPS format)
```

**Jitter plot**

```
. graph twoway scatter f1psepln bypared, name(jitterplot) ///
```

```
>       jitter(5) ///
>       msize(vtiny)

. graph export ${plotdir}jitterplot.$gtype, name(jitterplot) replace
(file ../plots/jitterplot.eps written in EPS format)
```

**Heatmap**

```
. tddens bypared f1psepln, title("") ///
>       xtitle("Parent's Level of Education") ///
>       ytitle("PS PLans")

.
. graph export ${plotdir}heatmap.$gtype, replace
(file ../plots/heatmap.eps written in EPS format)
```

**Tables**

When checking crosstabulations, we can produce two-way tables that include survey weights in the command itself.

```
. table byrace2 f1psepln [pw = bystuwt], by(bysex) contents(freq) row

--------------------------------------------------------------------------------
sex-composi |
te and      |
RECODE of   |
byrace      |
(student^s  |
race/ethnic |
ity-composi |          f1 post-secondary plans right after high school
te)         | don^t plan to contin  don^t know or planni  vocational, technica
------------+-------------------------------------------------------------------
male        |
    Am.Ind. |          959.9478            2,165.51             1,850.23
   Asian/PI |          955.4428            3,010.54             4,576.58
      Black |          3,280.14           13,926.5             22,419.1
   Hispanic |          2,481.69           23,987.1             25,572.5
Multiracial |          2,885.66           7,204.85             6,076.28
      White |            27,037             74,634               91,041
            |
      Total |          37,599.8           124,929              151,536
------------+-------------------------------------------------------------------
female      |
    Am.Ind. |                               961.4319            1,422.26
   Asian/PI |          182.5247            1,706.65             1,229.69
      Black |          1,682.59           9,165.85             12,928
   Hispanic |          1,098.54           12,487.4             16,983.4
Multiracial |          916.7742           2,516.75             2,778.86
      White |          5,139.84           37,381.5             44,545.3
            |
      Total |          9,020.27           64,219.6             79,887.6
```

```
--------------------------------------------------------------------------------


--------------------------------------------------------------------------------
sex-composi |
te and      |
RECODE of   |
byrace      |
(student^s  |
race/ethnic |
ity-composi |           f1 post-secondary plans right after high school
te)         | two-year community c  four-year college or  early hs grad attend
------------+-------------------------------------------------------------------
male        |
   Am.Ind.  |            1,168.03              8,487.8              509.1148
  Asian/PI  |            13,906.7              41,635.3             1,415.14
     Black  |            33,505.1              121,834              8,161.93
  Hispanic  |            68,736.3              82,651.8             7,870
Multiracial |            12,485.4              31,158.9             1,156.67
     White  |            172,672               531,000              18,210.1
            |
     Total  |            302,473               816,768              37,323
------------+-------------------------------------------------------------------
female      |
   Am.Ind.  |            2,538.59              6,913.7              198.1191
  Asian/PI  |            9,878.74              47,395.6             1,796.86
     Black  |            48,116.9              124,183              5,876.88
  Hispanic  |            76,304.8              114,582              7,564.25
Multiracial |            14,070.6              36,070.2             3,290.11
     White  |            195,002               606,845              20,989.9
            |
     Total  |            345,912               935,989              39,716.1
--------------------------------------------------------------------------------
```

Of course, if we want to use a table in a paper, we should use `esttab`.

```
. estpost svy: tabulate byrace2 newpln, row percent se
(running tabulate on estimation sample)

Number of strata    =         361        Number of obs     =       13,055
Number of PSUs      =         750        Population size    =    2,908,622
                                         Design df          =          389


----------------------------------------------------
RECODE of  |
byrace     |
(student^  |
s          |
race/ethn  |
icity-com  |                  PS Plans
posite)    | No plans  VoTech/C      4 yr      Total
-----------+----------------------------------------
  Am,Ind,  |    16.17     26.93     56.89        100
           |  (3.647)   (4.454)   (5.413)
           |
```

```
  Asian/PI |    4.746     23.67     71.59        100
           |  (.8931)   (1.933)   (2.107)
           |
     Black |     7.31      29.9     62.79        100
           |  (.7736)   (1.368)   (1.487)
           |
  Hispanic |    9.732     43.95     46.32        100
           |  (.8412)   (1.494)   (1.699)
           |
  Multirac |    12.27     30.13      57.6        100
           |   (1.72)   (2.436)   (2.591)
           |
     White |    8.147     28.26     63.59        100
           |  (.3851)   (.8152)   (.9036)
           |
     Total |    8.369      30.6     61.04        100
           |  (.3292)   (.6266)   (.7225)
-----------------------------------------------------
  Key:  row percentage
        (linearized standard error of row percentage)

  Pearson:
    Uncorrected   chi2(10)          =   253.9129
    Design-based  F(9.11, 3544.12)=   17.6793     P = 0.0000

Note: Variance scaled to handle strata with a single sampling unit.

saved vectors:
            e(b) =  row percentages
           e(se) =  standard errors of row percentages
           e(lb) =  lower 95% confidence bounds for row percentages
           e(ub) =  upper 95% confidence bounds for row percentages
         e(deff) =  deff for variances of row percentages
         e(deft) =  deft for variances of row percentages
         e(cell) =  cell percentages
          e(row) =  row percentages
          e(col) =  column percentages
        e(count) =  weighted counts
          e(obs) =  number of observations

row labels saved in macro e(labels)

. eststo racetab

. esttab racetab using ${tabsdir}race_tab.$ttype, ///
>     replace ///
>     nostar ///
>     nostar ///
>     unstack ///
>     nonotes ///
>     varlabels(`e(labels)') ///
>     eqlabels(`e(eqlabels)')
(output written to ../tables/race_tab.rtf)
```

| | (1) | | | |
| --- | --- | --- | --- | --- |
| | Ratio | | | |
| | No plans | VoTech/CC | 4 yr | Total |
| Am.Ind. | 16.17 | 26.93 | 56.89 | 100 |
| | (4.44) | (6.05) | (10.51) | (.) |
| Asian/PI | 4.746 | 23.67 | 71.59 | 100 |
| | (5.31) | (12.24) | (33.98) | (.) |
| Black | 7.310 | 29.90 | 62.79 | 100 |
| | (9.45) | (21.86) | (42.22) | (.) |
| Hispanic | 9.732 | 43.95 | 46.32 | 100 |
| | (11.57) | (29.42) | (27.26) | (.) |
| Multiracial | 12.27 | 30.13 | 57.60 | 100 |
| | (7.14) | (12.37) | (22.23) | (.) |
| White | 8.147 | 28.26 | 63.59 | 100 |
| | (21.16) | (34.67) | (70.37) | (.) |
| Total | 8.369 | 30.60 | 61.04 | 100 |
| | (25.42) | (48.83) | (84.48) | (.) |
| N | 13055 | | | |

*Init: 28 August 2015; Updated: 28 August 2015*