

Lab 1 (100 points): Introduction to Statistical Packages

Objectives: Loading files, cleaning and manipulating the data.

A. (10 pts.) Online Prelab

B. (90 points) US Demographic, Crime, and Test-Score Data. This semester, we are going to be exploring some Demographic, Crime, and Test-Score data for counties across the United States. The data we will analyze are in the data set "USData.txt". The variable names and definitions are listed in the file "US_Data set_Definition.pdf". In this lab, we are going to explore what is included in the data set, load it into the software package, and do some basic manipulations.

1. (10 points) How many variables does this data set contain? Which are [categorical or qualitative variables](#) and which are [quantitative or numeric variables](#)? Besides looking at the documentation file provided, you might want to look at the data file itself in a spreadsheet, notepad or the software package (R only).

Solution:

This data set contains 20 variables.

- Categorical. There are 5: State, Region, CountyIndex, UrbanIndicator, and IncomeCategory.
 - Quantitative. There are 15: Population, LandArea, PopulationDensity, PercentMaleDivorce, PercentFemaleDivorce, MedianIncome, PercentCollegeGraduates, MedianHouseAge, RobberiesPerPopulation, AssaultsPerPopulation, BurglariesPerPopulation, LarceniesPerPopulation, EducationSpending, EducationSpendingP2, and TestScore
2. (16 pts.) Write two analysis questions that can be answered from the data provided. In the project due at the end of the semester, your group will have to pose general questions that can be answered by three different statistical methods. You will be allowed to change the questions when you start the project, but this will get you thinking of possibilities.

Solution:

There are many possible answers, here are three examples:

- Does the number of robberies per 100,000 people differ significantly among the different regions of the US? Similar questions can be asked for the other crime variables.
- Is more Education Spending associated with higher Average Test Scores?
- Is a larger Percentage of College Graduates positively associated with Median Income?

3. (20 points) Load the data into your software package, and provide the programming code used to do so. If you used menu options to load the data, rather than code, please describe the procedure you followed. No output is required.

Solution:

Method 1: GUI Interface

File → Import Data → Tab Delimited File (.txt) → Next → browse for USData → Library: Work, Member: USData (or an appropriate name) → Finish

Method 2: Command line

```
data USData;
  infile "W:\USData.txt" delimiter = '09'x firstobs = 2 ;
  length IncomeCategory $11.; ** needed because the income category in the
                               first row is shorter than for other rows;
  input State $ Region $ CountyIndex $ UrbanIndicator $Population
        LandArea PopulationDensity PercentMaleDivorce
        PercentFemaleDivorce MedianIncome IncomeCategory
        PercentCollegeGraduates MedianHouseAge
        RobberiesPerPopulation AssaultsPerPopulation
        BurglariesPerPopulation LarceniesPerPopulation
        EducationSpending EducationSpendingP2 TestScore;
run;
```

Be sure that the number of variables matches what is in Part 1. The input statement will be the same for the rest of the semester. If you get this part wrong in Lab 1, please look at the key to be sure that you use the correct statements.

4. (19 points) Are there missing values (NA) in the data set? If so, please create a new data set by removing any rows that contain one or more NAs from the original data set. Please save this new data set to your computer and/or ITaP folder; this will be the data set that you will be using for the rest of the semester.

- a. (5 pts.) Code

Solution:

```
* Remove missing values;
data USData_cleaned;
  set USData;
  if nmiss( of _NUMERIC_ ) = 0 AND State ^= "NA" AND Region ^= "NA"
    AND CountyIndex ^= "NA" AND UrbanIndicator ^= "NA"
    AND IncomeCategory ^= "NA";
  /* only outputs the data if the numeric data is not text,
    and categorical variables are not NA) */
run;
```

To save the data, I followed the tutorial and used the following menu and wizard options:

File → Export Data → Select USData_cleaned in the Member drop-down box → Check “Write variable labels as column names” → Next → Selected Tab-Delimiter (*.txt) in the Data Source Drop Down Box → Saved as "W:\Labs\USData_cleaned.txt"

- b. (9 pts.) We want to know how many rows were removed, so please answer:

Solution:

- i. How many observations are there in the original data set? (The output is all that is required.)

Solution:

The NOTE can be found in “LOG” page.

NOTE: There were 1103 observations read from the data set WORK.USDATA.

There were 1,103 observations in the original data set.

- ii. How many observations are there after removing the incomplete data? (The output is all that is required.)

Solution:

NOTE: The data set WORK.USDATA_CLEANED has 1098 observations and 20 variables.
There are 1,098 observations after removing the incomplete data.

- iii. How many rows were removed (show the work, even though it is a quick calculation)?

Solution:

$$1103 - 1098 = 5$$

- c. (5 pt.) In which directory did you save your cleaned data set?

Solution:

W:\

It doesn't matter what this answer here. I just asked the question so that you can refer to it for the rest of the semester.

5. (10 points) For readability, we want to transform the values of "UrbanIndicator" from a number to what the number represents. That is, please create a new variable called "UrbanNew" such that:

If UrbanIndicator is "1", UrbanNew is "Urban" and

If UrbanIndicator is "0", UrbanNew is "Rural"

- a. (5 pts.) Code. Remember that all code needed to answer part b) needs to be included in this part.

Solution:

```
data USData_cleaned;
    set USData_cleaned;
    length UrbanNew $5;
    if UrbanIndicator = "1" then UrbanNew = "Urban";
    if UrbanIndicator = "0" then UrbanNew = "Rural";
run;

* Print ;
data printme;
    set USData_cleaned;
    if _n_ in (5, 55, 355, 555);
run;

proc print data = printme;
    var State CountyIndex UrbanIndicator UrbanNew;
run;
```

- b. (5 pts.) Print or display the data set (on the computer, not to physical paper), and take screen clippings which demonstrate the following rows: 5, 55, 355, and 555. Please highlight or somehow indicate the changes. These rows will prove that your code worked correctly. To save space, you are permitted to restrict the data set to show only the relevant columns and the columns for "State" and "CountyIndex."

Solution:

The SAS System				
Obs	State	CountyIndex	UrbanIndicator	UrbanNew
1	Alabama	5	0	Rural
2	Califor	11	1	Urban
3	Iowa	9	0	Rural
4	Missour	7	1	Urban

6. (15 points) We are going to show that "PopulationDensity" can be calculated from other variables in the data set.
- a. (5 pts.) Write down the equation relating "PopulationDensity" to "Population" and "LandArea."

Solution:

$$\text{PopulationDensity} = \frac{\text{Population}}{\text{LandArea}}$$

- b. (5 pts.) Write code (and provide it here) to create a new variable called "PopulationDensityNew" which implements the calculation described in part a). Remember that all code needed to answer part c) needs "to be included in this part."

Solution:

```
data USData_cleaned;
  set USData_cleaned;
  PopulationDensityNew = Population / LandArea;
run;
proc print data = USData_cleaned (obs = 6);
  var PopulationDensity PopulationDensityNew;
run;
```

- c. (5 pts.) Show that your code is correct by displaying the original variable "PopulationDensity" and "PopulationDensityNew". Please only print out the first 6 rows. To save space, you are permitted to restrict the data set to show only the relevant columns.

Solution:

The SAS System		
Obs	PopulationDensity	PopulationDensityNew
1	1650.83	1650.83
2	1273.83	1273.83
3	1199.62	1199.62
4	626.88	626.88
5	532.69	532.69
6	3485.71	3485.71

Our calculation matches the true value.