



I'm not robot



Continue

## Descriptives package in r

Packages used in this chapter include: • psych • DescTools • Rmisc • FSA • plyr • boot The following commands will install these packages if they are not yet installed: `if(!require(psych)){install.packages(psych)} if(!require(DescTools)){install.packages(DescTools)} if(!require(Rmisc)){install.packages(Rmisc)} if(!require(FSA){install.packages(FSA)} if(!require(plyr)){install.packages(plyr)} if(!require(boot)){install.packages(boot)}` Descriptive statistics Descriptive statistics are used to summarize data in a way that provides insight into the information contained in the data. This may include a study of the mean or median of the numerical data or the frequency of observation of the nominal data. You can create charts that show data and indicate summary statistics. The choice of summary statistics depends on the type of variable being tested. Different statistics should be used for interval/ratio, ordinal, and nominal data. When describing or examining data, you will typically deal with measures of location, variability, and shape. The location is also called the central trend. This is a measure of the value of the data. For example, are values close to 10 or 100 or 1000? Location measures include average and median, as well as slightly more exotic statistics such as M-estimators or winsorized means. The variety is also called dispersion. This is a measure of how far data points lie apart. Common statistics include standard deviation and coefficient of variability. For data that is not normally distributed, the percentile or inter quartile range can be used. A shape refers to a distribution of values. Histograms and related charts are the best tools for evaluating the shape of your data. Statistics include skewness and kurtosis, although they are less useful than visual inspection. We can describe the shape of the data as usually distributed, log-normal, uniform, oblique, bimodal, and others. Descriptive statistics for intervals/data ratio In this example, imagine that Ren and Stimpmy each held eight workshops educating the public about saving water at home. They are interested in how many people showed up at the workshops. Because the data is kept in a data frame, we can use the `Data$Participants` convention to access the `Participants` variable in the `Data` frame. `Entry = (Instructor Location Participants Ren North 7 Ren North 22 Ren North 6 Ren North 15 Ren South 12 Ren South 13 Ren South 14 Ren South 16 Stimpmy North 18 Stimpmy North 17 Stimpmy North 15 Stimpmy North 9 Stimpmy 15 Stimpmy South 11 Stimpmy South 19 Stimpmy South 23 )` `Data = read.table(textConnection(Input),header=TRUE) Date ###` This will cause that a data frame named `Data` (`Data`) `###` will be discussed later, `###` but shows the structure of the data frame `summary(Data) ###` Will be discussed later, `###` but summarizes the variables in the data frame The sum of functions and the length of the sum of the variable can be found using the `sum` function, and the number of observations can be found using the `length` function. `sum(Data$Participants) 232 length(Data$Participants) 16` Location statistics for interval/ratio data The mean is the arithmetic mean and is a common statistic used with interval/ratio data. This is simply the sum of the values divided by the number of values. The average function in R will return the average. `sum(Data$Attendees) / Length(Data$Attendees) average 14.5 (Data$Attendees) 14.5` Caution should be used when reporting average values with skewed data because the average may not be representative of the data center. Imagine, for example, a city with 10 families, nine of whom have an income of less than \$50,000 a year, but with one family with an income of \$2,000,000 a year. The average income for families in the city will be \$233,000, but this may not be a reasonable way to summarize the city's income. `income = c(49000, 44000, 25000, 18000, 32000, 47000, 37000, 45000, 36000, 200000) Average(Income) 233300` Median is defined as the value below, which represents 50% of observations. To find this value manually, order observations and separate the lowest 50% from the highest 50%. For datasets with odd number of observations, the median drops halfway between the two middle values. The median is a solid statistic because it is not affected by the addition of extreme values. For example, if we changed the value of the last Stimpmy participants from 23 to 1000, the median would not be affected. `median(Data$Attendees) 15 ###` Note that in this case the mean and median are close to the `###` value relative to each other. The average and median will be more different `###` the more data is warped. The median is appropriate for skewed or unseparated data. The average income for the city discussed above is \$40,500. Half of the families in the city have an income above this amount, and half have an income below that amount. `Income = c(49000, 44000, 25000, 18000, 32000, 47000, 37000, 45000, 36000, 200000) median(Income) 40500` It should be noted that median are sometimes reported as the average person or typical family. Saying: The average American family earned \$54,000 last year means that the average income for families The average family is the one with the median income. Mode mode is a summary statistic that is rarely used in practice, but is usually included in any average and median discussion. When there are discrete values for a variable, mode is simply the value that occurs most often. For example, in the Video Science Center statistics in the required readings below, Dr. Nic gives an example of counting the number of pairs of shoes each student has. The most common response was 10, so 10 is the mode for this dataset. For our example, The Rhine and Stimpmy value of 15 occurs three times and so is the mode. The Mode feature can be found in the DescTools package. `library(DescTools) Mode(Data$Attendees) 15` Variability statistics for interval/ratio data Standard deviation Standard deviation is a measure of variability that is commonly used with interval/ratio data. This is a measurement of how close the observations in the dataset are to the average. A useful rule is that, for normally distributed data, 68% of the data points are within the average deviation of  $\pm 1$  of the standard, 95% of the data points are in the average  $\pm 2$  standard deviations, and 99.7% of the data points are within the average  $\pm 3$  standard deviations. Since the average is often represented with the letter  $\mu$ , and the standard deviation is represented with the letter  $\sigma$ , saying that someone is a few sigmas away from  $\mu$  indicates that they are rather rare in nature. (Initially, I heard this joke on an episode of Car Talk for which I can't find a reference or transcript.) `sd(Data$Attendees) 4.830459` Standard deviation may not be suitable for skewed data. The standard error of the average standard average error is a measure that estimates how close the calculated average can be to the actual average of that population. It is commonly used in tables or charts where multiple measures are presented together. For example, we might want to present the average participants for the Rhine with a standard error for that average and the average stimpmy participants with the standard error that it means. A standard error is the standard deviation of a dataset divided by the square root of the number of observations. It can also be found in the output for the description function in the psych package, marked `se.sd(Data$Attendees) / sqrt(length(Data$Attendees)) 1.207615` `library(psych) describe(Data$Attendees) vars n mean sd median trimmed mad min max range skew kurtosis se1 1 16 14.5 4.83 15 14.5 4.45 6 23 17 -0.04 -0.88 1.21 ###` `se` indicates that the standard error of the average standard average error may not be appropriate for skewed data. A summary of five numbers, quartiles, median percentile is the same as the 50th percentile, because 50% of the value falls below this value. Other for the dataset can be identified to provide more information. Typically, 0, 25, 50, 75 and 100 percentiles are reported. This is sometimes referred to as a summary of five numbers. These values can also be called minimum, 1 quartile, second quartile, third quartile, and maximum. A summary of five numbers is a useful measure of variability for skewed interval/ratio or ordinal data. 25% of the value falls below the 1st quartile and 25% of the value exceeds the third quartile. This leaves the middle 50% of the value between 1 and 3 quartiles, giving a sense of the mid-range of the data. This range is called an inter quartile range (IQR). Percentiles and quartiles are relatively solid because they are not much affected by a few extreme values. They are suitable for both oblique and unenforced data. `summary(Data$Participants) Min. 1st Qu. Median Average 3 Qu. Max. 6.00 11.75 15.00 14.50 17.25 23.00 ###` Summary of five numbers and average optional technical note regarding the calculation of percentiles May occur as strange that 3. quartile for participants was reported as 17.25. After all, if you want to order the values of participants, the 75th percentile will fall between 17 and 18. But why does the R go from 5.25pm rather than 5.5pm? `sort(Data$Participants) 6 7 9 11 12 13 14 15 15 15 16 17 18 19 22 23` The answer is that there are several different methods for calculating the percentiles and they may give slightly different answers. Detailed information on the calculation can be found under ?quantiles. For participants, the default type 7 calculation is 75 percentile 17.25, while type 2 calculation simply divides the difference between 17 and 18 and yields 17.5. The type 1 calculation does not average two values, so it returns only 17. `quantile(Data$Attendees, 0.75, type=7) 75% 17.25` `quantil(Data$Attendees, 0.75, type=2) 75% 17.5` `quantile(Data$Attendees, 0.75, type=1) 75% 17` Percentiles other than 25, 50, and 75 can be calculated using the quantum function. For example, to calculate the 95th percentile: `quantile(Data$Attendees, .95) 95% 22.25` Confidence intervals are discussed in the next chapter. Statistics for grouped interval/ratio data In many cases, you will want to examine summary statistics for a variable in groups. For example, we may want to examine statistics for workshops conducted by the Rhine and those that run Stimpmy. Summarize in FSA Sum function in FSA returns the number of observations, average, standard deviation, minimum, 1 quartile, median, 3 quartile, and maximum for grouped data. Note the use of formula notation: Participants are a dependent variable (the variable for which you want to get statistics), and Instructor is an independent variable (grouping variable). The summary allows you to summarize the combination of multiple independent variables by listing them on the right side - separated by a plus sign (+). `library(FSA) summarize(Participants ~ Instructor, data=Data) Instructor n valid mean sd min Q1 median Q3 max percZero 1 Ren 8 8 13.125 5.083236 6 10 75 13.5 5 15 25 22 02 Stimpmy 8 8 15.875 4.454131 9 14 00 16 0 18 25 23 0 Summary (Participants ~ Instructor + Location, data=Data) Instructor location n valid mean sd min Q1 median Q3 max percZero 1 Ren North 4 4 12.50 7.505554 6 6 7 5 11 0 16 75 22 02 Stimpmy North 4 4 14.75 4.031129 9 13 50 16 0 17 25 18 03 Ren South 4 4 13.75 1.707825 12 12 75 13 5 14 50 16 04 Stimp South 4 4 4 17 00 5.163978 11 14 00 17 0 20 00 23 0 summarySE in Rmisc SummarySE function in the Rmisc package derives the number of observations, average, standard deviation, standard average error and confidence interval for grouped data. SummarySE summarizes a combination of multiple independent variables by listing them as a vector, such c(Instructor, Student), library(Rmisc) summarySE(data=Data, Attendees, groupvars=Instructor, conf.interval=0.95) Instructor N Attendees sd se ci1 Ren 8 13.125 5.083236 1.797195 4.2496912 Stimpmy 8 15.875 4.454131 1.574773 3.723747 summarySE(data=Data, Attendees, groupvars = c(Instructor, Location), conf.interval = 0.95) Instructor Location N Attendees sd se ci1 Ren North 4 12.50 7.505553 3.7527767 11.9430112 Ren South 4 13.75 1.707825 0.8539126 2.7175313 Stimpmy North 4 14.75 4.031129 2.0155644 6.414426 4 Stimpmy South 4 17.00 5.163978 2.5819889 8.217041 describeBy The function in the psych package returns the number of observations, average, median, cropped measures, minimum, maximum, range, tilt, kurtosis, and standard average error for grouped data. describeBy summarizes the combination of multiple independent variables by combining terms with a colon (:) library(psych) describeBy(Data$Attendees, group = Data$Instructor, digits= 4) group: Ren vars n mean sd median trimmed mad min max range skew kurtosis se1 1 8 13.12 13.5 13.12 2.97 6 22 16 0 13 -1.08 1.8----- group: Stimpmy vars n mean sd median trimmed mad min max range skew kurtosis se1 1 8 15.88 4.45 16 15.88 3.71 9 23 14 -0.06 -1.26 1.57 describeBy(Data$Attendees, group = Data$Instructor : Data$Location, digits= 4) group: Ren:North vars n mean sd median trimmed mad min max range skew kurtosis se1 1 4 13.75 1.71 13.5 13.75 1.48 12 16 4 0.28 -1.96 0.85----- group: Ren:South vars n mean sd median trimmed mad min max range skew kurtosis se1 1 4 13.75 1.71 13.5 13.75 1.48 12 16 4 0.28 -1.96 0.85----- group: Stimpmy:North vars n mean sd median trimmed mad min max range skew kurtosis se1 1 4 14.75 4.03 16 14.75 4.03 16 14.75 2.22 9 18 9 -0.55 -1.84 2.02----- group: Stimpmy:South vars n mean sd median trimmed mad min max range skew kurtosis se1 1 4 17 5.16 17 17 5.93 11 23 12 0 -2.08 2.58 Data frame summaries We will often want to summarize variables throughout the data frame, to get some summary statistics for each variable, or to verify that the variables have the values that we expect when entering data to make sure that no error has occurred. The str function of the str function in the native utils package will list variables for the data frame, with their types and levels. Data is the name of the data frame you created above. str(Date) 'data.frame': 16 obs. with 3 variables: $ Instructor: Factor w/ 2 levels Rhine,Stimpmy: 1 1 1 1 1 1 2 2 . $ Location : Factor w/ 2 levels North,South: 1 1 1 2 2 2 2 2 1 1 . $ Participants : int 7 22 6 15 12 13 14 16 18 17 . ### Instructor is a variable coefficient (nominal) with two levels. ### Location is a (nominal) factor variable with two levels. ### Participants is an integer variable. Summary function The summary function in the native base package summarizes all variables in the data frame by listing the frequency of nominal variable levels; and for interval/aspect ratio, minimum, 1 quartile, median, mean, third quartile and maximum. Summary(Data) Instructor Location Participants Rhine :8 North:8 Min. : 6.00 Stimpmy:8 Noon:8 1.1.11.75 Median :15.00 Average :14.50 3rd Max. :23.00 HeadTail function in psych The headTail function in the psych package reports the first and last observations for the data frame. library(psych) headTail(Data) Instructor Location Participants 1 Ren North 72 Ren North 23 Rhine North 64 Rhine North 15... &lt;!&NA&gt; &lt;!&NA&gt;. ... 13 Stimpmy South 1514 Stimpmy South 1115 Stimpmy South 1916 Stimpmy South 23 Description function in psych The function described in the psych package reports the number of observations, averages, standard deviation, trimmed measures, minimum, maximum, range, tilt, courtesy and standard error for variables in the data frame. Note that factor variables are marked with an asterisk (*) and factor levels are encoded as 1, 2, 3, etc. library(psych) describe(Data) vars n mean sd median cropped mad min max bevel range kurtosis seInstructor* 1 1 6 1.5 0.52 1.5 1 .0 5 74 1 2 1 0.00 -2.12 0.13 Location* 2 16 1.5 0.52 1.5 1.5 0.74 1 2 1 0.00 -2.12 0.13 Participants 3 16 14.5`

