# Translation Syntax
# (SPSS, Stata, SAS, R, and python)

## The Basics

The following conventions are used in this document:

• Bold font indicates code or other text that should be typed literally.

• Un-bolded font shows code or text that should be replaced with user-supplied values (i.e., your own variable names and other environment details).

*Calling in a data set*

| | |
|---|---|
| **SPSS** | **GET FILE='P:\QAC\qac201\Studies\study name\filename.sav'.** |
| **STATA** | **use "**P:\QAC\qac201\Studies\study name\filename**"** |
| **SAS** | **LIBNAME** mydata **"**P:\QAC\QAC201\study name**";**<br>**DATA new; set** mydata**.**filename**;** |
| **R** | **load ("**filename-including-path**.Rdata")**<br>myData_orig **<-** name-of-object-that-loaded-in-your-workspace<br><br>If calling in from a tab-delimited text file:<br>myData_orig **<- read.table(file = "**filename-including-path.txt**", sep = "\t",**<br>**header = TRUE**) |
| **PYTHON** | **import pandas**<br>**import numpy**<br>myData = **pandas.read_csv(**'nesarc_pds.csv'**)** |

*Selecting variables you want to examine*

| SPSS | **/KEEP** VAR1 VAR2 VAR3 VAR4 VAR5 VAR6 VAR7 VAR8. (Must follow the **SAVE OUTFILE=**'dataname' command) |
|---|---|
| STATA | **use** VAR1 VAR2 VAR3 VAR4 VAR5 VAR6 VAR7 VAR8 ///<br>**using "**P:\QAC\qac201\Studies\study name\filename**", clear** |
| SAS | **KEEP** VAR1 VAR2 VAR3 VAR4 VAR5 VAR6 VAR7 VAR8**;** |
| R (base) | var.keep **<- c(**"VAR1**", "**VAR2**", "**VAR3**", "**VAR4**", "**VAR5**", "**VAR6**", "**VAR7**",** "VAR8**")**<br><br>myData **<-** myData_orig **[ ,**var.keep**]** |
| R (Tidyverse) | myData <- myData_orig **%>% select(**VAR1, VAR2, VAR3, VAR4, VAR5, VAR6, VAR7, VAR8**)**<br><br>or, if these variables are in consecutive columns:<br>myData <- myData_orig **%>% select(VAR1:VAR8)** |
| PYTHON | myData = myData_orig**[[**'VAR1',VAR2','VAR3','VAR4','VAR5','VAR6','VAR7', 'VAR8'**]]** |

*Outputting your abbreviated data set*

| SPSS | **SAVE OUTFILE=** **'**Drive:\folder\folder\title_of_new_data_set**'.** |
|---|---|
| STATA | **save** filename |
| SAS | **Data libname**.title_of_new_data_set; **set** dataname; **by** unique_id**;** |
| R | To open in excel:<br>**write.table(**myData**, file = "**filename.txt**", sep = "\t", row.names = FALSE)**<br><br>To use in R:<br>**save (**myData**, file = "**filename.Rdata**")** |

| | |
|---|---|
| **PYTHON** | **pandas.DataFrame.to_csv(**myData,'filename.csv'**)** |

*Sorting the data*

| | |
|---|---|
| **SPSS** | **SORT CASES BY** unique_id. |
| **STATA** | **sort** unique_id |
| **SAS** | **proc sort; by** unique_id**;** |
| **R (base)** | myData **<-** myData**[order(**myData**$**unique_id**, decreasing = FALSE), ]** |
| **R (Tidyverse)** | myData <- myData %>% **arrange**(unique_id)<br><br>*or in descending order*<br><br>myData <- myData %>% arrange(**desc**(unique_id)) |
| **PYTHON** | myData = myData.sort_values(by='unique_id') |

*Displaying frequency tables*

| | |
|---|---|
| **SPSS** | **FREQUENCIES VARIABLES**=VAR1 VAR2 VAR3<br>    **/ORDER=ANALYSIS.** |
| **STATA** | **tab1** VAR1 VAR2 VAR3 |
| **SAS** | **PROC FREQ; tables** VAR1 VAR2 VAR3**;** |
| **R (base)** | **library(descr)**<br>**freq(as.ordered(**myData**$**VAR1**))**<br>**freq(as.ordered(**myData**$**VAR2**))**<br>**freq(as.ordered(**myData**$**VAR3**))** |
| **R (Tidyverse)** | myData %>% **group_by**(VAR1) %>% **tally()**<br><br>myData %>% **group_by**(VAR2) %>% **tally()** |

| | |
|---|---|
| | or<br>myData **%>%**<br>  **group_by**(VAR1) **%>%**<br>  **tally() %>%**<br>  **mutate(pct = n/sum(n)\*100,**<br>    **cumpct = cumsum(pct))** |
| **PYTHON** | c1 = **myData['**VAR1**'].value_counts(sort=False, dropna=False)**<br>**print(**c1**)** |

## Data Management

*Basic operations*

| SPSS | EQ or = | >= or GE | <= or LE | > or GT | < or LT | NE |
|---|---|---|---|---|---|---|
| **STATA** | == | >= | <= | > | < | != |
| **SAS** | EQ or = | >= or GE | <= or LE | > or GT | < or LT | != or NE |
| **R** | == | >= | <= | > | < | != |
| **PYTHON** | == | >= | <= | > | < | != or <> |

## Examples

*1. Need to identify missing data*

Often, you must define the response categories that represent missing data. For example, if the number 9 is used to represent a missing value, you must either designate in your program that this value represents missingness or else you must recode the variable into a missing data character that your statistical software recognizes. If you do not, the 9 will be treated as a real/meaningful value and will be included in each of your analyses.

| SPSS | **RECODE** VAR1 **(**9=**SYSMIS).** |
|---|---|

| | |
|---|---|
| **STATA** | **replace** VAR1**=. if** VAR1**==9** |
| **SAS** | **if** VAR1=9 **then** VAR1=.**;** |
| **R (base)** | myData**$**VAR1**[**myData**$**VAR1 **== 9 ] <- NA**<br><br>*if multiple values need to be set to missing*<br>myData**$**VAR1**[**myData**$**VAR1 **%in%** c(8, 9) **] <- NA** |
| **R (Tidyverse)** | *Set single value in a single variable to missing*<br>myData <- myData %>% **mutate**(VAR1 = **na_if**(VAR1, 9))<br><br>*Set multiple values in a single variable to missing*<br>myData <- myData %>% **mutate**(VAR1 = **ifelse**(VAR1 **%in%** c(8, 9), NA, VAR1)<br><br>*Set multiple values to missing across multiple variables*<br>**vlist <- c("**VAR1", "VAR2", "VAR3**")**<br>**set.to.na <- function(x) ifelse(x  %in% c(8,9), NA, x)**<br>**myData <- myData %>% mutate(across(vlist, set.to.na ))**<br><br>*Change a 999 for ALL numeric variables into missing*<br>myData <- myData %>% mutate(across(where(is.numeric), ~na_if(., 999)))<br><br>*General method to recode values of a single variable:*<br>***Syntax: recode(df$variable, "old" = "new")***<br><br>#For variable OUTCOME convert PASS to 0 and FAIL to 1<br>Df$OUTCOME <- recode(df$OUTCOME,<br>                         "PASS" = 0,<br>                         "FAIL = 1)<br><br>#For variable SMOKER convert y to Smoker and n to Non-Smoker<br>df$SMOKER <- recode(df$SMOKER,<br>                    "y" = "Smoker",<br>                    "n" = "Non-Smoker") |
| **PYTHON** | myData['VAR1']= myData['VAR1'].replace(9, numpy.nan) |

## 2. Need to recode responses to "no" based on skip patterns

There are a number of skip outs in some data sets. For example, if we ask someone whether or not they have ever used marijuana, and they say "no", it would not make sense to ask them more detailed questions about their marijuana use (e.g. quantity, frequency, onset, impairment, etc.). When analyzing more detailed questions regarding marijuana (e.g. have you ever smoked marijuana daily for a month or more?), those individuals that never used the substance may show up as missing data. Since they have never used marijuana, we can assume that their answer is "no", they have never smoked marijuana daily. This would need to be explicitly recoded. Note that we commonly code a no as 0 and a yes as 1.

| | |
|---|---|
| **SPSS** | **RECODE** VAR1 **(SYSMIS**=7**).** |
| **STATA** | **replace** VAR1=**7 if** VAR1==**.** |
| **SAS** | **if** VAR1=**. then** VAR1=7**;** |
| **R** (base) | myData$VAR1**[is.na(**myData$VAR1**)] <-** 7 |
| **R** (Tidyverse) | myData <- myData %>% **mutate**(VAR1 = ifelse(**is.na**(VAR1), 7, VAR1) |
| **PYTHON** | myData**[**'VAR1'**].fillna**(7, **inplace=True**) |

## 3. Creating a new variable by recoding a string variables into numeric variable

It is important when preparing to run statistical analyses in most software packages, that all variables have response categories that are numeric rather than "string" or "character" (i.e. response categories are actual strings of characters and/or symbols). All variables with string responses must therefore be recoded into numeric values. These numeric values are known as dummy codes in that they carry no direct numeric meaning.

| SPSS | **RECODE** TREE **(**'Maple'=1**) (**'Oak'=2**) INTO** TREE_N**.** |
|---|---|
| STATA | **generate** TREE_N**=.**<br>**replace** TREE_N=1 **if** TREE=="Maple"<br>**replace** TREE_N=2 if TREE=="Oak"<br>OR by using the encode command<br>**encode** TREE**, gen(**TREE_N**)** |
| SAS | **IF** TREE='Maple' **then** TREE_N=1**;**<br>**else if** TREE= 'Oak' **then** TREE_N=2**;** |
| R | Not necessary in R. |
| PYTHON | def TREE_N (row) :<br>      if row['TREE'] == 'Maple' :<br>            return 1<br>      if row['TREE'] == 'Oak' :<br>            return 2<br><br>myData['TREE_N'] = myData.**apply(lambda row**: TREE_N **(row), axis = 1)** |

## 4. Creating a new variable by collapsing response categories

If a variable has many response categories, it can be difficult to interpret the statistical analyses in which it is used. Alternatively, there may be too few subjects or observations identified by one or more response categories to allow for a successful analysis. In these cases, you would need to collapse across categories. For example, if you have the following categories for geographic region, you may want to collapse some of these categories:

Region: New England=1, Middle Atlantic=2, East North Central=3, West North Central=4, South Atlantic=5, East South Central=6, West South Central=7, Mountain=8, Pacific=9.
New_Region: East=1, West=2.

| | |
|---|---|
| **SPSS** | **COMPUTE** new_region**=2.**<br>**IF (**region=1\| region=2\|region=3\| region=5\|region=6**)** new_region=1**.** |
| **STATA** | **generate** new_region =2<br>**replace** new_region=1 **if** region==1\| region==2\|region==3\| region==5\|region==6<br><br>OR by using the recode command<br>**recode region** (1/3 5 6=2) **gen(**new_region**)** |
| **SAS** | **if r**egion=1 or region=2 or region=3 or region=5 or region=6 **then** new_region=1**;**<br>**else if** region=4 or region=7 or region=8 or region=9 **then** new_region=2**;** |
| **R** **(base)** | myData$new_region <- NA<br><br>myData**$**new_region**[**myData**$**region %in% c(1,2,3,5,6)**] <-** 1<br>myData**$**new_region**[**myData**$**region %in% c(4,7,8,9)**] <-2**<br><br>or<br>myData$new_region <- **ifelse(**region %in% c(1,2,3,5,6), **1**, **2)** |
| **R**<br>**(Tidyverse)** | *If the new variable only has 2 categories:*<br>myData %>% **mutate(**new_region =**ifelse(**region %in% c(1,2,3,5,6), **1**, **2))**<br><br>*If the new variable has more than 2 categories:*<br>myData %>% **mutate(**new_region = **recode(**region,<br>        "2" = "1", "3" = "1", "5" = "1", "6"="1",<br>        "4"="2", "7" = "2",<br>        "8"="3", "9"="3")) |
| **PYTHON** | def new_region (row) :<br>       if row['region'] == 1 or row['region'] == 2 or row['region'] == 3 or row['region'] |

```
                    row['region'] ==5 or row['region'] == 6 :
                        return 1
                elif row['region'] == 4 or row['region'] == 7 or row['region'] == 7 or row['region']
                        row['region']== 8 or row['region'] == 9 :
                        return 2

myData['new_region'] = myData.apply(lambda row: new_region (row), axis = 1)
```

## 5. Creating a new variable by aggregating across variables

In many cases, you will want to combine multiple variables into one. For example, while NESARC assesses several individual anxiety disorders, I may be interested in anxiety more generally. In this case I would create a general anxiety variable in which those individuals who received a diagnosis of social phobia, generalized anxiety disorder, specific phobia, panic disorder, agoraphobia, or obsessive compulsive disorder would be coded "yes" and those who were free from all of these diagnoses would be coded "no".

Syntax shown on next page.

## 5. Creating a new variable by aggregating across variables (continued)

| | |
|---|---|
| **SPSS** | **IF (**socphob=1\|gad=1\|specphob=1\| panic=1\|agora=1\|ocd=1**) anxiety=1.**<br>**RECODE** anxiety **(SYSMIS**=0**).** |
| **STATA** | **gen** anxiety=1 if socphob==1\|gad==1\|specphob==1\|<br>panic==1\|agora==1\|ocd==1<br>**replace** anxiety=0 if anxiety==. |
| **SAS** | **if** socphob=1 or gad=1 or specphob=1 or panic=1 or agora=1 or ocd=1 **then** anxiety=1**;**<br>**else** anxiety=0**;** |
| **R (base)** | myData**$**anxiety **<- rep(**0, **nrow(**myData**))**<br>myData**$**anxiety**[**myData**$**socphob == 1 \| myData**$**gad==1 \| myData**$**panic == 1 \|<br>          myData**$**agora==1 \| myData**$**ocd == 1**] <- 1**<br>myData**$**anxiety**[is.na(**myData**$**socphob**) & is.na(**myData**$**gad**) &**<br>          **is.na(**myData**$**panic**) & is.na(**myData**$**agora**) &**<br>          **is.na(**myData**$**ocd**)] <- NA** |
| **R**<br>**(Tidyverse)** | *Tidyverse way: specify to do operations by row (**rowwise**), count number of non-missing symptoms (**n.symptoms**), count number of missing values per row (**n.miss**), create **anxiety** using **ifelse**, overwrite anxiety if all specific symptoms were missing.*<br><br>myData <- myData %>%<br>   rowwise() %>% |

| | |
|---|---|
| | ```
  mutate(n.symptoms = sum(c(socphob==1, gad==1, panic==1, agora==1, ocd==1),
na.rm=TRUE),
            n.miss = sum(is.na(socphob), is.na(gad), is.na(panic), is.na(agora), is.na(ocd)),
            anxiety = ifelse(n.symptoms > 0, 1, 0),

anxiety = ifelse(nmiss==5, NA, anxiety)
``` |
| **PYTHON** | ```python
def anxiety (row) :
        if row['socphob'] == 1 or row['gad'] == 1 or row['panic'] == 1 or row['agora'] == 1 or
                row['ocd'] == 1 :
                return 1
        else :
                return 0

myData['anxiety'] = data.apply(lambda row: anxiety (row), axis = 1)
``` |

## 6. Need to create quantitative variables

If you are working with a number of items that represent a single construct, it may be useful to create a composite variable/score. For example, I want to use a list of nicotine dependence symptoms meant to address the presence or absence of nicotine dependence (e.g. tolerance, withdrawal, craving, etc.). Rather than using a dichotomous variable (i.e. nicotine dependence present/absent), I want to examine the construct as a dimensional scale (i.e. number of nicotine dependence symptoms). In this case, I would want to recode each symptom variable so that yes=1 and no=0 and then sum the items so that they represent one composite score.

| | |
|---|---|
| **SPSS** | **COMPUTE** nd_sum=**sum(**nd_symptom1 nd_symptom2 nd_symptom3 nd_symptom4)**.** |
| **STATA** | **egen** nd_sum=**rsum(**nd_symptom1 nd_symptom2 nd_symptom3 nd_symptom4**)** |
| **SAS** | nd_sum=**sum (of** nd_symptom1 nd_symptom2 nd_symptom3 nd_symptom4**);** |
| **R (base)** | ```
myData$nd_sum <- NA

myData$nd_sum <- myData$nd_symptom1 + myData$nd_symptom2 +
                    myData$nd_symptom3 + myData$nd_symptom4
``` |
| **R (Tidyverse)** | ```
myData <- myData %>% rowwise() %>%
mutate(nd_sum = sum(nd_symptom1, nd_symptom2, nd_symptom3, nd_symptom4))
```<br><br>*Can be shortened if the variables are in consecutive columns:* |

|  |  |
|---|---|
|  | myData **<-** myData %>% rowwise() %>% <br> mutate(nd_sum **= sum(**nd_symptom1: nd_symptom4)) |
| **PYTHON** | myData['nd_sum'] = myData['nd_symptom1'] **+** myData['nd_symptom2'] **+** <br> myData['nd_symptom3'] **+** myData['nd_symptom4'] |

## 7. Labeling variables

Given the often cryptic names that variables are given, it can sometimes be useful to label them.

| SPSS | **VARIABLE LABELS** VAR1 **'**label**'.** |
|---|---|
| **STATA** | **label variable** VAR1 **"**label**"** |
| **SAS** | **LABEL** VAR1**='**label**';** |
| **R** | For frequency tables: <br> **library (**Hmisc**)** <br> **label(**myData**$**VAR1**) <- "**label**"** |
| **PYTHON** | **N/A** |

## 8. Renaming variables

Given the often cryptic names that variables are given, it can sometimes be useful to give a variable a new name (something that is easier for you to remember or recognize).

| SPSS | **COMPUTE** newvarname**=**VAR1**.** |
|---|---|
| **STATA** | **rename** VAR1 newvarname |
| **SAS** | **RENAME** VAR1=newvarname**;** |
| **R (base)** | **names(**myData**)[names(**myData**)== "**VAR1**"] <- "**newvarname**"** |
| **R (Tidyverse)** | myData <- myData %>% **rename**(newvar == oldvar) |
| **PYTHON** | myData = myData**.rename(**'oldvar'**:**'newvar'**, axis='columns')** |

## 9. Labeling variable responses/values

Given that nominal and ordinal variables have, or are given numeric response values (i.e. dummy codes), it can be useful to label those values so that the labels are displayed in your output.

| | |
|---|---|
| **SPSS** | **VALUE LABELS** VAR1 0 **'**value0label**'** 1 **'**value1label**'** 2 **'**value2label**'** 3 **'**value3label**'.** |
| **STATA** | **label define** VAR1 0 "value0label" 1 "value1label" 2 "value2label" 3 "value3label"<br>**label values** VAR1 newvarname |
| **SAS** | Set up format before the data step.<br>**proc format; VALUE** FORMATNAME 0**=**"value0label" 1**=**"value1label" 2**=**"value2label" 3**=**"value3label"**;**<br><br>Before the end of the data step, tell SAS which variables you would like to format with these values.<br>**format** VAR1 FORMATNAME**.** |
| **R (base)** | Because the function doesn't rearrange the order of the levels (the default is alphabetical), make sure you write the new labels in the order that they currently appear in the data set.<br><br>This shows you the existing levels in the current order. If VAR1 is a character variable, this will return NULL<br>**levels(**myData**$**VAR1**)**<br><br>If VAR1 is already a factor, you can apply the labels like this:<br>**levels(**myData**$**VAR1**) <- c(**"value0label**", "**value1label**", "**value2label**", "**value3label**")**<br><br>If VAR1 is not yet a factor, you can convert it to a factor, and apply the labels like this. Again, the order you write the labels must match the order that is currently in the data.<br><br>myData$VAR1 <- factor(myData$VAR1, levels = **c(**"value0label**", "**value1label**", "**value2label**", "**value3label**"))** |
| **R (Tidyverse)** | *To reorder factor levels:*<br>myData$VAR1 <- myData$VAR1 %>% fct_relevel("value1label", "value2label", "value3label", "value0label")<br><br>*To entirely reverse the label ordering*<br>myData$Var1 <- fct_rev(myData$Var1)<br><br>*To reorder the factor levels based on the frequency they show up in the data*<br><br>myData$Var1 <- fct_infreq(myData$Var1) |
| **PYTHON** | Because the function doesn't name the existing levels, make sure you have them all in the right order. |

| | |
|---|---|
| | myData**['VAR1']=** myData**['VAR1'].**astype('category')<br>myData**['VAR1']=** myData**['VAR1'].cat.rename_categories([**"value0label"**,** "value1label"**,**<br>"value2label"**,** "value3label"**])** |

## 10. Need to further subset the sample

When using large data sets, it is often necessary to subset the data so that you are including only those observations that can assist in answering your particular research question. In these cases, you may want to select your own sample from within the survey's sampling frame. For example, if you are interested in identifying demographic predictors of depression among Type II diabetes patients, you would plan to subset the data to subjects endorsing Type II Diabetes.

| | |
|---|---|
| **SPSS** | **/SELECT=**diabetes2 EQ 1 (must be added as a command option) |
| **STATA** | **if** diabetes2==1 (put this after the command) |
| **SAS** | **if** diabetes2=1**;** (put in the data step before sorting the data) |
| **R (base)** | title_of_subsetted_data **<-** myData**[**myData**$**diabetes2 **== 1, ]** |
| **R**<br>**(Tidyverse)** | title_of_subsetted_data <- myData %>% **filter**(diabetes2 ==1) |
| **PYTHON** | title_of_subsetted_data **=** myData[myData.diabetes2 == 1] |

## 11. Need to create groups that will be compared to one another

Often, you will need to create groups or sub-samples from the data set for the purpose of making comparisons. It is important to be certain that the groups that you would like to compare are of adequate size and number. For example, if you were interested in comparing complications of depression in parents who had lost a child through miscarriage vs. parents who had lost a child in the first year of life, it would be important to have large enough groups of each. It would not be appropriate to attempt to compare 5000 observations in the miscarriage group to only 9 observations in the first year group.

Refer to other data management syntax examples.

# Graphing and Data Visualization

## 1.    Univariate

Code for Univariate Output (Categorical):

| SPSS | **FREQUENCIES VARIABLES=**CategVar1 CategVar2 CategVar3 **/ORDER=ANALYSIS.** |
|---|---|
| **STATA** | **tab1** CategVar1 CategVar2 CategVar3 |
| **SAS** | **PROC FREQ; tables** CategVar1 CategVar2 CategVar3**;** |
| **R (base)** | table(myData$CategVar1)<br><br>**or**<br><br>**library(descr)**<br>**freq(as.ordered(**myData$CategVar1**))**<br>**freq(as.ordered(**myData$CategVar2**))**<br>**freq(as.ordered(**myData$CategVar3**))** |
| **R (Tidyverse)** | myData %>% group_by(CategVar1) %>% tally() |
| **PYTHON** | c1 **=** myData**[**'CategVar'**].value_counts(sort=False)**<br>**print (**c1**)** |

Code for Univariate Graph (Categorical):

| SPSS | Use graphical user interface (GUI) |
|---|---|
| **STATA** | **histogram** BinaryVar |
| **SAS** | **Proc GCHART; VBAR** CategVar **/ Discrete type=PCT Width=**30**;** |
| **R** | **library(ggplot2)**<br>**ggplot(data=**myData**)+**<br>  **geom_bar(aes(x=**CategVar**))+**<br>  **ggtitle(**"Descriptive Title Here"**)** |

| PYTHON | `import seaborn`<br>`seaborn.countplot(x="CategVar", data=myData)`<br>`plt.xlabel("Label for CategVar")`<br>`plt.title("Descriptive Title")` |
| --- | --- |

Code for Univariate Output (Quantitative):

| SPSS | `DESCRIPTIVES VARIABLES=QuantVar1 QuantVar2 QuantVar3`<br>`/STATISTICS=MEAN STDDEV.` |
| --- | --- |
| STATA | `summarize QuantVar1 QuantVar2 QuantVar3` |
| SAS | `proc means; var QuantVar1 QuantVar2 QuantVar3;` |
| R (base) | Repeat for each variable.<br>`summary(myData$QuantVar1)`<br>`mean(myData$QuantVar1, na.rm = TRUE)`<br>`sd(myData$QuantVar1, na.rm = TRUE)` |
| R (Tidyverse) | summarize can calculate multiple summary stats at the same time.<br>`myData %>% summarize(mean_var1 = mean(QuantVar1, na.rm=TRUE),`<br>`                sd_var1 = sd(QuantVar1, na.rm=TRUE))` |
| PYTHON | Repeat for each variable.<br>`desc1 = myData['QuantVar1'].dropna().describe()`<br>`print (desc1)` |

Code for Univariate Graph (Quantitative):

| SPSS | Use graphical user interface (GUI) |
| --- | --- |
| STATA | `histogram QuantVar` |
| SAS | `Proc GCHART; VBAR QuantVar;` |
| R | `ggplot(data=myData)+`<br>`   geom_histogram(aes(x=QuantVar))+`<br>`   ggtitle("Descriptive Title Here")` |

| | appropriate geometries: **geom_histogram, geom_density, geom_boxplot** |
|---|---|
| **PYTHON** | **import seaborn**<br>**seaborn.distplot(**myData**["**QuantVar**"].dropna(), kde=False)**<br>**plt.xlabel("**Label for QuantVar**")**<br>**plt.title("**Descriptive Title Here**")** |

## 2. Bivariate

Code for Bivariate Output (Categorical Explanatory Variable and Categorical Response Variable):

| | |
|---|---|
| **SPSS** | **CROSSTABS**<br>**/TABLES=**CategResponseVar **by** CategExplanatoryVar<br>**/CELLS=COUNT ROW COLUMN TOTAL.** |
| **STATA** | **tab** CategResponseVar CategExplanatoryVar**, row column cell** |
| **SAS** | **Proc freq; tables** CategResponseVar**\***CategExplanatoryVar**;** |
| **R (base)** | tab1 **<- table (**myData**$**CategResponseVar**,** myData**$**CategExplanatoryVar**)**<br>tab1 # to output the table<br><br>tab1_colProp **<- prop.table(**tab1**, 2)** # column proportions<br>tab1_rowProp **<- prop.table(**tab1**, 1)** # row proportions<br>tab1_cellProp **<- prop.table(**tab1**)** # cell proportions<br><br>tab1_colProp<br>tab1_rowProp<br>tab1_cellProp |
| **R (Tidyverse)** | % of response var within levels of explanatory var. Corresponds to row proportions<br>myData %>% group_by(CategResponseVar, CategExplanatoryVar) %>%<br>    summarise(n=n()) %>% na.omit() %>%<br>    group_by(CategExplanatoryVar) %>%<br>    mutate(prop = n/sum(n))<br><br>% of explanatory var within levels of response var. Corresponds to column proportions<br>myData %>% group_by(CategResponseVar, CategExplanatoryVar) %>% |

```
    summarise(n=n()) %>% na.omit() %>%
    group_by(CategResponseVar) %>%
    mutate(prop = n/sum(n))

Cell proportions
myData %>% group_by(CategResponseVar, CategExplanatoryVar) %>%
    summarise(n=n()) %>% na.omit() %>%
    ungroup() %>%

    mutate(prop = n/sum(n))
```

**PYTHON**

```python
print (pandas.crosstab(myData['CategResponseVar'], myData['CategExplanatoryVar'],
margins=True))
# get column proportions
print (pandas.crosstab(myData['CategResponseVar'], myData['CategExplanatoryVar'], margins=True,
normalize='columns'))
# get row proportions
print (pandas.crosstab(myData['CategResponseVar'], myData['CategExplanatoryVar'], margins=True,
normalize='index'))
# get cell proportions
print (pandas.crosstab(myData['CategResponseVar'], myData['CategExplanatoryVar'], margins=True,
normalize='all'))
```

Code for Bivariate Bar Graph (Categorical Explanatory Variable and Bivariate Categorical Response Variable – should be dummy coded 0/1):

| SPSS | Use graphical user interface (GUI) |
|---|---|
| **STATA** | **graph bar (mean)** CategResponseVar**, over(**CategExplanatoryVar**)** |
| **SAS** | **Proc GCHART; vbar** CategExplanatoryVar **/discrete type=mean sumvar=**CategResponseVar**;** |
| **R (base)** | **ggplot(data=**myData**)+**<br>   **stat_summary(aes(x=**CategExplanatoryVar**, y=**CategResponseVar**),**<br>                          **fun=mean, geom="bar")+**<br>   **ggtitle(**"Descriptive Title Here"**)** |
| **R (Tidyverse)** | Alternative if you do not code your response variable as 0/1.<br><br>Data.to.plot <- myData %>%<br>   group_by(CategResponseVar, CategExplanatoryVar) %>%<br>   summarise(n=n()) %>% na.omit() %>%<br>   group_by(CategExplanatoryVar) %>%<br>   mutate(prop = n/sum(n))<br><br>ggplot(Data.to.plot) +<br>   geom_col(aes(x= CategExplanatoryVar, fill= **Categ**ResponseVar, y=prop),<br>position="dodge") |
| **PYTHON** | **seaborn.catplot(x=**"CategExplanatoryVar"**, y=**"QuantResponseVar"**, data=**myData**, kind="bar", ci=None)**<br>**plt.xlabel('**Label for CategExplanatoryVar**')**<br>**plt.ylabel('**Label for QuantResponseVar**')**<br>**plt.title('**Descriptive Title Here**')** |

Note: Your graph will display the proportion of observational units within a particular category who exhibit the indicated level of the response variable.

Code for Bivariate Output (Categorical Explanatory Variable and Quantitative Response Variable):

| | |
|---|---|
| **SPSS** | **MEANS TABLES=** CategExplanatoryVar **by** QuantResponseVar<br>**/CELLS MEAN COUNT STDDEV.** |
| **STATA** | **bys** CategExplanatoryVar**: su** QuantResponseVar |
| **SAS** | **proc sort; by** CategExplanatoryVar**;**<br>**proc means; var** QuantResponseVar**; by** CategExplanatoryVar**;** |
| **R (base)**<br><br>**R (Tidyverse)** | **by(**myData**$**QuantResponseVar**,** myData**$**CategExplanatoryVar**, mean, na.rm = TRUE)**<br>**by(**myData**$**QuantResponseVar**,** myData**$**CategExplanatoryVar**, sd, na.rm = TRUE)**<br>**by(**myData**$**QuantResponseVar**,** myData**$**CategExplanatoryVar**, length)**<br><br>myData %>% **group_by**(CategExplanatoryVar) %>%<br>   **summarize**(mean_var1 = mean(QuantVar1, na.rm=TRUE),<br>           sd_var1 = sd(QuantVar1, na.rm=TRUE),<br>           n_var1 = n()) |
| **PYTHON** | desc1 **=** myData**[**'QuantResponseVar'**]**.groupby(myData**[**'CategExplanatoryVar'**]**).describe()<br>print (desc1) |

Code for Bivariate Graphs (Categorical Explanatory Variable and Quantitative Response Variable):

| | |
|---|---|
| **SPSS** | Use graphical user interface (GUI) |
| **STATA** | **graph box** QuantResponseVar**, over(**CategExplanatoryVar**)** |
| **SAS** | **Proc GCHART; vbar** CategExplanatoryVar **/discrete**<br>**type=mean sumvar=**QuantResponseVar**;** |
| **R** | ##Below is code for bar graph<br>**ggplot(data=**myData**)+**<br>  **stat_summary(aes(x=**CategExplanatoryVar**, y=**QuantResponseVar**),**<br>                 **fun.y=mean, geom="bar")**<br>as.character(myData$ CategExplanatoryVar)<br><br>##Below is code for boxplots<br>**ggplot(data=**myData**)+**<br>  **geom_boxplot(aes(x=**CategExplanatoryVar**, y=**QuantResponseVar**))+**<br>  **ggtitle(**"Descriptive Title Here"**)**<br><br>##Below is code for density plots<br>**ggplot(data=**myData**)+**<br>  **geom_density(aes(x=**QuantResponseVar, **color=**CategExplanatoryVar**))+**<br>  **ggtitle(**"Descriptive Title Here) |
| **PYTHON** | scat1 **= seaborn.catplot(x=** 'CategExplanatoryVar ', **y=** 'QuantResponseVar ',<br>**data=**myData, **kind="bar", ci=None)**<br>**print(**scat1**)** |

Code for Bivariate Scatterplot (Quantitative Explanatory Variable and Quantitative Response Variable):

| | |
|---|---|
| **SPSS** | **GRAPH /scatterplot(bivar)=**QuantExplanatoryVar **with** QuantResponseVar**.** |
| **STATA** | **twoway (scatter** QuantResponseVar QuantExplanatoryVar**) (lfit** QuantResponseVar QuantExplanatoryVar**)** |
| **SAS** | **Proc GPLOT; Plot** QuantResponseVar ***QuantExplanatoryVar;** |
| **R** | **ggplot(data=**myData**)+**<br>  **geom_point(aes(x=**QuantExplanatoryVar**, y=**QuantResponseVar**))+**<br>  **geom_smooth(aes(x=**QuantExplanatoryVar**, y=**QuantResponseVar**), method="lm")**<br><br>### Note, you can also add a 3<sup>rd</sup> variable to a scatterplot by using the color argument:<br>**ggplot(data=**myData**)+**<br>  **geom_point(aes(x=**QuantExplanatoryVar**, y=**QuantResponseVar**, ,**<br>           **color=**CategThirdVar**))+**<br>  **geom_smooth(aes(x=**QuantExplanatoryVar**, y=**QuantResponseVar,<br>           **color=**CategThirdVar**)), method="lm")** |
| **PYTHON** | **seaborn.regplot(x=**"QuantExplanatoryVar", **y=**"QuantResponseVar", **fit_reg=False,**<br>**data**=myData**)**<br>**plt.xlabel(**'Label for QuantExplanatoryVar'**)**<br>**plt.ylabel(**'Label for QuantResponseVar'**)**<br>**plt.title(**'Descriptive Title Here'**)** |

Note: In the R section, "3<sup>rd</sup>" uses the superscript: 3$^{rd}$ variable.

## 3. Adding a Third Variable

Code for Output with a Third Variable (Categorical Explanatory Variable, Quantitative Response Variable, Categorical 3rd VAR):

| SPSS | **MEANS TABLES=** QuantResponseVar **BY** CategExplanatoryVar **BY** CategThirdVar **/CELLS MEAN COUNT STDDEV.** |
|---|---|
| **STATA** | **bys** CategExplanatoryVar CategThirdVar**: su** QuantResponseVar |
| **SAS** | **proc sort; by** CategExplanatoryVar CategThirdVar**;**<br>**proc means; var** QuantResponseVar**; by** CategExplanatoryVar CategThirdVar**;** |
| **R** (base) | **ftable(by(**myData**$**QuantResponseVar**,**<br>　　　**list(**myData**$**CategExplanatoryVar**,** myData**$**CategThirdVar**),**<br>　　　**mean, na.rm = TRUE))** |
| **R** (Tidyverse) | Summary statistics of quantitative response variable within cross-combination of two categorical variables<br><br>myData %>% **group_by**(CategExplanatoryVar, CategThirdVar) %>%<br>　**summarize**(mean = mean(QuantResponseVar, na.rm=TRUE),<br>　　　　　sd = sd(QuantResponseVar, na.rm=TRUE),<br>　　　　　n = n()) |
| **PYTHON** | #graphing code<br>**seaborn.catplot(x=**"CategExplanatoryVar"**, y=**"QuantResponseVar"**,**<br>　　　**hue=**"CategThirdVar"**, data=**myData**, kind='bar', ci=None)**<br>**plt.xlabel(**'Label for CategExplanatoryVar'**)**<br>**plt.ylabel(**'Label for QuantResponseVar'**)**<br>**plt.title(**'Descriptive Title Here'**)** |

Code for Output with a Third Variable (Categorical Explanatory Variable and Categorical Response Variable, Categorical 3rd VAR):

| SPSS | **CROSSTABS** **/TABLES=**CategResponseVar **BY** CategExplanatoryVar **BY** CategThirdVar**.** |
|---|---|
| STATA | **bys** CategExplanatoryVar CategThirdVar**: tab** CategResponseVar |
| SAS | **proc sort; by** CategThirdVar**;** **proc freq; tables** CategResponseVar*CategExplanatoryVar**; by** CategThirdVar**;** |
| R (base) | tab1 **<- ftable(**myData**$**CategResponseVar**,** myData**$**CategExplanatoryVar**,** myData**$**CategThirdVar**)** tab1 tab1_colProp **<- prop.table(**tab1**, 2)** tab1_colProp |
| R (Tidyverse) | myData %>% **group_by**(CategThirdVar , CategResponseVar , CategExplanatoryVar) %>% tally() |
| PYTHON | seaborn.catplot(x="CategExplanatoryVar", y="CategResponseVar", hue="CategThirdVar", data=myData, kind="bar", ci=None) **plt.xlabel(**'Label for CategExplanatoryVar'**)** **plt.ylabel(**'Label for CategResponseVar'**)** **plt.title(**'Descriptive Title Here'**)** |

Note: If your 3rd variable is quantitative, for graphing purposes, create meaningful categories and then use the code above.

# Bivariate Analysis

*ANOVA*

| | |
|---|---|
| **SPSS** | **UNIANOVA** QuantResponseVar **BY** CategExplanatoryVar**.** |
| **STATA** | **oneway** QuantResponseVar CategExplanatoryVar**, tabulate** |
| **SAS** | **proc anova;**<br>**class** CategExplanatoryVar**;**<br>**model** QuantResponseVar **=** CategExplanatoryVar**;**<br>**means** CategExplanatoryVar**;** |
| **R** | myAnovaResults **<- aov(**QuantResponseVar **~**<br>CategExplanatoryVar**, data =** myData**)**<br>**summary(**myAnovaResults**)** |
| **PYTHON** | import statsmodels.formula.api as smf<br>import statsmodels.stats.multicomp as multi<br><br>model1 = smf.ols(formula='QuantResponseVar ~ C(CategExplanatoryVar)', data=myData)<br>results1 = model1.fit()<br>print (results1.summary()) |

*Pearson Correlation*

| | |
|---|---|
| **SPSS** | **CORRELATIONS**<br>**/VARIABLES=** QuantResponseVar QuantExplanatoryVar<br>**/STATISTICS DESCRIPTIVES.** |
| **STATA** | **corr** QuantResponseVar QuantExplanatoryVar r<br>OR<br>**pwcorr** QuantResponseVar QuantExplanatoryVar**, sig** |
| **SAS** | **Proc corr; var** QuantResponseVar QuantExplanatoryVar**;** |
| **R** | **cor.test(**myData**$**QuantResponseVar, myData**$**QuantExplanatoryVar**)** |
| **PYTHON** | **import scipy** |

```
sub1= myData[['QuantResponseVar', 'QuantExplanatoryVar']].dropna()
print ('association between QuantExplanatoryVar and QuantResponseVar')
print (scipy.stats.pearsonr(sub1['QuantResponseVar'], sub1['QuantExplanatoryVar']))
```

*Chi-Square Test*

| SPSS | `CROSSTABS`<br>`/TABLES= CategResponseVar by CategExplanatoryVar`<br>`/STATISTICS=CHISQ.` |
|---|---|
| STATA | `tab CategResponseVar CategExplanatoryVar , chi2 row col` |
| SAS | `Proc freq; tables CategResponseVar*CategExplanatoryVar/ chisq;` |
| R | `myChi <- chisq.test(myData$CategResponseVar, myData$CategExplanatoryVar)`<br><br>`myChi`<br>`myChi$observed # for actual, observed cell counts`<br>`prop.table(myChi$observed, 2) # for column percentages`<br>`prop.table(myChi$observed, 1) # for row percentages` |
| PYTHON | `import scipy.stats`<br>`ct1=pandas.crosstab(myData['CategResponseVar'], myData['CategExplanatoryVar'])`<br>`print ('chi-square value, p value, degrees of freedom, expected counts')`<br>`cs1= scipy.stats.chi2_contingency(ct1)`<br>`print (cs1)`<br>`# column percentages`<br>`colsum=ct1.sum(axis=0)`<br>`colpct=ct1/colsum`<br>`print(colpct)` |

*POST HOC TESTS WITHIN ANOVA*

| SPSS | **UNIANOVA** QuantResponseVar **BY** CategExplanatoryVar<br>**/POSTHOC=**CategExplanatoryVar **(TUKEY)**<br>**/PRINT=ETASQ DESCRIPTIVE.** |
|---|---|
| STATA | **oneway** QuantResponseVar CategExplanatoryVar**, sidak** |
| SAS | **Proc anova;**<br>**class** CategExplanatoryVar**;**<br>**model** QuantResponseVar**=**CategExplanatoryVar**;**<br>**means** CategExplanatoryVar**/duncan;** |
| R | myAnovaResults **<- aov(**QuantResponseVar **~** CategExplanatoryVar, **data =** myData**)**<br>**TukeyHSD(**myAnovaResults**)** |
| PYTHON | model1 = **smf.ols(formula='**QuantResponseVar ~ **C(**CategExplanatoryVar**)'**, **data=**myData**)**<br>results1 = model1**.fit()**<br>print (results1**.summary()**)<br><br>sub1 = myData[['QuantResponseVar', 'CategExplanatoryVar']]**.dropna()**<br>mc1 = **multi.MultiComparison(**sub1['QuantResponseVar'], sub1['CategExplanatoryVar']**)**<br>res1 = mc1**.tukeyhsd()**<br>print(res1**.summary()**) |

*POST HOC TESTS FOR CHI SQUARE (must subset data in order to conduct 2X2 comparisons)*

| SPSS | **TEMPORARY.**<br>**SELECT IF** CategExplanatoryVar**=1 OR** CategExplanatoryVar **=3.**<br>**CROSSTABS**<br>**/TABLES=** CategResponseVar CategExplanatoryVar<br>**/STATISTICS=CHISQ.** |
|---|---|
| STATA | **tab** CategResponseVar CategExplanatoryVar **If**<br>CategExplanatoryVar**==1 |** CategExplanatoryVar**==3, chi2** |
| SAS | **IF (**CategExplanatoryVar **= 1) OR (**CategExplanatoryVar **= 3);** (in data step)<br>**Proc freq; tables** CategResponseVar*CategExplanatoryVar**/ chisq;** |

| R | ## You do not need to subset the data in R<br><br>pairwise.prop.test(table(myData$CategResponseVar, myData$CategExplanatoryVar), p.adjust.method="bonferroni") |
|---|---|
| PYTHON | ```python
#for each Chi Sq pair data subset (code below compares group 1 to group 2)
recode1 = {1: 1, 2: 2}
myData['COMP1v2']= myData['CategExplanatoryVar'].map(recode1)
ct1=pandas.crosstab(myData['CategResponseVar], myData['COMP1v2']
cs1= scipy.stats.chi2_contingency(ct1)
print (cs1)
``` |

# Statistical Interactions: Testing for Moderation

*Moderation: ANOVA*

In these analyses, the third variable must be categorical.

| SPSS | **SORT CASES BY** CategThirdVar**.**<br>**SPLIT FILE LAYERED BY** CategThirdVar**.**<br>**ONEWAY** QuantResponseVar **BY** CategExplanatoryVar<br>**/ STATISTICS DESCRIPTIVES**<br>**/ POSTHOC = BONFERRONI ALPHA (0.05).**<br>**SPLIT FILE OFF.** |
|---|---|
| STATA | **bys** CategThirdVar**: oneway** QuantResponseVar CategExplanatoryVar**, tab** |
| SAS | **Proc sort; by** CategThirdVar**;**<br>**Proc anova; class** CategExplanatoryVar**;**<br>**model** QuantResponseVar**=**CategExplanatoryVar**;**<br>**means** CategExplanatoryVar**; by** CategThirdVar**;** |
| R | **by(**myData**,** myData**$**CategThirdVar**, function(x)**<br>    **list(**<br>        **aov(**QuantResponseVar **~** CategExplanatoryVar**, data = x),**<br>        **summary(aov(** QuantResponseVar **~** CategExplanatoryVar**, data = x))))** |

| | |
|---|---|
| **PYTHON** | `#subset by categorical 3rd variable`<br>`sub2=myData[(myData['CategThirdVar']=='Group 1')]`<br>`sub3=myData[(myData['CategThirdVar']=='Group 2')]`<br><br>`import statsmodels.api`<br>`import statsmodels.formula.api as smf`<br>`model2 = smf.ols(formula='QuantResponseVar ~ C(CategExplanatoryVar)',`<br>`data=sub2).fit()`<br>`print (model2.summary())`<br><br>`model3 = smf.ols(formula= formula='QuantResponseVar ~ C(CategExplanatoryVar)',`<br>`data=sub3).fit()`<br>`print (model3.summary())` |

*Moderation: PEARSON CORRELATION*

In these analyses, the third variable must be categorical

| | |
|---|---|
| **SPSS** | `SORT CASES BY CategThirdVar.`<br>`SPLIT FILE LAYERED BY CategThirdVar.`<br>`CORRELATIONS`<br>`/VARIABLES= QuantResponseVar QuantExplanatoryVar`<br>`/STATISTICS DESCRIPTIVES.`<br>`SPLIT FILE OFF.` |
| **STATA** | `bys CategThirdVar: corr QuantResponseVar QuantExplanatoryVar`<br>OR<br>`bys CategThirdVar: pwcorr QuantResponseVar QuantExplanatoryVar, sig` |
| **SAS** | `Proc sort; by CategThirdVar;`<br>`Proc corr; var QuantResponseVar QuantExplanatoryVar; by CategThirdVar;` |
| **R** | `by(myData, myData$CategThirdVar, function(x)`<br>`    cor.test(x$QuantResponseVar, x$QuantExplanatoryVar))` |
| **PYTHON** | `#subset by categorical 3rd variable`<br>`sub1=myData[(myData['CategThirdVar']== 1)]`<br>`sub2=myData[(myData['CategThirdVar']== 2)]`<br>`sub3=myData[(myData['CategThirdVar']== 3)]`<br><br>`print (scipy.stats.pearsonr(sub1['QuantResponseVar'], sub1['QuantExplanatoryVar']))`<br>`print (scipy.stats.pearsonr(sub2['QuantResponseVar'], sub2['QuantExplanatoryVar']))`<br>`print (scipy.stats.pearsonr(sub3['QuantResponseVar'], sub3['QuantExplanatoryVar']))` |

*Moderation: CHI-SQUARE TEST*

In these analyses, the third variable must be categorical.

| SPSS | **CROSSTABS**<br>**/TABLES =** CategResponseVar **by** CategExplanatoryVar<br>**by** CategThirdVar<br>**/CELLS = COUNT ROW**<br>**/STATISTICS = CHISQ.** |
|---|---|
| **STATA** | **bys** CategThirdVar**: tab** CategResponseVar CategExplanatoryVar**,**<br>**chi2 row** |
| **SAS** | **Proc sort; by** CategThirdVar**;**<br>**Proc freq; tables** CategResponseVar**\*CategExplanatoryVar/chisq; by** CategThirdVar**;** |
| **R** | **by(**myData**,** myData**$CategThirdVar, function(x)**<br>   **list(**<br>      **chisq.test(x$**CategResponseVar**, x$**CategExplanatoryVar**),**<br>      **chisq.test(x$**CategResponseVar**, x$**CategExplanatoryVar**)$observed,**<br>      **prop.table(chisq.test(x$**CategResponseVar**, x$**CategExplanatoryVar**)$observed,**<br>      **2)))** # column %s |
| **PYTHON** | #subset by categorical 3rd variable<br>sub2=myData**[(**myData**['CategThirdVar']==**'Group 1**')]**<br>sub3=myData**[(**myData**['CategThirdVar']==**'Group 2**')]**<br><br>ct2**=pandas.crosstab(**sub2**['CategResponseVar'],** sub2**['CategExplanatoryVar'])**<br>**print (**ct2**)**<br>ct3**=pandas.crosstab(**sub3**['CategResponseVar'],** sub3**['CategExplanatoryVar'])**<br>**print (**ct3**)** |

# Regression Analyses

In these analyses, treat ordinal variables (e.g., variables scaled Strongly Disagree (1) to Strongly Agree (5)) as quantitative. Dummy code (0 = no, 1 = yes) all multi-level categorical variables.

Adding variables other than your response variable and your explanatory variable, as shown by the use of ThirdVar in the following syntax, lets you test for confounding.

# Multivariate Regression

*MULTIPLE REGRESSION*

| SPSS | **REGRESSION**<br>**/DEPENDENT** QuantResponseVar<br>**/METHOD ENTER** ExplanatoryVar ThirdVar1 ThirdVar2**.** |
|---|---|
| **STATA** | **reg** QuantResponseVar ExplanatoryVar ThirdVar1 ThirdVar2 |
| **SAS** | Code binary variables as yes = 1 and no = 2.<br>**Proc glm; class** CategExplanatoryVar CategThirdVar;<br>**model** QuantResponseVar**=** CategExplanatoryVar CategThirdVar QuantThirdVar<br>**/solution;** |
| **R** | my.lm **<- lm(**QuantResponseVar **~** ExplanatoryVar **+** ThirdVar1 **+** ThirdVar2**, data =** myData**)**<br>**summary(**my.lm**)** |
| **PYTHON** | **import statsmodels.api**<br>**import statsmodels.formula.api as smf**<br>#note that categorical explanatory/third variables have to be entered as C(CategVar)<br>lm1 = **smf.ols('**QuantResponseVar **~** ExplanatoryVar + **C(**CategThirdVar1**)** +<br>QuantThirdVar**', data=**myData**).fit()**<br>**print (**lm1**.summary())** |

*LOGISTIC REGRESSION*

| SPSS | **LOGISTIC REGRESSION** BinaryResponseVar **with** ExplanatoryVar ThirdVar1 ThirdVar2**.** |
|---|---|
| **STATA** | **logistic** BinaryResponseVar ExplanatoryVar ThirdVar1 ThirdVar2 |
| **SAS** | Code your binary variables as yes = 1 and no = 2.<br><br>**Proc logistic;**<br>**class** CategExplanatoryVar CategThirdVar**;**<br>**model** BinaryResponseVar**=**CategExplanatoryVar CategThirdVar QuantThirdVar**;** |
| **R** | my.logreg **<- glm(**BinaryResponseVar **~** ExplanatoryVar **+** ThirdVar1 +<br>ThirdVar2**, data =** myData**, family = "binomial")**<br><br>**summary(**my.logreg**)**        # for p-values<br>**exp(my.logreg$coefficients)**   # for odds ratios<br>**exp(confint(**my.logreg**))**      # for confidence intervals on the odds ratios |
| **PYTHON** | **import statsmodels.api**<br>**import statsmodels.formula.api as smf**<br># logistic regression |

```python
lreg1 = smf.logit(formula = 'BinaryResponseVar ~ ExplanatoryVar + C(CategThirdVar) +
QuantThirdVar', data = myData).fit()
print (lreg1.summary())

# odd ratios with 95% confidence intervals
params = lreg1.params
conf = lreg1.conf_int()
conf['OR'] = params
conf.columns = ['Lower CI', 'Upper CI', 'OR']
print (numpy.exp(conf))
```