**STAT 506: Homework 3**

For these problems you will need to access the data in the PG1/data folder. Use the `libname` statement we learned to load this each time you work on your assignments. You should call it 'pg1' to be consistent with the SAS materials.

I tried to *italicize* the parts where I expect you to actually show me something in your homework solutions if it is not obvious.

1. **DATA step processing and filtering**

   Write a DATA step to do the following:

   - Read in the table **pg1.eu_occ**.

   - Add a WHERE statement to select only the stays that were reported in the year 2014. [Note that **YearMon** is a character column, and the first four characters represent the year.]

   - Assign the COMMA17. format to the **Hotel**, **ShortStay**, and **Camp** columns.

   - Save the new table as **eu_occ2014**, but exclude the column **Geo**.

   Print the first 6 observations of **eu_occ2014**. *Show your code and the output.*

2. **Creating New Columns**

   Write a DATA step to do the following:

   - Read in the table **pg1.np_summary**.

   - Create a new column named **SqMiles** by dividing the column **Acres** by 640.

   - Create a new column named **Camping** as the sum of **OtherCamping**, **TentCampers**, **RVCampers**, and **BackcountryCampers**.

   - Format **SqMiles** and **Camping** to include commas and zero decimal places.

   - Save the new table as **np_summary_update**, but only include the new columns created above and **ParkName**.

   Print the first 10 observations of **np_summary_update**. *Show your code and the output.*

3. **Using Conditional Processing to Re-Categorize and Clean Data**

   a. As we've seen previously, the table **pg1.np_summary** is using some inconsistent codes for the column **Type**. Create a frequency table for **Type**. *Show your code.*

   b. Write a DATA step to create a new table named **park_type** that includes everything from **pg1.np_summary**. Also use IF-THEN/ELSE statements to create a new column named **ParkType** based on the value of **Type**:

      - NP → Park

      - NM → Monument

      - NS → Seashore

      - RVR or RIVERWAYS → River

      - PRE, NPRE, or PRESERVE → Preserve

      *Show your code and the corresponding log notes.*

   c. Create a frequency table for **ParkType**. *Show your code and the output.*

4. **Two-Way Frequency Reports**

Make a two-way frequency report for the columns **sex** and **birthdate** in **pg1.class_birthdate**.

- Use **birthdate** as the row variable.

- Use a format to group the values of **birthdate** by year instead of individual date; if done properly, this should result in a table with 6 rows.

- Add the label "Year" to **birthdate**.

- Add the titles "Class Overview" on the first line and "Birth Year versus Sex" on the third line.

- Add your name as a footnote.

- Use an option in the TABLES statement to suppress the column percentages.

- Add code to clear the titles and footnote after the report is generated.

*Show your code and the output.*

5. **Using Labels in PROC PRINT**

a. Write a PROC CONTENTS step to display the descriptor portion of **pg1.eu_occ** to see the permanent labels assigned to the columns. *Show the relevant part of the output (the part that shows the labels).*

b. Print the first 6 observations from **pg1.eu_occ**. All the columns should be displayed with their permanent labels, except for **ShortStay**, which should have the temporarily assigned label "Nights Spent at Short Stays" displayed instead. *Show your code and output.*

6. **Creating an Output Summary Table**

a. Write a PROC MEANS step that will calculate summary statistics for the variable **hotel** in **pg1.eu_occ** using **country** as the class variable. Save the output as a new temporary table named **med_hotel** which includes the median values for the **hotel** variable as a variable named **MedianHotel**. Use the NOPRINT option. *Show your code and the corresponding log notes.*

b. Write a PROC SORT step to sort **med_hotel** by **MedianHotel** in descending order. In this step, also filter out the row that summarizes the entire table (the row with a blank Country) if you didn't already do part a in a way that automatically removes that row. *Show your code and the corresponding log notes.*

c. Write a DATA step to update **med_hotel** by eliminating the columns **_TYPE_** and **_FREQ_**. In this step, also assign **MedianHotel** the permanent label "Median of Nights Spent at Hotels". *Show your code and the corresponding log notes.*

d. Finally, print the first 6 observations from **med_hotel** and display the labels for the variables. *Show your code and output.*