

# STSCI 5060 Final Project

(Fall 2018)

**Important: This project must be submitted by 4:30 PM, December 13, 2018, and any overdue project will not be accepted. It is a one-person project. Cornell's Code of Academic Integrity is strictly enforced. Do your own work!**

## General Instructions

### 1. Coding:

- A. At the very beginning, start your Oracle and SAS code with the following comment block.  

```
/* Fall 2018 STSCI 5060 Final Project */  
/* Name:                               */  
/* NetID:                              */
```
- B. Use ORACLE SQL DEVELOPER for SQL coding unless otherwise specified.
- C. If a task can be done with Oracle SQL\*PLUS or with SAS PROC SQL, use Oracle SQL\*PLUS. You use SAS PROC SQL only when you are explicitly told to do so.
- D. Always leave a blank line between two blocks of code.
- E. When you code in Oracle, start each project step with a title like:
  - i. `title '***** Step # *****' skip 2` (# is the step number).
- F. When you code in SAS, start your project step with a title like:
  - i. `title '***** Step # *****';` (# is the step number).If this guideline is not followed, up to 5 points will be taken off.
- G. Set your line size (e.g., 5000) and page size (e.g., 1000) properly so that you produce well-organized outputs, e.g., your Oracle output should not wrap to next line(s). If your output is messy without following this guideline, up to 5 points will be taken off.
- H. Insert a brief comment at the beginning of a code block to briefly explain what the code does. Up to 5 points will be taken off if these comments are missing.

### 2. Submission:

- A. Your Oracle code/script, named LastName\_FirstName\_Oracle\_Fall2018STSCI5060FP.sql (Note: use your real name to replace LastName\_FirstName, the same below).
- B. Your SAS code, named LastName\_FirstName\_SAS\_Fall2018STSCI5060FP.sas.
- C. All the Oracle script outputs, i.e., your Oracle query results. (LastName\_FirstName\_SQL\_output.txt). You may produce this file in a new, single run after all your code works correctly, but you will need to first drop State\_t, School\_Finance\_2010\_t, School\_Finance\_2014\_t, and School\_t, etc. (hint: drop table state\_t CASCADE CONSTRAINTS;) and then create them again when you run your SQL script.
- D. Your SAS log (LastName\_FirstName\_SAS\_log.log).
- E. An MS Word (or PDF) file containing all your screenshots, code entered in the Oracle command line interface, the ERD from Step 9, SAS analysis results, and answers to the questions (LastName\_FirstName\_FP\_Report.docx or LastName\_FirstName\_FP\_Report.PDF).
- F. At the end, compress all the above files with 7-ZIP (or a similar tool) into one final file named LastName\_FirstName\_Fall2018STSCI5060FP.7z and submit it to the course website.

## The Project Background and Steps

This project utilizes much knowledge and programming skills (both SQL and SAS) covered in the class. You will create a database in Oracle, based on two big SAS datasets called *School\_Finance\_2010.sas7bdat* and *School\_Finance\_2014.sas7bdat*, and then you use this database to analyze the data with sql\*plus and SAS. The *School\_Finance\_201X.sas7bdat* data files contain real data about the finance of US public schools in the years of 2010 and 2014. For our international students who are not familiar with the US school systems, here is some short explanation. In each state, the public schools are divided into different school districts. A school district normally includes many different levels of schools (elementary schools, middle schools, and high schools in a city/town/area). A school district is normally funded by three sources: the local revenue, the state revenue and the federal revenue. The datasets describe these revenues and expenditures. You need to read *Column Description.PDF* so that you understand the datasets and the meaning of column names. Please strictly follow the following steps as you do the project and your project report.

- ① 1. In Oracle command line, connect to the Oracle SYSTEM user and create a new user called **LastName\_FirstName\_STSCI5060FP** and grant all the privileges to this user. Then, in ORACLE SQL DEVELOPER create a new connection to this new user with the same name, LastName\_FirstName\_STSCI5060FP.
- ② 2. In SQL DEVELOPER, create a table, **State\_t**, by importing a csv file called State\_Code.csv. To import, right click "Tables (filtered)" and choose "Import Data ..." and then follow the instructions. In the Column Definition window (in the 4<sup>th</sup> Step of Data Import Wizard), you need to delete the trailing blank(s) for all the variables that are listed in the Source Data Columns window, in order to have valid column names. Use the default setting when you import the table, except setting the attributes of columns StCode, StName and StAbb to varchar2(2), varchar2(26) and char(2) respectively.

Attach the following screenshots of ORACLE SQL DEVELOPER: the **Columns** tab and **Data** tab of your State\_T table (just showing a portion of the data of the table is OK).

- ③ 3. In Oracle, programmatically (i.e., you have to code it with SQL) update your State\_t table by changing the single-digit values, 1-9, of state code to two-digit values, 01-09. This will make sure that these values are consistent with those in other tables. In this step, you are required to accomplish the task just with one update statement with the following requirements:
  - A. Use a WHERE clause to specify a condition.
  - B. Use the concatenation operator.
  - C. Use the SUBSTR() function.
  - D. Use the CAST function to explicitly convert the data type for your comparison in the WHERE clause (although in this case Oracle can do data type automatic conversion without using this function).

Query your updated State\_t table to confirm that you have successfully made the changes. Only display the 9 rows whose Stcode values are less than 10.

(4)

4. In SAS 9.4, create a libref called **"Final"** that references your file location of all your final project SAS files, and then connect SAS to the Oracle database user, "LastName\_FirstName\_STSCI5060FP," with the LIBNAME statement by creating a libref called **myoracle**. Create an Oracle database table, **School\_Finance\_2010\_t**, using the myoracle libref and PROC SQL by querying the School\_Finance\_2010.sas7bdat dataset (which is in SAS format).

Right after you ran your SAS code, refresh the connections in your ORACLE SQL Developer. In Oracle, describe the School\_Finance\_2010\_t table and display the first 10 rows of the table (hint: use rownum<=10).

(2)

5. Change table properties of the School\_Finance\_2010\_t table:
  - A. Change the property of IDCENSUS to varchar2(15).
  - B. Change the property of NAME to varchar2(60).

(2)

6. Change column names: rename the column "NAME" to "SD\_NAME" and the column "State" to "STCODE" in the School\_Finance\_2010\_t table.

(4)

7. The School\_Finance\_2010\_t table is big table with many functional dependencies. You will create four tables from this big table. The first three tables are Fedrev\_t, Strev\_t and Locrev\_t for the federal, state and local revenues respectively. The fourth table is called School\_t.
  - A. Your Fedrev\_t table should include the following columns: idcensus, stcode, and fed\_rev, which is the sum of columns c14, c15, c16, c17, c18, c19, b11, c20, c25, c36, b10, b12, and b13.
  - B. Your Strev\_t table should include the following columns: idcensus, stcode, and st\_rev, which is the sum of columns c01, c04, c05, c06, c07, c08, c09, c10, c11, c12, c13, c24, c35, c38, and c39.
  - C. Your Locrev\_t table should include the following columns: idcensus, stcode, and loc\_rev, which is the sum of columns t02, t06, t09, t15, t40, t99, d11, d23, a07, a08, a09, a11, a13, a15, a20, a40, u11, u22, u30, u50, and u97.
  - D. Your School\_t table should include these columns: idcensus, stcode, and sd\_name.

(6.5)

8. Programmatically set primary keys and foreign keys for the following tables:
  - A. Set the stcode column as the primary key of the State\_t table.
  - B. Set the idcensus column in the Fedrev\_t, Strev\_t, Locrev\_t, and School\_t tables as the primary key.
  - C. Set the idcensus column of the Fedrev\_t, Strev\_t, Locrev\_t tables as the foreign key that references the idcensus column of the School\_t table.
  - D. Set the stcode column of the School\_t table as the foreign key that references the stcode column of the State\_t table.

[Hint: To set the foreign key, use the following syntax: **ALTER TABLE** *table\_name* **add constraint** *foreign\_key\_name* **foreign KEY** (*column\_name*) **references** *name\_of\_the\_table\_referenced* (*name\_of\_the\_column\_referenced*);]

9. Each of the table you have created is correspondent to an entity. As a result, you have five entities: STATE, SCHOOL, FEDREV, STREV and LOCREV. Draw an ERD (using the book's HRT notation) to represent the relationships of these entities. You are required to list all the attributes. Then, convert your ERD to relations and use arrows to represent the primary key/foreign key constraints. If you draw the diagrams with your hand, you may scan it and then insert it into the MS Word (or PDF) file.

10. Based on Fedrev\_t, list all the school districts that received more than \$1,000,000 K from the federal source (note that the revenue values are expressed in thousands of dollars in the database tables). You should display three columns: idcensus, stcode and an alias, fed\_revenue for the federal revenue. Do the same for the Strev\_t and Locrev\_t tables and find out all the school districts that received more than \$1,000,000 K from the state or local sources. You should name the state and local total revenue aliases as st\_revenue and loc\_revenue. The revenue values you display should not be in scientific notation (**do the same for the rest of the project**). Hint: use the function **to\_char(value, '999999999.9')** to achieve this; the numbers of the digit "9" represent the scale and precision of the number.

11. Create a view called **sd#\_v** to calculate the total number of school districts (SD#) in each state. This view has two columns, SD# and stcode. Then,
- find the state(s) that with the highest number of school districts by using sd#\_v. In your output, list the state code, state name and the total number of school districts.
  - find the state(s) that with the lowest number of school districts by using sd#\_v. List the state code, state name and the total number of school districts.

12. In each state, among all the school districts, find out the highest local, state and federal revenues, and sort your result by state. Your final output should be listed in the following column order and format:

STCODE	MAX_FED_REV	MAX_ST_REV	MAX_LOC_REV	STATE_NAME
1	xxxxx.0	xxxxx.0	xxxxx.0	Alabama
2	xxxxx.0	xxxxx.0	xxxxx.0	Alaska
.	.	.	.	.
.	.	.	.	.

To achieve this, you need to do the following:

- Based on the tables you created in Step 7, create three views in Oracle called **mfr\_v**, **msr\_v**, and **mlr\_v** to calculate the maximum federal, state, and local revenues in each state. Each view should only contain two columns: stcode and the value of the

maximum revenue, sorted by stcode. Name the maximum revenues as MAX\_FED\_REV, MAX\_ST\_REV and MAX\_LOC\_REV, respectively.

- B. Connect SAS 9.4 to Oracle using the libname approach. In a SAS DATA step, use an appropriate dataset combining method to combine the data in the above three views to directly create a table called **mfslr\_t** in Oracle. This table should contain four columns: stcode, MAX\_FED\_REV, MAX\_ST\_REV and MAX\_LOC\_REV.
- C. Go back to Oracle and use the mfslr\_t table created by above SAS DATA Step to get your results. You should space your four columns as shown above, i.e., do not leave too much white space between the columns (hint: use the to\_char() function to achieve this).

5 13. This is an extension of above question. In addition to listing the state name, state code and the highest federal revenue (use aliases, state\_code for state code, state\_name for stname, and max\_fed\_rev for the highest total federal revenue of the school district in that state), you are required to list the name of the school district that received the highest federal revenue in that state, as the 4<sup>th</sup> column. Sort your result by the revenue in descending order.

1 14. Create a view called **Total\_Rev\_v** from fedrev\_t, strev\_t, and locrev\_t by including idcensus, state code, total federal revenue (named tfedrev), total state revenue (named tstrev), and total local revenue (named tlocrev) of each school district.

4 15. From Total\_Rev\_v, calculate the total revenues (use an alias, total\_revenue) of these three sources (tfedrev, tstrev, and tlocrev) for each school district in the US. With the information from other tables, display the columns in the order of stcode, stname, idcensus, total\_revenue and sd\_name. The result should be sorted by the total revenue in descending order. Just output the first 100 rows (Hint: use ROWNUM <= 100).

3 16. The total expenditure of a school district is indicated by the TOTALEXP column in the SCHOOL\_FINANCE\_2010\_T table. Find out the total school expenditure of each state. Include the following columns in your query output: stcode, stname, and the total school expenditure of the state. Sort your output with the total school expenditure in descending order.

2 17. Calculate the total amount of the money that the United State spent on the public school systems in that year. Your query output **must** read something like

"The total amount that the United States spent on the public school systems in 2010 was \$ XXXXXXXXXXXXX.X K."

Note:

- The integer part of the dollar amount is not necessary the same length as indicated by the number of X'es here, but you do need to keep one decimal point.

- There is no column heading displayed. (Hint: use the “SET HEADING OFF” command to suppress headings and use “SET HEADING ON” command to restore to the normal condition).

18. Based on fedrev\_t or strev\_t, or locrev\_t and other necessary tables, create 3 views, excluding any school districts that had no expenses (otherwise, you will have an error in calculation). You should check the value range of the ratios you are getting, which will affect how you choose the display format of your ratios (all the values must be displayed as required).

- A. fed\_contribution\_v, to calculate the federal revenue contribution to each school district (federal revenue/total expense), including these columns: idcensus, stcode, sd\_name and the fed\_pcmt (for the federal revenue/total expense ratio). Keep 4 decimal points for the ratio. Find out school districts that received federal revenues greater than the total expense, listing all the columns that exist in the fed\_contribution\_v and sorting in descending order by fed\_pcmt.
- B. st\_contribution\_v, to calculate the state revenue contribution to each school district (state revenue/total expense), including these columns: idcensus, stcode, sd\_name and the st\_pcmt (for the state revenue/total expense ratio). Keep 4 decimal points for the ratio. List all the columns that exist in the st\_contribution\_v, and sort in descending order by st\_pcmt. Find out school districts that received state revenues greater than the total expense, and sort in descending order by st\_pcmt.
- C. loc\_contribution\_v, to calculate the local revenue contribution to each school district (local revenue/total expense), including these columns: idcensus, stcode, sd\_name and the loc\_pcmt (for the revenue/total expense ratio). Keep 4 decimal points for the ratio. Find out school districts that received local revenues greater than the total expense. List all the columns that existed in the loc\_contribution\_v, and sort in descending order by loc\_pcmt.

19. Based on the three views that were created in Step 18, create another view called fsl\_contribution\_v, including these columns: idcensus, stcode, sd\_name and the fsl\_pcmt (for the total ratio, which is the sum of fed\_pcmt, st\_pcmt and loc\_pcmt). Keep 4 decimal points.

- A. Find out the school districts that received total revenues (federal+state+local) over 3 times of the total amount they actually spent in that year. List all the columns that exist in the fsl\_contribution\_v, and sort in descending order by fsl\_pcmt.
- B. Find out the school districts that received total revenues (federal+state+local) up to 30% of the total amount they actually spent in that year. List all the columns that exist in the fsl\_contribution\_v, and sort in descending order fsl\_pcmt.

20. In SASUSER, use PROC SQL to create a dataset in SAS called **Total\_Rev** by querying the

Total\_Rev\_v view you created in Step 14 in the Oracle database. This will create a SAS dataset that has the same contents as the view.

21. Do a correlation analysis (PROC CORR) on the three variables (tfedrev, tstrev, and tlocrev) of the Total\_Rev dataset. Put the analysis results in the MS Word (or PDF) file and comment on the results.

(Hint: use the following the syntax of PROC CORR

```
proc corr data=...  
plots(maxpoints=NONE)=matrix(histogram);  
var variable1 variable2 variable3;  
run;)
```

22. Do a regression analysis (PROC REG) on variables (tfedrev, tstrev, and tlocrev) of the Total\_Rev dataset. Use tfedrev as the dependent variable and tstrev and tlocrev as the independent variables. Put the analysis results in the MS Word (or PDF) file and comment on the results.

(Hint: use the following the syntax of PROC REG

```
proc reg data=...;  
model dependent-variable = independent-variable1 independent-variable2;  
run;)
```

23. Do the correlation analysis and regression analysis by directly getting data from the Total\_Rev\_v view in the Oracle database through the myoracle libref you created earlier. You should achieve the same results as above. Is this an instance of in-database processing? Why or why not?

24. Using PROC SQL and the librefs (myoracle and Final) created earlier in Step 4, create an Oracle database table, **School\_Finance\_2014\_t** by querying the School\_Finance\_2014.sas7bdat dataset.

25. From 2010 to 2014, some school districts total revenues increased, some decreased and a small number of them stayed the same. Using sql\*plus, find out

- A. The top 10 school districts that had increased total revenues.
- B. The top 10 school districts that had decreased total revenues.
- C. All the school districts whose total revenues stayed the same.

You are required to display the state code, state name, IDCENSUS, the name of school district, the difference of the total revenues of the two years (revdif=2014's totalrev – 2010's totalrev) and the change percentage (change\_percentage = 100\*revdif/2010's totalrev). The results are sorted by revdif. All the values should keep one decimal point.

26. Report a screenshot of your final ORACLE SQL DEVELOPER interface, showing **all** your tables, views and the connection that you have created in the whole project. Make sure that you expand the Tables and Views folder so that all your tables and Views are visible.