

Module 7 - Final Project

DAD215

March 3, 2019

Evan Leet

I. Introduction to the Software Interface

A. Use a libname to show the creation of a data library. This library name will be used at different times in the code to reference the location of the data sets.

CODE:

```
/*The following code will satisfy the Introduction to Software interface portion of the final project rubric.  
3/3/19  
Evan Leet  
  
input: final project folder and files  
output: library  
  
*/  
  
libname bikes 'C:\Users\evan.leet\Documents\bikesheets1';  
/*Use Libname statement. Name new library. Point SAS towards proper folder.*/  
RUN;
```

LOG:

```
46 /*The following code will satisfy the Introduction to Software interface portion of the final  
46 ! project rubric.  
47  
48 3/3/19  
49 Evan Leet  
50  
51 input: final project folder and files  
52 output: library  
53  
54  
55 */  
56  
57 libname bikes 'C:\Users\evan.leet\Documents\bikesheets1';  
NOTE: Libref BIKES was successfully assigned as follows:  
Engine: U9  
Physical Name: C:\Users\evan.leet\Documents\bikesheets1  
58 /*Use Libname statement. Name new library. Point SAS towards proper folder.*/  
59 RUN;
```

II. Working With the Data: In this section, you will utilize SAS methods and data-manipulation techniques to work with the provided data. At the end, you will submit your finalized code that demonstrates your ability to do the following:

A. Utilize an appropriate method for importing and converting the data from the provided CSV file to an Excel file.

CODE:

```
/*The folowing code will use PROC IMPORT to convert the excel files into SAS files and place them in the library
```

```
3/3/19
Evan Leet
```

```
input: final project folder and files
output: SAS files
```

```
*/
```

```
=PROC IMPORT out=bikes.bikes72013 datafile='C:\Users\evan.leet\Documents\bikesheets1\bikes72013.csv'
/*Use proc import to import data files. Use OUT to name the file and place it in the proper livrary
Point SAS towards the data file on the computer
*/
```

```
DBMS=csv replace; /*original file is a CSV file*/
getnames= yes;
```

```
RUN;
```

```
/*Run and repeat for the other files*/
```

```
=PROC IMPORT out=bikes.bikes72014 datafile='C:\Users\evan.leet\Documents\bikesheets1\bikes72014.csv'
DBMS=csv replace;
getnames= yes;
```

```
RUN;
```

```
=PROC IMPORT out=bikes.bikes72015 datafile='C:\Users\evan.leet\Documents\bikesheets1\bikes72015.csv'
DBMS=csv replace;
getnames= yes;
```

```
RUN;
```

LOG:

```
100 100% The following code will use PROC IMPORT to convert the excel files into SAS files and place them
101 in the library.
102
103 100% 3/3/19
104 100% Evan Leet
105
106 100% input: final project folder and files
107 100% output: SAS files
108
109 100% */
110
111 100% =PROC IMPORT out=bikes.bikes72013 datafile='C:\Users\evan.leet\Documents\bikesheets1\bikes72013.csv'
112 /*Use proc import to import data files. Use OUT to name the file and place it in the proper
113 livrary
114 Point SAS towards the data file on the computer
115 */
116
117 100% DBMS=csv replace; /*original file is a CSV file*/
118
119 100% getnames= yes;
120
121 100% RUN;
122
123 100% /*Run and repeat for the other files*/
124
125 100% =PROC IMPORT out=bikes.bikes72014 datafile='C:\Users\evan.leet\Documents\bikesheets1\bikes72014.csv'
126 DBMS=csv replace;
127
128 100% getnames= yes;
129
130 100% RUN;
131
132 100% =PROC IMPORT out=bikes.bikes72015 datafile='C:\Users\evan.leet\Documents\bikesheets1\bikes72015.csv'
133 DBMS=csv replace;
134
135 100% getnames= yes;
136
137 100% RUN;
138
139 100%
140
141 100%
142
143 100%
144
145 100%
146
147 100%
148
149 100%
150
151 100%
152
153 100%
154
155 100%
156
157 100%
158
159 100%
160
161 100%
162
163 100%
164
165 100%
166
167 100%
168
169 100%
170
171 100%
172
173 100%
174
175 100%
176
177 100%
178
179 100%
180
181 100%
182
183 100%
184
185 100%
186
187 100%
188
189 100%
190
191 100%
192
193 100%
194
195 100%
196
197 100%
198
199 100%
200
201 100%
202
203 100%
204
205 100%
206
207 100%
208
209 100%
210
211 100%
212
213 100%
214
215 100%
216
217 100%
218
219 100%
220
221 100%
222
223 100%
224
225 100%
226
227 100%
228
229 100%
230
231 100%
232
233 100%
234
235 100%
236
237 100%
238
239 100%
240
241 100%
242
243 100%
244
245 100%
246
247 100%
248
249 100%
250
251 100%
252
253 100%
254
255 100%
256
257 100%
258
259 100%
260
261 100%
262
263 100%
264
265 100%
266
267 100%
268
269 100%
270
271 100%
272
273 100%
274
275 100%
276
277 100%
278
279 100%
280
281 100%
282
283 100%
284
285 100%
286
287 100%
288
289 100%
290
291 100%
292
293 100%
294
295 100%
296
297 100%
298
299 100%
300
301 100%
302
303 100%
304
305 100%
306
307 100%
308
309 100%
310
311 100%
312
313 100%
314
315 100%
316
317 100%
318
319 100%
320
321 100%
322
323 100%
324
325 100%
326
327 100%
328
329 100%
330
331 100%
332
333 100%
334
335 100%
336
337 100%
338
339 100%
340
341 100%
342
343 100%
344
345 100%
346
347 100%
348
349 100%
350
351 100%
352
353 100%
354
355 100%
356
357 100%
358
359 100%
360
361 100%
362
363 100%
364
365 100%
366
367 100%
368
369 100%
370
371 100%
372
373 100%
374
375 100%
376
377 100%
378
379 100%
380
381 100%
382
383 100%
384
385 100%
386
387 100%
388
389 100%
390
391 100%
392
393 100%
394
395 100%
396
397 100%
398
399 100%
400
401 100%
402
403 100%
404
405 100%
406
407 100%
408
409 100%
410
411 100%
412
413 100%
414
415 100%
416
417 100%
418
419 100%
420
421 100%
422
423 100%
424
425 100%
426
427 100%
428
429 100%
430
431 100%
432
433 100%
434
435 100%
436
437 100%
438
439 100%
440
441 100%
442
443 100%
444
445 100%
446
447 100%
448
449 100%
450
451 100%
452
453 100%
454
455 100%
456
457 100%
458
459 100%
460
461 100%
462
463 100%
464
465 100%
466
467 100%
468
469 100%
470
471 100%
472
473 100%
474
475 100%
476
477 100%
478
479 100%
480
481 100%
482
483 100%
484
485 100%
486
487 100%
488
489 100%
490
491 100%
492
493 100%
494
495 100%
496
497 100%
498
499 100%
500
501 100%
502
503 100%
504
505 100%
506
507 100%
508
509 100%
510
511 100%
512
513 100%
514
515 100%
516
517 100%
518
519 100%
520
521 100%
522
523 100%
524
525 100%
526
527 100%
528
529 100%
530
531 100%
532
533 100%
534
535 100%
536
537 100%
538
539 100%
540
541 100%
542
543 100%
544
545 100%
546
547 100%
548
549 100%
550
551 100%
552
553 100%
554
555 100%
556
557 100%
558
559 100%
560
561 100%
562
563 100%
564
565 100%
566
567 100%
568
569 100%
570
571 100%
572
573 100%
574
575 100%
576
577 100%
578
579 100%
580
581 100%
582
583 100%
584
585 100%
586
587 100%
588
589 100%
590
591 100%
592
593 100%
594
595 100%
596
597 100%
598
599 100%
600
601 100%
602
603 100%
604
605 100%
606
607 100%
608
609 100%
610
611 100%
612
613 100%
614
615 100%
616
617 100%
618
619 100%
620
621 100%
622
623 100%
624
625 100%
626
627 100%
628
629 100%
630
631 100%
632
633 100%
634
635 100%
636
637 100%
638
639 100%
640
641 100%
642
643 100%
644
645 100%
646
647 100%
648
649 100%
650
651 100%
652
653 100%
654
655 100%
656
657 100%
658
659 100%
660
661 100%
662
663 100%
664
665 100%
666
667 100%
668
669 100%
670
671 100%
672
673 100%
674
675 100%
676
677 100%
678
679 100%
680
681 100%
682
683 100%
684
685 100%
686
687 100%
688
689 100%
690
691 100%
692
693 100%
694
695 100%
696
697 100%
698
699 100%
699
```

B. Create two new variables that represent the month and year of each file.

CODE:

```
/*The following code will use the same format for the date for each data set. It will also add a Month and Year variable.

3/3/19
Evan Leet

input: final project folder and files, SAS files created earlier
output: data sets with formatted dates

*/

%DATA bikes.Bikes72013FD; /*New data set name. FD stands for formatted date*/
  set bikes.bikes72013; /*Point to original data set*/

  dateformatted = input(starttime, ANYDTDTM21.); /*Create new variable for the formatted date
  Set the input as starttime variable, use ANYDTDTM21. to allow multiple date formats as input*/

  dateformatted = datepart(dateformatted); /*use datepart to keep only the date, use date9 format*/
  format dateformatted date9.;

  Month = Month(dateformatted); /*extract month from dateformatted, assign to new variable*/
  Year = Year(dateformatted); /*extract year*/

RUN;
/*Run and repeat for other data sets*/

%DATA bikes.Bikes72014FD;
  set bikes.bikes72014;

  dateformatted = input(starttime, ANYDTDTM21.);

  dateformatted = datepart(dateformatted);
  format dateformatted date9.;

  Month = Month(dateformatted);
  Year = Year(dateformatted);

RUN;

%DATA bikes.Bikes72015FD;
  set bikes.bikes72015;

  dateformatted = input(starttime, ANYDTDTM21.);

  dateformatted = datepart(dateformatted);
  format dateformatted date9.;

  Month = Month(dateformatted);
  Year = Year(dateformatted);

RUN;

/*verify with proc print*/

%PROC PRINT data=bikes.bikes72013fd (obs=5);
RUN;
%PROC PRINT data=bikes.bikes72014fd (obs=5);
RUN;
%PROC PRINT data=bikes.bikes72015fd (obs=5);
RUN;
```

LOG:

```

1524 /*The following code will use the same format for the date for each data set. It will also add a
1525 Date and Year.
1526
1527 3/3/19
1528 Even Last
1529
1530 input: final project folder and files
1531 output: SAS files
1532 */
1533
1534 NOTE: There were 5 observations read from the data set BIKES.BIKEST2015FD.
1535 NOTE: PROCEDURE PRINT used (Total process time):
1536    real time          0.19 seconds
1537    cpu time           0.02 seconds
1538
1539 DATA bikes.bikest2013fd; /*New data set name. FD stands for formatted dates/
1540    set bikes.bikest2013; /*Point to original data sets/
1541    dateformatted = input(starttime, ANYDTDTM2.); /*Create new variable for the formatted date
1542    Set the input as starttime variable, use
1543    ANYDTDTM2. to allow multiple date formats as inputs/
1544    dateformatted = datepart(dateformatted); /*use datepart to keep only the date, use dates
1545    format dateformatted dates.;
1546    Month = Month(dateformatted); /*extract month from dateformatted, assign to new variable/
1547    Year = Year(dateformatted); /*extract years/
1548 RUN;
1549
1550 NOTE: There were 843416 observations read from the data set BIKES.BIKEST2013.
1551 NOTE: The data set BIKES.BIKEST2013FD has 843416 observations and 18 variables.
1552 NOTE: DATA statement used (Total process time):
1553    real time          1.71 seconds
1554    cpu time           1.15 seconds
1555
1556 /*Run and repeat for other data sets/
1557
1558 DATA bikes.bikest2014fd;
1559    set bikes.bikest2014;
1560    dateformatted = input(starttime, ANYDTDTM2.);
1561    dateformatted = datepart(dateformatted);
1562    format dateformatted dates.;
1563    Month = Month(dateformatted);
1564    Year = Year(dateformatted);
1565 RUN;
1566
1567 NOTE: There were 868842 observations read from the data set BIKES.BIKEST2014.
1568 NOTE: The data set BIKES.BIKEST2014FD has 868842 observations and 18 variables.
1569 NOTE: DATA statement used (Total process time):
1570    real time          2.13 seconds
1571    cpu time           1.24 seconds
1572
1573
1574 DATA bikes.bikest2015fd;
1575    set bikes.bikest2015;
1576    dateformatted = input(starttime, ANYDTDTM2.);
1577    dateformatted = datepart(dateformatted);
1578    format dateformatted dates.;
1579    Month = Month(dateformatted);
1580    Year = Year(dateformatted);
1581 RUN;
1582
1583 NOTE: There were 1685876 observations read from the data set BIKES.BIKEST2015.
1584 NOTE: The data set BIKES.BIKEST2015FD has 1685876 observations and 18 variables.
1585 NOTE: DATA statement used (Total process time):
1586    real time          2.31 seconds
1587    cpu time           1.68 seconds
1588
1589
1590 /*verify with proc prints/
1591
1592 PROC PRINT data=bikes.bikest2013fd (obs=5);
1593 RUN;
1594
1595 NOTE: There were 5 observations read from the data set BIKES.BIKEST2013FD.
1596 NOTE: PROCEDURE PRINT used (Total process time):
1597    real time          0.09 seconds
1598    cpu time           0.01 seconds
1599
1600 PROC PRINT data=bikes.bikest2014fd (obs=5);
1601 RUN;
1602
1603 NOTE: There were 5 observations read from the data set BIKES.BIKEST2014FD.
1604 NOTE: PROCEDURE PRINT used (Total process time):
1605    real time          0.02 seconds
1606    cpu time           0.01 seconds
1607
1608 PROC PRINT data=bikes.bikest2015fd (obs=5);
1609 RUN;
1610
1611 NOTE: There were 5 observations read from the data set BIKES.BIKEST2015FD.
1612 NOTE: PROCEDURE PRINT used (Total process time):
1613    real time          0.01 seconds
1614    cpu time           0.01 seconds

```

OUTPUT:

ir	gender	dateformatted	Month	Year
0		01JUL2013	7	2013
0		01JUL2013	7	2013
2		01JUL2013	7	2013
0		01JUL2013	7	2013
1		01JUL2013	7	2013

ir	gender	dateformatted	Month	Year
2		01JUL2014	7	2014
1		01JUL2014	7	2014
2		01JUL2014	7	2014
1		01JUL2014	7	2014
2		01JUL2014	7	2014

ir	gender	dateformatted	Month	Year
1		01JUL2015	7	2015
1		01JUL2015	7	2015
1		01JUL2015	7	2015
1		01JUL2015	7	2015
1		01JUL2015	7	2015

C. Before merging your data, ensure that you have sorted your data appropriately. For example, both data sets must be sorted by the same ID variable.

CODE:

```

/*The following code will sort data sets by BIKEID

3/3/19
Evan Leet

input: final project folder and files, SAS files created earlier
output: sorted data sets

*/

PROC SORT data=bikes.bikes72013fd; /*Use PROC SORT. Point SAS to the relevant data file*/
  BY bikeid; /*Sort by bikeid*/
RUN;
/*Run and repeat for other datasets*/

PROC SORT data=bikes.bikes72014fd;
  BY bikeid;
RUN;

PROC SORT data=bikes.bikes72015fd;
  BY bikeid;
RUN;

```

LOG:

```

1628 /*The following code will sort data sets by BIKEID
1629
1630 3/3/19
1631 Evan Leet
1632
1633 input: final project folder and files, SAS files created earlier
1634 output: sorted data sets
1635
1636 */
1637
1638 PROC SORT data=bikes.bikes72013fd; /*Use PROC SORT. Point SAS to the relevant data file*/
1639   BY bikeid; /*Sort by bikeid*/
1640 RUN;
NOTE: Input data set is already sorted, no sorting done.
NOTE: PROCEDURE SORT used (Total process time):
      real time          0.00 seconds
      cpu time           0.01 seconds

1641 /*Run and repeat for other datasets*/
1642
1643 PROC SORT data=bikes.bikes72014fd;
1644   BY bikeid;
1645 RUN;
NOTE: Input data set is already sorted, no sorting done.
NOTE: PROCEDURE SORT used (Total process time):
      real time          0.00 seconds
      cpu time           0.00 seconds

1646
1647 PROC SORT data=bikes.bikes72015fd;
1648   BY bikeid;
1649 RUN;
NOTE: Input data set is already sorted, no sorting done.
NOTE: PROCEDURE SORT used (Total process time):
      real time          0.00 seconds
      cpu time           0.01 seconds

```

D. Merge the data from the data sets.

CODE:

```

/*The following code will merge the previous data sets

3/3/19
Evan Leet

input: final project folder and files, SAS files created earlier
output: one merged set
*/

DATA bikes.bikescomplete; /*create and name new dataset*/
  merge bikes.bikes72013fd bikes.bikes72014fd bikes.bikes72015fd; /*use merge, call on relevant datafiles*/
  BY bikeid; /*merge by bikeid since sets were sorted by this variable earlier*/
RUN;

```

LOG:

```

1650 /*The following code will merge the previous data sets
1651
1652 3/3/19
1653 Evan Leet
1654
1655 input: final project folder and files, SAS files created earlier
1656 output: one merged set
1657 */
1658
1659 DATA bikes.bikescomplete; /*create and name new dataset*/
1660     merge bikes.bikes72013fd bikes.bikes72014fd bikes.bikes72015fd; /*use merge, call on
1661     relevant datafiles*/
1662     BY bikeid; /*merge by bikeid since sets were sorted by this variable earlier*/
1663 RUN;

WARNING: Multiple lengths were specified for the variable start_station_name by input data set(s).
This can cause truncation of data.
WARNING: Multiple lengths were specified for the variable end_station_id by input data set(s). This
can cause truncation of data.
WARNING: Multiple lengths were specified for the variable end_station_name by input data set(s). This
can cause truncation of data.
WARNING: Multiple lengths were specified for the variable end_station_latitude by input data set(s).
This can cause truncation of data.
WARNING: Multiple lengths were specified for the variable end_station_longitude by input data set(s).
This can cause truncation of data.
NOTE: MERGE statement has more than one data set with repeats of BY values.
NOTE: There were 843416 observations read from the data set BIKES.BIKES72013FD.
NOTE: There were 968842 observations read from the data set BIKES.BIKES72014FD.
NOTE: There were 1085676 observations read from the data set BIKES.BIKES72015FD.
NOTE: The data set BIKES.BIKESCOMPLETE has 1554715 observations and 18 variables.
NOTE: DATA statement used (Total process time):
      real time          3.90 seconds
      cpu time           1.35 seconds

```

E. Review your work to this point and ensure that you identify and resolve any syntax or logic errors that you have encountered. Provide at least one example of an issue you debugged or justify why no changes were needed.

The work up to this point was making a usable data set that we can manipulate and ultimately use for the purposes set out in the assignment. I had trouble converting the given excel files into SAS data sets. At first, I was treating them as xlsx files but I then realized that they were CSV files. Once I made this important distinction, incorporating the files into SAS went more smoothly. I also had a hard time converting the dates into a single format. In the 3 original files, the dates were formatted differently and included the time. I did some research via google and found a solution that easily allowed me to convert all the dates to date9 format, and then extract the month and year from that.

F. Within your data, manipulate the text so that “Gender” is represented by numbers. (Male = 1; Female = 2; Unknown = 0).

Already Complete.

G. After formatting the data, utilize a function to merge two variables together. (Station A is one variable, and Station B is one variable) = Station A – Station B (new variable).

CODE:

```

/*The following create a new variable which is the combination of two existing variables

3/3/19
Evan Leet

input: bikescomplete set
output: bikes complete set with an additional variable
*/

DATA bikes.bikescomplete;
  set bikes.bikescomplete;

  monthyear = catx('/', month, year); /*create new variable, use catx, determine separator and existing variables to use*/

RUN;

```

LOG:

```

1732 /*The following create a new variable which is the combination of two existing variables
1733
1734 3/3/19
1735 Evan Leet
1736
1737 input: bikescomplete set
1738 output: bikes complete set with an additional variable
1739 */
1740
1741 DATA bikes.bikescomplete;
1742   set bikes.bikescomplete;
1743
1744   monthyear = catx('/', month, year); /*create new variable, use catx, determine separator and
1745   existing variables to use*/
1746 RUN;

NOTE: There were 1554715 observations read from the data set BIKES.BIKESCOMPLETE.
NOTE: The data set BIKES.BIKESCOMPLETE has 1554715 observations and 19 variables.
NOTE: DATA statement used (Total process time):
      real time           3.30 seconds
      cpu time            1.59 seconds

```

H. Utilize conditional logic to create a new variable based on trip duration. For example, if duration is less than 5 minutes, then trip_time = “short trip.”

CODE:

```

/*The following code will use if/then to define a variable based on given conditions
if trip duration is greater than 5 minutes, call it long. Otherwise, call it short

3/3/19
Evan Leet

input: bikescomplete set
output: bikescomplete set with an additional variable
*/

DATA bikes.bikescomplete;
  set bikes.bikescomplete;

  tripdur = 'short'; /*define new variable*/

  if tripduration > 300 then tripdur = 'long'; /*use if then to assign new value to variable based on given condition*/

RUN;

```

LOG:


```

NOTE: DATA statement used (Total process time):
      real time      5.86 seconds
      cpu time       1.45 seconds

1781 /*The following code will use if/then to define a variable based on given conditions
1782 if trip duration is greater than 5 minutes, call it long. Otherwise, call it short
1783
1784 3/3/19
1785 Evan Leet
1786
1787 input: bikescomplete set
1788 output: bikescomplete set with an additional variable
1789 */
1790
1791 DATA bikes.bikescomplete;
1792 set bikes.bikescomplete;
1793
1794 tripdur = 'short'; /*define new variable*/
1795
1796 if tripduration > 300 then tripdur = 'long'; /*use if then to assign new value to variable
1797 based on given conditions*/
1798
1799 RUN;

NOTE: Character values have been converted to numeric values at the places given by: (Line):(Column).
1796:8
NOTE: There were 1554715 observations read from the data set BIKES.BIKESCOMPLETE.
NOTE: The data set BIKES.BIKESCOMPLETE has 1554715 observations and 20 variables.
NOTE: DATA statement used (Total process time):
      real time      7.11 seconds
      cpu time       1.54 seconds

```

I. From this data set, you will create multiple data sets, applying different criteria to each. You will need to create one data set for each breakout of the city (boroughs) represented in the data.

```

/*The following code will combine the bikeboroughs data set with the bikescomplete dataset,
it will give each observation a borough.
Afterward, it will remove observations that are missing a stationid
It will also create a variable ageatuse which is the difference between birthyear and year of use

3/3/19
Evan Leet

input: bikescomplete set and nyc bikeboroughs set
output: bikescomplete set with borough added, other items as above
*/

PROC SORT data=bikes.nycbikeboroughs;
  by station_id;

RUN;

DATA bikes.bikescomplete;
  set bikes.bikescomplete;
  station_id = input(start_station_id, 5.);

RUN;

PROC SORT data=bikes.bikescomplete;
  by station_id;

RUN;

DATA bikes.bikescomplete;
  merge bikes.bikescomplete bikes.nycbikeboroughs;
  by station_id;

RUN;

DATA bikes.bikescomplete;
  set bikes.bikescomplete(where=(station_id ne .));

RUN;

DATA bikes.bikescomplete;
  set bikes.bikescomplete;

  ageatuse = year - birth_year;

RUN;

```

J. Review your work to this point and ensure that you identify and resolve any syntax or logic errors that you have encountered. Provide at least one example of an issue you debugged or justify why no changes were needed.

The code just above is representative of an issue encountered. It was sort of tricky to get the name of the borough alongside the other data. Luckily, we were provided the chart that provided the station ID and the borough each is located in on a separate sheet. I sorted the borough sheet

and the bike data set by station ID and merged them. This gave each observation in the bike data set a value for location. I also had to calculate the customer's age at use based on their birth year and the year they used the Citi Bike service. Luckily, the Year variable was calculated prior to merging.

K. Submit your finalized code that runs successfully without errors and yields appropriate outcomes. Submit your code and log files by copying and pasting into a Microsoft Word document from SAS.

See above and below

CODE:

```
/*The following code will further trim the bikescomplete set to include only the necessary data
| 3/3/19
Evan Leet
input: bikescomplete set and nyc bikeboroughs set
output: bikescomplete set with borough added, other items as above
*/

DATA bikes.bikescompleteref; /*create and name new data set*/
    set bikes.bikescomplete (keep=tripdur bikeid usertype birth_year gender
                             station_id location ageatuse dateformatted year);
    /*Use keep to keep only the required data*/

    if location = "Brooklyn" or location = "Manhattan" then output;
    /*output if location is brooklyn or manhattan only*/

RUN;
```

LOG:

```
2195
2196 /*The following code will further trim the bikescomplete set to include only the necessary data
2197 3/3/19
2198 Evan Leet
2199
2200 input: bikescomplete set and nyc bikeboroughs set
2201 output: bikescomplete set with borough added, other items as above
2202 */
2203
2204 DATA bikes.bikescompleteref; /*create and name new data set*/
2205     set bikes.bikescomplete (keep=tripdur bikeid usertype birth_year gender
2206                             station_id location ageatuse dateformatted year);
2207     /*Use keep to keep only the required data*/
2208
2209     if location = "Brooklyn" or location = "Manhattan" then output;
2210     /*output if location is brooklyn or manhattan only*/
2211
2212 RUN;
2213
NOTE: There were 1554872 observations read from the data set BIKES.BIKESCOMPLETE.
NOTE: The data set BIKES.BIKESCOMPLETEREF has 1551326 observations and 10 variables.
NOTE: DATA statement used (Total process time):
      real time           1.19 seconds
      cpu time            1.15 seconds
```

The CONTENTS Procedure

Data Set Name	BIKES.BIKESCOMPLETEREF	Observations	1551326
Member Type	DATA	Variables	10
Engine	V9	Indexes	0
Created	03/03/2019 20:51:26	Observation Length	80
Last Modified	03/03/2019 20:51:26	Deleted Observations	0
Protection		Compressed	NO
Data Set Type		Sorted	NO
Label			
Data Representation	WINDOWS_64		
Encoding	wlatin1 Western (Windows)		

Engine/Host Dependent Information	
Data Set Page Size	65536
Number of Data Set Pages	1899
First Data Page	1
Max Obs per Page	817
Obs in First Data Page	792
Number of Data Set Repairs	0
ExtendObsCounter	YES
Filename	C:\Users\evan.leet\Documents\bikesheets1\bikescompleteref.sas7bdat
Release Created	9.0401M3
Host Created	X64_BPRO

Alphabetic List of Variables and Attributes						
#	Variable	Type	Len	Format	Informat	Label
6	Year	Num	8			
10	ageatuse	Num	8			
1	bikeid	Char	7	\$7.	\$7.	
3	birth_year	Char	6	\$6.	\$6.	
5	dateformatted	Num	8	DATE9.		
4	gender	Char	3	\$3.	\$3.	
9	location	Char	9	\$9.	\$9.	location
8	station_id	Num	8	BEST.		station id
7	trpdur	Char	5			
2	usertype	Char	12	\$12.	\$12.	

CODE:

```
/*The folowing code will add age categories

Evan Leet
input: bikescompleteref
output: bikescompleteref with agecat
*/

DATA bikes.bikescompleteref;
    set bikes.bikescompleteref;

    agecat = 'unknown'; /*define new variable*/

    /*Set up if then statements to change agecat as necessary*/
    if ageatuse ge 0 and ageatuse le 25 then agecat = '0-25';
    else if ageatuse ge 26 and ageatuse le 35 then agecat = '26-35';
    else if ageatuse ge 36 and ageatuse le 49 then agecat = '36-49';
    else if ageatuse ge 50 then agecat = '50+';
    else if ageatuse = . then agecat = 'unknown'; /*assign unknown if no value*/

RUN;
```

LOG:

```

2240
2241      /*The following code will add age categories
2242
2243      Evan Leet
2244      input: bikescompleteref
2245      output: bikescompleteref with agecat
2246      */
2247
2248      DATA bikes.bikescompleteref;
2249          set bikes.bikescompleteref;
2250
2251          agecat = 'unknown';
2252
2253          if ageatuse ge 0 and ageatuse le 25 then agecat = '0-25';
2254          else if ageatuse ge 26 and ageatuse le 35 then agecat = '26-35';
2255          else if ageatuse ge 36 and ageatuse le 49 then agecat = '36-49';
2256          else if ageatuse ge 50 then agecat = '50+';
2257          else if ageatuse = . then agecat = 'unknown';
2258
2259      RUN;
NOTE: There were 1551326 observations read from the data set BIKES.BIKESCOMPLETEREF.
NOTE: The data set BIKES.BIKESCOMPLETEREF has 1551326 observations and 11 variables.
NOTE: DATA statement used (Total process time):
      real time           0.58 seconds
      cpu time            0.51 seconds

```

III. Understanding the Data In this section, you will write your data summary to explain your findings to stakeholders.

A. Explain how you imported the data into SAS, the id variables used to merge the data, the formats and functions used to transform the data, and the analysis done to summarize the data.

A library for the Citi Bike data was created in SAS. The excel files were converted into SAS files and were incorporated into this library. The date of use for each entry was created and the month and year of use were extracted from this date. Each of the 3 files was sorted by the same variable (bikeid) and then the three files were combined into one dataset. The age at use was calculated for each subscriber based on the year of use and the subscriber's year of birth – a new variable was created and each entry was the result of the calculation of the difference between age at use and birth year. Bike users were further placed in age groups. Unknown entries for age and gender come from non subscribers. Each entry included a station ID and this was used to merge the bike use data set with a key that included the borough in which each station ID is located. This merging process ensured that each entry would have the appropriate borough assigned to it. These three categories - age group, gender, and location -were the basis for the following analyses.

B. Use SAS graphing procedures to create graphs of your outcome data.

CODE:

```

/*The following code will create various graphs

Evan Leet
input: bikescompleteref
output: graphs
*/

title "Total Frequency of Gender"; /*assign title*/

❑ PROC GCHART data=bikes.bikescompleteref; /*Use PROC GCHART, point SAS to the proper file*/
    vbar gender; /*Use vertical bar graph for the variable gaender*/

    RUN;

title "Total Frequency of Age by group";

❑ PROC GCHART data=bikes.bikescompleteref;
    vbar agecat; /*Use vertical bar graph for the variable agecat*/

    RUN;

title "Frequency of Gender By Year and Location";

❑ PROC GCHART data=bikes.bikescompleteref;
    vbar gender /
    subgroup=location /*create a verticle chart of gender, subgrouped by location*/
    GROUP=year;      /*further group results by year*/

    RUN;

title "Frequency of Age Group By Year and Location";

❑ PROC GCHART data=bikes.bikescompleteref;
    vbar agecat /
    subgroup=location
    GROUP=year;

    RUN;

```

LOG:

```

2411 title "Total Frequency of Gender";
2412
NOTE: There were 1551191 observations read from the data set BIKES.BIKESCOMPLETEREF.
NOTE: PROCEDURE GCHART used (Total process time):
      real time      3:35.41
      cpu time       9.85 seconds

2413 PROC GCHART data=bikes.bikescompleteref;
2414 vbar gender;
2415
2416 RUN;
NOTE: 12773 bytes written to C:\Users\EVAN\1.LEE\AppData\Local\Temp\SAS Temporary
Files\TD10132_STEM-SAS199\_gchart10.png.
2417 title "Total Frequency of Age by group";
2418
NOTE: There were 1551191 observations read from the data set BIKES.BIKESCOMPLETEREF.
NOTE: PROCEDURE GCHART used (Total process time):
      real time      10.85 seconds
      cpu time       2.09 seconds

2419 PROC GCHART data=bikes.bikescompleteref;
2420 vbar agecat;
2421
2422 RUN;
NOTE: 13570 bytes written to C:\Users\EVAN\1.LEE\AppData\Local\Temp\SAS Temporary
Files\TD10132_STEM-SAS199\_gchart11.png.
2423 title "Frequency of Gender By Year and Location";
2424
NOTE: There were 1551191 observations read from the data set BIKES.BIKESCOMPLETEREF.
NOTE: PROCEDURE GCHART used (Total process time):
      real time      10.94 seconds
      cpu time       2.04 seconds

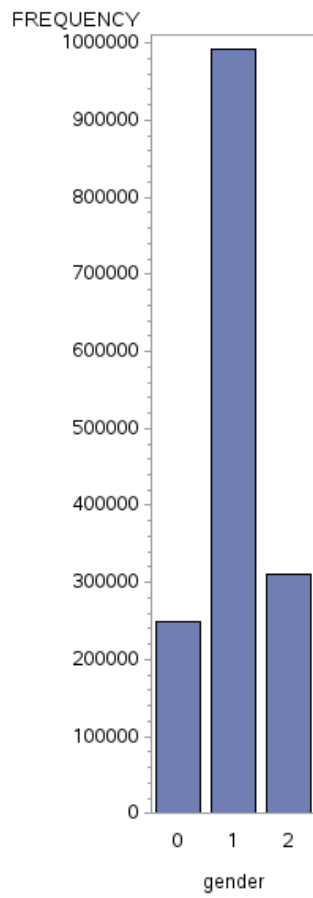
2425 PROC GCHART data=bikes.bikescompleteref;
2426 vbar gender /
2427 subgroup=location
2428 GROUP=year;
2429
2430 RUN;
NOTE: 16195 bytes written to C:\Users\EVAN\1.LEE\AppData\Local\Temp\SAS Temporary
Files\TD10132_STEM-SAS199\_gchart12.png.
2431 title "Frequency of Age Group By Year and Location";
2432
NOTE: There were 1551191 observations read from the data set BIKES.BIKESCOMPLETEREF.
NOTE: PROCEDURE GCHART used (Total process time):
      real time      11.09 seconds
      cpu time       3.17 seconds

2433 PROC GCHART data=bikes.bikescompleteref;
2434 vbar agecat /
2435 subgroup=location
2436 GROUP=year;
2437
2438 RUN;
NOTE: 17499 bytes written to C:\Users\EVAN\1.LEE\AppData\Local\Temp\SAS Temporary
Files\TD10132_STEM-SAS199\_gchart13.png.

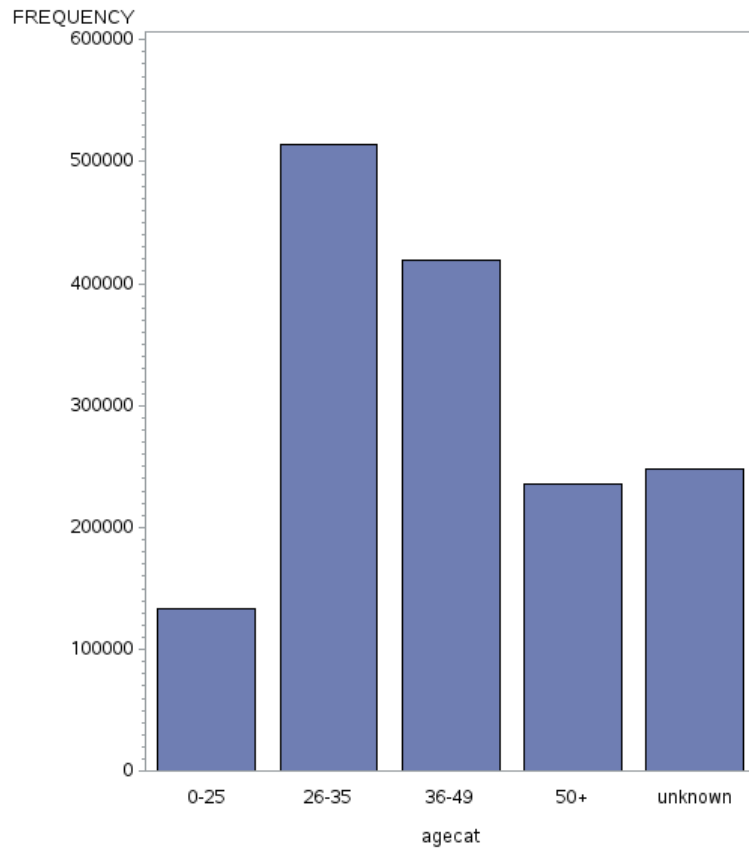
```

RESULTS:

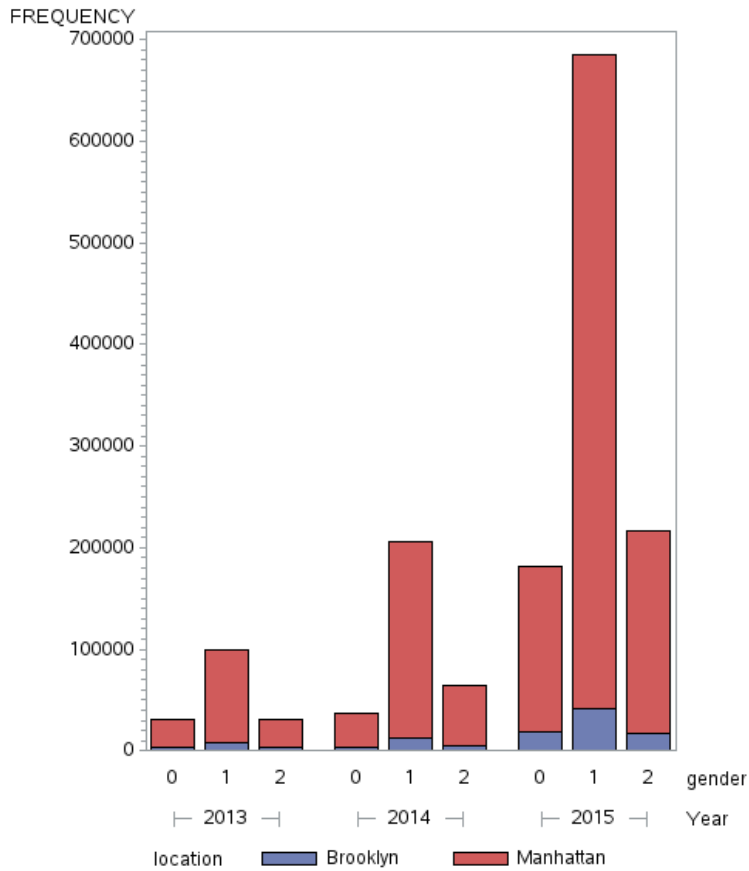
Total Frequency of Gender



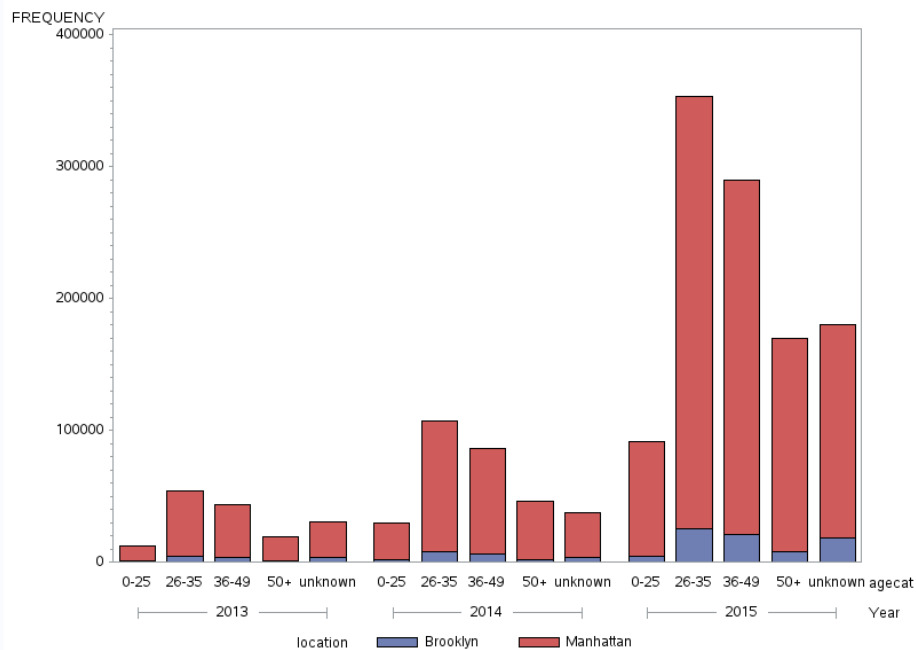
Total Frequency of Age by group

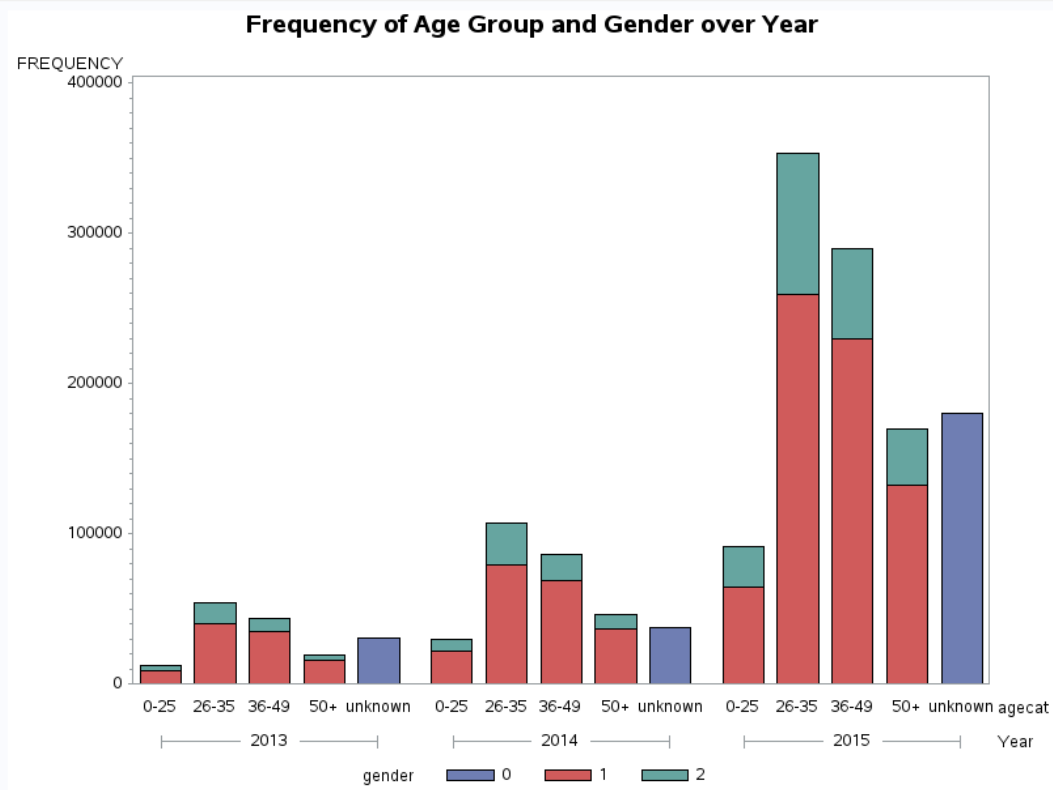


Frequency of Gender By Year and Location

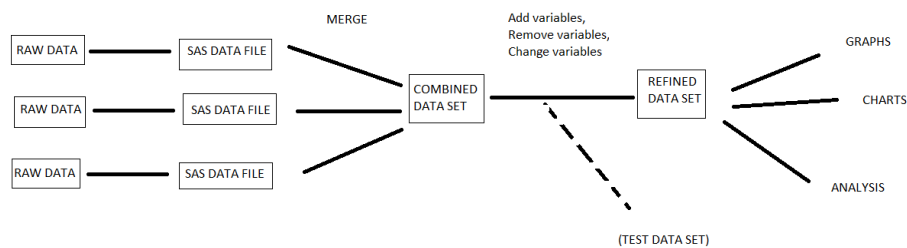


Frequency of Age Group By Year and Location





C. Create a visual illustration of the data workflow, from import to multiple data set creation.



D. Discuss the errors or difficulties encountered during the coding process, including how you resolved these issues.

The majority of my difficulties came from creating and manipulating the data sets. I frequently created duplicate data sets so that I could experiment with them and not worry about ruining something important. In terms of graphs, I tried to learn from my classmates and included a graph which I learned from the discussion assignment in previous weeks.

E. Explain what the data tells you about the utilization of the bike-sharing program. Support your response with examples from the data. This support could be presented in table form but should come from the descriptive analysis.

F. Summarize what you learned about the bike-sharing program and its users. Support your response with examples from the data. (Remember this should be the final summary of the outcomes discovered from the analysis and should be directed toward a professional audience.)

In looking at the graphs created by the data set, it looks like use of the Citi Bikes has been increasing year to year. This is true among genders and across the age categories. Usage in Manhattan is greater than usage in Brooklyn, and usage in Manhattan has also been increasing at a much greater rate. It seems that males are more commonly users than females, and the most frequent age group is the 25-36 group. Usage among males in the 25-36 year old age group is increasing, and usage among females in this category is increasing as well. In fact, usage among females in general appears to have tripled or more from 2014 to 2015.

In terms of location, Manhattan locations are clearly more frequently used. Growth in this location is also more apparent. In both locations and across all years, males are the most common gender group and 25-36 is the most common age group. Usage in Brooklyn is increasing over the years (21,496 users in 2014 to 77,754 users in 2015), but the numbers from Manhattan are far greater – 285,752 users in 2014 to 1,005,661 in 2015. The following charts further demonstrate these findings:

3 Way Chart of Gender and Location by Year			
The FREQ Procedure			
Frequency Percent Row Pct Col Pct	Table 1 of gender by location		
	Controlling for Year=2013		
	location(location)		
gender	Brooklyn	Manhattan	Total
0	3473 2.16 11.36 25.53	27101 16.87 88.64 18.43	30574 19.03
1	7368 4.59 7.40 54.16	92168 57.38 92.60 62.69	99536 61.97
2	2762 1.72 9.05 20.30	27756 17.28 90.95 18.88	30518 19.00
Total	13603 8.47	147025 91.53	160628 100.00

Frequency Percent Row Pct Col Pct	Table 2 of gender by location		
	Controlling for Year=2014		
	location(location)		
gender	Brooklyn	Manhattan	Total
0	4094 1.33 10.97 19.05	33227 10.81 89.03 11.63	37321 12.15
1	12498 4.07 6.06 58.14	193828 63.09 93.94 67.83	206326 67.15
2	4904 1.60 7.71 22.81	58697 19.10 92.29 20.54	63601 20.70
Total	21496 7.00	285752 93.00	307248 100.00

Frequency Percent Row Pct Col Pct	Table 3 of gender by location		
	Controlling for Year=2015		
	location(location)		
gender	Brooklyn	Manhattan	Total
0	18669 1.72 10.31 24.01	162354 14.99 89.69 16.15	181023 16.71
1	42022 3.88 6.13 54.04	643760 59.43 93.87 64.02	685782 63.30
2	17063 1.58 7.88 21.94	199447 18.41 92.12 19.83	216510 19.99
Total	77754 7.18	1005561 92.82	1083315 100.00

3 Way Chart of Age Category and Location by Year

The FREQ Procedure

Frequency Percent Row Pct Col Pct	Table 1 of agecat by location			
	Controlling for Year=2013			
	location(location)			
	agecat	Brooklyn	Manhattan	Total
0-25		795	11489	12284
		0.49	7.15	7.65
		6.47	93.53	
		5.84	7.81	
26-35		4705	49446	54151
		2.93	30.78	33.71
		8.69	91.31	
		34.59	33.63	
36-49		3427	40614	44041
		2.13	25.28	27.42
		7.78	92.22	
		25.19	27.62	
50+		1204	18378	19582
		0.75	11.44	12.19
		6.15	93.85	
		8.85	12.50	
unknown		3472	27098	30570
		2.16	16.87	19.03
		11.36	88.64	
		25.52	18.43	
Total		13603	147025	160628
		8.47	91.53	100.00

Frequency Percent Row Pct Col Pct	Table 2 of agecat by location			
	Controlling for Year=2014			
	location(location)			
	agecat	Brooklyn	Manhattan	Total
0-25		1568	28608	30176
		0.51	9.31	9.82
		5.20	94.80	
		7.29	10.01	
26-35		7796	99614	107410
		2.54	32.42	34.96
		7.26	92.74	
		36.27	34.86	
36-49		5879	80270	86149
		1.91	26.13	28.04
		6.82	93.18	
		27.35	28.09	
50+		2173	44059	46232
		0.71	14.34	15.05
		4.70	95.30	
		10.11	15.42	
unknown		4080	33201	37281
		1.33	10.81	12.13
		10.94	89.06	
		18.98	11.62	
Total		21496	285752	307248
		7.00	93.00	100.00

Frequency Percent Row Pct Col Pct	Table 3 of agecat by location			
	Controlling for Year=2015			
	location(location)			
	agecat	Brooklyn	Manhattan	Total
0-25		4626	86472	91098
		0.43	7.98	8.41
		5.08	94.92	
		5.95	8.60	
26-35		25527	327436	352963
		2.36	30.23	32.58
		7.23	92.77	
		32.83	32.56	
36-49		21100	268418	289518
		1.95	24.78	26.73
		7.29	92.71	
		27.14	26.69	
50+		7939	161817	169756
		0.73	14.94	15.67
		4.68	95.32	
		10.21	16.09	
unknown		18562	161418	179980
		1.71	14.90	16.61
		10.31	89.69	
		23.87	16.05	
Total		77754	1005561	1083315
		7.18	92.82	100.00