```sas
1   ******************************
2       Exam 2
3       Name: Gavin Frias
4       Version: 1
5   ******************************;
6
7
8   ******************************
9   ******** Task 1: DATA *********
10  ******************************;
11
12  /* Question 1: Import Data */
13  TITLE 'Task1 Q1: Import Data';
14
15  %web_drop_table(WORK.IMPORT);
16
17
18  FILENAME REFFILE '/home/u61397358/sasuser.v94/coaster1.csv';
19
20  PROC IMPORT DATAFILE=REFFILE
21      DBMS=CSV
22      OUT=coaster1;
23      GETNAMES=YES;
24  RUN;
25
26  PROC CONTENTS DATA=coaster1; RUN;
27
28
29  %web_open_table(WORK.IMPORT);
30
31  /* Question 2: Remove the rows that contain missing data (see PDF for column) */
32  TITLE 'Task1 Q2: Remove Missing Data';
33
34  DATA Coaster1_Task1;
35      SET Coaster1;
36      IF Drop = . THEN DELETE;
37  RUN;
38
39  /* Question 3: Create a new character variable */
40  TITLE 'Task1 Q3: Create Character Variable';
41
42  DATA Coaster1_Task1;
43      SET Coaster1;
44      LENGTH LengthGroup $6.;
45      IF Length<2500 THEN LengthGroup="Short";
46      IF Length>=2500 AND Length<4000 THEN LengthGroup="Medium";
47      IF Length>=4000 THEN LengthGroup="Long";
48  RUN;
49
50  /* Question 4: Create a new variable called Ratio */
51  TITLE 'Task1 Q4: Create Ratio';
52
53  DATA Coaster1_Task1;
54      SET Coaster1;
55      Ratio=Height/Drop;
56  RUN;
57
58
59  /* Question 5: Create a New Dataset called High_Ratio and Print it */
60  TITLE 'Task1 Q5: Create Dataset High_Ratio';
61
62  DATA High_Ratio;
63      SET Coaster1_Task1;
64      WHERE Ratio>1.15;
65      KEEP Track Height Drop Length;
66  RUN;
67
68
69  ********************************************
70  ******** Task 2: INTRODUCTORY ANALYSIS *********
71  ********************************************;
72
73  /* Question 6: Compute values of sample mean / median / std dev / IQR
74      / # Observations / # Missing */
75  TITLE 'Task2 Q6: Summary Statistics';
76
```

```
77  PROC SORT DATA=Coaster1; by Duration; RUN;
78  PROC MEANS DATA=Coaster1 MEAN MEDIAN STDDEV QRANGE N NMISS;
79  by Duration;
80  RUN;
81
82
83  /* Question 7: Histogram with density kernel */
84  TITLE 'Task2 Q7: Histogram with Density Kernel';
85
86  PROC SGPLOT DATA=Coaster1;
87      HISTOGRAM Height;
88      DENSITY Height / type=kernel;
89  RUN;
90
91
92  /* Question 8: Bar Chart */
93  TITLE 'Task2 Q8: Bar Chart';
94
95  PROC SGPLOT DATA=Coaster1;
96      VBAR SpeedGroup;
97  RUN;
98
99
100 /* Question 9: Boxplot */
101 TITLE 'Task2 Q9: Boxplot';
102 /* CODE */
103
104 PROC SGPLOT DATA=Coaster1;
105     HBOX Drop;
106 RUN;
107
108
109 /*
110 Are there outliers?
111 Yes, there is at least one outlier to the far right based on the box plot.
112 */
113
114
115 **********************************************
116 ******** Task 3: INFERENCE ********************
117 **********************************************;
118
119 TITLE 'Task3 Q10, Q11: Inference';
120     /* CODE */
121 proc ttest data=Coaster1 h0=0 sides=2 ALPHA=0.017 plots;
122 var Length;
123 run;
124
125 proc ttest data=Coaster1 ho=-800 sides=u ALPHA=0.017 plots;
126 var Length;
127 RUN;
128
129     /* Question 10: Equal Variance Test */
130     /*  Hypotheses
131             H0: Steel Tracks - Wood Tracks = 0
132             H1: Steel Tracks - Wood Tracks != 0
133         Test Statistic: 23.31
134         P-Value: <0.0001
135         Decision: Reject H0
136         Conclusion: There is enough evidence to suggest a difference in length between wood and steel track roller coast
137     */
138
139
140     /* Question 11: Mean Testing */
141     /*  Hypotheses
142             H0: Steel Tracks - Wood Tracks = -800
143             H1: Steel Tracks - Wood Tracks < -800
144         Test Statistic: 30.23
145         P-Value: <0.0001
146         Decision: Reject H0
147         Conclusion: There is enough evidence to suggest that the mean length of Steel Tracks - Wood Tracks is less than
148     */
149
150
151 **********************************************
152 ******** Task 4: REGRESSION *******************
153 **********************************************;
```

```
154  TITLE 'Task4 Q12: Multiple Linear Regression';
155  /* CODE */
156
157  PROC REG DATA=Coaster1 ALPHA=0.04 ;
158      MODEL Duration = Length Type / corrb;
159  RUN;
160  /*
161
162
163  Part a - Check model assumptions
164      Linearity
165          Graph / results looked at: Plot of residuals vs Length and Type.
166          Is the linearity condition met or not? Yes.
167
168      Normality
169          Graph / results looked at: Plots of residual vs quantile and percent vs residual.
170          Is the normality of residuals condition met or not? Yes
171
172      Equal Variance
173          Graph / results looked at: Plot of residual vs predicted value
174          Is the equal variance of residuals condition met or not? Yes.
175
176
177  Part b - Give the equation of the Multiple Linear Regression line
178
179  Duration = B0+B1Length+B2Type
180  Y = 45.15060 + 0.02386Length + 12.28970Type
181
182
183  Part c - Does the model in total explain variability in Duration?
184      Hypotheses
185          H0: beta_length = beta_type = 0
186          H1: beta_length = beta_type != 0
187      Test Statistic: 146.85
188      P-Value: <0.0001
189      Decision: Reject H0
190      Conclusion: There is enough evidence to suggest that at least one variable explains the variability in Duration.
191
192
193  Part d (If needed. If not needed, state why.)
194
195      Testing Individual Variables (Variable 1)
196          Hypotheses
197              H0: beta_length = 0
198              H1: beta_length != 0
199          Test Statistic: 16.89
200          P-Value: <0.001
201          Decision: Reject H0
202          Conclusion: There is enough evidence to suggest that Length explains some variability in Duration.
203
204
205      Testing Individual Variables (Variable 2)
206          Hypotheses
207              H0: beta_type = 0
208              H1: beta_type != 0
209          Test Statistic: 1.95
210          P-Value: 0.0532
211          Decision: Do Not Reject H0
212          Conclusion: There is not enough evidence to suggest that Type explains some variability in Duration.
213
214
215  Part e - Value of R^2 and interpretation
216      R^2: 0.6835
217      Interpretation: We can interpret this as 68.35% of the variability observed in Duration is explained by the model.
218  */
219
220
221  **********************************************
222  ******** Task 5: 1-way ANOVA ******************
223  ***********************************************;
224  TITLE 'Task5 Q13: 1-Way ANOVA';
225
226  TITLE2 'Part a: Mean Duration for each Group';
227      /* CODE */
228  PROC MEANS; CLASS SpeedGroup;
229  RUN;
230
```

```sas
231      /* Detail any difference by group.
232 Some differences to note are that the Fast rollercoasters have a higher mean duration
233 and also have the highest Duration of any rollercoaster.
234      */
235
236
237 TITLE2 'Part b: Side by Side Boxpots';
238      /* CODE */
239 PROC SGPLOT DATA=Coaster1;
240     HBOX Duration / Category=SpeedGroup;
241 RUN;
242
243
244      /* Detail any difference by group.
245 The boxplot for the Fast variable has several outliers, while the Middle variable had the widest interval.
246 Another thing to note is that the Small variable boxplot seemed to be the most normal.
247      */
248
249
250 TITLE2 'Part c: Run a 1-way ANOVA model';
251
252 PROC GLM DATA=Coaster1 ALPHA=0.015;
253 CLASS SpeedGroup;
254 MODEL Duration = SpeedGroup;
255 MEANS SpeedGroup / BON CLDIFF HOVTEST=LEVENE;
256 OUTPUT OUT = ANOVA13 r = residual;
257 RUN;
258
259 TITLE2 'Part d: Normality Test';
260 /* Will you test the normality assumption using the overall dataset, or for each group individually?
261 The overall dataset. */
262
263
264 /* CODE, if needed */
265 PROC UNIVARIATE NORMAL PLOT DATA=Coaster1 ALPHA=0.015;
266 VAR Duration;
267 RUN;
268
269
270 /* Conclusion(s): The data passes the normality check. Shapiro-Wilk = 0.0955 which is greater than 0.05. */
271
272
273 TITLE2 'Part e: Equal Variance Assumption Check';
274 /* Conclusion: The data passes the equal variance check.*/
275
276
277 TITLE2 'Part f: Is there a significant evidence of an effect?';
278 /*  Hypotheses
279        H0: = 0
280        H1: != 0
281     Test Statistic: 23.22
282     P-Value: <0.0001
283     Decision: Reject H0
284     Conclusion: There is enough evidence to suggest that Speed Group explains some variability in Duration.
285 */
286
287
288 TITLE2 'Part g: Bonerroni or Tukey';
289 /* Are you providing Bonferroni or Tukey Intervals?
290 Bonferroni Intervals */
291
292
293 /* Provide confidence intervals for each difference
294     (make sure to indicate the difference you are writing a confidence interval for):
295 Fast - Middle (5.348,55.575)
296 Fast - Slow (33.885,83.587)
297 Middle - Fast (-55.575,-5.348)
298 Middle - Slow (6.655,49.894)
299 Slow - Fast (-83.587,-33.885)
300 Slow - Middle (-49.894,-6.655)
301 */
302
303
304 /* For each pair, state whether the difference is significant or not
305 According to my output the difference of each of these confidence intervals are significant.
306 */
307
```

```
308
309 TITLE;
```