

Appendix A

Exploratory Data Analysis

Raw distributions before preprocessing.

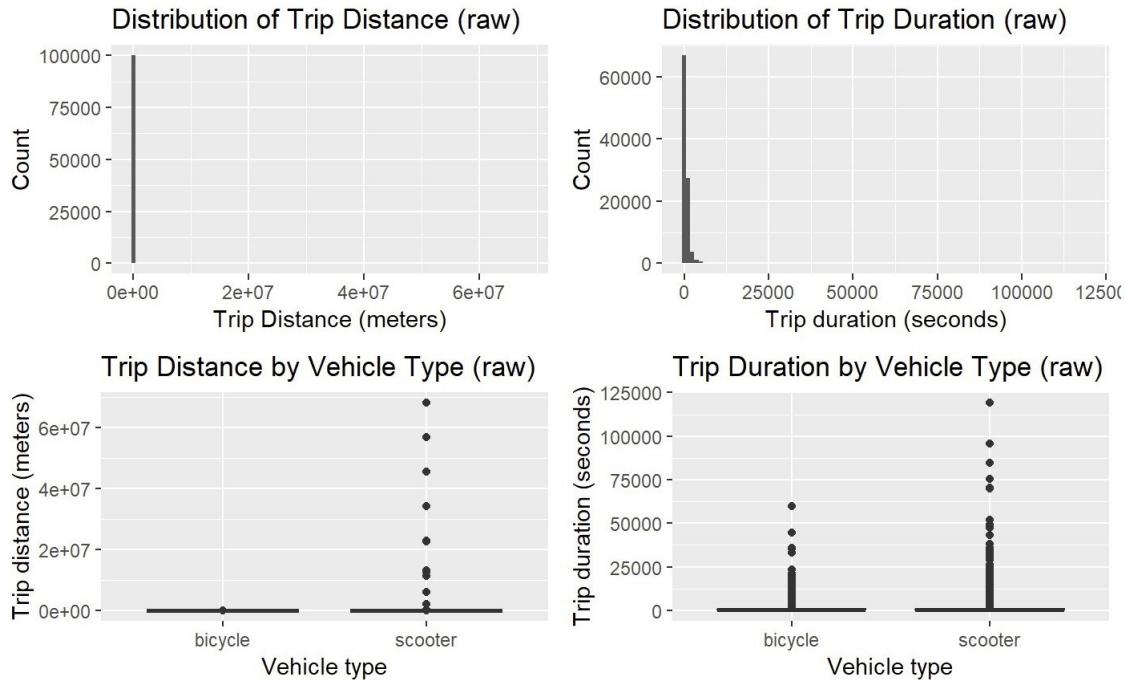


Figure A1

Raw-scale exploratory plots for trip distance and duration (with outliers).

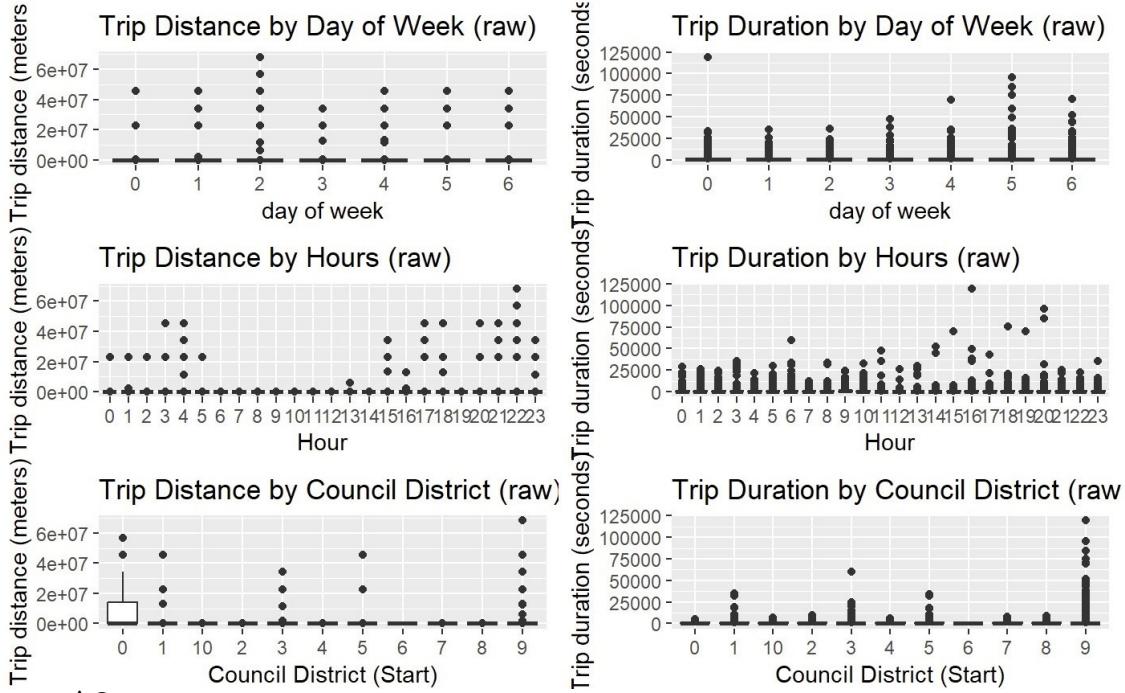


Figure A2

Additional raw-scale exploratory plots (with outliers).

Raw distributions after preprocessing (raw scale).

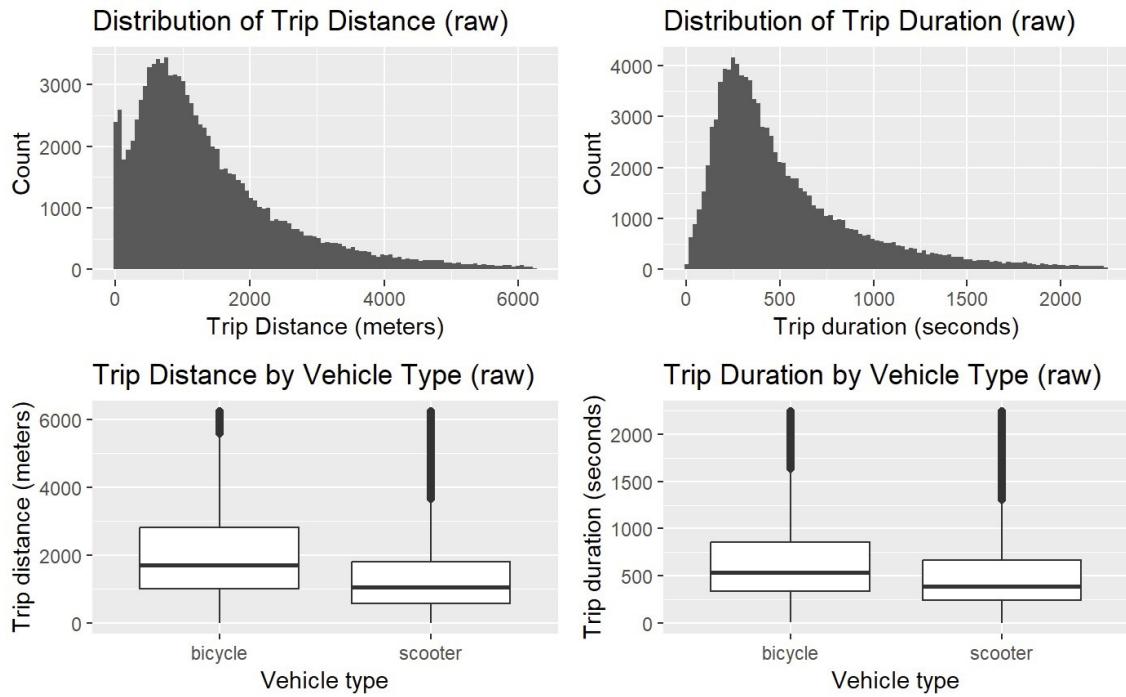


Figure A3

Raw-scale distribution plots after IQR-based trimming.

n	mean_dist	median_dist	max_dist	mean_dur	median_dur	max_dur
<int>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>
5486004	1383.348	1078	6243	525.632	400	2243

Figure A4

Summary of data for raw-scale variables after IQR-based trimming.

Measure	Q1	Q3	IQR	Lower_1_5	Upper_1_5	Lower_3	Upper_3
<chr>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>
Trip Distance	611	2019	1408	0	4131	0	6243
Trip Duration	251	749	498	0	1496	0	2243

Figure A5

IQR-based trimming thresholds (range visualization).

Distributions after transformation (log scale).

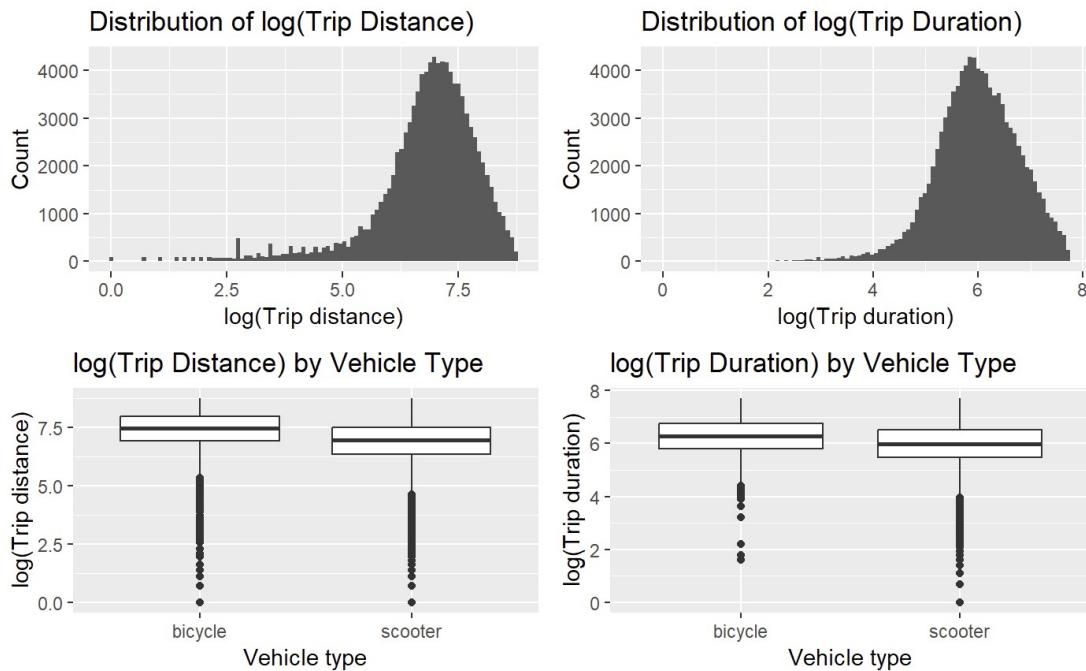


Figure A6

Log-scale exploratory plots for transformed trip distance and duration.

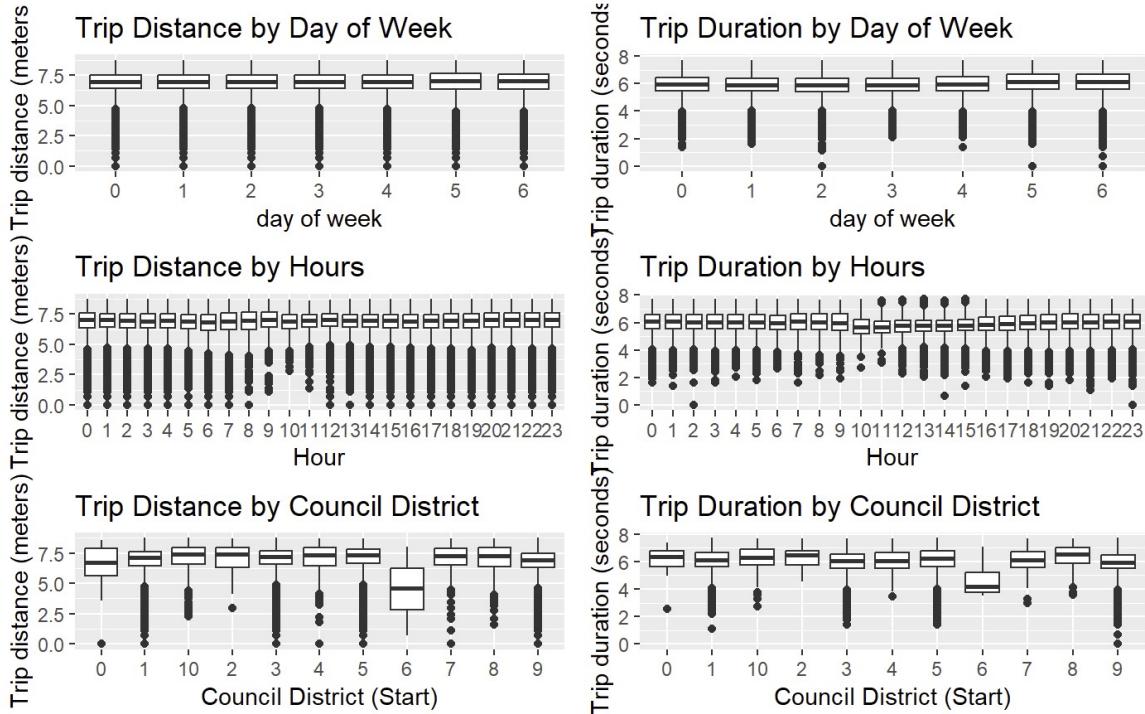


Figure A7

Additional log-scale exploratory plots (log scale).

Appendix B

Model Diagnostics: Question 1 (Trip Distance)

Raw distance model (no transformation).

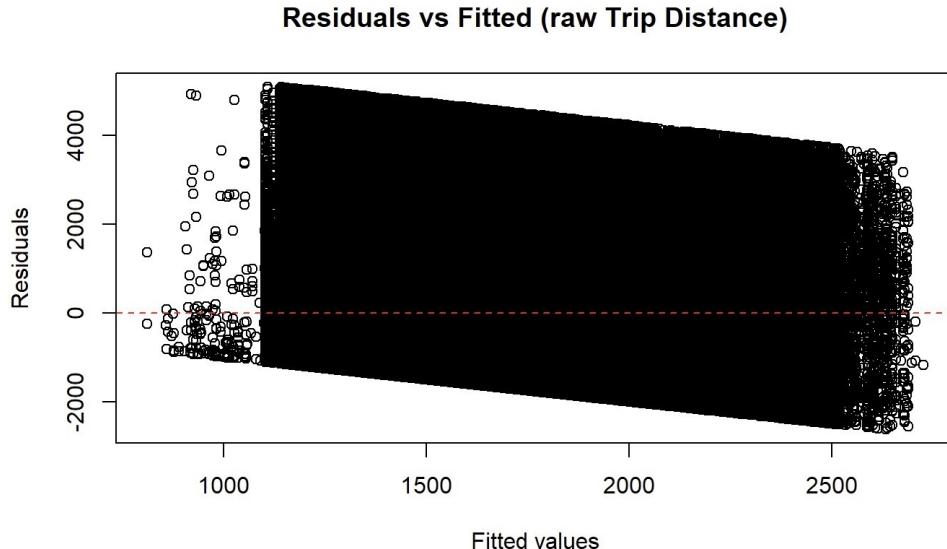


Figure B1

Residuals vs. fitted values for the raw-distance OLS model.

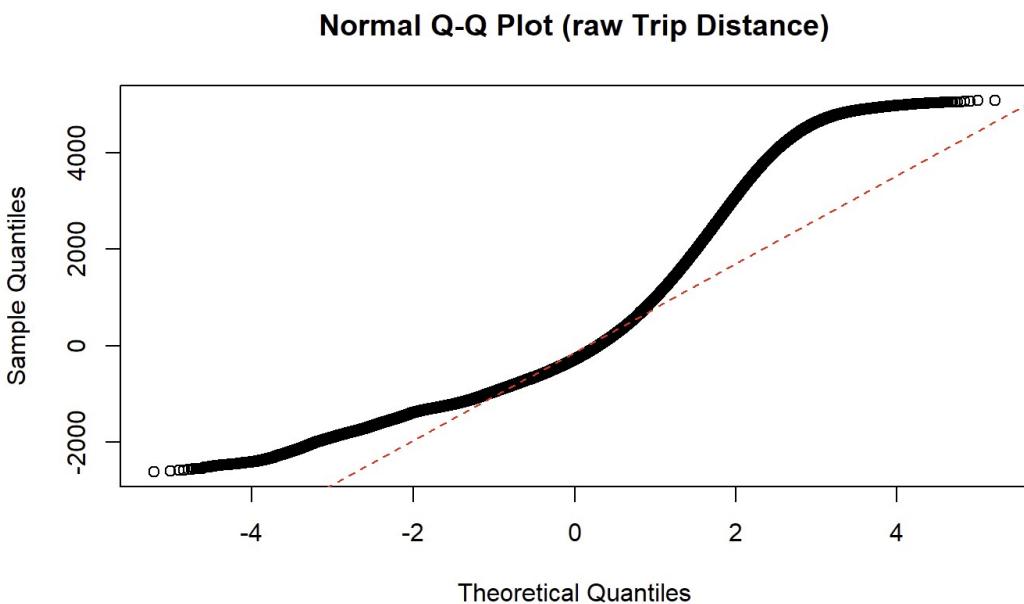


Figure B2

Normal Q-Q plot for the raw-distance OLS model residuals.

Log-transformed distance model.

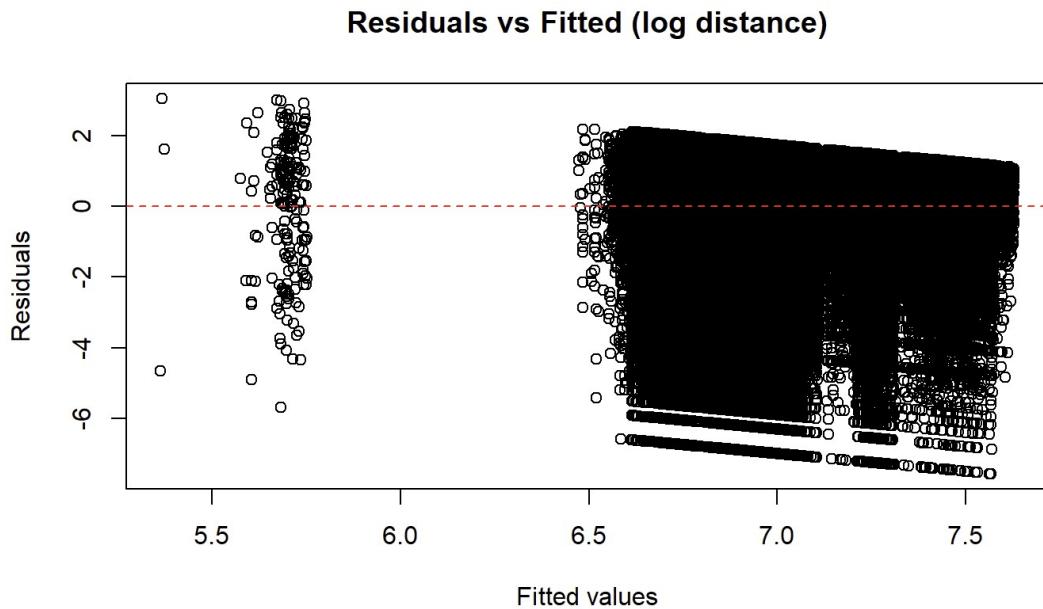


Figure B3

Residuals vs. fitted values for the log-distance OLS model.

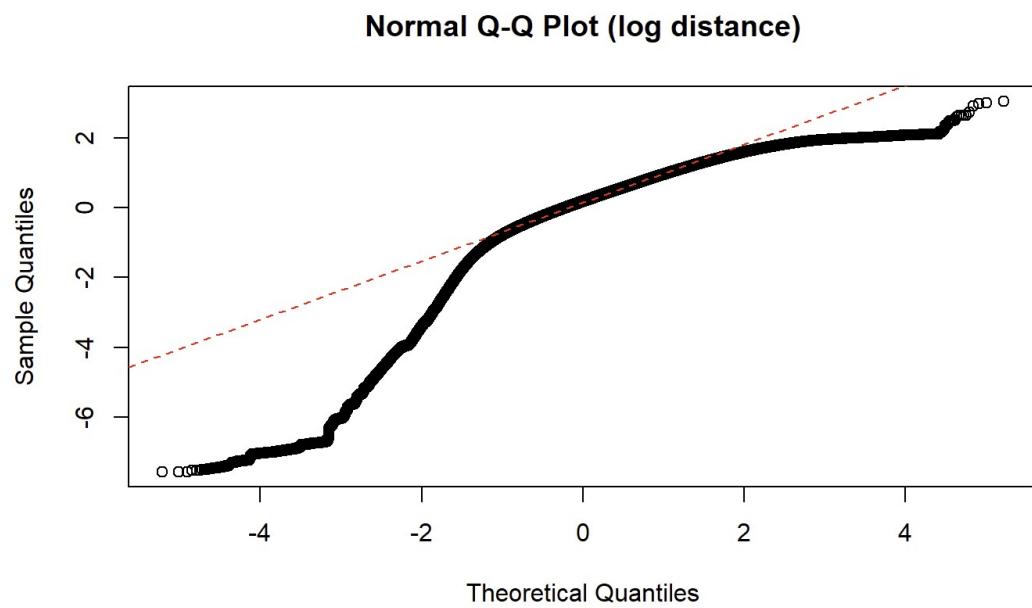


Figure B4

Normal Q-Q plot for the log-distance OLS model residuals.

Robustness checks.

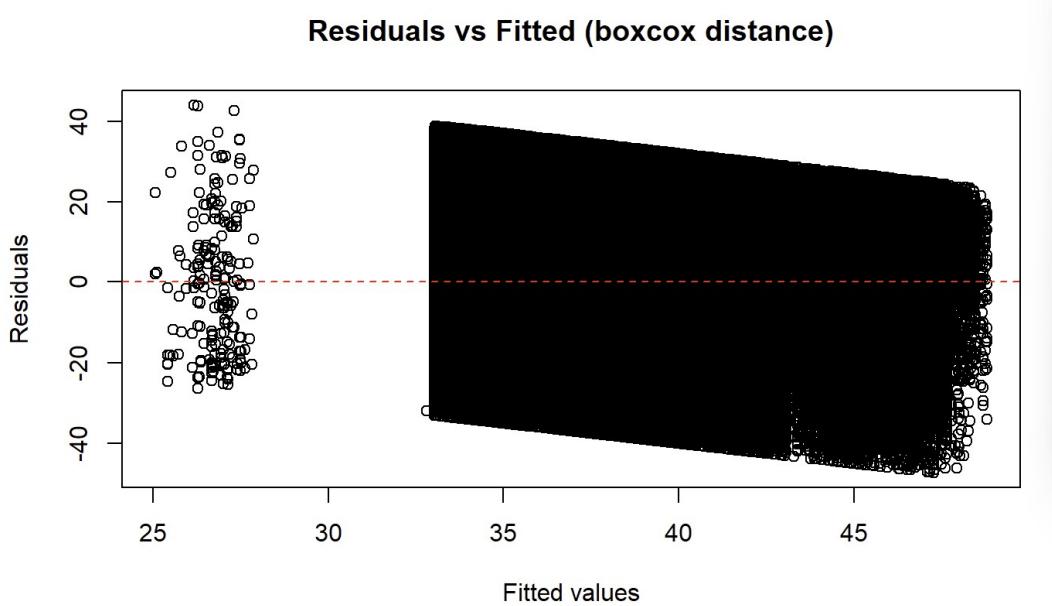


Figure B5

Residuals vs. fitted values for the Box–Cox distance model.

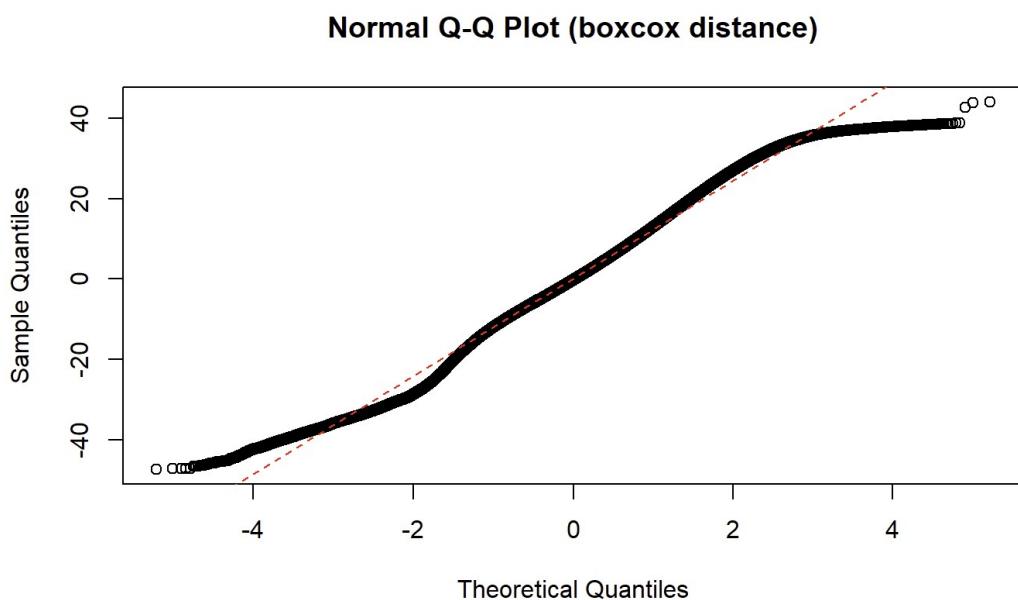


Figure B6

Normal Q–Q plot for the Box–Cox distance model residuals.

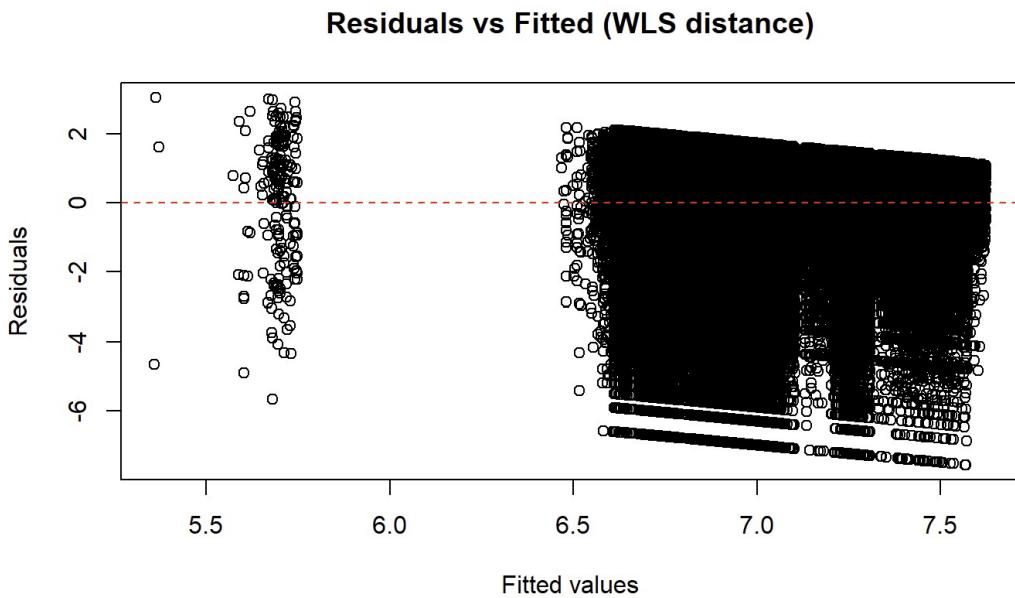


Figure B7

Residuals vs. fitted values for the WLS distance model.

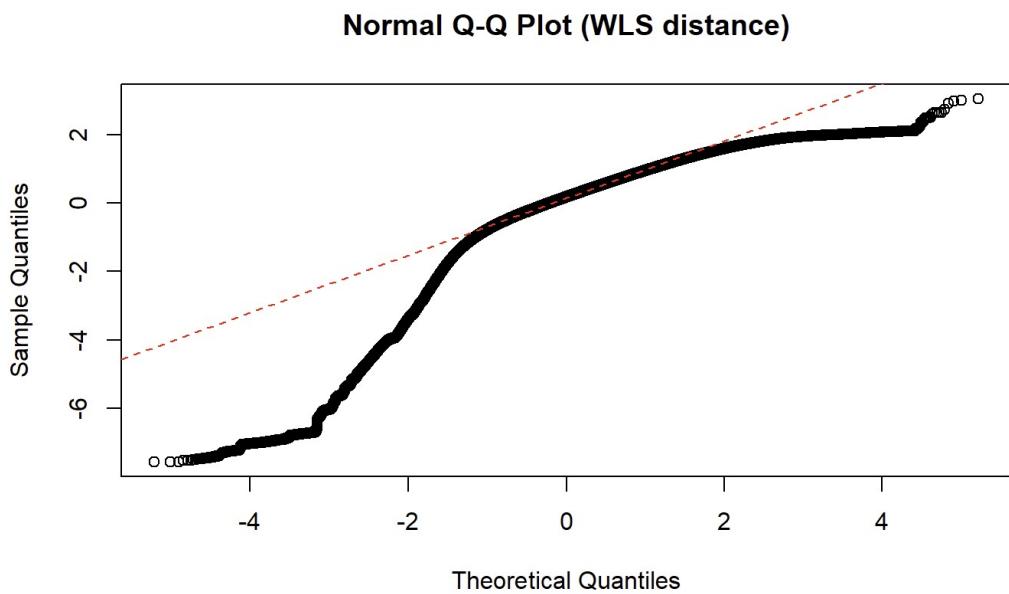


Figure B8

Normal Q-Q plot for the WLS distance model residuals.

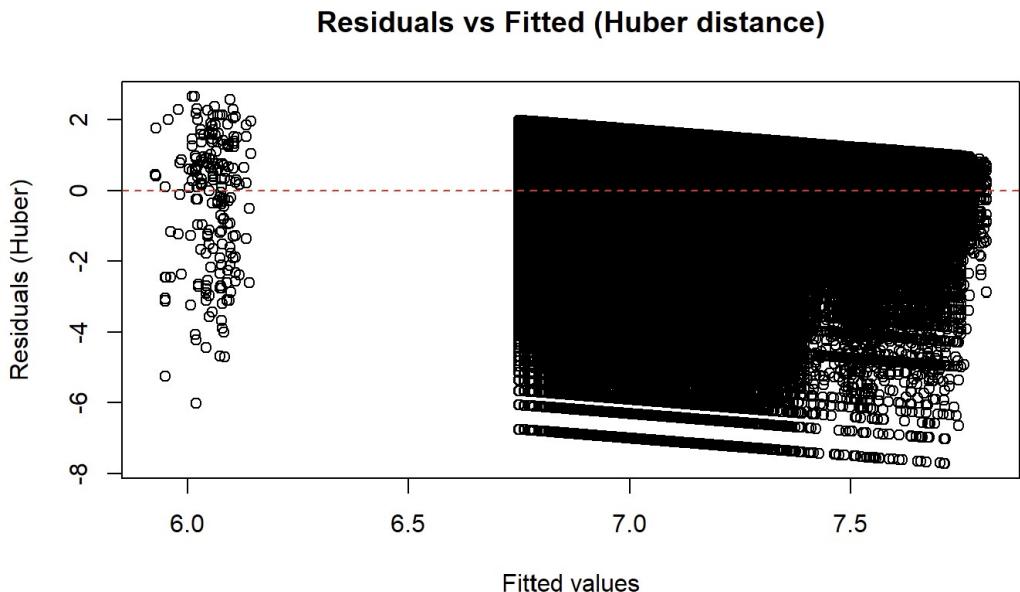


Figure B9

Residuals vs. fitted values for the Huber distance model.

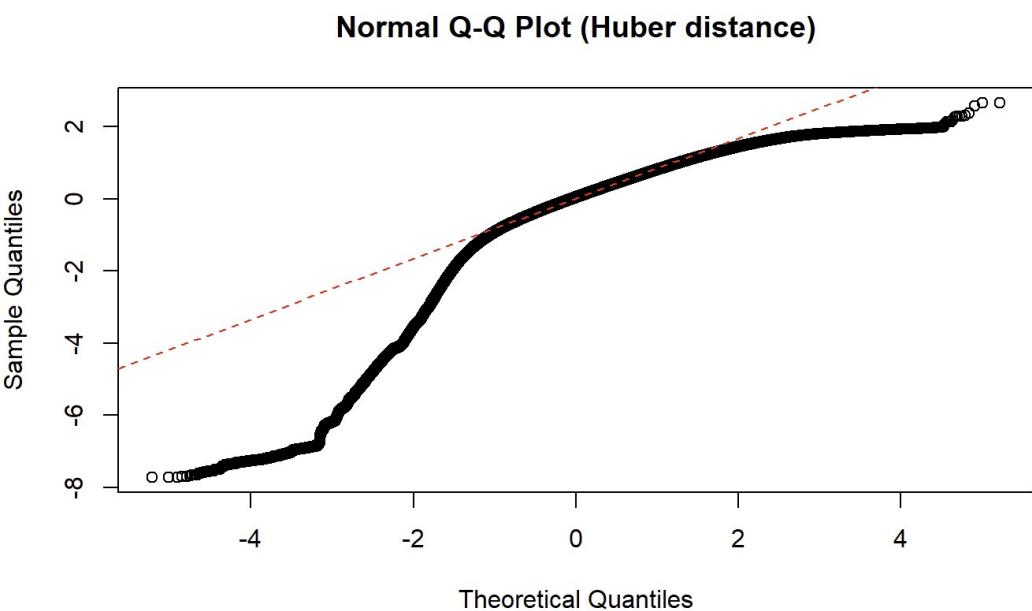


Figure B10

Normal Q-Q plot for the Huber distance model residuals.

	Vehicle <chr>	OLS <dbl>	Boxcox <dbl>	WLS <dbl>	Huber <dbl>
vehicle_type	bicycle (baseline)	0.0000000	0.0000000	0.0000000	0.0000000
scooter	-0.5183198	-7.217821	-0.5190153	-0.4542413	

Figure B11

Key coefficient comparison across distance models.

	GVIF	Df	GVIF^(1/(2*Df))
vehicle_type	1.009107	1	1.004543
hour_f	1.058705	23	1.001241
dow_f	1.054343	6	1.004420
district_start_f	1.018104	11	1.000816

Figure B12

GVIF diagnostics for the distance model predictors.

Model <chr>	BIC <dbl>
OLS_raw	92487424
OLS_log	17005040
BoxCox_log	43809476
WLS_log	17019218
Huber_log	17099228

Figure B13

BIC comparison across candidate distance models.

term <chr>	estimate <dbl>	std_error <dbl>	ratio <dbl>	pct_change <dbl>
Scooter (vs bicycle)	-0.5183198	0.002148842	0.5955203	-40.44797

Figure B14

Results of Question 1.

Appendix C

Model Diagnostics: Question 2 (Trip Duration)

Raw duration model and log-transformed model.

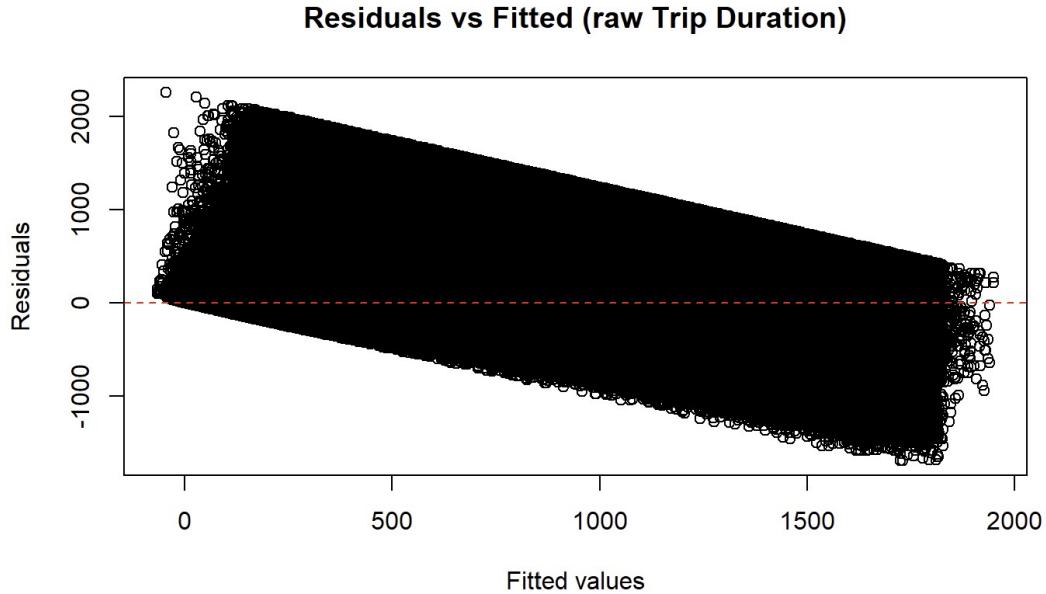


Figure C1

Residuals vs. fitted values for the raw-duration OLS model.

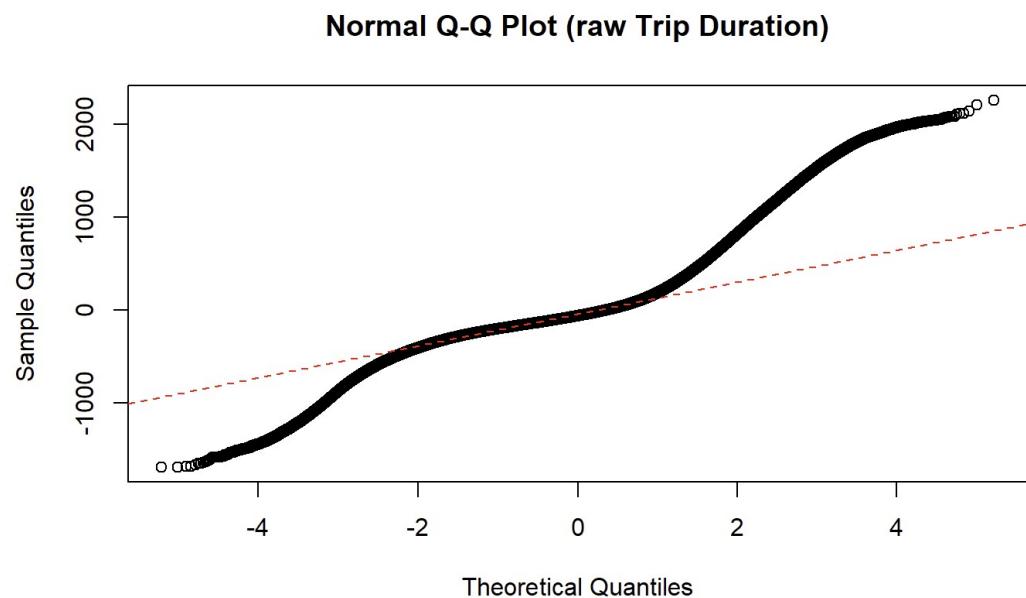


Figure C2

Normal Q-Q plot for the raw-duration OLS model residuals.

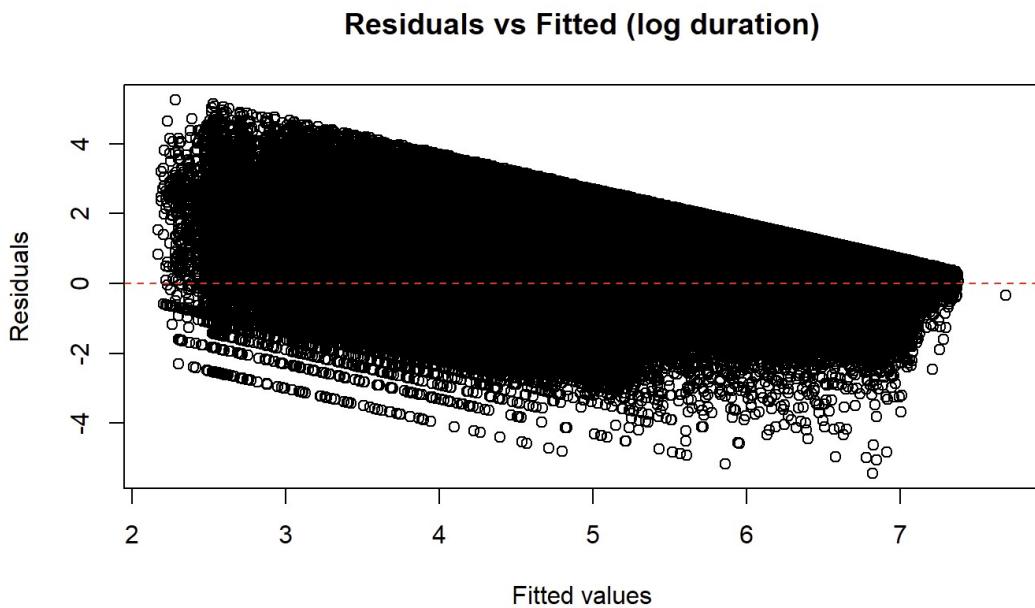


Figure C3

Residuals vs. fitted values for the log-duration OLS model.

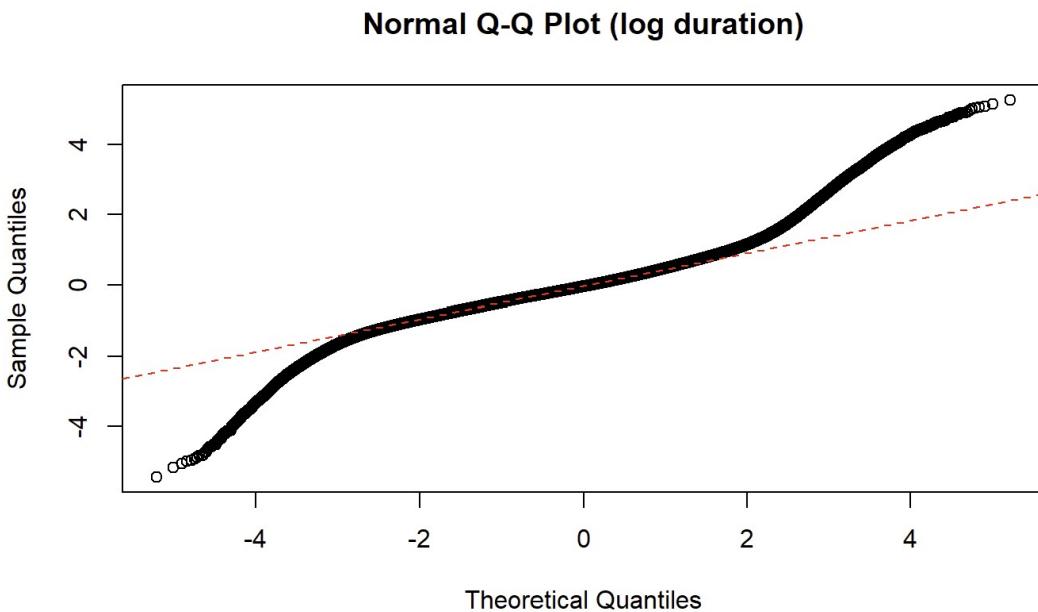


Figure C4

Normal Q-Q plot for the log-duration OLS model residuals.

Robustness checks and interaction model.

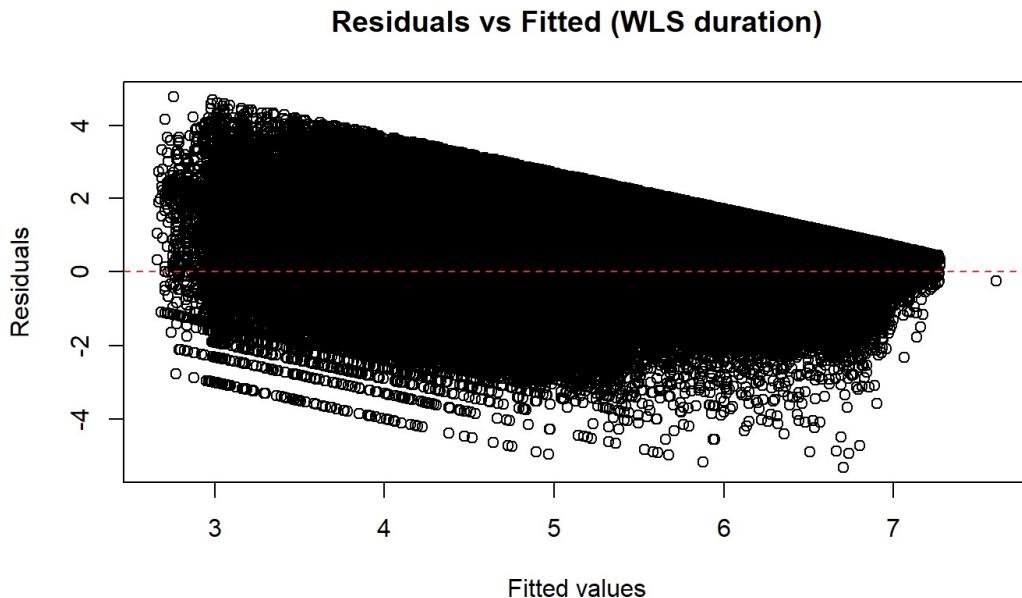


Figure C5

Residuals vs. fitted values for the WLS duration model.

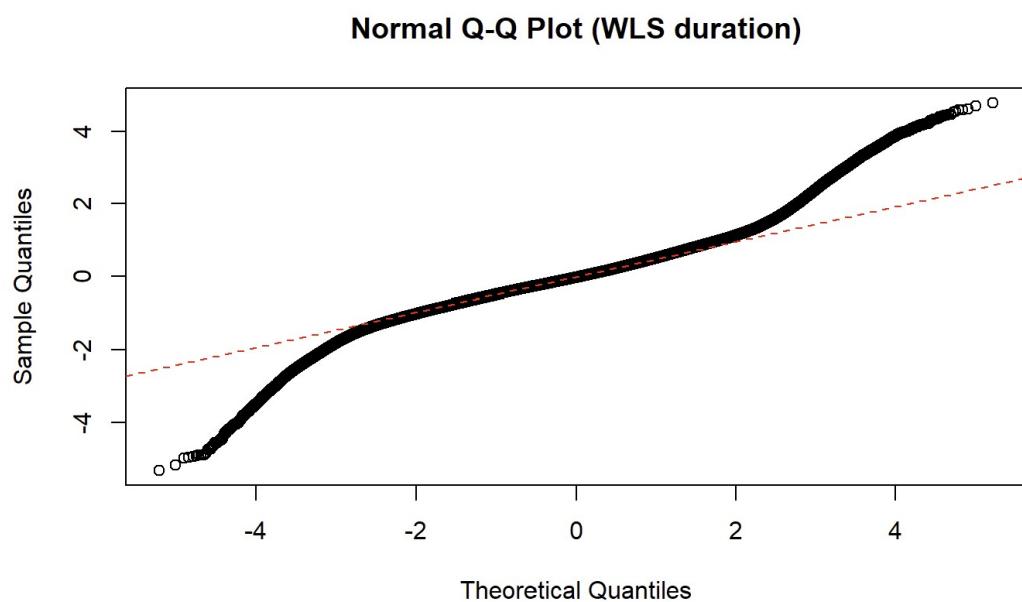


Figure C6

Normal Q-Q plot for the WLS duration model residuals.

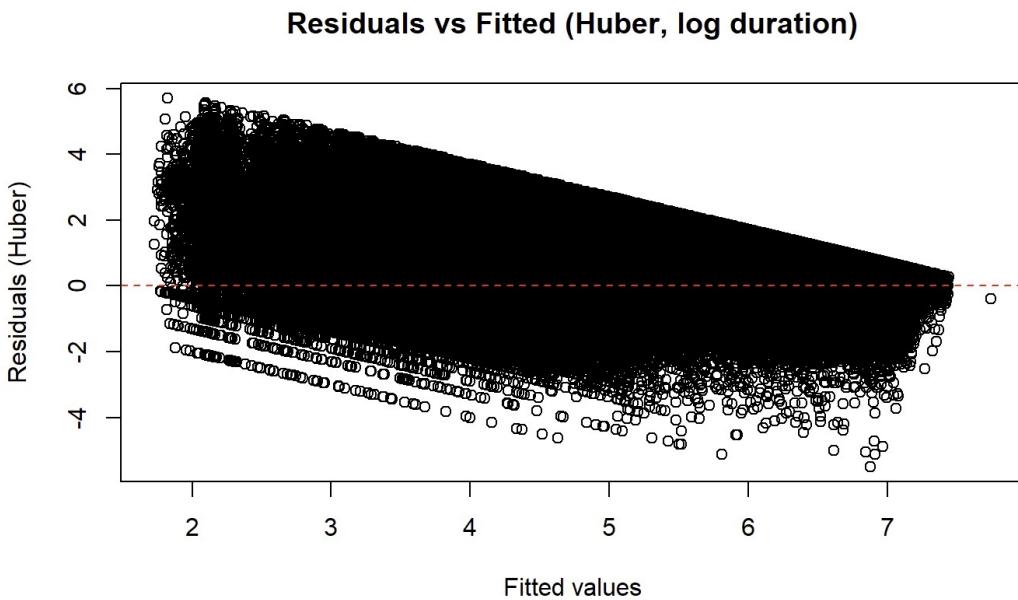


Figure C7

Residuals vs. fitted values for the Huber (log-duration) model.

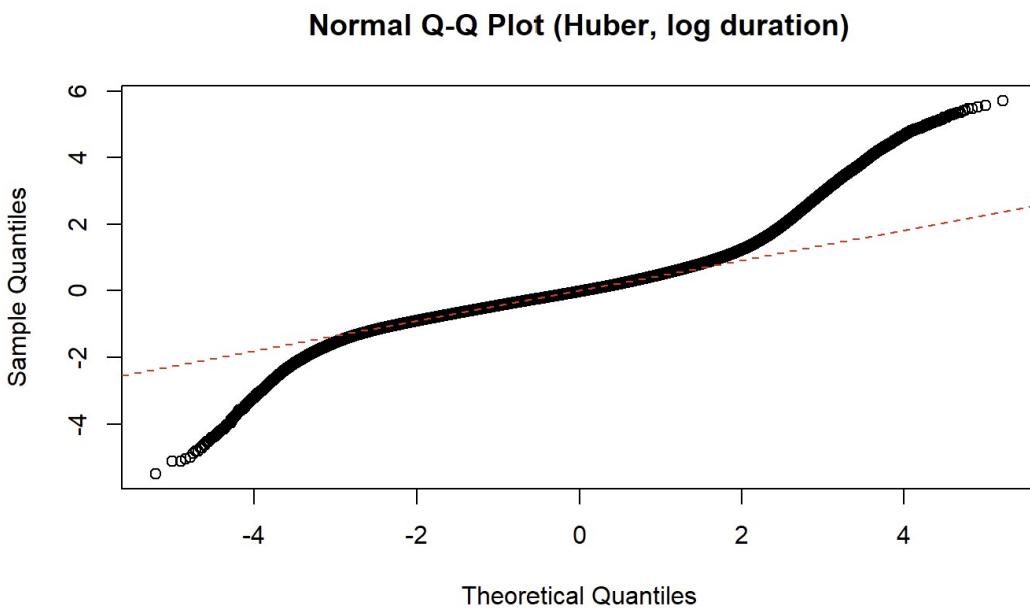


Figure C8

Normal Q-Q plot for the Huber (log-duration) model residuals.

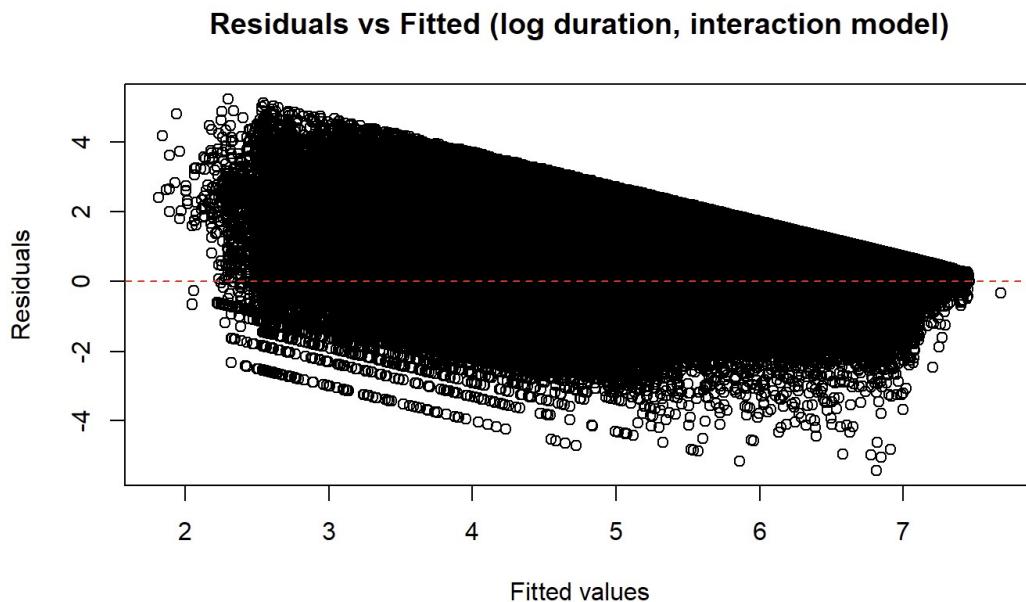


Figure C9

Residuals vs. fitted values for the log-duration interaction model.

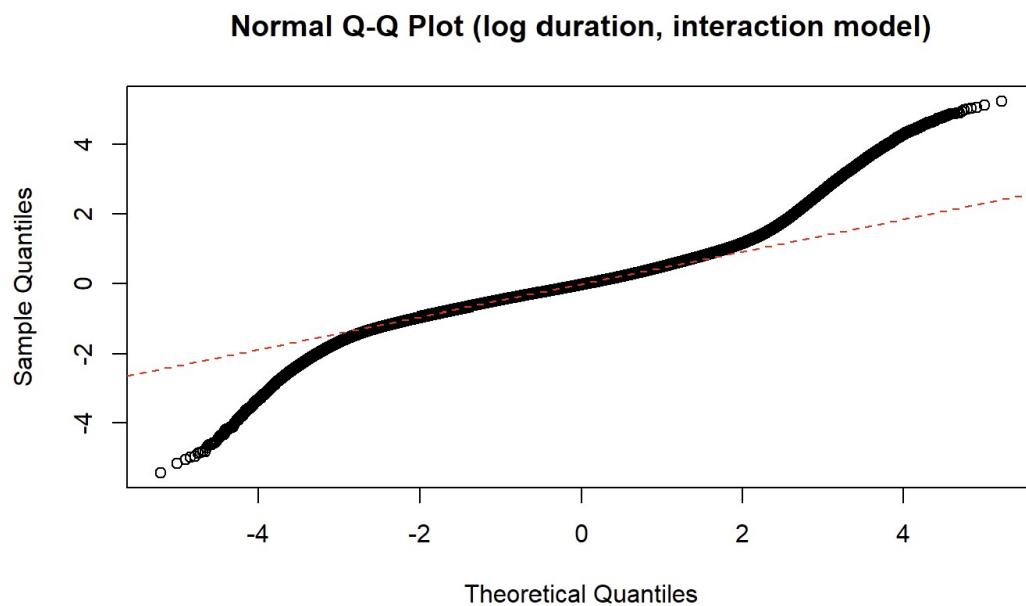


Figure C10

Normal Q-Q plot for the interaction model residuals.

Model <chr>	BIC <dbl>
OLS_raw	77480492
OLS_log	8621600
BoxCox_log	25917405
WLS_log	9992741
Huber_log	8728249

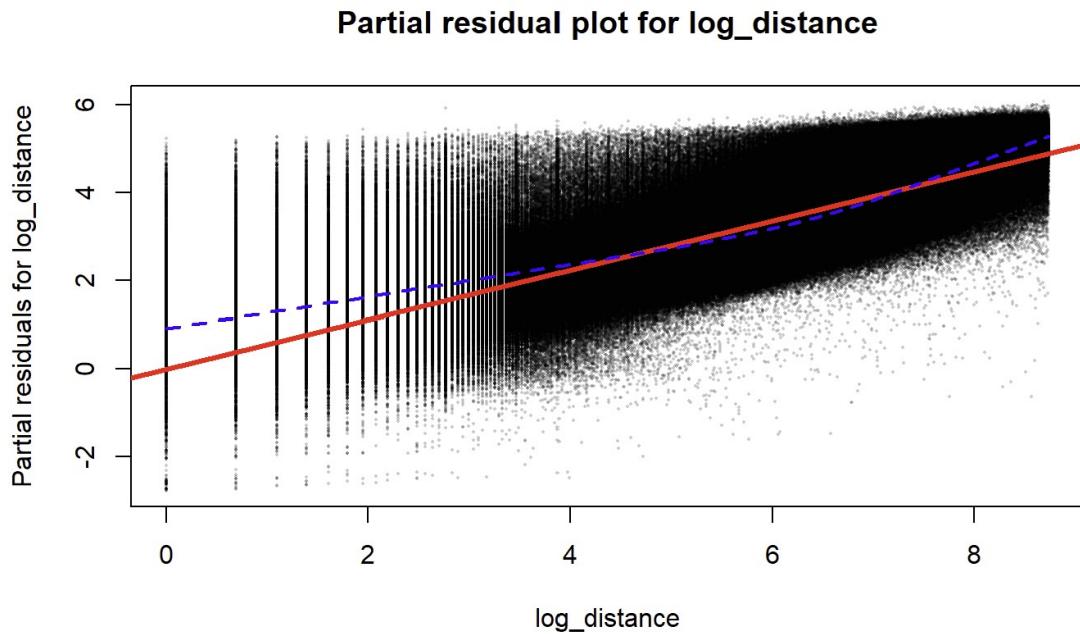
Figure C11

BIC comparison across candidate duration models.

Model <chr>	BIC <dbl>
main_log	8621600
Veh_dis_log	8618148
Veh_hour_log	8619729
Veh_district_log	8621300

Figure C12

BIC comparison across interaction candidates (log-duration).

**Figure C13**

Partial residual plot for log_distance in the duration model.

	Vehicle <chr>	OLS <dbl>	BoxCox <dbl>	WLS <dbl>	Huber <dbl>
vehicle_type	bicycle (baseline)	0.00000000	0.0000000	0.00000000	0.000000000
scooter	-0.02805409	0.1982758	-0.05877432	-0.006723275	

Figure C14

Key coefficient comparison across duration models.

	GVIF	DF	GVIF^(1/(2*DF))
vehicle_type	61.715967	1	7.855951
log_distance	28.067034	1	5.297833
hour_f	1.059670	23	1.001261
dow_f	1.054380	6	1.004422
district_start_f	1.023551	11	1.001059
vehicle_type:log_distance	79.985775	1	8.943477

Figure C15

GVIF diagnostics for the duration model predictors.

Model <chr>	a1_scooter <dbl>	a3_interaction <dbl>	log_distance_used <dbl>	delta_scooter_log <dbl>
Main-effects (no interaction)	-0.028054	NA	NA	-0.028054
Interaction (vehicle_type × log_distance)	0.428973	-0.062602	6.804167	0.003020

Figure C16

Results of Q2.

Appendix D

Full Summary of Selected Models

Model from Question 1: Full summary of OLS Model (log scale).

```

Call:
lm(formula = log_distance ~ vehicle_type + hour_f + dow_f + district_start_f,
    data = trips_2019)

Residuals:
    Min      1Q  Median      3Q     Max 
-7.5680 -0.4148  0.1839  0.7148  3.0537 

Coefficients:
            Estimate Std. Error t value Pr(>|t|)    
(Intercept) 7.135650  0.038454 185.564 < 2e-16 ***
vehicle_typescooter -0.518320  0.002149 -241.209 < 2e-16 ***
hour_f1      -0.017469  0.002631 -6.640 3.15e-11 ***
hour_f2      -0.022280  0.002740 -8.130 4.28e-16 ***
hour_f3      -0.023225  0.002835 -8.191 2.58e-16 ***
hour_f4      -0.023485  0.002974 -7.898 2.84e-15 ***
hour_f5      -0.099775  0.003381 -29.508 < 2e-16 ***
hour_f6      -0.129893  0.003929 -33.064 < 2e-16 ***
hour_f7      -0.063489  0.004337 -14.639 < 2e-16 ***
hour_f8      -0.038733  0.005957 -6.502 7.95e-11 ***
hour_f9      -0.006326  0.008843 -0.715 0.474384  
hour_f10     0.011898  0.009962  1.194 0.232368  
hour_f11     -0.094305  0.007863 -11.993 < 2e-16 ***
hour_f12     0.005546  0.005068  1.094 0.273805  
hour_f13     0.031317  0.003508  8.926 < 2e-16 ***
hour_f14     -0.017633  0.003193 -5.522 3.34e-08 ***
hour_f15     -0.057582  0.003138 -18.350 < 2e-16 ***
hour_f16     -0.048138  0.002936 -16.396 < 2e-16 ***
hour_f17     -0.049691  0.002717 -18.290 < 2e-16 ***
hour_f18     -0.030394  0.002644 -11.497 < 2e-16 ***
hour_f19     -0.017133  0.002627 -6.523 6.90e-11 ***
hour_f20     -0.008002  0.002593 -3.086 0.002026 **  
hour_f21     0.002844  0.002556  1.113 0.265876  
hour_f22     0.028587  0.002503  11.423 < 2e-16 ***
hour_f23     0.025952  0.002500  10.379 < 2e-16 ***
dow_f1      -0.016064  0.001986 -8.086 6.14e-16 ***
dow_f2      -0.014945  0.001988 -7.516 5.64e-14 ***
dow_f3      -0.004832  0.001943 -2.487 0.012887 *  
dow_f4      0.001306  0.001864  0.701 0.483476  
dow_f5      -0.003959  0.001788 -2.214 0.026853 *  
dow_f6      -0.010944  0.001818 -6.021 1.73e-09 *** 
district_start_f1 0.276188  0.038390  7.194 6.28e-13 ***
district_start_f10 0.460076  0.039137 11.755 < 2e-16 ***
district_start_f2 0.232989  0.042787  5.445 5.17e-08 ***
district_start_f3 0.336295  0.038366  8.765 < 2e-16 ***
district_start_f4 0.420622  0.040031 10.508 < 2e-16 ***
district_start_f5 0.405061  0.038426 10.541 < 2e-16 ***
district_start_f6 -0.896498  0.089802 -9.983 < 2e-16 ***
district_start_f7 0.416205  0.039506 10.535 < 2e-16 ***
district_start_f8 0.307449  0.038909  7.902 2.75e-15 ***
district_start_f9 0.141769  0.038340  3.698 0.000218 *** 
district_start_fNone -1.224563  0.659186 -1.858 0.063213 . 
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 1.14 on 5485962 degrees of freedom
Multiple R-squared:  0.01775, Adjusted R-squared:  0.01774 
F-statistic: 2418 on 41 and 5485962 DF, p-value: < 2.2e-16

```

Figure D1

Model from Question 2: Full summary of Main-Effect Model (log scale).

```

Call:
lm(formula = log_duration ~ vehicle_type + log_distance + hour_f +
dow_f + district_start_f, data = trips_2019)

Residuals:
    Min      1Q  Median      3Q     Max 
-5.4317 -0.3336 -0.0367  0.2937  5.2592 

Coefficients:
            Estimate Std. Error t value Pr(>|t|)    
(Intercept) 2.8845026  0.0179665 160.549 < 2e-16 ***
vehicle_typescooter -0.0280541  0.0010061 -27.883 < 2e-16 ***
log_distance   0.4999024  0.0001989 2513.887 < 2e-16 ***
hour_f1        0.0072145  0.0012254  5.887 3.93e-09 ***
hour_f2        0.0012516  0.0012764  0.981 0.326814  
hour_f3        -0.0040595  0.0013206 -3.074 0.002112 ** 
hour_f4        -0.0200463  0.0013850 -14.474 < 2e-16 ***
hour_f5        0.0158552  0.0015750 10.067 < 2e-16 ***
hour_f6        0.0122678  0.0018300  6.704 2.03e-11 *** 
hour_f7        0.0267895  0.0020201 13.262 < 2e-16 *** 
hour_f8        -0.0116477  0.0027748 -4.198 2.70e-05 *** 
hour_f9        -0.1128507  0.0041186 -27.400 < 2e-16 *** 
hour_f10       -0.2525382  0.0046400 -54.426 < 2e-16 *** 
hour_f11       -0.3411018  0.0036626 -93.132 < 2e-16 *** 
hour_f12       -0.3014430  0.0023604 -127.709 < 2e-16 *** 
hour_f13       -0.2645194  0.0016341 -161.870 < 2e-16 *** 
hour_f14       -0.2293047  0.0014872 -154.188 < 2e-16 *** 
hour_f15       -0.1783729  0.0014616 -122.042 < 2e-16 *** 
hour_f16       -0.1182519  0.0013675 -86.475 < 2e-16 *** 
hour_f17       -0.0840287  0.0012654 -66.403 < 2e-16 *** 
hour_f18       -0.0468970  0.0012313 -38.086 < 2e-16 *** 
hour_f19       -0.0143503  0.0012234 -11.730 < 2e-16 *** 
hour_f20       -0.0081430  0.0012076 -6.743 1.55e-11 *** 
hour_f21       -0.0031135  0.0011906 -2.615 0.008919 ** 
hour_f22       -0.0087702  0.0011656 -7.524 5.31e-14 *** 
hour_f23       -0.0042086  0.0011646 -3.614 0.000302 *** 
dow_f1         -0.0475681  0.0009252 -51.412 < 2e-16 *** 
dow_f2         -0.0613278  0.0009261 -66.221 < 2e-16 *** 
dow_f3         -0.0407633  0.0009050 -45.044 < 2e-16 *** 
dow_f4         0.0112928  0.0008684 13.004 < 2e-16 *** 
dow_f5         0.1147309  0.0008330 137.730 < 2e-16 *** 
dow_f6         0.1219175  0.0008466 144.012 < 2e-16 *** 
district_start_f1 -0.2085172  0.0178808 -11.662 < 2e-16 *** 
district_start_f10 -0.2159013  0.0182290 -11.844 < 2e-16 *** 
district_start_f2 -0.0959542  0.0199289 -4.815 1.47e-06 *** 
district_start_f3 -0.3238657  0.0178696 -18.124 < 2e-16 *** 
district_start_f4 -0.2260694  0.0186450 -12.125 < 2e-16 *** 
district_start_f5 -0.2250908  0.0178978 -12.576 < 2e-16 *** 
district_start_f6  0.0220634  0.0418271  0.527 0.597853  
district_start_f7 -0.2870359  0.0184006 -15.599 < 2e-16 *** 
district_start_f8  0.0094983  0.0181223  0.524 0.600193  
district_start_f9 -0.2617185  0.0178576 -14.656 < 2e-16 *** 
district_start_fNone 0.5041681  0.3070251  1.642 0.100568 
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 0.5309 on 5485961 degrees of freedom
Multiple R-squared:  0.5505,    Adjusted R-squared:  0.5505 
F-statistic: 1.6e+05 on 42 and 5485961 DF,  p-value: < 2.2e-16

```

Figure D2

Model from Question 2: Full summary of Interaction Model (log scale).

```

Call:
lm(formula = log_duration ~ vehicle_type * log_distance + hour_f +
    dow_f + district_start_f, data = trips_2019)

Residuals:
    Min      1Q  Median      3Q     Max 
-5.4275 -0.3332 -0.0371  0.2932  5.2435 

Coefficients:
                                         Estimate Std. Error t value Pr(>|t|)    
(Intercept)                      2.4415454  0.0194719 125.388 < 2e-16 ***
vehicle_typescooter               0.4289727  0.0078246  54.823 < 2e-16 ***
log_distance                       0.5602536  0.0010438  536.751 < 2e-16 ***
hour_f1                            0.0073281  0.0012251   5.982 2.21e-09 ***
hour_f2                            0.0014422  0.0012760   1.130 0.258382  
hour_f3                            -0.0038601  0.0013202  -2.924 0.003456 ** 
hour_f4                            -0.0198047  0.0013846  -14.304 < 2e-16 ***
hour_f5                            0.0160625  0.0015745   10.201 < 2e-16 ***
hour_f6                            0.0124417  0.0018294   6.801 1.04e-11 ***
hour_f7                            0.0268699  0.0020195   13.306 < 2e-16 ***
hour_f8                            -0.0116280  0.0027739  -4.192 2.77e-05 *** 
hour_f9                            -0.1125968  0.0041173  -27.347 < 2e-16 ***
hour_f10                           -0.2520008  0.0046386  -54.327 < 2e-16 ***
hour_f11                           -0.3409569  0.0036614  -93.122 < 2e-16 ***
hour_f12                           -0.3013432  0.0023596  -127.707 < 2e-16 ***
hour_f13                           -0.2645329  0.0016336  -161.930 < 2e-16 ***
hour_f14                           -0.2293100  0.0014867  -154.241 < 2e-16 ***
hour_f15                           -0.1782284  0.0014611  -121.981 < 2e-16 ***
hour_f16                           -0.1180901  0.0013670  -86.384 < 2e-16 ***
hour_f17                           -0.0838511  0.0012650  -66.283 < 2e-16 ***
hour_f18                           -0.0466565  0.0012310  -37.903 < 2e-16 ***
hour_f19                           -0.0142044  0.0012230  -11.614 < 2e-16 ***
hour_f20                           -0.0079734  0.0012072  -6.605 3.98e-11 *** 
hour_f21                           -0.0030407  0.0011902  -2.555 0.010623 *  
hour_f22                           -0.0088022  0.0011653  -7.554 4.23e-14 *** 
hour_f23                           -0.0042743  0.0011642  -3.671 0.000241 *** 
dow_f1                            -0.0476582  0.0009249  -51.526 < 2e-16 ***
dow_f2                            -0.0613926  0.0009258  -66.312 < 2e-16 ***
dow_f3                            -0.0408346  0.0009047  -45.137 < 2e-16 *** 
dow_f4                            0.0112013  0.0008681  12.903 < 2e-16 *** 
dow_f5                            0.1147095  0.0008328  137.748 < 2e-16 *** 
dow_f6                            0.1218560  0.0008463  143.985 < 2e-16 *** 
district_start_f1                 -0.2073884  0.0178752  -11.602 < 2e-16 ***
district_start_f10                -0.2172245  0.0182232  -11.920 < 2e-16 *** 
district_start_f2                 -0.0934031  0.0199226  -4.688 2.75e-06 *** 
district_start_f3                 -0.3227818  0.0178640  -18.069 < 2e-16 *** 
district_start_f4                 -0.2252413  0.0186391  -12.084 < 2e-16 *** 
district_start_f5                 -0.2243001  0.0178921  -12.536 < 2e-16 *** 
district_start_f6                 0.0207705  0.0418139   0.497 0.619375  
district_start_f7                 -0.2860714  0.0183947  -15.552 < 2e-16 *** 
district_start_f8                 0.0090387  0.0181166   0.499 0.617838  
district_start_f9                 -0.2605560  0.0178519  -14.595 < 2e-16 *** 
district_start_fNone              0.5020058  0.3069281   1.636 0.101927  
vehicle_typescooter:log_distance -0.0626017  0.0010629  -58.897 < 2e-16 *** 
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 0.5307 on 5485960 degrees of freedom
Multiple R-squared:  0.5508,    Adjusted R-squared:  0.5508 
F-statistic: 1.564e+05 on 43 and 5485960 DF,  p-value: < 2.2e-16

```

Figure D3

Appendix E

Full R code

```

1 # Load packages
2
3 ``'{r, load packages}
4 # load packages
5 library(readr)
6 library(dplyr)
7 library(ggplot2)
8 library(gridExtra)
9 library(car)
10 library(MASS)
11 library(broom)
12 '''
13 ``'{r, read data}
14 trips_raw <-
15   read_csv("https://data.austintexas.gov/api/views/7d8e-dm7r/rows.csv?accessType=DOWNLOAD")
16 '''
17
18 # Preprocessing and EDA
19
20 ## Preprocessing data##
21 #####
22 ## Cleaning data
23 # Filter for data of 2019 (year), get rid of NA's and negative values,
24 # select essential predictors that will be utilized
25 #####
26
27 trips_2019 <- trips_raw %>%
28   filter(
29     'Year (US/Central)' == 2019,
30     'Vehicle Type' %in% c("bicycle", "scooter"),
31     !is.na('Trip Distance'),
32     !is.na('Trip Duration'),
33     'Trip Duration' > 0)
34 
```

```

32     'Trip Distance' > 0,
33     'Trip Duration' > 0
34 ) %>%
35 select(
36   'Vehicle Type', 'Trip Distance', 'Trip Duration',
37   'Hour', 'Day of Week', 'Council District (Start)'
38 )
39
40
41 saveRDS(trips_2019, "trips_2019.rds")
42 ' '
43
44 ``'{r, preprocessing2}
#####
# Change predictors from characters to factors
#####
45 trips_2019_nolog <- trips_2019 %>%
46   mutate(
47     vehicle_type = factor('Vehicle Type'),
48     hour_f = factor(Hour),
49     dow_f = factor('Day of Week'),
50     district_start_f = factor('Council District (Start)')
51   )
52
53
54 saveRDS(trips_2019_nolog, "trips_2019_nolog.rds")
55 ' '
56
57
58 ``'{r, echo=FALSE}
59
60 trips_2019 <- readRDS("trips_2019.rds")
61 trips_2019_nolog <- readRDS("trips_2019_nolog.rds")
62 ' '
63
64 ** Raw EDA**
65 ``'{r, 1st raw EDA}

```

```

66 #####
67 ## EDA of preprocessed raw data (no transformation applied)
68 # The preprocessed dataset is large, so take a random subset and explore the data
69 # Check predictor (vehicle type) with response (trip distance, trip duration)
70 #####
71
72 # Create sample of 100,000
73 set.seed(500)
74 trips_sample_nolog <- trips_2019_nolog %>%
75   sample_n(100000)
76
77 # Trip distance histogram
78 dis_hist <-
79 ggplot(trips_sample_nolog, aes(x = 'Trip Distance')) +
80   geom_histogram(bins = 100) +
81   labs(x = "Trip\u00d7Distance\u00d7(meters)", y = "Count",
82        title = "Distribution\u00d7of\u00d7Trip\u00d7Distance\u00d7(raw)")
83
84 # Trip duration histogram
85 dur_hist <-
86 ggplot(trips_sample_nolog, aes(x = 'Trip Duration')) +
87   geom_histogram(bins = 100) +
88   labs(x = "Trip\u00d7duration\u00d7(seconds)", y = "Count",
89        title = "Distribution\u00d7of\u00d7Trip\u00d7Duration\u00d7(raw)")
90
91 # boxplot (x:vehicle_type / y:log_distance)
92 dis_box <-
93 ggplot(trips_sample_nolog, aes(x = vehicle_type, y = 'Trip Distance')) +
94   geom_boxplot() +
95   labs(x = "Vehicle\u00d7type", y = "Trip\u00d7distance\u00d7(meters)",
96        title = "Trip\u00d7Distance\u00d7by\u00d7Vehicle\u00d7Type\u00d7(raw)")
97
98 # boxplot (x:vehicle_type / y:log_duration)
99 dur_box <

```

```

100 ggplot(trips_sample_nolog, aes(x = vehicle_type, y = 'Trip Duration')) +
101   geom_boxplot() +
102   labs(x = "Vehicle\u00e1type", y = "Trip\u00d7duration\u00d7(seconds)", 
103         title = "Trip\u00d7Duration\u00d7by\u00d7Vehicle\u00d7Type\u00d7(raw)")
104 
105 grid.arrange(dis_hist, dur_hist, dis_box, dur_box, nrow=2, ncol=2)
106 ''
107 
108 '''{r, additional 1st EDA}
109 #####
110 ## EDA of preprocessed raw data (no transformation applied)
111 # Check other predictors with response predictor (trip distance, trip duration)
112 #####
113 
114 # boxplot (x:day of week / y:log_distance)
115 dis_box_dof <-
116   ggplot(trips_sample_nolog, aes(x = dow_f, y = 'Trip Distance')) +
117   geom_boxplot() +
118   labs(x = "day\u00d7of\u00d7week", y = "Trip\u00d7distance\u00d7(meters)", 
119         title = "Trip\u00d7Distance\u00d7by\u00d7Day\u00d7of\u00d7Week\u00d7(raw)")
120 
121 # boxplot (x:day of week / y:log_duration)
122 dur_box_dof <-
123   ggplot(trips_sample_nolog, aes(x = dow_f, y = 'Trip Duration')) +
124   geom_boxplot() +
125   labs(x = "day\u00d7of\u00d7week", y = "Trip\u00d7duration\u00d7(seconds)", 
126         title = "Trip\u00d7Duration\u00d7by\u00d7Day\u00d7of\u00d7Week\u00d7(raw)")
127 
128 # boxplot (x:hour / y:log_distance)
129 dis_box_hour <-
130   ggplot(trips_sample_nolog, aes(x = hour_f, y = 'Trip Distance')) +
131   geom_boxplot() +
132   labs(x = "Hour", y = "Trip\u00d7distance\u00d7(meters)", 
133         title = "Trip\u00d7Distance\u00d7by\u00d7Hours\u00d7(raw)")

```

```

134
135 # boxplot (x:hour / y:log_duration)
136 dur_box_hour <- ggplot(trips_sample_nolog, aes(x = hour_f, y = 'Trip Duration')) +
137   geom_boxplot() +
138   labs(x = "Hour", y = "Trip\u00d7duration\u207b(seconds)", 
139         title = "Trip\u00d7Duration\u207bby\u207bHours\u207b(raw)")
140
141 # boxplot (x: Council District (Start) / y:log_distance)
142 dis_box_district <-
143   ggplot(trips_sample_nolog, aes(x = district_start_f, y = 'Trip Distance')) +
144   geom_boxplot() +
145   labs(x = "Council\u00d7District\u207b(Start)", y = "Trip\u00d7distance\u207b(meters)", 
146         title = "Trip\u00d7Distance\u207bby\u207bCouncil\u00d7District\u207b(raw)")
147
148 # boxplot (x: Council District (Start) / y:log_duration)
149 dur_box_district <-
150   ggplot(trips_sample_nolog, aes(x = district_start_f, y = 'Trip Duration')) +
151   geom_boxplot() +
152   labs(x = "Council\u00d7District\u207b(Start)", y = "Trip\u00d7duration\u207b(seconds)", 
153         title = "Trip\u00d7Duration\u207bby\u207bCouncil\u00d7District\u207b(raw)")
154
155 grid.arrange(dis_box_dof, dur_box_dof, dis_box_hour, dur_box_hour, dis_box_district,
156               dur_box_district, nrow=3, ncol=2)
157
158 """
159 {r, Calculate IQR}
160 #####
161 ## Calculate IQR
162 # 1.5IQR for checking potential outliers, 3IQR for extreme outliers
163 #####
164 # Trip Distance(d), Trip Duration(t)
165 qd <- quantile(trips_2019$'Trip Distance', probs = c(0.25, 0.75), na.rm = TRUE)
166 IQR_d <- qd[2] - qd[1]

```

```

167
168 qt <- quantile(trips_2019$'Trip Duration', probs = c(0.25, 0.75), na.rm = TRUE)
169 IQR_t <- qt[2] - qt[1]
170
171 # 1.5*IQR standard (for potential outliers)
172 lower_d_1_5 <- qd[1] - 1.5 * IQR_d
173 upper_d_1_5 <- qd[2] + 1.5 * IQR_d
174 lower_t_1_5 <- qt[1] - 1.5 * IQR_t
175 upper_t_1_5 <- qt[2] + 1.5 * IQR_t
176
177 # 3*IQR standard (for extreme outliers)
178 lower_d_3 <- qd[1] - 3 * IQR_d
179 upper_d_3 <- qd[2] + 3 * IQR_d
180 lower_t_3 <- qt[1] - 3 * IQR_t
181 upper_t_3 <- qt[2] + 3 * IQR_t
182
183 # Distance must be positive(cutoff > 0)
184 lower_d_1_5 <- max(lower_d_1_5, 0)
185 lower_d_3 <- max(lower_d_3, 0)
186 lower_t_1_5 <- max(lower_t_1_5, 0)
187 lower_t_3 <- max(lower_t_3, 0)
188
189 qd; qt
190 IQR_d; IQR_t
191 lower_d_1_5; upper_d_1_5
192 lower_d_3; upper_d_3
193 lower_t_1_5; upper_t_1_5
194 lower_t_3; upper_t_3
195
196 outlier_table <- data.frame(
197   Measure = c("Trip_Distance", "Trip_Duration"),
198   Q1 = c(unname(qd[1]), unname(qt[1])),
199   Q3 = c(unname(qd[2]), unname(qt[2])),
200   IQR = c(unname(IQR_d), unname(IQR_t)),

```

```

201 Lower_1_5 = c(unname(lower_d_1_5), unname(lower_t_1_5)),
202 Upper_1_5 = c(unname(upper_d_1_5), unname(upper_t_1_5)),
203 Lower_3 = c(unname(lower_d_3), unname(lower_t_3)),
204 Upper_3 = c(unname(upper_d_3), unname(upper_t_3)),
205 check.names = FALSE
206 )
207
208 outlier_table
209 """
210
211 """{r, remove outlier by 3 X IQR}
212
213 ## Remove extreme outliers by 3*IQR
214 trips_2019_nolog <- trips_2019_nolog %>%
215   filter(
216     `Trip Distance` >= lower_d_3,
217     `Trip Distance` <= upper_d_3,
218     `Trip Duration` >= lower_t_3,
219     `Trip Duration` <= upper_t_3
220   )
221
222 trips_2019_nolog %>%
223   summarise(
224     n = n(),
225     mean_dist = mean(`Trip Distance`),
226     median_dist= median(`Trip Distance`),
227     max_dist = max(`Trip Distance`),
228     mean_dur = mean(`Trip Duration`),
229     median_dur = median(`Trip Duration`),
230     max_dur = max(`Trip Duration`)
231   )
232
233 saveRDS(trips_2019_nolog, "trips_2019_nolog.rds")
234 """

```

```

235
236 ## Raw data EDA without outliers (IQR)**
237 ## {r, Raw data EDA without outliers (IQR)}
238 #####
239 ## Remove extreme outliers by 3*IQR and check dataset
240 #####
241
242 # Create sample of 100,000
243 set.seed(500)
244 trips_sample_nolog <- trips_2019_nolog %>%
245   sample_n(100000)
246
247 # Trip distance histogram
248 dis_hist <-
249   ggplot(trips_sample_nolog, aes(x = 'Trip Distance')) +
250     geom_histogram(bins = 100) +
251     labs(x = "Trip\u00d7Distance\u20d7(meters)", y = "Count",
252           title = "Distribution\u20d7of\u20d7Trip\u00d7Distance\u20d7(raw)")
253
254 # Trip duration histogram
255 dur_hist <-
256   ggplot(trips_sample_nolog, aes(x = 'Trip Duration')) +
257     geom_histogram(bins = 100) +
258     labs(x = "Trip\u00d7duration\u20d7(seconds)", y = "Count",
259           title = "Distribution\u20d7of\u20d7Trip\u00d7Duration\u20d7(raw)")
260
261 # boxplot (x:vehicle_type / y:log_distance)
262 dis_box <-
263   ggplot(trips_sample_nolog, aes(x = vehicle_type, y = 'Trip Distance')) +
264     geom_boxplot() +
265     labs(x = "Vehicle\u00d7type", y = "Trip\u00d7distance\u20d7(meters)",
266           title = "Trip\u00d7Distance\u20d7by\u20d7Vehicle\u00d7Type\u20d7(raw)")
267
268 # boxplot (x:vehicle_type / y:log_duration)

```

```

269 dur_box <-
270   ggplot(trips_sample_nolog, aes(x = vehicle_type, y = 'Trip Duration')) +
271   geom_boxplot() +
272   labs(x = "Vehicle\u2014type", y = "Trip\u2014duration\u2014(seconds)",
273        title = "Trip\u2014Duration\u2014by\u2014Vehicle\u2014Type\u2014(raw)")
274
275 grid.arrange(dis_hist, dur_hist, dis_box, dur_box, nrow=2, ncol=2)
276 ''
277
278 # Question 1.
279
280 ** Raw model**
281 '''{r, raw model with no transformation}
282
283 ## Check raw model results with no transformation
284 model_raw_dis <- lm('Trip Distance' ~ vehicle_type + hour_f + dow_f +
285   district_start_f, data=trips_2019_nolog)
286 summary(model_raw_dis)
287 ''
288
289 ** Raw model diagnostics**
290 '''{r, Diagnostics for raw model}
291
292 ## Check nonlinearity, heteroscedasticity, normality of raw model
293
294 fitted_raw_dis <- fitted(model_raw_dis)
295 resid_raw_dis <- resid(model_raw_dis)
296
297 # par(mfrow = c(2, 1))
298
299 # Residual vs Fitted
300 plot(fitted_raw_dis, resid_raw_dis,
301       xlab = "Fitted\u2014values", ylab = "Residuals",
302       main = "Residuals\u2014vs\u2014Fitted\u2014(raw\u2014Trip\u2014Distance)")
```

```

302 abline(h = 0, lty = 2, col = "red")
303
304 # QQ plot
305 qqnorm(resid_raw_dis,
306         main = "Normal_Q-Q_Plot_(raw_Trip_Distance)")
307 qqline(resid_raw_dis, col = "red", lty = 2)
308
309 ''
310
311 **-1-Log transformation of response variables**
312 '''{r, log-transform response variables}
313 #####
314 # In order to overcome nonlinearity, heteroscedasticity, normality
315 # try log transformation
316 #####
317
318 ## Mutate log-transform response (distance & duration)
319 trips_2019 <- trips_2019_nolog %>%
320   mutate(
321     log_distance = log('Trip Distance'),
322     log_duration = log('Trip Duration')
323   )
324
325 saveRDS(trips_2019, "trips_2019.rds")
326 ''
327
328 **-1-EDA: Exploration of Data (Log transformed)**
329 '''{r, 2nd EDA}
330 # Create sample of 100,000
331 set.seed(500)
332 trips_sample_log <- trips_2019 %>%
333   sample_n(100000)
334
335 # log distance histogram

```

```

336 dis_hist_log <-
337   ggplot(trips_sample_log, aes(x = log_distance)) +
338     geom_histogram(bins = 100) +
339     labs(x = "log(Trip\u2022distance)", y = "Count",
340           title = "Distribution\u2022of\u2022log(Trip\u2022Distance)")
341
342 # log duration histogram
343 dur_hist_log <-
344   ggplot(trips_sample_log, aes(x = log_duration)) +
345     geom_histogram(bins = 100) +
346     labs(x = "log(Trip\u2022duration)", y = "Count",
347           title = "Distribution\u2022of\u2022log(Trip\u2022Duration)")
348
349 # boxplot (x:vehicle_type / y:log_distance)
350 dis_box_log <-
351   ggplot(trips_sample_log, aes(x = vehicle_type, y = log_distance)) +
352     geom_boxplot() +
353     labs(x = "Vehicle\u2022type", y = "log(Trip\u2022distance)",
354           title = "log(Trip\u2022Distance)\u2022by\u2022Vehicle\u2022Type")
355
356 # boxplot (x:vehicle_type / y:log_duration)
357 dur_box_log <-
358   ggplot(trips_sample_log, aes(x = vehicle_type, y = log_duration)) +
359     geom_boxplot() +
360     labs(x = "Vehicle\u2022type", y = "log(Trip\u2022duration)",
361           title = "log(Trip\u2022Duration)\u2022by\u2022Vehicle\u2022Type")
362
363
364 grid.arrange(dis_hist_log, dur_hist_log, dis_box_log, dur_box_log, nrow=2, ncol=2)
365 """
366
367 """{r, additional 2nd EDA}
368 # boxplot (x:day of week / y:log_distance)
369 dis_box_dof_log <

```

```

370 ggplot(trips_sample_log, aes(x = dow_f, y = log_distance)) +
371   geom_boxplot() +
372   labs(x = "day\u00d7of\u00d7week", y = "Trip\u00d7distance\u00d7(meters)",
373         title = "Trip\u00d7Distance\u00d7by\u00d7Day\u00d7of\u00d7Week")
374
375 # boxplot (x:day of week / y:log_duration)
376 dur_box_dof_log <-
377   ggplot(trips_sample_log, aes(x = dow_f, y = log_duration)) +
378   geom_boxplot() +
379   labs(x = "day\u00d7of\u00d7week", y = "Trip\u00d7duration\u00d7(seconds)",
380         title = "Trip\u00d7Duration\u00d7by\u00d7Day\u00d7of\u00d7Week")
381
382 # boxplot (x:hour / y:log_distance)
383 dis_box_hour_log <-
384   ggplot(trips_sample_log, aes(x = hour_f, y = log_distance)) +
385   geom_boxplot() +
386   labs(x = "Hour", y = "Trip\u00d7distance\u00d7(meters)",
387         title = "Trip\u00d7Distance\u00d7by\u00d7Hours")
388
389 # boxplot (x:hour / y:log_duration)
390 dur_box_hour_log <-
391   ggplot(trips_sample_log, aes(x = hour_f, y = log_duration)) +
392   geom_boxplot() +
393   labs(x = "Hour", y = "Trip\u00d7duration\u00d7(seconds)",
394         title = "Trip\u00d7Duration\u00d7by\u00d7Hours")
395
396 # boxplot (x: Council District (Start) / y:log_distance)
397 dis_box_district_log <-
398   ggplot(trips_sample_log, aes(x = district_start_f, y = log_distance)) +
399   geom_boxplot() +
400   labs(x = "Council\u00d7District\u00d7(Start)", y = "Trip\u00d7distance\u00d7(meters)",
401         title = "Trip\u00d7Distance\u00d7by\u00d7Council\u00d7District")
402
403 # boxplot (x: Council District (Start) / y:log_duration)

```

```

404 dur_box_district_log <-
405   ggplot(trips_sample_log, aes(x = district_start_f, y = log_duration)) +
406   geom_boxplot() +
407   labs(x = "Council\u2022District\u2022(Start)", y = "Trip\u2022duration\u2022(seconds)",
408         title = "Trip\u2022Duration\u2022by\u2022Council\u2022District")
409
410 grid.arrange(dis_box_dof_log, dur_box_dof_log, dis_box_hour_log, dur_box_hour_log,
411               dis_box_district_log, dur_box_district_log, nrow=3, ncol=2)
412
413 ***-1-Modeling & Diagnostics: Log transformed model***
414 ***{r, log transformed model diagnostics}
415 #####
416 ## Q1. Log transformed Model
417 # Modeling Diagnose
418 #####
419
420 # Modeling Log transformed model
421 model_log_dis <- lm(log_distance ~ vehicle_type + hour_f + dow_f + district_start_f,
422                       data = trips_2019)
423
424 # Residual vs Fitted
425 plot(fitted(model_log_dis), resid(model_log_dis),
426       xlab = "Fitted\u2022values",
427             ylab = "Residuals",
428             main = "Residuals\u2022vs\u2022Fitted\u2022(log\u2022distance)")
429 abline(h = 0, lty = 2, col = "red")
430
431 # QQ plot
432 qqnorm(resid(model_log_dis),
433         main = "Normal\u2022Q-Q\u2022Plot\u2022(log\u2022distance)")
434 qqline(resid(model_log_dis), col = "red", lty = 2)
435
436

```

```

437  **-2-Modeling & Diagnostics: Box-Cox model**
438  ``'{r, Boxcox Diagnostics}
439 #####
440 ## Q1. Boxcox Model
441 # Modeling Diagnose
442 #####
443
444 # Modeling Boxcox
445 boxcox_lambda_dis <- boxcox(model_raw_dis, plotit = TRUE)
446
447 lambda_hat_dis <- boxcox_lambda_dis$x[which.max(boxcox_lambda_dis$y)]
448 lambda_hat_dis
449
450 trips_2019_nolog$y_boxcox_dis <- if (abs(lambda_hat_dis) < 1e-3) {
451   log(trips_2019_nolog$'Trip Distance')
452 } else {
453   (trips_2019_nolog$'Trip Distance'^lambda_hat_dis - 1) / lambda_hat_dis
454 }
455
456 model_boxcox_dis <- lm(y_boxcox_dis ~ vehicle_type + hour_f + dow_f +
457                           district_start_f,
458                           data = trips_2019_nolog)
459
460 # Residual vs Fitted
461 # par(mfrow = c(2,1))
462 plot(fitted(model_boxcox_dis), resid(model_boxcox_dis),
463       xlab = "Fitted values",
464       ylab = "Residuals",
465       main = "Residuals vs Fitted (boxcox distance)")
466 abline(h = 0, lty = 2, col = "red")
467
468 # QQ plot
469 qqnorm(resid(model_boxcox_dis)),

```

```

470     main = "Normal_Q-Q_Plot_(boxcox_distance)")
471 qqline(resid(model_boxcox_dis), col = "red", lty = 2)
472 ' '
473
474 **-3-Modeling & Diagnostics: Weighted Least Square (WLS) model**
475 ' '{r, WLS Diagnostics}
476 #####
477 ## Q1. WLS Model
478 # Modeling Diagnose
479 #####
480
481 # Modeling WLS
482 f_hat <- fitted(model_log_dis)
483 e_hat <- resid(model_log_dis)
484
485 w_simple <- 1 / (f_hat^2 + 1e-8)
486
487 model_wls_dis <- lm(log_distance ~ vehicle_type + hour_f + dow_f + district_start_f,
488                         data = trips_2019,
489                         weights = w_simple)
490
491
492 # Residual vs Fitted
493 # par(mfrow = c(2,1))
494 plot(fitted(model_wls_dis), resid(model_wls_dis),
495       xlab = "Fitted_values",
496       ylab = "Residuals",
497       main = "Residuals_vs_Fitted_(WLS_distance)")
498 abline(h = 0, lty = 2, col = "red")
499
500 # QQ plot
501 qqnorm(resid(model_wls_dis),
502         main = "Normal_Q-Q_Plot_(WLS_distance)")
503 qqline(na.omit(resid(model_wls_dis)), col = "red", lty = 2)

```

```

504   ' '
505
506   ' ' '{r}
507 #####
508 ## Q1. Huber Model
509 # Modeling Diagnose
510 #####
511
512 # Modeling Huber
513 model_huber_dis <- rlm(log_distance ~ vehicle_type + hour_f + dow_f +
514   district_start_f,
515   data = trips_2019)
516
517 # Residual vs Fitted
518 plot(fitted(model_huber_dis), resid(model_huber_dis),
519   xlab = "Fitted_values", ylab = "Residuals_(Huber)",
520   main = "Residuals_vs_Fitted_(Huber_distance)")
521 abline(h = 0, col = "red", lty = 2)
522
523 # QQ plot
524 qqnorm(resid(model_huber_dis),
525   main = "Normal_Q-Q_Plot_(Huber_distance)")
526 qqline(resid(model_huber_dis), col = "red", lty = 2)
527 ' '
528 **-5 Check BIC for candidate models**
529 ' ' '{r}
530 #####
531 ## Q1. Compare BIC of candidate models
532 # - model_raw_dis : distance (raw)
533 # - model_log_dis : log_distance
534 # - model_boxcox_dis : boxCox
535 # - model_wls_dis : WLS (log_distance)
536 # - model_huber_dis : Huber (log_distance, rlm)

```

```

537 #####
538
539 ## BIC Comparison table
540 info_dis <- data.frame(
541   Model = c("OLS_raw", "OLS_log", "BoxCox_log", "WLS_log", "Huber_log"),
542   BIC = c(BIC(model_raw_dis),
543           BIC(model_log_dis),
544           BIC(model_boxcox_dis),
545           BIC(model_wls_dis),
546           BIC(model_huber_dis)))
547 )
548 info_dis
549 ''
550
551 **-6 Compare Coefficient table**
552 '''
553 # Compare vehicle_type coefficients
554 coef_table_dis <- data.frame(
555   Vehicle = c("bicycle_(baseline)", "scooter"),
556   OLS = c(0, coef(model_log_dis)["vehicle_typescooter"]),
557   Boxcox = c(0, coef(model_boxcox_dis)["vehicle_typescooter"]),
558   WLS = c(0, coef(model_wls_dis)["vehicle_typescooter"]),
559   Huber = c(0, coef(model_huber_dis)["vehicle_typescooter"]))
560 )
561
562 coef_table_dis
563
564 '''
565
566 ** Model-based outlier, influencial point check**
567 '''
568 #####
569 ## Q1. Check Model-based outliers & influencial points
570 # Apply studentized residuals, leverage, Cook's Distance

```

```

571 #####
572
573 ## Define numbers of observations and coefficients
574 n_dis <- nrow(trips_2019)
575 p_dis <- length(coef(model_log_dis))
576 alpha <- 0.05
577
578
579 ## Calculate Diagnostic measures
580 student_dis <- rstudent(model_log_dis) # externally studentized residuals
581 lev_dis <- hatvalues(model_log_dis) # leverage (hat values)
582 cook_dis <- cooks.distance(model_log_dis) # Cook's distance
583
584
585 ## Set Cutoff
586 # studentized residual Bonferroni cutoff
587 cutoff_student_dis <- qt(1 - alpha/(2*n_dis), df = model_log_dis$df.residual)
588
589 # leverage rule of thumb:  $2*(p+1)/n$ 
590 cutoff_lev_dis <- 2 * (p_dis + 1) / n_dis
591
592 # Cook's distance rule of thumb:  $4/n$ 
593 cutoff_cook_dis <- 4 / n_dis
594
595
596 ## Visualize results as a Table
597 diag_tbl_dis <- data.frame(
598   id = 1:n_dis,
599   student = student_dis,
600   leverage = lev_dis,
601   cook = cook_dis,
602   flag_student = abs(student_dis) > cutoff_student_dis,
603   flag_leverage = lev_dis > cutoff_lev_dis,
604   flag_cook = cook_dis > cutoff_cook_dis

```

```

605 )
606
607 # filtering only the points that might need to be considered
608 diag_tbl_dis$flag_any <- with(diag_tbl_dis,
609                         flag_student | flag_leverage | flag_cook)
610 suspects_dis <- subset(diag_tbl_dis, flag_any)
611 suspects_dis[order(-suspects_dis$cook), ][1:20, ]
612 ```

613
614 ** Check Collinearity**
615 `'{r, check collinearity (VIF)}
616 vif(model_log_dis)
617 ```

618
619 ** Log transformed model interpretation**
620 `'{r}
621 #####
622 ## Q1: scooter effect summary table (log-distance model)
623 #####
624 summary(model_log_dis)

625
626 q1_scooter_tab <- tidy(model_log_dis) %>%
627   filter(term == "vehicle_typescooter") %>%
628   transmute(
629     term = "Scooter $\sqcup$ bicycle",
630     estimate = estimate,
631     std_error = std.error,
632     ratio = exp(estimate),
633     pct_change = 100 * (exp(estimate) - 1)
634   )
635 q1_scooter_tab
636 ```
637
638 # Question 2.

```

```

639
640 ** Raw model**
641
642 '{r, raw model with no transformation}
643
644 ## Check raw model results with no transformation
645 model_raw_dur <- lm(
646   'Trip Duration' ~ vehicle_type + 'Trip Distance' + hour_f + dow_f +
647     district_start_f,
648   data = trips_2019_nolog
649 )
650
651 ** Raw model diagnostics**
652
653 '{r, Diagnostics for raw model}
654
655 ## Check nonlinearity, heteroscedasticity, normality of raw model
656
657 fitted_raw_dur <- fitted(model_raw_dur)
658 resid_raw_dur <- resid(model_raw_dur)
659
660 # par(mfrow = c(2, 1))
661
662 # Residual vs Fitted
663 plot(fitted_raw_dur, resid_raw_dur,
664   xlab = "Fitted values", ylab = "Residuals",
665   main = "Residuals_vs_Fitted_(raw_Trip_Duration)")
666 abline(h = 0, lty = 2, col = "red")
667
668 # QQ plot
669 qqnorm(resid_raw_dur,
670   main = "Normal_Q-Q_Plot_(raw_Trip_Duration)")
671 qqline(resid_raw_dur, col = "red", lty = 2)

```

```

672
673   ''
674
675 **-1-Modeling & Diagnostics: Log transformed model**
676   ``'{r, log transformed model diagnostics}
677 #####
678 ## Q2: Log transformed model for duration
679 ## Response: log_duration
680 ## Predictors: vehicle_type + log_distance + hour_f + dow_f + district_start_f
681 #####
682
683 # Modeling Log transformed model (duration)
684 model_log_dur <- lm(
685   log_duration ~ vehicle_type + log_distance + hour_f + dow_f + district_start_f,
686   data = trips_2019
687 )
688
689 # Residual vs Fitted
690 plot(fitted(model_log_dur), resid(model_log_dur),
691       xlab = "Fitted_values",
692       ylab = "Residuals",
693       main = "Residuals_vs_Fitted_(log_duration)")
694 abline(h = 0, lty = 2, col = "red")
695
696 # QQ plot
697 qqnorm(resid(model_log_dur),
698         main = "Normal_Q-Q_Plot_(log_duration)")
699 qqline(resid(model_log_dur), col = "red", lty = 2)
700   ''
701
702 **-2-Modeling & Diagnostics: Box-Cox model**
703   ``'{r, Boxcox Diagnostics}
704 #####
705 ## Q2: Box-Cox model for duration

```

```

706 #####
707
708 # Box-Cox for duration
709 boxcox_lambda_dur <- boxcox(model_raw_dur, plotit = TRUE)
710
711 lambda_hat_dur <- boxcox_lambda_dur$x[which.max(boxcox_lambda_dur$y)]
712 lambda_hat_dur
713
714 # Box-Cox transformed response for duration
715 trips_2019_nolog$y_boxcox_dur <- if (abs(lambda_hat_dur) < 1e-3) {
716   log(trips_2019_nolog$'Trip Duration')
717 } else {
718   (trips_2019_nolog$'Trip Duration'^lambda_hat_dur - 1) / lambda_hat_dur
719 }
720
721 model_boxcox_dur <- lm(
722   y_boxcox_dur ~ vehicle_type + 'Trip Distance' + hour_f + dow_f + district_start_f,
723   data = trips_2019_nolog
724 )
725
726 # Residual vs Fitted (Box-Cox duration)
727 plot(fitted(model_boxcox_dur), resid(model_boxcox_dur),
728       xlab = "Fitted values",
729       ylab = "Residuals",
730       main = "Residuals_vs_Fitted_(Box-Cox_duration)")
731 abline(h = 0, lty = 2, col = "red")
732
733 # QQ plot
734 qqnorm(resid(model_boxcox_dur),
735         main = "Normal_Q-Q_Plot_(Box-Cox_duration)")
736 qqline(resid(model_boxcox_dur), col = "red", lty = 2)
737 ''
738
739 **-3-Modeling & Diagnostics: Weighted Least Square (WLS) model**

```

```

740  ``'{r}
741 #####
742 ## Q2: WLS model for log_duration
743 #####
744
745 # Fitted values & residuals from log-duration model
746 f_hat_dur <- fitted(model_log_dur)
747 e_hat_dur <- resid(model_log_dur)
748
749 # Simple variance function: var(e) fitted^2 (same idea as distance)
750 w_simple_dur <- 1 / (f_hat_dur^2 + 1e-8)
751
752 model_wls_dur <- lm(
753   log_duration ~ vehicle_type + log_distance + hour_f + dow_f + district_start_f,
754   data = trips_2019,
755   weights = w_simple_dur
756 )
757
758 # Residual vs Fitted (WLS duration)
759 plot(fitted(model_wls_dur), resid(model_wls_dur),
760       xlab = "Fitted_values",
761       ylab = "Residuals",
762       main = "Residuals_vs_Fitted_(WLS_duration)")
763 abline(h = 0, lty = 2, col = "red")
764
765 # QQ plot
766 qqnorm(resid(model_wls_dur),
767         main = "Normal_Q-Q_Plot_(WLS_duration)")
768 qqline(resid(model_wls_dur), col = "red", lty = 2)
769 ``
770
771 **-4-Modeling & Diagnostics: Huber model**
772 ``'{r}
773 #####

```

```

774 ## Q2: Huber robust regression for log_duration
775 #####
776
777 library(MASS)
778
779 model_hubер_dur <- rlm(
780   log_duration ~ vehicle_type + log_distance + hour_f + dow_f + district_start_f,
781   data = trips_2019
782 )
783
784 # Residual vs Fitted (Huber duration)
785 plot(fitted(model_hubер_dur), resid(model_hubер_dur),
786       xlab = "Fitted_values",
787       ylab = "Residuals_(Huber)",
788       main = "Residuals_vs_Fitted_(Huber,_log_duration)")
789 abline(h = 0, col = "red", lty = 2)
790
791 # QQ plot
792 qqnorm(resid(model_hubер_dur),
793         main = "Normal_Q-Q_Plot_(Huber,_log_duration)")
794 qqline(resid(model_hubер_dur), col = "red", lty = 2)
795
796
797 **-6 Check BIC for candidate models**
798 {r}
799 #####
800 ## Compare BIC of candidate models
801 # - model_raw_dur : duration (raw)
802 # - model_log_dur : log_duration
803 # - model_boxcox_dur : boxCox
804 # - model_wls_dur : WLS (log_duration)
805 # - model_hubер_dur : Huber (log_duration, rlm)
806 #####
807 ## BIC Comparison table

```

```

808 info_dur <- data.frame(
809   Model = c("OLS_raw", "OLS_log", "BoxCox_log", "WLS_log", "Huber_log"),
810   BIC = c(BIC(model_raw_dur),
811           BIC(model_log_dur),
812           BIC(model_boxcox_dur),
813           BIC(model_wls_dur),
814           BIC(model_huber_dur)))
815 )
816 info_dur
817 """
818
819 **-7 Compare Coefficient table (duration, scooter effect)**
820 '''
821 #####
822 ## Q2: Compare vehicle_type coefficients across OLS / WLS / Huber (duration)
823 #####
824
825 coef_table_dur <- data.frame(
826   Vehicle = c("bicycle_(baseline)", "scooter"),
827   OLS = c(0, coef(model_log_dur)[["vehicle_typescooter"]]),
828   BoxCox = c(0, coef(model_boxcox_dur)[["vehicle_typescooter"]]),
829   WLS = c(0, coef(model_wls_dur)[["vehicle_typescooter"]]),
830   Huber = c(0, coef(model_huber_dur)[["vehicle_typescooter"]]))
831 )
832
833 coef_table_dur
834 """
835
836 ** Check Plausible Interaction models (optional candidates)**
837 '''
838 #####
839 # model_int1: log_duration ~ log_distance * vehicle_type + hour_f + dow_f +
     district_start_f

```

```

840 # model_int2: log_duration ~ log_distance + vehicle_type * hour_f + dow_f +
841   district_start_f
842 # model_int3: log_duration ~ log_distance + hour_f + dow_f + vehicle_type *
843   district_start_f
844 #####
845
846
847 # Vehicle_type x log_distance
848 model_int1 <- update(model_log_dur, . ~ . + vehicle_type:log_distance)
849
850 # Vehicle_type x hour_f
851 model_int2 <- update(model_log_dur, . ~ . + vehicle_type:hour_f)
852
853 # Vehicle_type x district_start_f
854 model_int3 <- update(model_log_dur, . ~ . + vehicle_type:district_start_f)
855
856
857
858
859 )
860 info_dur_interaction
861 """
862
863 ** Fit log interaction model**
864 """
865 #####
866 ## Log duration model with interaction
867 ## Response : log_duration
868 ## Predictors : vehicle_type * log_distance + hour_f + dow_f + district_start_f
869 #####
870 model_log_interaction <- lm(log_duration ~ vehicle_type * log_distance +
871   hour_f + dow_f + district_start_f, data = trips_2019)

```

```

872   ' '
873
874 ** Fit log interaction model**
875 ' '{r, log interaction model diagnostics}
876
877 # Residual vs Fitted
878 plot(fitted(model_log_interaction), resid(model_log_interaction),
879       xlab = "Fitted_values",
880       ylab = "Residuals",
881       main = "Residuals_vs_Fitted_(log_duration,interaction_model)")
882 abline(h = 0, lty = 2, col = "red")
883
884 # Q-Q plot
885 qqnorm(resid(model_log_interaction),
886         main = "Normal_Q-Q_Plot_(log_duration,interaction_model)")
887 qqline(resid(model_log_interaction), col = "red", lty = 2)
888 ' '
889
890 ** Check structure of log transformed model**
891 ' '{r, Check structure of log transformed model}
892
893 ## partial residual plot
894 partial_resid <- residuals(model_log_interaction) +
895           coef(model_log_interaction)["log_distance"] * trips_2019$log_distance
896
897 plot(trips_2019$log_distance, partial_resid,
898       xlab = "log_distance",
899       ylab = "Partial_residuals_for_log_distance",
900       main = "Partial_residual_plot_for_log_distance",
901       pch = 16,
902       cex = 0.3,
903       col = adjustcolor("black", alpha.f = 0.2))
904 abline(a = 0,
905        b = coef(model_log_interaction)["log_distance"],

```

```

906     lwd = 3, col ="red")
907 lines(lowess(trips_2019$log_distance, partial_resid),lwd = 2, lty = 2, col = "blue")
908 ''
909
910 ** Model-based outlier, influencial point check**
911 ``-{r}
912 #####
913 ## Q2: Check Model-based outliers & influential points (duration model)
914 #####
915
916 ## Number of observations and coefficients
917 n_dur <- nrow(trips_2019)
918 p_dur <- length(coef(model_log_interaction))
919 alpha <- 0.05
920
921 ## Diagnostic measures
922 student_dur <- rstudent(model_log_interaction) # externally studentized residuals
923 lev_dur <- hatvalues(model_log_interaction) # leverage (hat values)
924 cook_dur <- cooks.distance(model_log_interaction) # Cook's distance
925
926 ## Cutoffs
927 cutoff_student_dur <- qt(1 - alpha/(2*n_dur), df = model_log_interaction$df.residual)
928 cutoff_lev_dur <- 2 * (p_dur + 1) / n_dur
929 cutoff_cook_dur <- 4 / n_dur
930
931 ## Table of diagnostics
932 diag_tbl_dur <- data.frame(
933   id = 1:n_dur,
934   student = student_dur,
935   leverage = lev_dur,
936   cook = cook_dur,
937   flag_student = abs(student_dur) > cutoff_student_dur,
938   flag_leverage = lev_dur > cutoff_lev_dur,
939   flag_cook = cook_dur > cutoff_cook_dur

```

```

940 )
941
942 diag_tbl_dur$flag_any <- with(diag_tbl_dur,
943                               flag_student | flag_leverage | flag_cook)
944
945 suspects_dur <- subset(diag_tbl_dur, flag_any)
946 suspects_dur[order(-suspects_dur$cook), ][1:20, ]
947 ' '
948
949 ** Check Collinearity**
950 ' '{r}
951 #####
952 ## Q2: Check collinearity (VIF) for duration model
953 #####
954
955 library(car)
956 vif(model_log_interaction)
957 ' '
958
959 ** Log transformed model interpretation**
960 ' '{r}
961 #####
962 ## Q2: Scooter vs Bicycle duration effect
963 ## (1) Main-effects model: constant effect across distance
964 ## (2) Interaction model: effect depends on log_distance via 1 + 3 x
965 #####
966
967 # Summary of main-effects & interaction models
968 summary(model_log_dur) # main-effects
969 summary(model_log_interaction) # interaction
970
971 # mean log_distance
972 mean_log_dist <- mean(trips_2019$log_distance, na.rm = TRUE)
973

```

```

974 # Scooter effects from main-effects model: = 1
975 coef_main <- coef(model_log_dur)
976 a1_main <- unname(coef_main["vehicle_typescooter"])
977
978 delta_main <- a1_main
979 ratio_main <- exp(delta_main)
980
981 # Scooter effects from Interaction model: (x) = 1 + 3 x
982 coef_int <- coef(model_log_interaction)
983 a1_int <- unname(coef_int["vehicle_typescooter"]) # 1
984 a3_int <- unname(coef_int["vehicle_typescooter:log_distance"]) # 3
985
986 delta_int_at_mean <- a1_int + a3_int * mean_log_dist
987 ratio_int_at_mean <- exp(delta_int_at_mean)
988
989 # Table of results
990 results_tbl <- data.frame(
991   Model = c("Main-effects_(no_interaction)",
992             "Interaction_(vehicle_type*log_distance)"),
993   a1_scooter = c(a1_main, a1_int),
994   a3_interaction = c(NA, a3_int),
995   log_distance_used = c(NA, mean_log_dist),
996   delta_scooter_log = c(delta_main, delta_int_at_mean),
997   ratio_exp_delta = c(ratio_main, ratio_int_at_mean),
998   pct_diff = 100 * (c(ratio_main, ratio_int_at_mean) - 1)
999 )
1000
1001
1002 results_tbl_round <- within(results_tbl, {
1003   a1_scooter <- round(a1_scooter, 6)
1004   a3_interaction <- round(a3_interaction, 6)
1005   log_distance_used <- round(log_distance_used, 6)
1006   delta_scooter_log <- round(delta_scooter_log, 6)
1007   ratio_exp_delta <- round(ratio_exp_delta, 6)

```

```
1008     pct_diff <- round(pct_diff, 3)
1009 }
1010
1011 results_tbl_round
1012   ````
```