

IMPERIAL COLLEGE OF SCIENCE, TECHNOLOGY AND
MEDICINE

DEPARTMENT OF MECHANICAL ENGINEERING

ME3 Statistics Coursework

Author:

Arohan SUBRAMONIA

CID:

01054062

Course Leader:

Dr. Nicola FITZ-SIMON

March 13, 2018

Introduction

The fuel efficiencies of vehicles are dependent on a number of parameters. With the unique dataset given, the aim of this coursework is to construct linear regression models to aim to estimate fuel efficiencies based on a combination of these parameters. The analysis of this data is given in the following report.

1 Question 1

Exploratory Data Analysis

The following parameters were calculated for the fuel efficiency (litres/100km) data given:

Arithmetic Mean	4.7
Geometric Mean	4.6
Median	4.4
10% Trimmed mean	4.7
Arithmetic standard deviation	1.30
Geometric standard deviation	1.33

Table 1: Exploratory data analysis values

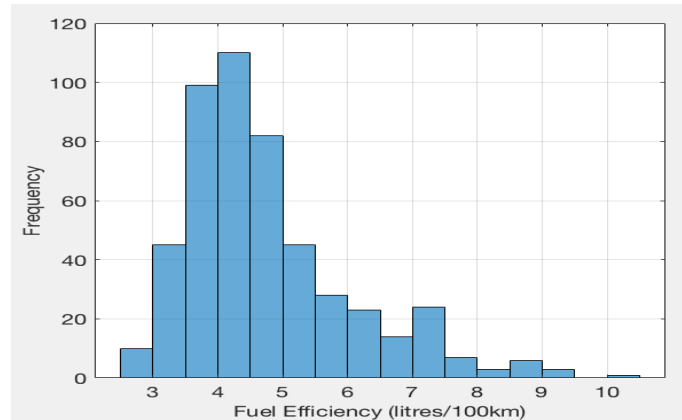


Figure 1: A histogram of the fuel efficiency distribution

The 10% trimmed mean is a better measure than the means or the median, as it is less affected by skewed data and outliers.

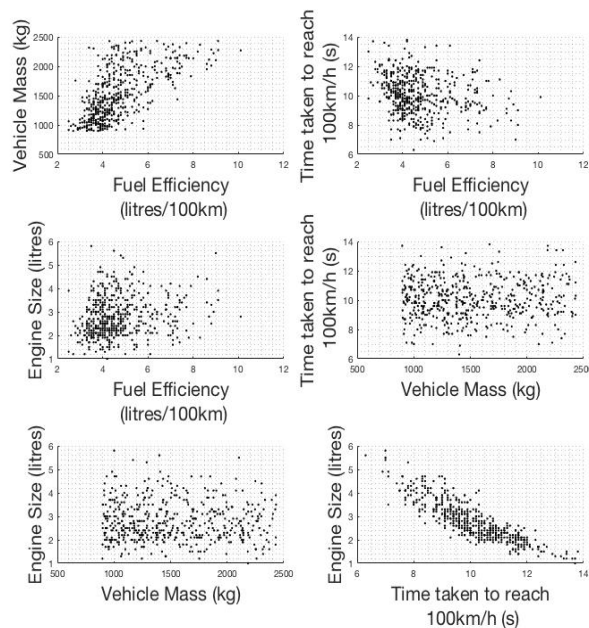


Figure 2: Scatter plots of all continuous variables

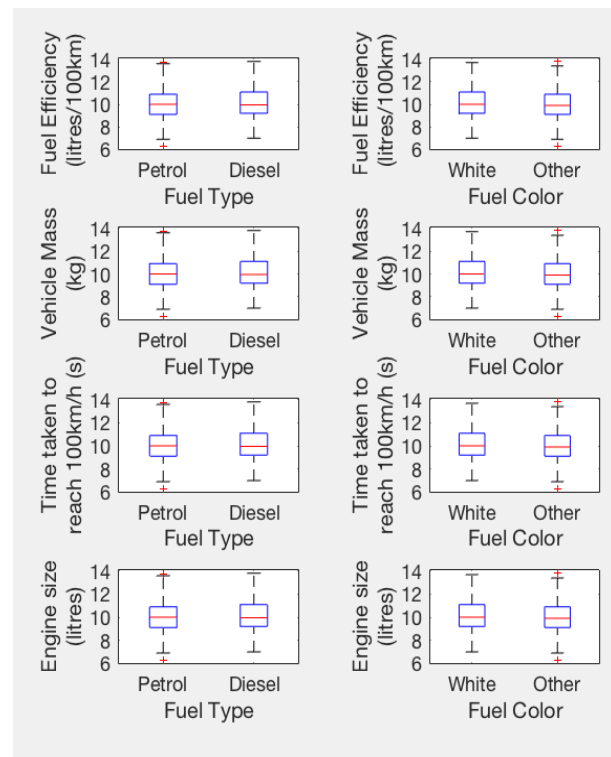


Figure 3: Boxplots of continuous by categorical variables

The scatter plots show that there may be some form of linear relationship between acceleration time and engine size; and potentially some relationship between vehicle mass and fuel efficiency. However, for the other pairs of data, there are not many easily discernible linear (or otherwise) relationships. To satisfy the assumptions of linearity and additivity, all continuous variables should be standardized by their arithmetic means and standard deviations. Therefore, the data was standardized in this way before moving onto the next section.

2 Question 2

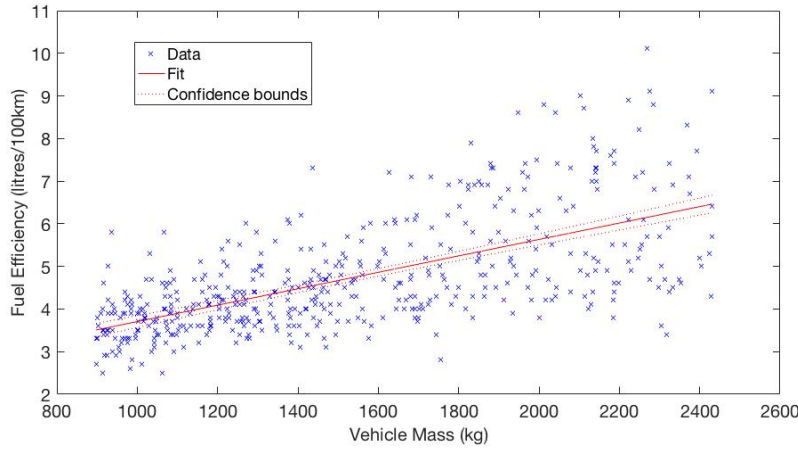


Figure 4: Scatter plot showing linear regression model of fuel efficiency in terms of mass.

This model is a passable initial linear approximation for fuel efficiency in terms of mass, as the two variables seem to have a positive relationship on the graph. This relationship could also be modeled by a quadratic relationship - however, this gives the exact same R squared estimate and larger values for the MSE and AIC.

Prove that R^2 is the square of the sample correlation coefficient $r(x, y)$:

$$\hat{y}_i = \hat{\alpha} + \hat{\beta}x_i \quad \bar{y} = \hat{\alpha} + \hat{\beta}\bar{x} \quad (1)$$

$$R^2 = 1 - \frac{\sum_i^n (y_i - \hat{y}_i)^2}{\sum_i^n (y_i - \bar{y})^2} = \frac{\sum_i^n (y_i - \bar{y})^2 - \sum_i^n (y_i - \hat{y}_i)^2}{\sum_i^n (y_i - \bar{y})^2} \quad (2)$$

For simple linear regression, $\sum (y_i - \bar{y})^2 = \sum (y_i - \hat{y}_i)^2 + \sum (\hat{y}_i - \bar{y})^2$ which means:

$$R^2 = \frac{\sum (\hat{y}_i - \bar{y})^2}{\sum (y_i - \bar{y})^2} = \frac{\sum (\hat{\alpha} + \hat{\beta}x_i - \bar{y})^2}{\sum (y_i - \bar{y})^2} \quad (3)$$

From (1.b),

$$R^2 = \frac{\Sigma(\bar{y} - \hat{\beta}\bar{x} + \hat{\beta}x_i - \bar{y})^2}{\Sigma(y_i - \bar{y})^2} = \frac{\hat{\beta}^2 \Sigma(x_i - \bar{x})^2}{\Sigma(y_i - \bar{y})^2} \quad (4)$$

$$R^2 = \frac{[\Sigma(x_i - \bar{x})(y_i - \bar{y})]^2 \Sigma(x_i - \bar{x})^2}{[\Sigma(x_i - \bar{x})^2 \Sigma(y_i - \bar{y})^2]} \quad (5)$$

Cancelling the $\Sigma(x_i - \bar{x})^2$ terms above, we get:

$$R^2 = \frac{[\Sigma(x_i - \bar{x})(y_i - \bar{y})]^2}{\Sigma(x_i - \bar{x})^2 \Sigma(y_i - \bar{y})^2} = r(x, y)^2 \quad (6)$$

As $r(x, y) = \frac{\Sigma(x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\Sigma(x_i - \bar{x})^2 \Sigma(y_i - \bar{y})^2}}$ which is the correlation coefficient.

Measures of accuracy - a comparison

The **R squared model** is very widely used, and a common predictor of the percentage variance of the model from the actual values.

The **(root) mean squared estimate** is a good parameter of a model as it gives an absolute value for the average residual or error of the data sample. By squaring the errors, we ensure that the mean error does not come out to be zero.

The **Akaike Information Criterion** is similar to the mean squared estimate, except it accounts for the number of parameters in the model, and therefore prevents over-fitting.

First we try and predict values of fuel efficiency with single parameter models for vehicle mass, acceleration time, and engine size. After obtaining R squared, MSE and AIC estimates for these, we try out linear and quadratic multivariate models with and without interactions, with the aim of maximizing the R squared value and minimizing the AIC value for each parameter. The overall results are plotted in the table below:

				Linear regression model: Nt100 ~ 1 + Nmass*Nt100 + Nmass*type + Nt100*type				
Type	Rsquared	MSE	AIC	Estimated Coefficients:				
					Estimate	SE	tStat	pValue
'mass (linear)'	0.40815	0.59304	1159.7	(Intercept)	-0.46943	0.031081	-15.104	1.2352e-42
'mass (quad)'	0.40815	0.59423	1161.7	Nmass	0.30604	0.030579	10.008	1.3695e-21
'time (linear)'	0.043828	0.95809	1399.5	Nt100	-0.1912	0.030667	-6.2348	9.7257e-10
'time (quad)'	0.043843	0.96	1401.5	type	0.88475	0.043975	20.12	3.5955e-66
'engine size (linear)'	0.034128	0.96781	1404.6	Nmass:Nt100	-0.064931	0.021255	-3.0549	0.0023731
'engine size (quad)'	0.039556	0.96431	1403.8	Nmass:type	0.63732	0.044088	14.456	9.6989e-40
'Linear Multivariate'	0.65966	0.34379	891.04	Nt100:type	-0.12849	0.04396	-2.9228	0.003628
'Pure Quadratic Multivariate'	0.66309	0.3424	895.97	Number of observations: 500, Error degrees of freedom: 493				
'Linear Interactions'	0.76933	0.23782	716.55	Root Mean Squared Error: 0.489				
'Quadratic Interactions'	0.77371	0.23476	716.96	R-squared: 0.764, Adjusted R-Squared 0.761				
				F-statistic vs. constant model: 266, p-value = 6.05e-151				

Figure 5: Parameters for various models

Figure 6: Final linear model to predict fuel efficiency

The multivariate quadratic interactions model has 21 terms, which should be refined. Therefore, smallest (absolute) coefficients have been taken out, and using the 'step' function in MATLAB, coefficients are shuffled around to minimize the AIC value of the model. A trade off is made between number of coefficients to involve and R squared value. The final model is given in Figure 6 above.

3 Question 3

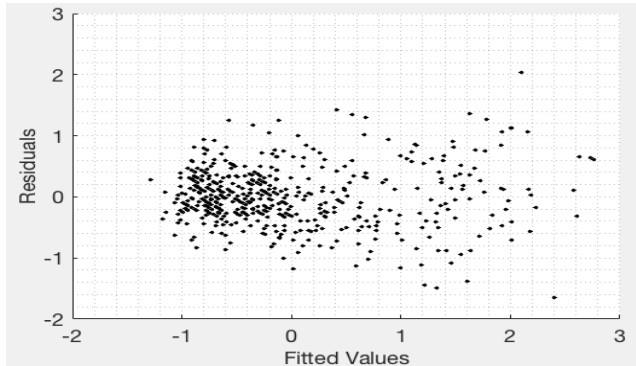


Figure 7: Fitted values vs Residuals

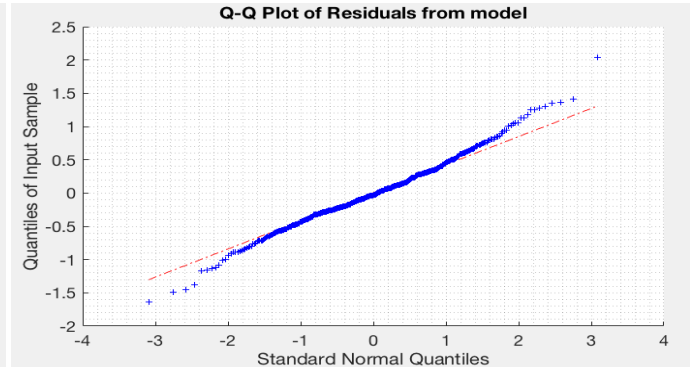


Figure 8: Q-Q plot of residuals

The residuals have no particular relationship with the fitted values, which shows that there is no systematic error throughout the model. Furthermore, it is assumed that the error term in a linear regression model has a distribution of $N(0, 1)$.

Interpreting regression coefficients

The regression coefficients for each of the predictor terms shows how much change there is in fuel efficiency when there is one unit of change in the predictor. Furthermore, the 'intercept' value is the value of fuel efficiency when all the predictor terms are set to zero. To make these coefficients more meaningful, all the values could be arithmetically standardized¹ which would mean the coefficients could be directly compared. This has already been done for the model in Question 2, and so is not necessary to do again.

4 Question 4

For the bootstrapping part of the assignment, the instructions (as given in the coursework handout) were followed. A new response variable was calculated using bootstrapped residuals from the model in Question 2 the same model was used to predict this new data. The aforementioned model is stated below:

$$l100 \sim 1 + mass + t100 + type + mass : t100 + mass : type + t100 : type$$

This process was repeated 10 times, and each time an array of fuel efficiency estimates (predicted using the linear model) for each acceleration time value, was calculated. Then, the mean value and the 0.025

¹To standardize a term x_i , one can do

$$z_i = \frac{x_i - \mu_x}{\sigma_x}$$

and 0.975 quantiles for each acceleration time were calculated, and are plotted on the following page (Figures 9 and 10).

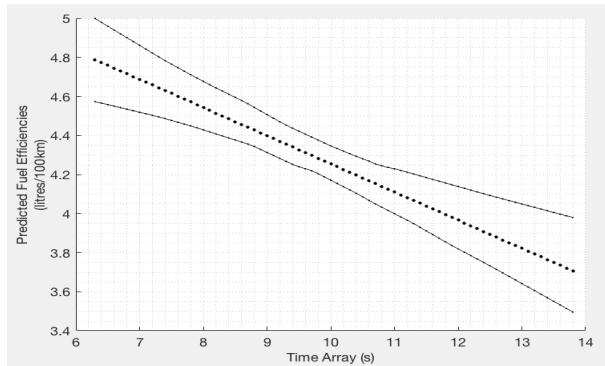


Figure 9: Predicted fuel efficiency by acceleration time, and bootstrapped 95% confidence band.

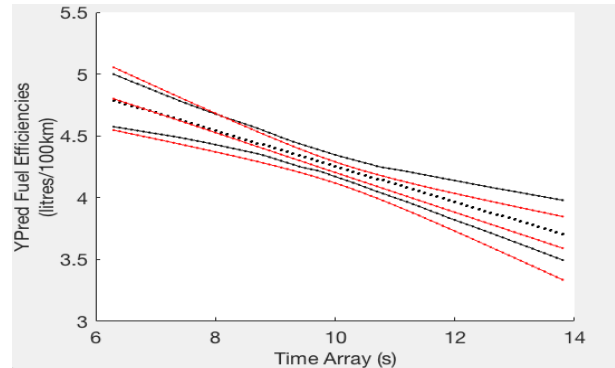


Figure 10: Model using 'predict' function (red) plotted alongside bootstrapped model.

It was found that, as the number of repeats for the process increased, the bootstrapped data would be more similar to the model obtained by the 'predict' function on MATLAB.

5 Question 5

Fitting the model from Question 2 to the test data was yielded a range of results. These are given below"

Training Data			Test Data		
R Squared	MSE	AIC	R Squared	MSE	AIC
0.764	0.2391	710.4619	0.7638	0.4061	975.3545

Therefore, the test data did not perform as well as the training data did; however, the values are not far off each other. Furthermore, the mean squared error value for the test data is larger than the same value for the training data, which suggests that there is little to no overfitting. Based on this evidence, it may be worth considering a slightly more complex model with more interaction terms.

By creating a linear model using this data, with quadratic terms and interactions (and no refining), we are able to achieve an R squared value of 0.774. However, the MSE is still larger than the simplified training data value (test data MSE is 0.3987) and the AIC value is 981.855, which is also higher than the corresponding value for the training data

6 Appendix - Code

```

1
2 %% ME3 STATS COURSEWORK
3 %AROHAN SUBRAMONIA
4 rng(01054062);
5
6 clear all
7 %% Reading the data from the csv file
8 AS10415 = readtable('as10415.csv');
9 DATA=table2cell(AS10415);
10 num=cell2mat(DATA(:,1:4));
11 num(:,5) = string(DATA(:,5)) == 'petrol'; %Fuel type logical (0 if
    Diesel, 1 if Petrol)
12 num(:,6) = string(DATA(:,6)) == 'white'; %Colour logical (1 if white,
    0 if other)
13
14 % Splitting the data up into four sections
15 l100=num(:,1);
16 mass=num(:,2);
17 t100=num(:,3);
18 disp=num(:,4);
19
20 %% Question 1
21 %% Exploratory Data Analysis
22 A_mean(1:4) = mean(num(:,1:4)); %Arithmetic mean
23 G_mean(1:4) = geomean(num(:,1:4)); %Geometric mean
24 median(1:4) = median(num(:,1:4)); %Median
25 T_mean(1:4) = mean(num(51:450,1:4)); %10% trimmed mean
26 A_std(1:4) = std(num(:,1:4)); %Arithmetic standard deviation
27 top(1:4) = sum(log(num(:,1:4)/G_mean).^2); %Geometric standard
    deviation
28 G_std(1:4) = exp((top/499).^(0.5));
29
30 %% Histogram of Fuel Efficiency
31 histogram(num(:,1));
32 grid ON
33 fs=20; %FontSize
34 %title('Histogram of Fuel Efficiency per 100km','FontSize',18);
35 set(gca,'FontSize',fs)
36 xlabel('Fuel Efficiency (litres/100km)','FontSize',fs);
37 ylabel('Frequency','FontSize',fs);
38
39 %% Scatterplots
40 subplot(3,2,1), scatter(l100,mass,'.','k'), xlabel({'Fuel Efficiency';'
    (litres/100km)'},'FontSize',fs), ylabel('Vehicle Mass (kg)',
    'FontSize',fs), grid MINOR

```

```

41 subplot(3,2,2), scatter(l100,t100, '.', 'k'), xlabel({'Fuel Efficiency';'
    (litres/100km)'}), 'FontSize', fs), ylabel({'Time taken to reach';'
    100km/h (s)'}), 'FontSize', fs), grid MINOR
42 subplot(3,2,3), scatter(l100, disp, '.', 'k'), xlabel({'Fuel Efficiency';'
    (litres/100km)'}), 'FontSize', fs), ylabel('Engine Size (litres)', '
    FontSize', fs), grid MINOR
43 subplot(3,2,4), scatter(mass,t100, '.', 'k'), xlabel('Vehicle Mass (kg)',
    'FontSize', fs), ylabel({'Time taken to reach';' 100km/h (s)'}), '
    FontSize', fs), grid MINOR
44 subplot(3,2,5), scatter(mass, disp, '.', 'k'), xlabel('Vehicle Mass (kg)',
    'FontSize', fs), ylabel('Engine Size (litres)', 'FontSize', fs), grid
    MINOR
45 subplot(3,2,6), scatter(t100, disp, '.', 'k'), xlabel({'Time taken to
    reach';' 100km/h (s)'}), 'FontSize', fs), ylabel('Engine Size (litres)'
    , 'FontSize', fs), grid MINOR
46
47 %% Boxplot drawing routine
48 fueltype=[['Petrol'] ['Diesel']]; {'White' ['Other']]; %Box labels
49
50 XLabel=[{'Fuel Type'} {'Fuel Color'}];
51 YLabel=[{'Fuel Efficiency';' (litres/100km)'};{'Vehicle Mass';'(kg)
    '};{'Time taken to';'reach 100km/h (s)'};{'Engine size';'(litres)'
    }];
52 yl=[1,1,3,3,5,5,7,7];
53 var=[1,1,2,2,3,3,4,4];
54 fs=18;
55 b=[1,2,1,2,1,2,1,2];
56 for a=1:8
57     subplot(4,2,a), boxplot(num(:, var(a)), num(:, b(a)+4), 'labels',
        fueltype(b(a), :));
58     xlabel(XLabel(b(a)), 'FontSize', fs);
59     ylabel(YLabel(yl(a):yl(a)+1), 'FontSize', fs);
60     set(gca, 'FontSize', fs)
61 end
62
63 %% Question 2
64 %% Standardizing and Organising data
65 NI100=zscore(num(:,1));
66 Nmass=zscore(num(:,2));
67 Nt100=zscore(num(:,3));
68 Ndisp=zscore(num(:,4));
69 type = num(:,5);
70 colour=num(:,6);
71 %Creating a table of standardized variables
72 Nvarnames(1:6)=[{'NI100'} {'Nmass'} {'Nt100'} {'Ndisp'} {'type'} {'
    colour'}];
73 Nnum=table(NI100, Nmass, Nt100, Ndisp, num(:,5), num(:,6), 'VariableNames',

```



```

Nvarnames);
74
75 %% Linear Models – single variables without interaction
76
77 %Linear model of l100 vs mass
78 LMmass1=fitlm(mass,l100,'linear');
79 plot(LMmass1);
80 set(gca,'FontSize',fs)
81 ylabel('Fuel Efficiency (litres/100km)','FontSize',fs);
82 xlabel('Vehicle Mass (kg)','FontSize',fs);
83 title(' ');
84
85 %Quadratic model of l100 vs mass
86 LMmass2=fitlm(mass,l100,'purequadratic');
87 subplot(), plot(LMmass2);
88 set(gca,'FontSize',fs)
89 ylabel('Fuel Efficiency (litres/100km)','FontSize',fs);
90 xlabel('Vehicle Mass (kg)','FontSize',fs);
91
92 %Linear model of l100 vs t100
93 LMtime1=fitlm(t100,l100,'linear');
94 plot(LMtime1);
95 set(gca,'FontSize',fs)
96 ylabel('Fuel Efficiency (litres/100km)','FontSize',fs);
97 xlabel('Time taken to reach 100km/h (s)','FontSize',fs);
98
99 %Quadratic model of l100 vs t100
100 LMtime2=fitlm(t100,l100,'purequadratic');
101 plot(LMtime2);
102 set(gca,'FontSize',fs)
103 ylabel('Fuel Efficiency (litres/100km)','FontSize',fs);
104 xlabel('Time taken to reach 100km/h (s)','FontSize',fs);
105
106 %Linear model of l100 vs disp
107 LMdisp1=fitlm(disp,l100,'linear');
108 plot(LMdisp1);
109 set(gca,'FontSize',fs)
110 ylabel('Fuel Efficiency (litres/100km)','FontSize',fs);
111 xlabel('Engine Size (litres)','FontSize',fs);
112
113 %Quadratic model of l100 vs disp
114 LMdisp2=fitlm(disp,l100,'purequadratic');
115 plot(LMdisp2);
116 set(gca,'FontSize',fs)
117 ylabel('Fuel Efficiency (litres/100km)','FontSize',fs);
118 xlabel('Engine Size (litres)','FontSize',fs);
119

```

```

120 %% Parameters for single variable models without interaction (not
    standardized)
121 %Parameters for l100 vs mass – linear
122 Mass_linear=fitlm(Nmass,NI100,'linear');
123 RSm1=Mass_linear.Rsquared.Ordinary;
124 MSEM1=Mass_linear.MSE;
125 AICm1=Mass_linear.ModelCriterion.AIC;
126
127 %Parameters for l100 vs mass – quadratic
128 Mass_quad=fitlm(Nmass,NI100,'purequadratic');
129 RSm2=Mass_quad.Rsquared.Ordinary;
130 MSEM2=Mass_quad.MSE;
131 AICm2=Mass_quad.ModelCriterion.AIC;
132
133 %Parameters for l100 vs time – linear
134 Time_linear=fitlm(Nt100,NI100,'linear');
135 RSt1=Time_linear.Rsquared.Ordinary;
136 MSET1=Time_linear.MSE;
137 AICt1=Time_linear.ModelCriterion.AIC;
138
139 %Parameters for l100 vs time – quadratic
140 Time_quad=fitlm(Nt100,NI100,'purequadratic');
141 RSt2=Time_quad.Rsquared.Ordinary;
142 MSET2=Time_quad.MSE;
143 AICt2=Time_quad.ModelCriterion.AIC;
144
145 %Parameters for l100 vs disp – linear
146 Disp_linear=fitlm(Ndisp,NI100,'linear');
147 RSd1=Disp_linear.Rsquared.Ordinary;
148 MSED1=Disp_linear.MSE;
149 AICd1=Disp_linear.ModelCriterion.AIC;
150
151 %Parameters for l100 vs disp – quadratic
152 Disp_quad=fitlm(Ndisp,NI100,'purequadratic');
153 RSd2=Disp_quad.Rsquared.Ordinary;
154 MSED2=Disp_quad.MSE;
155 AICd2=Disp_quad.ModelCriterion.AIC;
156
157 RES_RS=[RSm1;RSm2;RSt1;RSt2;RSd1;RSd2];
158 RES_MSE=[MSEM1;MSEM2;MSET1;MSET2;MSED1;MSED2];
159 RES_AIC=[AICm1;AICm2;AICt1;AICt2;AICd1;AICd2];
160
161 %Final table of single variable parameters
162 modelname={'mass (linear)','mass (quad)','time (linear)','time (quad)',
    'engine size (linear)','engine size (quad)'}';
163 estnames={'Type' 'Rsquared' 'MSE' 'AIC'};
164 statistics=table(modelname,RES_RS,RES_MSE,RES_AIC,'VariableNames',

```

```

    estnames);
165
166 %% Linear Models – multivariable, some with interactions
167
168 %Linear multivariate model, no interactions
169 Linear_Model = fitlm(Nnum, 'linear', 'PredictorVars', {'Nmass', 'Nt100', '
    Ndisp', 'type', 'colour'}, 'ResponseVar', 'NI100');
170 lintab=table({'Linear Multivariate'}, Linear_Model.Rsquared.Ordinary,
    Linear_Model.MSE, Linear_Model.ModelCriterion.AIC, 'VariableNames',
    estnames);
171
172 %Quadratic multivariate model, no interactions
173 Quad_Model = fitlm(Nnum, 'purequadratic', 'PredictorVars', {'Nmass', 'Nt100
    ', 'Ndisp', 'type', 'colour'}, 'ResponseVar', 'NI100');
174 quadtab=table({'Pure Quadratic Multivariate'}, Quad_Model.Rsquared.
    Ordinary, Quad_Model.MSE, Quad_Model.ModelCriterion.AIC, 'VariableNames
    ', estnames);
175
176 %Linear multivariate model, with interactions
177 Inter_Model1 = fitlm(Nnum, 'interactions', 'PredictorVars', {'Nmass', '
    Nt100', 'Ndisp', 'type', 'colour'}, 'ResponseVar', 'NI100');
178
179 %Quadratic multivariate model, with interactions
180 Inter_Model2 = fitlm(Nnum, 'quadratic', 'PredictorVars', {'Nmass', 'Nt100',
    'Ndisp', 'type', 'colour'}, 'ResponseVar', 'NI100');
181
182 inter1tab=table({'Linear Interactions'}, Inter_Model1.Rsquared.Ordinary,
    Inter_Model1.MSE, Inter_Model1.ModelCriterion.AIC, 'VariableNames',
    estnames);
183 inter2tab=table({'Quadratic Interactions'}, Inter_Model2.Rsquared.
    Ordinary, Inter_Model2.MSE, Inter_Model2.ModelCriterion.AIC, '
    VariableNames', estnames);
184
185 %Final table of all models
186 statistics(7:10, 1:4)=[lintab; quadtab; inter1tab; inter2tab];
187
188 %% Refining the best performing (quadratic) model
189
190 LM2=Inter_Model2;
191 %Removing the least significant coefficients
192 for a=1:(0.5*length(LM2.Coefficients.Estimate))
193     [X, I]=min(abs(LM2.Coefficients.Estimate));
194     if I ~= 1
195         LM2=removeTerms(LM2, char(LM2.CoefficientNames(I)));
196     end
197 end
198

```

```

199 for a=1:10
200     LM2=step(LM2);
201 end
202
203 %% Adding terms in to finalise model
204 mdl_eqn2 = ('NI100~1+Nmass+Nt100+type+Nmass:Nt100+Nmass:type+Nt100:type
    ');
205 LM2=fitlm(Nnum, mdl_eqn2);
206
207 %% Question 3
208 %% Residual and Fitted Plot
209 res=LM2.Residuals.Raw;
210 fit=LM2.Fitted;
211
212 scatter(fit,res,10,'k','MarkerFaceColor','k','MarkerEdgeColor','k');
213 grid MINOR;
214 set(gca,'FontSize',fs);
215 xlabel('Fitted Values','FontSize',fs);
216 ylabel('Residuals','FontSize',fs);
217
218 %% QQ Plot of residuals
219 qqplot(res);
220 grid MINOR;
221 title('Q-Q Plot of Residuals from model','FontSize',fs);
222 set(gca,'FontSize',fs);
223
224 %% Coefficient Plot
225 coeff=LM2.Coefficients.Estimate;
226
227 plot(coeff,'k.-','MarkerSize',15);
228 grid MINOR;
229 set(gca,'FontSize',fs);
230 xlabel('Coefficients index','FontSize',fs);
231 ylabel('Coefficients values','FontSize',fs);
232
233 %% Question 4
234 %% Bootstrapping
235 mint=min(t100);
236 maxt=max(t100);
237 time=[mint:0.1:maxt];
238 res1=LM2.Residuals.Raw;
239 l100hat=LM2.Fitted;
240 res=l100-l100hat;
241
242 %Using lmstar, predict
243 for j=1:10
244     res_star = randsample(res,length(res),'true');

```

```

245 l100star=LM2.Fitted+res_star;
246
247 % Same model as Q2 fitted to Nl100star
248 mdl_eqnstar = ('l100star~1+mass+t100+type+mass:t100+mass:type+t100:
    type');
249 numstar=table(l100star, mass, t100, disp, num(:,5), num(:,6), '
    VariableNames', {'l100star' 'mass' 't100' 'disp' 'type' 'colour'
    });
250 LM2star=fitlm(numstar, mdl_eqnstar);
251
252 for i=1:length(time)
253     l100starhat(i)=LM2star.Coefficients.Estimate(1)+...
254         LM2star.Coefficients.Estimate(2).*mean(mass)+...
255         LM2star.Coefficients.Estimate(3).*time(i)+...
256         LM2star.Coefficients.Estimate(4).*min(type)+...
257         LM2star.Coefficients.Estimate(5).*(mean(mass).*time(i))+...
258         LM2star.Coefficients.Estimate(6).*(mean(mass).*min(type))
259         +...
260         LM2star.Coefficients.Estimate(5).*(time(i).*min(type));
261
262     l100stardata(:,j)=l100starhat';
263 end
264
265 l100starmean=mean(l100stardata');
266 lower = quantile(l100stardata,0.025,2);
267 upper = quantile(l100stardata,0.975,2);
268
269 scatter(time, l100starmean, 5, 'MarkerFaceColor', 'k', 'MarkerEdgeColor', 'k'
    );
270 hold;
271 plot(time, lower, 'k.-', time, upper, 'k.-');
272 grid MINOR
273 set(gca, 'FontSize', fs)
274 xlabel('Time Array (s)', 'FontSize', fs);
275 ylabel({'Predicted Fuel Efficiencies'; '(litres/100km)'}, 'FontSize', fs);
276
277 %% Using a separate model to predict confidence intervals
278
279 xnew=zeros(length(time),5);
280 xnew(:,1)=mean(mass);
281 xnew(:,2)=time;
282 xnew(:,3)=mean(disp);
283
284 [ypred, yci]=predict(LM2star, xnew);
285
286 plot(time, ypred, 'r.-', time, yci, 'r.-');

```

```

287 grid MINOR
288 set(gca, 'FontSize', fs)
289 xlabel('Time Array (s)', 'FontSize', fs);
290 ylabel({'YPred Fuel Efficiencies'; '(litres/100km)'}, 'FontSize', fs);
291
292 %% Question 5
293 %% Reading the data from the csv file
294 testdata = readtable('testdata.csv');
295 TDATA=table2cell(AS10415);
296 tnum=cell2mat(TDATA(:,1:4));
297 tnum(:,5) = string(TDATA(:,5)) == 'petrol'; %Fuel type logical (0 if
        Diesel, 1 if Petrol)
298 tnum(:,6) = string(TDATA(:,6)) == 'white'; %Colour logical (1 if
        white, 0 if other)
299
300 % Splitting the data up into six sections
301 TI100=tnum(:,1);
302 Tmass=tnum(:,2);
303 Tt100=tnum(:,3);
304 Tdisp=tnum(:,4);
305 Ttype=tnum(:,5);
306 Tcolour=tnum(:,6);
307
308 %% Using Question 2 model on Test data
309 Tvarnames(1:6)=[{'TI100'} {'Tmass'} {'Tt100'} {'Tdisp'} {'Ttype'} {'
        Tcolour'}];
310 Tnum=table(TI100, Tmass, Tt100, Tdisp, Ttype, Tcolour, 'VariableNames',
        Tvarnames);
311
312 mdl_eqnT = ('TI100~1+Tmass+Tt100+Ttype+Tmass:Tt100+Tmass:Ttype+Tt100:
        Ttype');
313 LMT=fitlm(Tnum, mdl_eqnT);
314
315 %Quadratic Interactions model on test data
316 Test_Inter_Quad_Model=fitlm(Tnum, 'quadratic', 'PredictorVars', {'Tmass' '
        Ttype' 'Tt100' 'Tdisp' 'Ttype' 'Tcolour'}, 'ResponseVar', 'TI100');

```