

Factor Analysis on CO2 Emissions from Vehicles

Author: Zhizhi Wei

Discussants:

Websites used include **International Energy Agency's Official Website** and **dplyr Recode Reference**. I have discussed my analyses part with Prof. Meyers.

Introduction

As global concern over climate change grows, there is an increasing focus on the carbon footprint generated by human activities. Among these, CO2 emissions from vehicles stand out as a major contributor to global warming. This brings to a pivotal question: What are the primary factors influencing CO2 emissions by vehicles? Could it be related to engine characteristics, fuel choice, or vehicle's internal structure? Bearing this question in mind, this report aims to quantify and analyze the key factors that would significantly affect CO2 emissions from vehicles, and based on these, to develop strategies for effectively reducing the emissions.

This factor analysis is important because understanding the determinants of vehicle CO2 emissions could provide us a scientific basis for informed policy design and technical innovation. Identifying those key factors also helps in directing consumer choices regarding vehicles and fuels. Meanwhile, despite the rising popularity of electric vehicles, the global market is still dominated by traditional internal combustion engine vehicles. As of last year, these traditional vehicles accounted for approximately 85% of total car sales, which emphasizes the importance of targeting them in CO2 emission reduction strategies.

The dataset originates from **Canada Government Official open data website**. The version I will be using in this report is a compiled one available on **Kaggle**. This contains data over a period of 7 years and includes 7385 observations, each representing a vehicle. It captures the details of how CO2 emissions by a vehicle can vary with different features. This is my first time working with this dataset, and it is not being used for another research project or class. Hundreds of analyses have been done with the data, mostly concentrating on EDA or applying machine learning techniques, one example is **CO2 Emission EDA & Visualization & ML**. I haven't delved deeply into these analyses; I just briefly skimmed through some of them to gain a basic understanding of the dataset and identify potential limitations. I noticed that none of these analyses had performed comprehensive model selection for multiple regression models or adequately quantified the significance of influencing factors. These aspects are what my report will be focusing on.

Results

Data wrangling:

```
emission_df <- emission |>
  select(make, model, vehicle_class, engine_size, cylinders, transmission, fuel_type,
         fuel_consump_comb, co2_emission) |>
  filter(fuel_type %in% c("D", "E", "X", "Z")) |>
  mutate(transmission = recode(transmission, A4 = "A", A5 = "A", A6 = "A", A7 = "A",
```

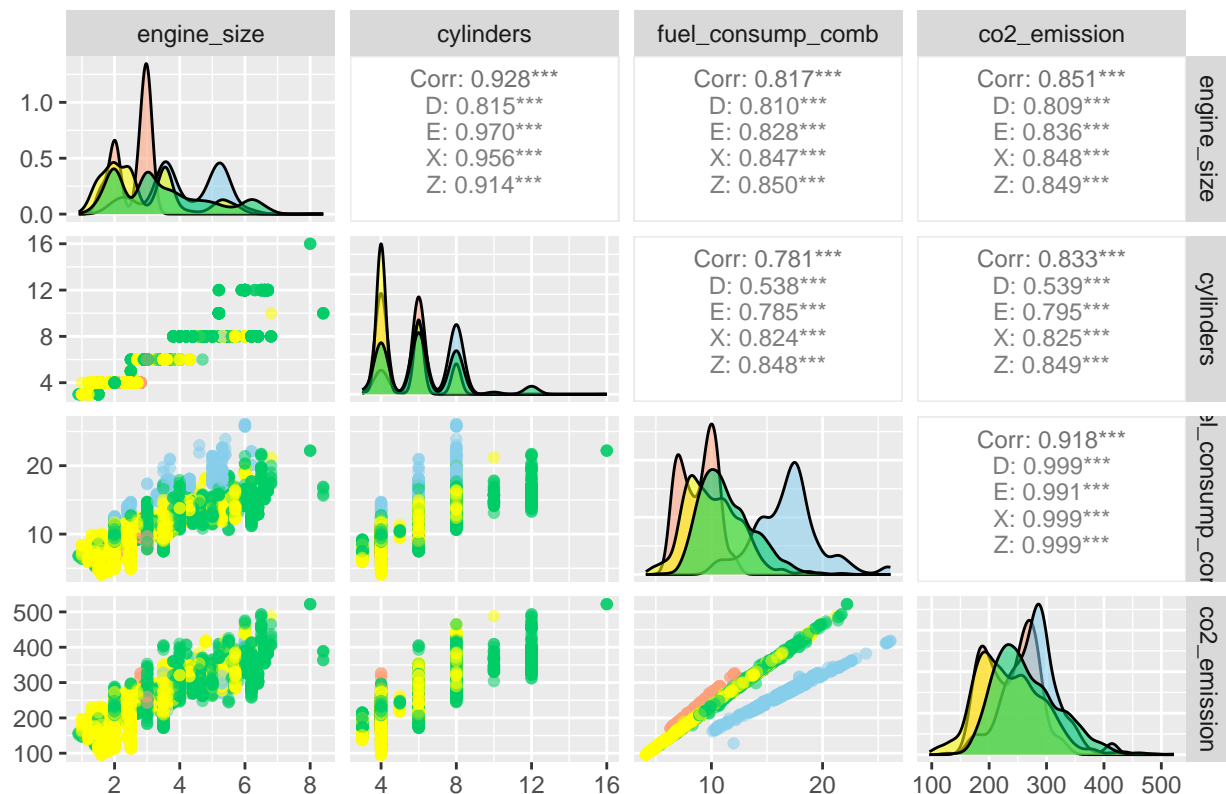
```
A8 = "A", A9 = "A", A10 = "A", AM5 = "AM", AM6 = "AM", AM7 = "AM", AM8 = "AM",
AM9 = "AM", AS4 = "AS", AS5 = "AS", AS6 = "AS", AS7 = "AS", AS8 = "AS", AS9 = "AS",
AS10 = "AS", AV6 = "AV", AV7 = "AV", AV8 = "AV", AV10 = "AV", M5 = "M", M6 = "M",
M7 = "M"))
```

The data is relatively clean for analysis. I checked that there are no missing values (NA) or erroneous outliers, so what I did was mainly for more efficient data processing. I first renamed columns to avoid processing errors (see appendix i). Next, I removed the variables `fuel_consump_city`, `fuel_consump_hwy`, `fuel_consump_comb_mpg`, they are redundant because the variable `fuel_consump_comb` adequately represents fuel consumption. According to dataset description: `fuel_consump_comb` = 55% `fuel_consump_city` + 45% `fuel_consump_hwy`; and `fuel_consump_comb_mpg` merely converts `fuel_consump_comb` into a different unit, with both displaying nearly identical distribution shapes (see appendix ii). Then, I filtered out the fuel type N, as it contained only one single observation, lacking sufficient generality for analysis. Lastly, I simplified the variable `transmission` by categorizing its detailed entries into five primary groups: A, AM, AS, AV, and M, following the categorization outlined in dataset description.

Visualize the Data: Correlation between Quantitative Variables, Grouped by Fuel Type

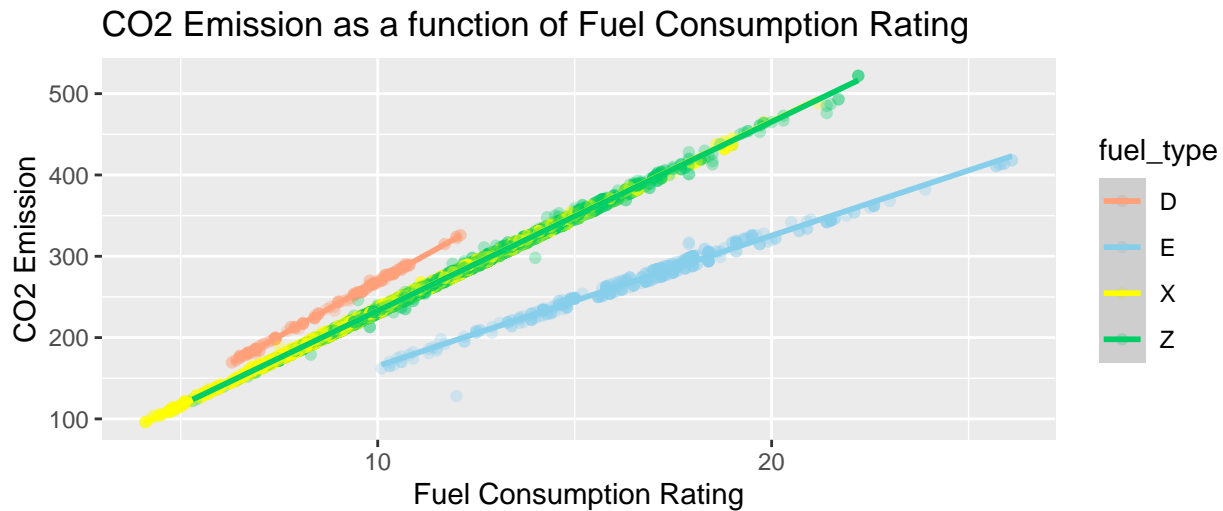
```
emission_df |>
  select(engine_size, cylinders, fuel_consump_comb, co2_emission,
         fuel_type) |>
  ggpairs(columns = 1:4, aes(color = fuel_type, alpha = 0.1),
         upper = list(continuous = wrap("cor", size = 3))) + scale_color_manual(values = color) +
  scale_fill_manual(values = color) + ggtitle("Pairs Plot")
```

Pairs Plot



Firstly, I generated a pairs plot to visualize the relationship among quantitative variables in the dataset. This plot provides three key insights for my subsequent analysis: (1) The response variable `co2_emissions` exhibits strong positive correlations with the explanatory variables `engine_size`, `cylinders`, and `fuel_consump_comb`, suggesting that they are appropriate to be added as predictors in regression model; (2) There are high correlations between the pairs `engine_size` & `cylinders`, `engine_size` & `fuel_consump_comb`, indicating potential multicollinearity issues in model fitting; (3) The scatterplot for `fuel_consump_comb` versus `co2_emission` displays several types of linear patterns with different slopes, implying that the relationship between CO2 emissions and fuel consumption rating may vary depending on `fuel_type`.

```
emission_df |>
  ggplot(aes(x = fuel_consump_comb, y = co2_emission, col = fuel_type)) +
  geom_point(alpha = 0.3) + geom_smooth(method = "lm", se = TRUE,
    level = 0.95) + xlab("Fuel Consumption Rating") + ylab("CO2 Emission") +
  ggtitle("CO2 Emission as a function of Fuel Consumption Rating") +
  scale_colour_manual(values = color) + theme(aspect.ratio = 0.4)
```



To further explore the relationship between `fuel_consump_comb` and `co2_emission`, I enlarged the scatter-plot specifically for these two variables, and added smooth lines to the data in each `fuel_type` category. This plot clearly shows varying slopes in the relationship between CO2 emissions and fuel consumption rating across different fuel types. Therefore, incorporating an interaction term between `fuel_consump_comb` and `fuel_type` in the regression model could be taken into consideration.

Analyses: Model Selection and Influencing Factor Quantification in Multiple Regression

In order to determine and quantify the key factors influencing vehicle CO2 emissions, I build 4 nested multiple regression models. Based on the key insights from previous plot analyses, my initial model incorporates three quantitative variables `engine_size`, `cylinders`, `fuel_consump_comb`. Then, I introduce `fuel_type`, `fuel_type * fuel_consump_comb` and `transmission` into the model one by one:

```
lm_fit_1 <- lm(co2_emission ~ engine_size + cylinders + fuel_consump_comb, data = emission_df)
lm_fit_2 <- lm(co2_emission ~ engine_size + cylinders + fuel_consump_comb + fuel_type,
  data = emission_df)
lm_fit_3 <- lm(co2_emission ~ engine_size + cylinders + fuel_consump_comb + fuel_type +
```

```

    fuel_type * fuel_consump_comb, data = emission_df)
lm_fit_4 <- lm(co2_emission ~ engine_size + cylinders + fuel_consump_comb + fuel_type +
    fuel_type * fuel_consump_comb + transmission, data = emission_df)
models <- list(lm_fit_1, lm_fit_2, lm_fit_3, lm_fit_4)
sapply(models, summary)[8:9, 1:4]

```

```

##           [,1]      [,2]      [,3]      [,4]
## r.squared    0.8793584 0.9913684 0.9974744 0.9974842
## adj.r.squared 0.8793094 0.9913614 0.9974713 0.9974798

```

```
sapply(models, AIC)
```

```
## [1] 65442.41 45973.91 36905.42 36884.47
```

```
sapply(models, BIC)
```

```
## [1] 65476.94 46029.17 36981.40 36988.08
```

Comparing the statistics for model fit R -squared, adjusted R -squared, AIC, BIC, it is clear that Model 3 and 4 better fit the data than Model 1 and 2. When comparing Model 3 with Model 4, I find that they both have very similar values for R -squared and adjusted R -squared, which is around 0.997; AIC indicates a preference for Model 4, while BIC favors model 3.

In order to choose between model 3 and model 4, I will run an ANOVA test based on the F-statistic:

$$H_0 : \beta_{transmissionAM} = \beta_{transmissionAS} = \beta_{transmissionAV} = \beta_{transmissionM} = 0$$

H_A : At least one of the β_i above is not 0

```
anova(lm_fit_3, lm_fit_4)
```

```

## Analysis of Variance Table
##
## Model 1: co2_emission ~ engine_size + cylinders + fuel_consump_comb +
##   fuel_type + fuel_type * fuel_consump_comb
## Model 2: co2_emission ~ engine_size + cylinders + fuel_consump_comb +
##   fuel_type + fuel_type * fuel_consump_comb + transmission
##   Res.Df  RSS Df Sum of Sq    F    Pr(>F)
## 1    7374 63847
## 2    7370 63597   4    249.82 7.2378 8.199e-06 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

```

The ANOVA table shows p -value is very small and close to 0, indicating that adding the variable `transmission` leads to a statistically significant increase in the amount of variability explained. Therefore, Model 4 gives a better fit to the data compared to Model 3.

The coefficients in Model 4 is shown as follows:

```
summary(lm_fit_4)$coefficients
```

	Estimate	Std. Error	t value	Pr(> t)
## (Intercept)	-0.71903279	1.29596876	-0.5548226	5.790328e-01
## engine_size	-0.01966442	0.07789278	-0.2524549	8.006965e-01
## cylinders	0.29652205	0.05613882	5.2819431	1.314767e-07
## fuel_consump_comb	26.79510749	0.14426436	185.7361588	0.000000e+00
## fuel_typeE	5.20707795	1.57256613	3.3111981	9.334421e-04
## fuel_typeX	1.31049373	1.30884974	1.0012561	3.167359e-01
## fuel_typeZ	1.91647245	1.31946634	1.4524603	1.464162e-01
## transmissionAM	0.12346974	0.14663481	0.8420220	3.998029e-01
## transmissionAS	0.22395587	0.09258554	2.4189076	1.559126e-02
## transmissionAV	-0.32956838	0.16045352	-2.0539804	4.001279e-02
## transmissionM	-0.23615784	0.12079509	-1.9550285	5.061745e-02
## fuel_consump_comb:fuel_typeE	-10.86097519	0.15290828	-71.0293464	0.000000e+00
## fuel_consump_comb:fuel_typeX	-3.68201485	0.14519528	-25.3590540	4.281307e-136
## fuel_consump_comb:fuel_typeZ	-3.77113677	0.14591795	-25.8442277	4.679170e-141

Based on the summary table of Model 4 above, the p -values for coefficients of `engine_size`, `fuel_typeX`, `fuel_typeZ`, `transmissionAM` are above 0.05, meaning that these coefficients are not significantly different from 0 at 5% level. `Fuel_typeX`, `fuel_typeZ` and `transmissionAM` are levels of categorical variables, their zero coefficients imply that their model intercepts are equivalent to the intercept of their reference category, which we would discuss later. The statistical insignificance of the `engine_size` coefficient is, however, a matter of concern; it suggests we need to reevaluate whether to include `engine_size` or not.

To address this problem, I will compute the 95% confidence interval for $\beta_{engine-size}$.

- 95% Confidence Interval for $\beta_{engine-size}$ using t-distribution:

```
confint(lm_fit_4)[2, ]
```

```
##      2.5 %      97.5 %
## -0.1723565  0.1330277
```

- 95% Confidence Interval for $\beta_{engine-size}$ using bootstrap to double check:

```
num_cases <- dim(emission_df)[1]
boot_dist <- NULL
for (i in 1:1000) {
  boot_df <- sample_n(emission_df, size = num_cases, replace = TRUE)
  boot_fit <- lm(co2_emission ~ engine_size + cylinders + fuel_consump_comb + fuel_type +
    fuel_type * fuel_consump_comb + transmission, data = boot_df)
  boot_coefs <- coef(boot_fit)
  boot_dist[i] <- boot_coefs[2]
}
(boot_CI <- quantile(boot_dist, c(0.025, 0.975)))
```

```
##      2.5%      97.5%
## -0.1873428  0.1404631
```

The 95% confidence interval for $\beta_{engine-size}$ using t-distribution and using bootstrap both contain 0, indicating that CO2 emissions is not significantly dependent on engine size when the other explanatory variables in this model are included, and thus `engine_size` should be dropped from Model 4. This is resulted from the multicollinearity issue as discussed in the visualization part. Engine size is highly correlated with the

number of cylinders and fuel consumption rating. Typically, larger engines with more cylinders tend to consume more fuel, leading to higher CO2 emissions. Indeed, a simple linear regression model with CO2 emissions as a function of engine size shows a strong positive correlation between them and a relatively high *R*-squared of 72.4%, suggesting that engine size is a potential key factor influencing CO2 emissions (see appendix iii). However, considering the observed multicollinearity between `engine_size`, `cylinders` and `fuel_consump_comb`, I will proceed to remove `engine_size` from Model 4.

My final model after excluding `engine_size` from Model 4 is as follows:

```
# Final model
lm_fit <- lm(co2_emission ~ cylinders + fuel_consump_comb + fuel_type + fuel_type *
  fuel_consump_comb + transmission, data = emission_df)
summary(lm_fit)
```

```
##
## Call:
## lm(formula = co2_emission ~ cylinders + fuel_consump_comb + fuel_type +
##     fuel_type * fuel_consump_comb + transmission, data = emission_df)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -69.066  -2.508   0.650   1.831  25.019
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    -0.70663    1.29495  -0.546 0.585304
## cylinders         0.28579    0.03665   7.797 7.22e-15 ***
## fuel_consump_comb 26.79382    0.14417 185.855 < 2e-16 ***
## fuel_typeE        5.21725    1.57195   3.319 0.000908 ***
## fuel_typeX        1.32536    1.30744   1.014 0.310757
## fuel_typeZ        1.92392    1.31905   1.459 0.144729
## transmissionAM    0.12867    0.14517   0.886 0.375444
## transmissionAS    0.22733    0.09161   2.481 0.013107 *
## transmissionAV   -0.33049    0.16040  -2.060 0.039396 *
## transmissionM    -0.23339    0.12029  -1.940 0.052388 .
## fuel_consump_comb:fuel_typeE -10.86174    0.15287 -71.053 < 2e-16 ***
## fuel_consump_comb:fuel_typeX  -3.68377    0.14502 -25.402 < 2e-16 ***
## fuel_consump_comb:fuel_typeZ  -3.77196    0.14587 -25.858 < 2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 2.937 on 7371 degrees of freedom
## Multiple R-squared:  0.9975, Adjusted R-squared:  0.9975
## F-statistic: 2.435e+05 on 12 and 7371 DF, p-value: < 2.2e-16
```

The main model assumptions are homoscedasticity, normality and linearity. By checking the diagnostic plots for the model (see appendix iv), it is clear that these assumptions are adequately satisfied.

The final model has lower AIC and BIC compared with the previous four models:

```
AIC(lm_fit); BIC(lm_fit)
```

```
## [1] 36882.53
```

```
## [1] 36979.23
```

Cross-validation further corroborates that the final model achieves lowest mean squared prediction error (MSPE):

```
# Cross Validation:
total_num_points <- dim(emission_df)[1]
num_training_points <- floor(total_num_points/2)
training_data <- emission_df[1:num_training_points, ]
test_data <- emission_df[(num_training_points + 1):total_num_points, ]
lm_train_1 <- lm(co2_emission ~ engine_size + cylinders + fuel_consump_comb, data = training_data)
lm_train_2 <- lm(co2_emission ~ engine_size + cylinders + fuel_consump_comb + fuel_type, data = training_data)
lm_train_3 <- lm(co2_emission ~ engine_size + cylinders + fuel_consump_comb + fuel_type + fuel_type *
  fuel_consump_comb, data = training_data)
lm_train_4 <- lm(co2_emission ~ engine_size + cylinders + fuel_consump_comb + fuel_type + fuel_type *
  fuel_consump_comb + transmission, data = training_data)
lm_train <- lm(co2_emission ~ cylinders + fuel_consump_comb + fuel_type + fuel_type * fuel_consump_comb
  + transmission, data = training_data)
model_train <- list(lm_train_1, lm_train_2, lm_train_3, lm_train_4, lm_train)
all_MSPE <- NULL
for (i in 1:5) {
  curr_model <- model_train[[i]]
  curr_test_predicted <- predict(curr_model, newdata = test_data)
  all_MSPE[i] <- mean((test_data$co2_emission - curr_test_predicted)^2)
}
all_MSPE
```

```
## [1] 390.39944 24.82576 14.08427 14.07625 14.06829
```

To interpret the coefficients in our final model: (1) There is a positive relationship between number of cylinders and CO2 emissions, vehicles with more cylinders tend to emit more CO2; (2) There is a strong positive relationship between fuel consumption rating and CO2 emissions; (3) The relationship between fuel consumption rating and CO2 emissions varies depending on fuel type; (4) While the ANOVA test comparing Model 3 and 4 indicates that transmission types do have an impact on CO2 emissions, it is worth noting that the coefficients for transmission categories are quite small, suggesting a minimal variation in CO2 emissions across different transmission types, so **transmission** will not be considered as a major factor in CO2 emissions.

To further investigate the effect of fuel type on CO2 emissions, I have computed CO2 emissions for each fuel type based on model coefficients, average number of cylinders, average fuel consumption rating, and the reference transmission category Automatic (A). The results, presented in the table below, show that: Ethanol E85 (type E) emits the least CO2, followed by premium gasoline (type Z), regular gasoline (type X), and diesel (type D) which emits the most (see appendix v for computation details).

Fuel Type	CO2 Emissions
E = ethanol (E85)	181.08 g/km
Z = premium gasoline	255.59 g/km
X = regular gasoline	255.96 g/km
D = diesel	295.07 g/km

Note: CO2 Emissions is computed based on avg. fuel consumption rating & avg. number of cylinders.

Conclusion

In summary, fuel consumption rating, type of fuel and number of cylinders are the primary factors with the most significant impact on vehicle CO2 emissions. Based on the analysis of my fitted regression model, fuel indeed appears to be the most direct contributor to CO2 emissions. A lower fuel consumption rating signifies higher fuel efficiency and consequently lower emissions. In the relationship between fuel consumption rating and CO2 emissions, fuel type plays a critical role. For vehicles with identical features, including fuel consumption rating, the choice of fuel type can result in varying levels of CO2 emissions. Among the four common fuel types, Ethanol (E85) is associated with notably lower CO2 emissions, whereas diesel tends to result in the highest emissions. Another influencing factor is transmission type. As is validated through ANOVA test, different transmission types can lead to changes in CO2 emissions level, though its impact is relatively minor. These results are strongly supported by model's reliability and high explanatory power, which is confirmed through cross-validation and a high adjusted R -squared score of 99.75%. During the model selection process, the variable `engine_size` is excluded from the model due to multicollinearity concerns; however, it is still worth noting that engine size is highly correlated with number of cylinders and fuel consumption rating, and thus indirectly affects CO2 emissions level.

These findings are important, as they provide valuable insights to policy design and technical innovation in the automotive industry. Given the direct link between fuel and CO2 emissions, it would be helpful if fuel efficiency can be improved through advanced engine designs. The effect of fuel types on emission levels underscores the need for policies that promote cleaner diesel technology and encourage greener alternatives such as Ethanol (E85). For example, implementing subsidies for E85-compatible vehicles could steer consumer preferences towards more environmental-friendly vehicle and fuel choices. Moreover, engine downsizing and cylinder reduction could be considered as a key aspect in technical improvement for further decreasing CO2 emissions. For future research, I would suggest to explore the link between vehicle weight and CO2 emissions. Although this dataset contains a categorical variable called `vehicle_class` which offers some general size information, more detailed data is required to study the effect of vehicle weight. Considering that vehicle weight is associated with fuel efficiency, looking into this relationship could be significantly beneficial in advancing CO2 emissions reduction strategies.

Appendix

i.

```
# Rename columns:
emission <- read.csv("/Users/zhizhiwei/Downloads/archive (2)/CO2 Emissions_Canada.csv")
colnames(emission) <- c("make", "model", "vehicle_class", "engine_size", "cylinders",
  "transmission", "fuel_type", "fuel_consump_city", "fuel_consump_hwy", "fuel_consump_comb",
  "fuel_consump_comb_mpg", "co2_emission")
```

ii.


```

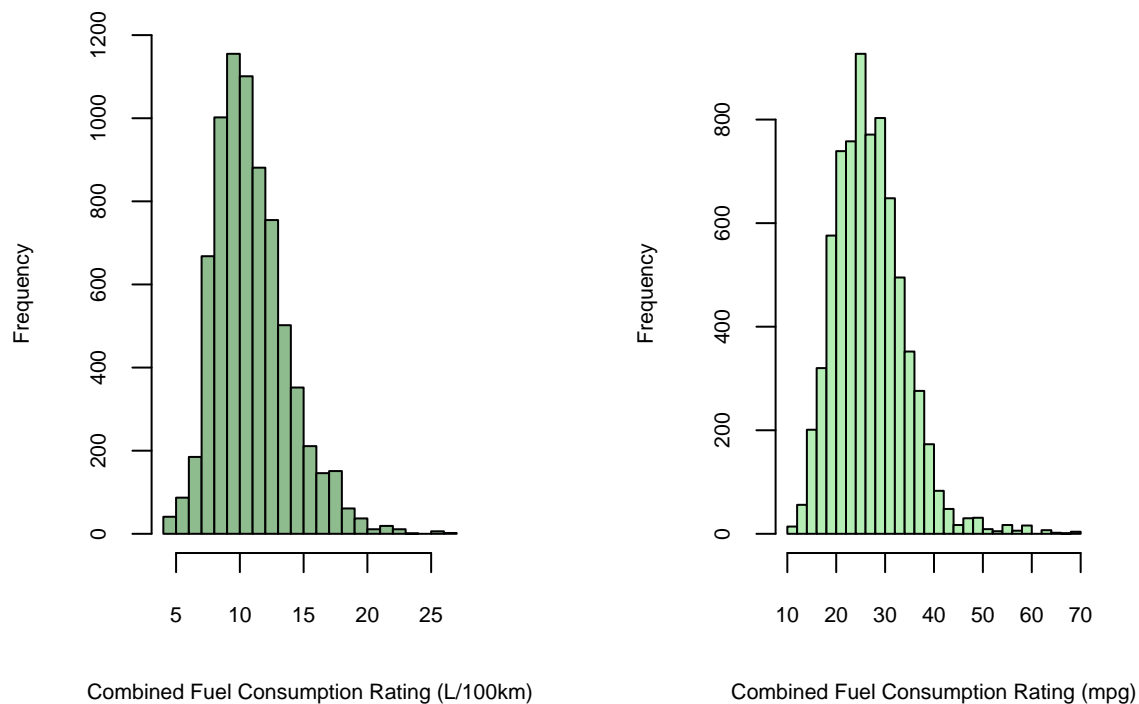
par(mfrow = c(1, 2))
par(mar=c(5, 4, 4, 4))

# Distribution of Combined Fuel Consumption Rating in L/100km
hist(emission$fuel_consump_comb, breaks=30,
     xlab = "Combined Fuel Consumption Rating (L/100km)",
     main = "Histogram of Combined Fuel Consumption Rating in L/100km",
     cex.lab=0.65, cex.main = 0.65, cex.axis = 0.7, col = "darkseagreen")

# Distribution of Combined Fuel Consumption Rating in mpg
hist(emission$fuel_consump_comb_mpg, breaks=30,
     xlab = "Combined Fuel Consumption Rating (mpg)",
     main = "Histogram of Combined Fuel Consumption Rating in mpg",
     cex.lab=0.65, cex.main = 0.65, cex.axis = 0.7, col = "darkseagreen2")

```

Histogram of Combined Fuel Consumption Rating in L/100km Histogram of Combined Fuel Consumption Rating in mpg



iii.

```

# Simple linear regression: CO2 emissions as a function of engine size
slm <- lm(co2_emission ~ engine_size, data = emission_df)
summary(slm)

```

```

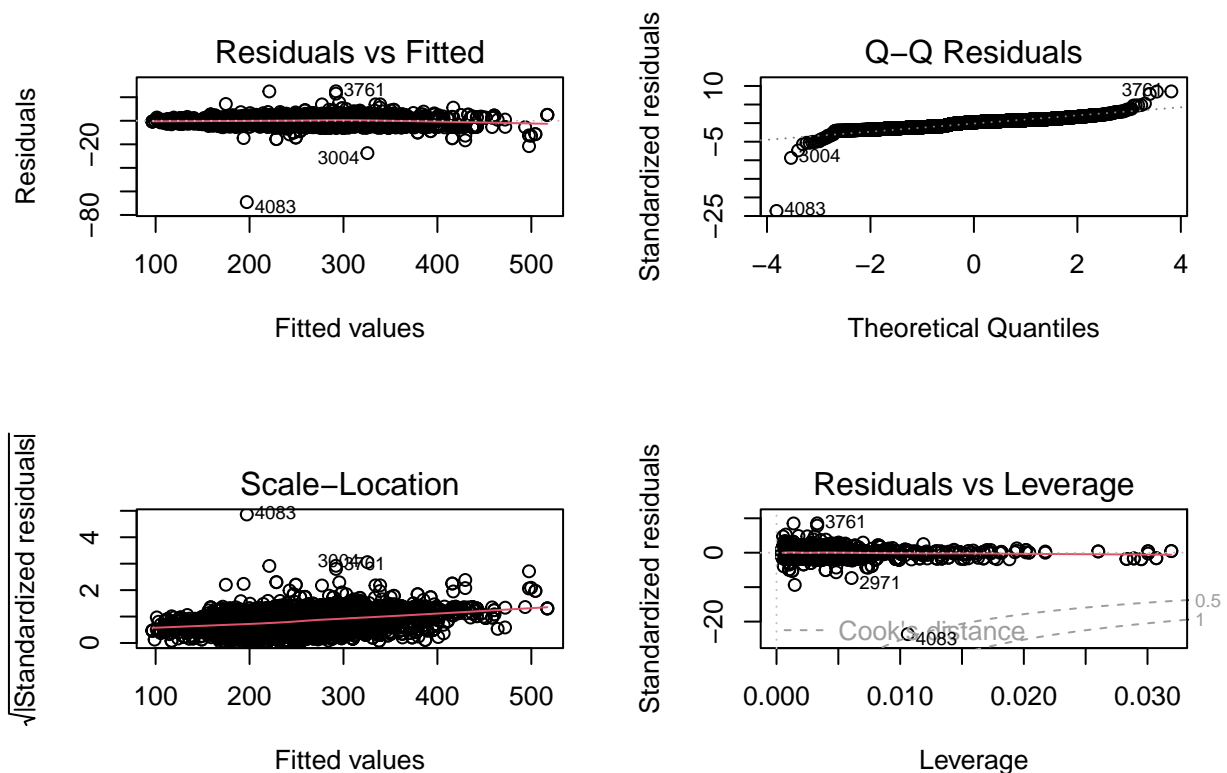
##
## Call:
## lm(formula = co2_emission ~ engine_size, data = emission_df)
##
## Residuals:

```

```
##      Min      1Q  Median      3Q      Max
## -113.31 -18.32   -1.45   19.07  142.91
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept) 134.3677     0.9074   148.1  <2e-16 ***
## engine_size  36.7791     0.2639   139.3  <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 30.71 on 7382 degrees of freedom
## Multiple R-squared:  0.7245, Adjusted R-squared:  0.7245
## F-statistic: 1.942e+04 on 1 and 7382 DF,  p-value: < 2.2e-16
```

iv.

```
# Diagnostic plots of the final model:
par(mfrow = c(2, 2))
plot(lm_fit)
```



v.

```
mean_fuel_consump <- mean(emission_df$fuel_consump_comb)
mean_cylinder <- round(mean(emission_df$cylinders))
```

```

b0 <- coef(lm_fit)[1]
b1 <- coef(lm_fit)[2]
b2 <- coef(lm_fit)[3]
b3 <- coef(lm_fit)[4]
b4 <- coef(lm_fit)[5]
b5 <- coef(lm_fit)[6]
b6 <- coef(lm_fit)[7]
b7 <- coef(lm_fit)[8]
b8 <- coef(lm_fit)[9]
b9 <- coef(lm_fit)[10]
b10 <- coef(lm_fit)[11]
b11 <- coef(lm_fit)[12]
b12 <- coef(lm_fit)[13]

# Diesel (type D):
(co2_D <- b0 + b1 * mean_cylinder + b2 * mean_fuel_consump)

```

```

## (Intercept)
##      295.0659

```

```

# Ethanol E85 (type E):
(co2_E <- b0 + b1 * mean_cylinder + (b2 + b10) * mean_fuel_consump + b3)

```

```

## (Intercept)
##      181.0773

```

```

# Regular gasoline (type X):
(co2_X <- b0 + b1 * mean_cylinder + (b2 + b11) * mean_fuel_consump + b4)

```

```

## (Intercept)
##      255.9625

```

```

# Premium gasoline (type Z):
(co2_Z <- b0 + b1 * mean_cylinder + (b2 + b12) * mean_fuel_consump + b5)

```

```

## (Intercept)
##      255.5932

```

vi.

```

# Side-by-side boxplot showing the relationship between vehicle class and CO2 emission:
emission_df |>
  ggplot(aes(x = factor(vehicle_class), y = co2_emission, fill = factor(vehicle_class))) +
  geom_boxplot() +
  xlab("Vehicle Class") +
  ylab("CO2 Emission") +
  ggtitle("Side by Side Boxplot") +
  theme(axis.text.x = element_text(angle = 90, vjust = 0.5, hjust = 1), text = element_text(size = 9))

```

