

Report 2

Lab2 GroupA

2025-02-07

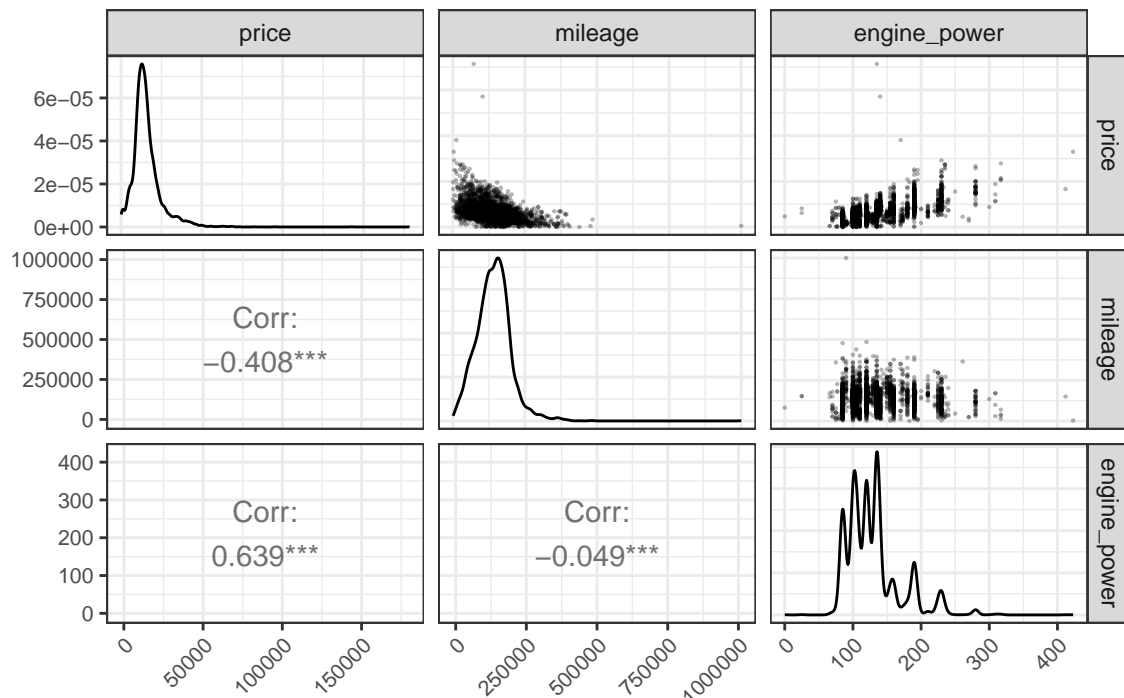
Loading the Library and Dataset

```
mylibrary <- c("dplyr", "ggplot2", "cowplot", "GGally", "MASS")
invisible(lapply(mylibrary, library, character.only = TRUE))
### load the data
dat_bmw <- read.csv("BMWpricing_updated.csv")
### Assigning NA to data points with negative mileage
dat_bmw[dat_bmw$mileage < 0, "mileage"] <- NA
### Converting mileage, engine power, and price into numeric
dat_bmw[,c(3,4,17)] <- apply(dat_bmw[,c(3,4,17)], 2, as.numeric)
### Create a data frame only has the variable of interests
dat_bmw2 <- dat_bmw[,c("price", "mileage", "engine_power")]
### Drop the data point with NA
dat_bmw2 <- dat_bmw2[complete.cases(dat_bmw2),]
```

scatterplot matrix

```
ggpairs(dat_bmw2,
        upper = list(continuous = wrap("points", alpha = 0.3, size = 0.1)),
        lower = list(continuous = wrap("cor", size = 4))) +
labs(title = "Distribution and Correlation between Price, Mileage, and Engine Power") +
theme_bw() +
theme(axis.text.x = element_text(angle = 45, hjust = 1, size = 8),
      axis.text.y = element_text(size = 8),
      title = element_text(size = 8))
```

Distribution and Correlation between Price, Mileage, and Engine Power



scatter plot of price and mileage We first set the negative mileage (in one sample) as the missing data and excluded it in our following analysis.

Our scatter plot shows a negative correlation between price and mileage, indicating that price tends to decrease as mileage increases. However, this relationship is not strongly linear and a straight-line regression model may not be the best option. We find some problems in our data: 1) Both price and mileage follow highly right-skewed distributions; 2) High-leverage points appear on the right side of the scatter plot, particularly the red point with approximately 1,000,000 miles; 3) Outliers with extreme high sold prices are observed on the top-left of the plot, especially these blue points with price over 100000.

To improve normality, we applied the Box-Cox transformation to the price variable. Though the λ with the greatest log-likelihood ($\lambda_{price} = 0.4, \lambda_{mileage} = 0.7$) and its 95% CI do not include 0, 0.5, and 1, there are slight differences in log-likelihood for price when $\lambda \in [0, 0.5]$ and for mileage when $\lambda \in [0.5, 1]$. For interpretability, we finally apply log-transformation on the price and keep the mileage in the data.

In the scatter plot of $\log(\text{price})$ and mileage, we observe some potential outliers with low prices and short mileage.

Scatter plot of price vs. mileage with highlighted outliers

```
scatterPlotPriceMileage <- ggplot(data = dat_bmw2, mapping = aes(x = mileage, y = price)) +
  geom_point(pch=19, cex=0.3) + # Default points
  geom_point(data = subset(dat_bmw2, mileage > 500000), mapping = aes(x = mileage, y = price), pch=19, cex=0.3) +
  geom_point(data = subset(dat_bmw2, price > 100000), mapping = aes(x = mileage, y = price), pch=19, cex=0.3) +
  labs(title = "Scatter plot between price and mileage") +
  theme_bw() + theme(title = element_text(size=8))
```

Box-Cox transformation for the price variable

```
boxcox_price <- boxcox(lm(dat_bmw2$price ~ 1), plotit = F)
boxcox_priceDF <- data.frame("lambda" = boxcox_price$x, "ll" = boxcox_price$y)
maxlabmda_price <- boxcox_price$x[which.max(boxcox_price$y)]
```

Box-cox Plot for the price variable

```
boxcoxPlotPrice <- ggplot(data = boxcox_priceDF, mapping = aes(x = lambda, y = ll)) +
```

```

geom_line(color = "black", linewidth = 0.5) +
geom_vline(xintercept = maxlabmda_price, linetype = "dashed", color = "red") +
geom_hline(yintercept = max(boxcox_price$y) - qchisq(0.95, df = 1)/2, linetype = "dashed", color = "black") +
labs(title = "Log-likelihood for the Box-Cox transformation",
      x = expression(lambda),
      y = "Log-likelihood for price") +
annotate("text", label = paste0("lambda ==", maxlabmda_price), x = 1, y = -30000, parse = T) +
theme_bw() + theme(title = element_text(size=8))

# Box-cox transformation for the mile variable
boxcox_mile <- boxcox(lm(dat_bmw2$mileage ~ 1), plotit = F)
boxcox_mileDF <- data.frame("lambda" = boxcox_mile$x, "ll" = boxcox_mile$y)
maxlabmda_mile <- boxcox_mile$x[which.max(boxcox_mile$y)]

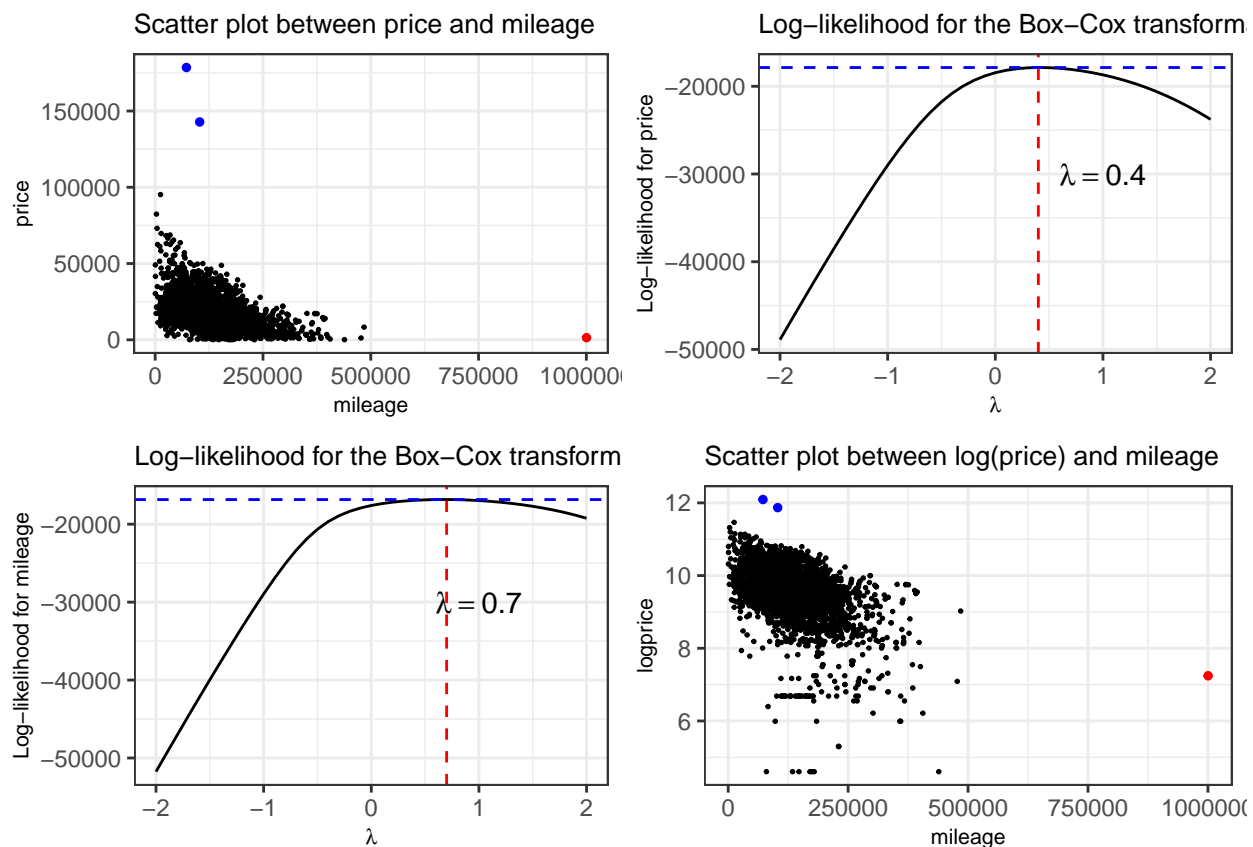
# Box-cox Plot for the mile variable
boxcoxPlotMile <- ggplot(data = boxcox_mileDF, mapping = aes(x = lambda, y = ll)) +
  geom_line(color = "black", linewidth = 0.5) +
  geom_vline(xintercept = maxlabmda_mile, linetype = "dashed", color = "red") +
  geom_hline(yintercept = max(boxcox_mile$y) - qchisq(0.95, df = 1)/2, linetype = "dashed", color = "black") +
  labs(title = "Log-likelihood for the Box-Cox transformation",
        x = expression(lambda),
        y = "Log-likelihood for mileage") +
  annotate("text", label = paste0("lambda ==", maxlabmda_mile), x = 1, y = -30000, parse = T) +
  theme_bw() + theme(title = element_text(size=8))

# Create a column of logprice in dat_bmw2
dat_bmw2 <- dat_bmw2 %>% mutate(logprice = log(price))

# Scatter plot of logprice vs. mileage with highlighted outliers
scatterPlotLogPriceMileage <- ggplot(data = dat_bmw2, mapping = aes(x = mileage, y = logprice)) +
  geom_point(pch=19, cex=0.3) +
  geom_point(data = subset(dat_bmw2, mileage > 500000), mapping = aes(x = mileage, y = logprice), pch=19, cex=0.3) +
  geom_point(data = subset(dat_bmw2, logprice > log(100000)), mapping = aes(x = mileage, y = logprice), pch=19, cex=0.3) +
  labs(title = "Scatter plot between log(price) and mileage") +
  theme_bw() + theme(title = element_text(size=8))

plot_grid(scatterPlotPriceMileage, boxcoxPlotPrice, boxcoxPlotMile, scatterPlotLogPriceMileage, align = "left")

```



Outliers

```
modelDF <- dat_bmw2
### Assigning NA to car with mileage > 50,000 miles
modelDF[modelDF$mileage > 500000, "mileage"] <- NA
### Assigning NA to car with price > 100,000 usd
modelDF[modelDF$price > 100000, "price"] <- NA
modelDF <- modelDF[complete.cases(modelDF),]
modelDF$logprice <- log(modelDF$price)
```

model fitting We applied an Ordinary Least Square regression model to predict the natural logarithm of BMW car prices from mileage. The model may be represented as:

$$\log(\text{price}) = \beta_0 + \beta_1 \times \text{mileage} + \epsilon$$

The residuals represent the difference between the model-predicted values and the actual $\log(\text{price})$ values. Ideally, they should be symmetrically distributed around zero for a well-behaved model. But the minimum of residual is -5.1907 . However, the minimum of residual is -5.1907 . That means that, for at least one observation, the model-predicted $\log(\text{price})$ is much larger than the actual value.

The “Coefficients” section provides us with detailed statistics for the intercept and the mileage coefficient. The intercept is estimated to be 10.17 with very small standard error (0.02175) and thus having an extremely large t-value of 467.7 and an extremely small p-value of less than $2e-16$. This is very strong evidence against the null hypothesis that the intercept is zero. Similarly, mileage coefficient is also estimated as $-4.746e-06$, indicating that for one unit increase in mileage, logprice decreases by about $4.746e-06$ units. Its standard

error is $1.425\text{e-}07$, which yields a t-value of -33.3 and a p-value similarly less than $2\text{e-}16$. These allow us to reject the null hypothesis for both coefficients, meaning mileage is a statistically significant logprice predictor.

Besides the individual coefficients, the fit of the entire model is also tested by means of the F-test. The F-statistic of 1109 with degrees of freedom 1 and 4837, and p-value smaller than 2.2×10^{-16} tests whether the model that includes mileage as a predictor explains significantly more variation in $\log(\text{price})$ than does a model with no predictors. This result confirms that our model as a whole is statistically significant although mileage by itself explains only about 18.65% of the variation in $\log(\text{price})$ (as indicated by the Multiple R-squared of 0.1865). The relatively low R-squared value suggests that while mileage is a significant predictor, other variables likely contribute to the variation in car prices.

```
mod1 <- lm(logprice ~ mileage, data = modelDF)
summary(mod1)

##
## Call:
## lm(formula = logprice ~ mileage, data = modelDF)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -5.1907 -0.2404  0.0496  0.3229  1.3558
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  1.017e+01  2.175e-02   467.7  <2e-16 ***
## mileage      -4.746e-06  1.425e-07   -33.3  <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.5838 on 4837 degrees of freedom
## Multiple R-squared:  0.1865, Adjusted R-squared:  0.1864
## F-statistic: 1109 on 1 and 4837 DF, p-value: < 2.2e-16

# Create a scatter plot with the best-fit line
# Add predicted values to the dataset
modelDF$predLinear <- predict(mod1)

# Create the scatter plot with the best-fit line
ggplot(modelDF, aes(x = mileage, y = logprice)) +
  geom_point(size = 0.15, color = 'blue') + # Scatter plot points
  geom_line(aes(x = mileage, y = predLinear, color = "blue")) + # Best-fit line from predictions
  scale_color_discrete(name = "Prediction", labels = c("Linear")) + # Legend for the line
  theme_bw()
```

