

# Assignment 3

Zheyang Liu

## Contents

<b>Introduction</b>	<b>1</b>
<b>Question (a)</b>	<b>2</b>
Target variable mpg_cat . . . . .	2
mpg_cat and categorical variables . . . . .	2
mpg_cat and continuous variables . . . . .	4
<b>Question (b)</b>	<b>6</b>
Build logistic regression model and get significant predictors . . . . .	6
Confusion matrix and fraction of correct predictions . . . . .	7

## Introduction

In this problem, you will develop a model to predict whether a given car gets high or low gas mileage based on the dataset “auto.csv”. The dataset contains 392 observations. The response variable is mpg cat, which indicates whether the miles per gallon of a car is high or low. The predictors are:

- cylinders: Number of cylinders between 4 and 8
- displacement: Engine displacement (cu. inches)
- horsepower: Engine horsepower
- weight: Vehicle weight (lbs.)
- acceleration: Time to accelerate from 0 to 60 mph (sec.)
- year: Model year (modulo 100)
- origin: Origin of car (1. American, 2. European, 3. Japanese)

Split the dataset into two parts: training data (70%) and test data (30%).

```
# read data
df =
  read_csv('data/auto.csv', show_col_types = FALSE) %>%
  janitor::clean_names() %>%
  mutate(cylinders = as.factor(cylinders),
         origin = as.character(origin),
         origin =
           case_when(origin == '1' ~ 'American',
                     origin == '2' ~ 'European',
```

```

        origin == '3' ~ 'Japanese'),
  origin = as.factor(origin),
  # target
  mpg_cat = as.factor(mpg_cat),
  mpg_cat = fct_relevel(mpg_cat, 'low'))

# split data
rowTrain <- createDataPartition(y = df$mpg_cat,
                                p = 0.7,
                                list = FALSE)

```

## Question (a)

Produce some graphical or numerical summaries of the data.

The model has 392 observations and 7 independent variables including 2 categorical variables (cylinders, origin) and 5 continuous variables (displacement, horsepower, weight, acceleration, year).

### Target variable mpg\_cat

Category high and low are balanced

```

df %>%
  group_by(mpg_cat) %>%
  summarise(cnt = n()) %>%
  knitr::kable()

```

mpg_cat	cnt
low	196
high	196

### mpg\_cat and categorical variables

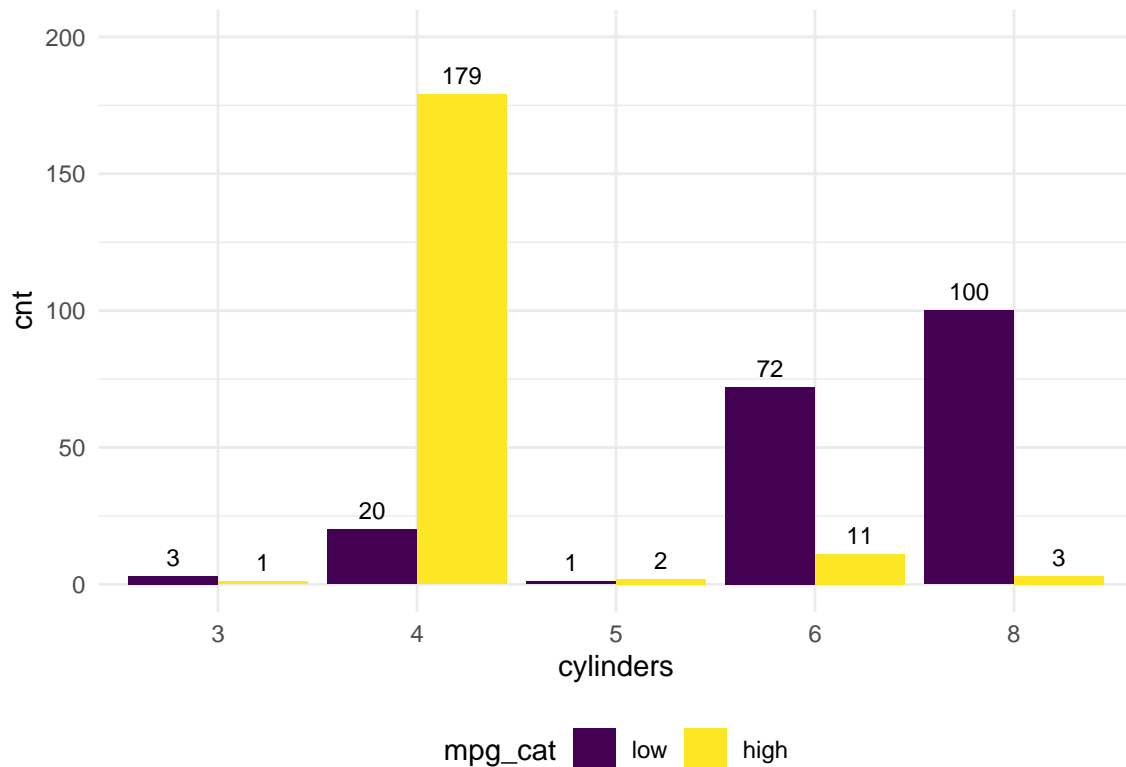
Cars with low mpg mostly has 6 or 8 cylinders while those with high mpg has 4 cylinders.

```

df %>% group_by(cylinders, mpg_cat) %>%
  summarise(cnt = n()) %>%
ggplot(aes(x = cylinders, y = cnt, fill = mpg_cat, label = cnt)) +
  geom_bar(stat = "identity", position = "dodge") +
  geom_text(
    aes(label = cnt),
    colour = "black", size = 3.2,
    vjust = -0.6, position = position_dodge(.9)
  ) + ylim(0, 200)

```

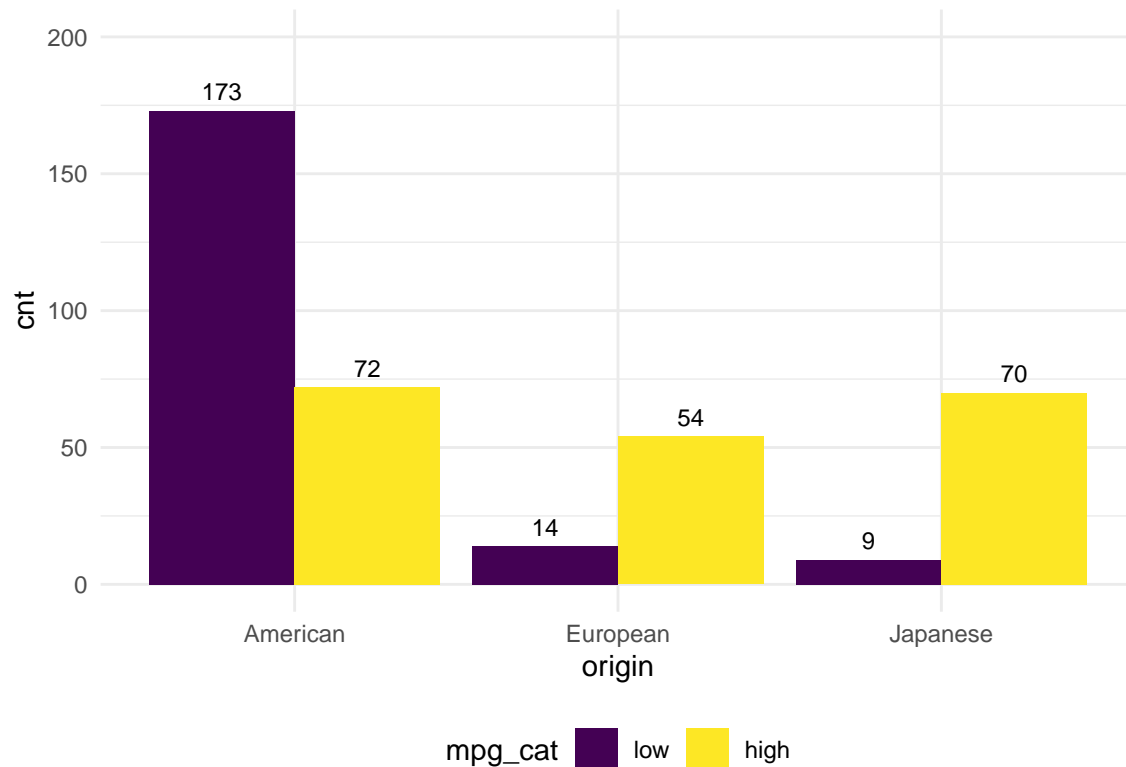
## 'summarise()' has grouped output by 'cylinders'. You can override using the '.groups' argument.



Cars originates in American are more likely to have low mpg (2.4 times more likely), while cars from European and Japanese are more likely to have high mpg.

```
df %>% group_by(origin, mpg_cat) %>%
  summarise(cnt = n()) %>%
  ggplot(aes(x = origin, y = cnt, fill = mpg_cat, label = cnt)) +
  geom_bar(stat = "identity", position = "dodge") +
  geom_text(
    aes(label = cnt),
    colour = "black", size = 3.2,
    vjust = -0.6, position = position_dodge(.9)
  ) + ylim(0, 200)
```

## 'summarise()' has grouped output by 'origin'. You can override using the '.groups' argument.

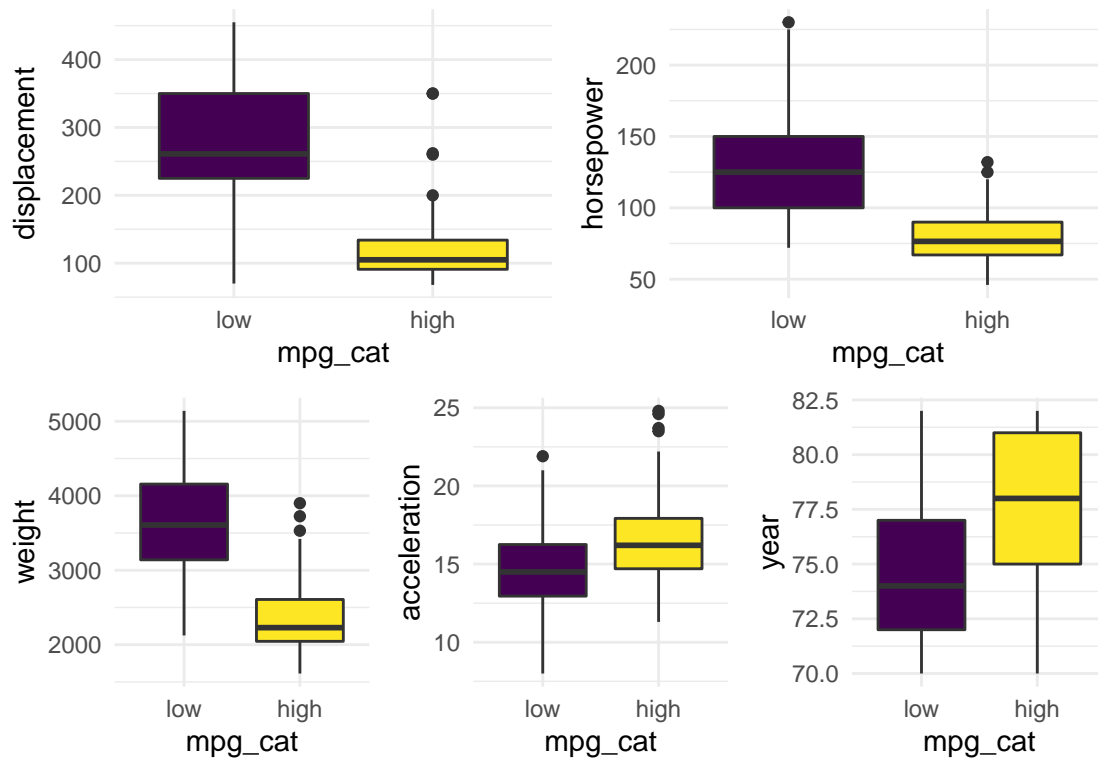


## mpg\_cat and continuous variables

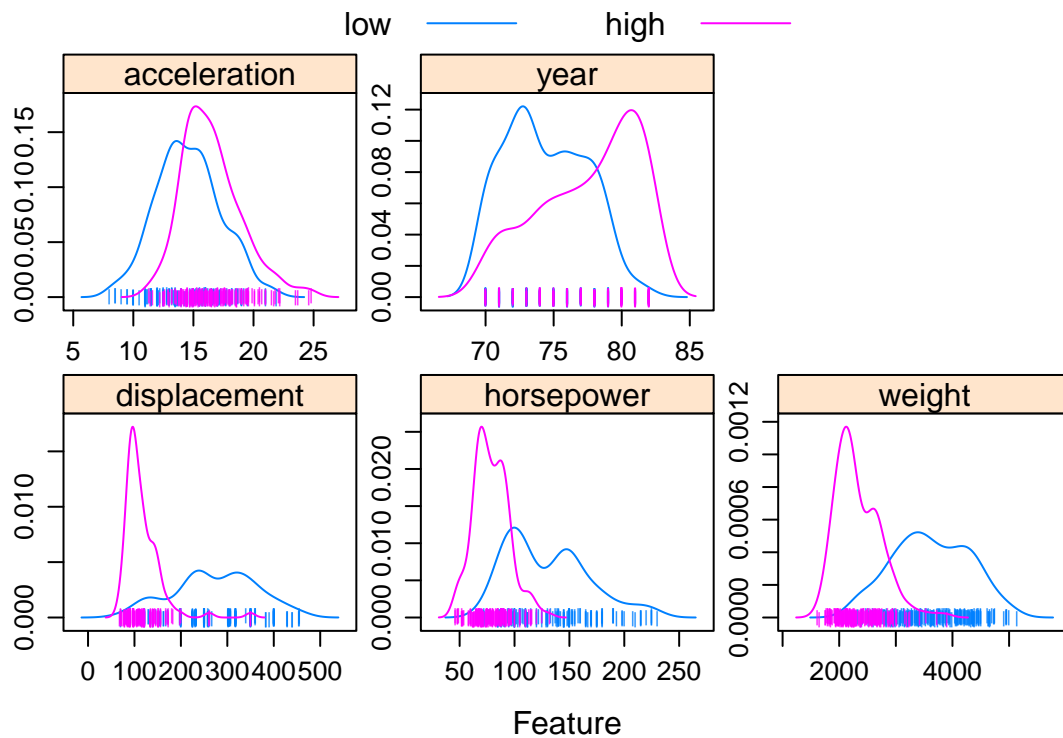
From the boxplot and feature plot, the median of displacement, horsepower, weight of the high mpg cars is lower than that of the low mpg cars, while the median of acceleration, year of the high mpg cars is higher than that of the low mpg cars

```
library(patchwork)
p1 = ggplot(df, aes(x=mpg_cat, y=displacement, fill = mpg_cat)) +
  geom_boxplot() + theme(legend.position = "none")
p2 = ggplot(df, aes(x=mpg_cat, y=horsepower, fill = mpg_cat)) +
  geom_boxplot() + theme(legend.position = "none")
p3 = ggplot(df, aes(x=mpg_cat, y=weight, fill = mpg_cat)) +
  geom_boxplot() + theme(legend.position = "none")
p4 = ggplot(df, aes(x=mpg_cat, y=acceleration, fill = mpg_cat)) +
  geom_boxplot() + theme(legend.position = "none")
p5 = ggplot(df, aes(x=mpg_cat, y=year, fill = mpg_cat)) +
  geom_boxplot() + theme(legend.position = "none")

(p1 + p2)/(p3 + p4 + p5)
```



```
featurePlot(x = df %>% select(displacement, horsepower, weight, acceleration, year),
            y = df$mpg_cat,
            scales = list(x = list(relation = "free"),
                          y = list(relation = "free")),
            plot = "density", pch = "|",
            auto.key = list(columns = 2))
```



## Question (b)

Perform a logistic regression using the training data. Do any of the predictors appear to be statistically significant? If so, which ones? Compute the confusion matrix and overall fraction of correct predictions using the test data. Briefly explain what the confusion matrix is telling you.

**Build logistic regression model and get significant predictors**

```
glm.fit <- glm(mpg_cat ~ .,
               data = df,
               subset = rowTrain,
               family = binomial(link = "logit"))

summary(glm.fit)
```

```
##
## Call:
## glm(formula = mpg_cat ~ ., family = binomial(link = "logit"),
##      data = df, subset = rowTrain)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -2.53245  -0.05528   0.00028   0.11302   2.65490
##
```

```
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)  -2.174e+01  8.887e+00  -2.446  0.01445 *
## cylinders4    6.014e+00  1.980e+00   3.037  0.00239 **
## cylinders5    1.684e+01  1.455e+03   0.012  0.99077
## cylinders6    6.034e+00  2.581e+00   2.338  0.01939 *
## cylinders8    1.050e+01  3.769e+00   2.785  0.00535 **
## displacement -3.243e-02  1.999e-02  -1.622  0.10483
## horsepower   -9.003e-02  3.691e-02  -2.439  0.01473 *
## weight       -3.017e-03  1.792e-03  -1.684  0.09226 .
## acceleration -3.781e-01  1.996e-01  -1.894  0.05818 .
## year          5.752e-01  1.299e-01   4.429 9.47e-06 ***
## originEuropean 8.191e-01  1.076e+00   0.761  0.44637
## originJapanese 1.005e-01  1.011e+00   0.099  0.92076
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##      Null deviance: 382.617  on 275  degrees of freedom
## Residual deviance:  80.343  on 264  degrees of freedom
## AIC: 104.34
##
## Number of Fisher Scoring iterations: 14
```

Under 0.05 significance level, The significant predictors are cylinders 4 (reference category is cylinders 2), weight and year.

## Confusion matrix and fraction of correct predictions

Confusion matrix

```
test.pred.prob <- predict(glm.fit, newdata = df[-rowTrain,],
                           type = "response")
test.pred <- rep("low", length(test.pred.prob))
test.pred[test.pred.prob>0.5] <- "high"

confusionMatrix(data = as.factor(test.pred),
                 reference = df$mpg_cat[-rowTrain],
                 positive = "high")
```

```
## Confusion Matrix and Statistics
##
##              Reference
## Prediction low high
##      low   51    5
##      high    7   53
##
##              Accuracy : 0.8966
##              95% CI : (0.8263, 0.9454)
##      No Information Rate : 0.5
##      P-Value [Acc > NIR] : <2e-16
```

```
##
##           Kappa : 0.7931
##
## Mcnemar's Test P-Value : 0.7728
##
##           Sensitivity : 0.9138
##           Specificity : 0.8793
##           Pos Pred Value : 0.8833
##           Neg Pred Value : 0.9107
##           Prevalence : 0.5000
##           Detection Rate : 0.4569
##           Detection Prevalence : 0.5172
##           Balanced Accuracy : 0.8966
##
##           'Positive' Class : high
##
```

Fraction of correct predictions is 0.8534483.

If we set the threshold to be 0.5, The confusion matrix is telling that

- Sensitivity = 0.8621, 0.8621 of the high mpg cars are detected by the model
- Specificity = 0.8448, 0.8448 of the low mpg cars are detected by the model
- PPV = 0.8475, 0.8475 of the predicted high are actually high
- NPV = 0.8596, 0.8596 of the predicted low are actually low