

Report



Project Phase I

Car Insurance Claim Prediction

Statistical Inference

Fatemeh Nadi

810101285

December, 2022

Contents

| | |
|---------------------|-----------|
| 1 Question 0 | 1 |
| 1.1 A | 1 |
| 1.2 B | 2 |
| 1.3 C | 4 |
| 1.4 D | 5 |
| 2 Question 1 | 6 |
| 2.1 A | 6 |
| 2.2 B | 9 |
| 2.3 C | 9 |
| 2.4 D | 11 |
| 2.5 E | 12 |
| 2.6 F | 13 |
| 2.7 G | 15 |
| 3 Question 2 | 17 |
| 3.1 A | 17 |
| 3.2 B | 18 |
| 3.3 C | 19 |
| 3.4 D | 20 |
| 4 Question 3 | 21 |
| 4.1 A | 21 |
| 4.2 B | 22 |
| 4.3 C | 23 |
| 4.4 D | 23 |
| 4.5 E | 26 |
| 5 Question 4 | 27 |
| 5.1 A | 27 |
| 5.2 B | 29 |
| 5.3 C | 30 |
| 6 Question 5 | 32 |
| 6.1 A | 32 |
| 6.2 B | 33 |
| 6.3 C | 34 |
| 6.4 D | 35 |
| 7 Question 6 | 36 |
| 7.1 A | 36 |
| 7.2 B | 36 |
| 8 Question 7 | 38 |
| 8.1 A | 39 |

| | | |
|-----------|-------------------|-----------|
| 8.2 | B | 40 |
| 8.3 | C | 40 |
| 8.4 | D | 42 |
| 8.5 | E | 43 |
| 8.6 | F | 44 |
| 8.7 | G | 45 |
| 9 | Question 8 | 46 |
| 9.1 | A | 47 |
| 9.2 | B | 48 |
| 9.3 | C | 49 |
| 10 | Question 9 | 51 |
| 10.1 | A | 52 |
| 10.2 | B | 54 |
| 11 | R codes | 56 |

1. Question 0

1.1 A

Briefly describe your dataset and why studying your dataset can be interesting?

Dataset: Car Insurance Claim Prediction

There are attributes in the dataset about policyholders, including tenure of the policy, age of the car, age of the owner, population density, car make and model, power, engine type, etc., and the target variable indicates whether a claim will be filed in the next six months.

| X | policy_id | policy_tenure | age_of_car | age_of_policyholder | area_cluster | population_density | segment | model | fuel_type | engine_type | |
|----|-----------|---------------|-------------|---------------------|--------------|--------------------|---------|-------|-----------|-------------|-----------------------|
| 1 | 5756 | ID05757 | 0.082997707 | 0.07 | 0.3269231 | C11 | 6108 | B2 | M6 | Petrol | K Series Dual jet |
| 2 | 21942 | ID21943 | 0.097829849 | 0.03 | 0.5769231 | C9 | 17804 | B2 | M7 | Petrol | 1.2 L K Series Engine |
| 3 | 53860 | ID53861 | 0.059199281 | 0.01 | 0.7019231 | C3 | 4076 | A | M1 | CNG | F8D Petrol Engine |
| 4 | 38217 | ID38218 | 0.541986698 | 0.17 | 0.6153846 | C2 | 27003 | B2 | M6 | Petrol | K Series Dual jet |
| 5 | 7724 | ID07725 | 0.093530064 | 0.09 | 0.5096154 | C5 | 34738 | C2 | M4 | Diesel | 1.5 L U2 CRDi |
| 6 | 31567 | ID31568 | 1.109789324 | 0.23 | 0.4038462 | C2 | 27003 | B2 | M6 | Petrol | K Series Dual jet |
| 7 | 42661 | ID42662 | 0.080624964 | 0.00 | 0.5000000 | C14 | 7788 | B2 | M7 | Petrol | 1.2 L K Series Engine |
| 8 | 2989 | ID02990 | 0.337275748 | 0.18 | 0.6730769 | C8 | 8794 | B2 | M6 | Petrol | K Series Dual jet |
| 9 | 39633 | ID39634 | 0.598877634 | 0.12 | 0.4807692 | C8 | 8794 | B2 | M6 | Petrol | K Series Dual jet |
| 10 | 42488 | ID42489 | 1.231427992 | 0.12 | 0.4134615 | C8 | 8794 | B2 | M7 | Petrol | 1.2 L K Series Engine |

Figure 1: First 10 rows of the dataset

Analyzing this data help the insurance companies detects some cases which haven't enough revenue.

If a case has a high risk of accident or destruction so reject insurance it.

And finally, make the best decision to prevent failure in it's market.

Studying and analyzing the information in this database lets us determine what influences consumers' claims most. It also allows us to study any two columns and see if there is an interesting statistical relationship between them.

Also exploring the statistical distribution of the factors can also help us discover the most meaningful factors.

1.2 B

How many variables (features) and cases does your dataset have?

We have a dataset of 30000 policyholders, each policyholder has 28 features which we explain in detail below:

- policy_id: Unique identifier of the policyholder.
- policy_tenure: Time period of the policy.
- age_of_car: Normalized age of the car in years.
- age_of_policyholder: Normalized age of policyholder in years.
- area_cluster: Area cluster of the policyholder.
- population_density: Population density of the city (Policyholder City).
- segment: Segment of the car (A/ B1/ B2/ C1/ C2).
- model: Encoded name of the car.
- fuel_type: Type of fuel used by the car.
- engine_type: Type of engine used in the car.
- airbags: Number of airbags installed in the car.
- is_parking_sensors: Boolean flag indicating whether parking sensors are present in the car or not.
- is_parking_camera: Boolean flag indicating whether the parking camera is present in the car or not.
- rear_brakes_type: Type of brakes used in the rear of the car
- displacement: Engine displacement of the car (cc).
- cylinder: Number of cylinders present in the engine of the car.
- transmission_type: Transmission type of the car.
- gear_box: Number of gears in the car.
- steering_type: Type of the power steering present in the car.
- turning_radius: The space a vehicle needs to make a certain turn (Meters).
- length: Length of the car (Millimetre).
- width: Width of the car (Millimetre).
- height: Height of the car (Millimetre).
- gross_weight: The maximum allowable weight of the fully-loaded car, including passengers, cargo and equipment (Kg).

- `is_speed_alert`: Boolean flag indicating whether the speed alert system is available in the car or not.
- `ncap_rating`: Safety rating given by NCAP (out of 5).
- `is_claim`: Outcome: Boolean flag indicating whether the policyholder file a claim in the next 6 months or not.

```
> summary(car)
      X      policy_id      policy_tenure      age_of_car      age_of_policyholder area_cluster
Min. : 0 Length:30000 Min. :0.00275 Min. :0.00000 Min. :0.2885 Length:30000
1st Qu.:14647 Class :character 1st Qu.:0.20760 1st Qu.:0.02000 1st Qu.:0.3654 Class :character
Median :29397 Mode :character Median :0.57231 Median :0.06000 Median :0.4519 Mode :character
Mean :29314 Mean :0.60925 Mean :0.06948 Mean :0.4697
3rd Qu.:43948 3rd Qu.:1.03817 3rd Qu.:0.11000 3rd Qu.:0.5481
Max. :58590 Max. :1.38652 Max. :0.82000 Max. :0.9808
population_density      segment      model      fuel_type      engine_type      airbags
Min. : 290 Length:30000 Length:30000 Length:30000 Length:30000 Min. :1.00
1st Qu.: 6112 Class :character Class :character Class :character Class :character 1st Qu.:2.00
Median : 8794 Mode :character Mode :character Mode :character Mode :character Median :2.00
Mean :18816 Mean :3.14
3rd Qu.:27003 3rd Qu.:6.00
Max. :73430 Max. :6.00
is_parking_sensors is_parking_camera rear_brakes_type displacement cylinder transmission_type
Length:30000 Length:30000 Length:30000 Min. : 796 Min. :3.000 Length:30000
Class :character Class :character Class :character 1st Qu.: 796 1st Qu.:3.000 Class :character
Mode :character Mode :character Mode :character Median :1197 Median :4.000 Mode :character
Mean :1163 Mean :3.628
3rd Qu.:1493 3rd Qu.:4.000
Max. :1498 Max. :4.000
gear_box steering_type turning_radius length width height gross_weight
Min. :5.000 Length:30000 Min. :4.500 Min. :3445 Min. :1475 Min. :1475 Min. :1051
1st Qu.:5.000 Class :character 1st Qu.:4.600 1st Qu.:3445 1st Qu.:1515 1st Qu.:1475 1st Qu.:1185
Median :5.000 Mode :character Median :4.800 Median :3845 Median :1735 Median :1530 Median :1335
Mean :5.245 Mean :4.853 Mean :3851 Mean :1672 Mean :1553 Mean :1385
3rd Qu.:5.000 3rd Qu.:5.000 3rd Qu.:3995 3rd Qu.:1755 3rd Qu.:1635 3rd Qu.:1510
Max. :6.000 Max. :5.200 Max. :4300 Max. :1811 Max. :1825 Max. :1720
is_speed_alert ncap_rating is_claim
Length:30000 Min. :0.000 Min. :0.00000
Class :character 1st Qu.:0.000 1st Qu.:0.00000
Mode :character Median :2.000 Median :0.00000
Mean :1.762 Mean :0.06377
3rd Qu.:3.000 3rd Qu.:0.00000
Max. :5.000 Max. :1.00000
```

Figure 2: Summary of the dataset

Except for `X` as index and `policy_id`, we have 15 numerical and 11 categorical variables.

1.3 C

Is there any missing value in your data? Provide a summary of a portion of missing values for each variable (feature) and describe how you handle these missing values for each variable (on what basis).

There are no missing values in this dataset.

```
> any(is.na(car))
[1] FALSE
```

Figure 3: Detecting missing values

Portion of missing values for each variable (feature):

| | x | 0 |
|---------------------|---|---|
| policy_id | | 0 |
| policy_tenure | | 0 |
| age_of_car | | 0 |
| age_of_policyholder | | 0 |
| area_cluster | | 0 |
| population_density | | 0 |
| segment | | 0 |
| model | | 0 |
| fuel_type | | 0 |
| engine_type | | 0 |
| airbags | | 0 |
| is_parking_sensors | | 0 |
| is_parking_camera | | 0 |
| rear_brakes_type | | 0 |
| displacement | | 0 |
| cylinder | | 0 |
| transmission_type | | 0 |
| gear_box | | 0 |
| steering_type | | 0 |
| turning_radius | | 0 |
| length | | 0 |
| width | | 0 |
| height | | 0 |
| gross_weight | | 0 |
| is_speed_alert | | 0 |
| ncap_rating | | 0 |
| is_claim | | 0 |

Figure 4: Missing values for each feature

The dataset shows all the values haven't any missing values. I also looked for special characters indicating missing values, such as ' ', '-' , 'NA' , and 'NAN' , by unique on each column.

If the dataset has missing values, these could be fixed by deleting, replacing with mean, median, mod, etc., or using some algorithm like fill forward or estimating or backward eliminating.

1.4 D

Using this elementary view of your dataset, which variables do you think may be the most relevant (contain some important information)? Why?

Age of the car: if a car is too old, it's maybe destroyed as soon as possible, and generally, people care about new cars more than old cars.

Displacement: High displacement (compared to the car's age) shows this policyholder driving more, so the probability of accidents increases.

Policy tenure: one or two months is an alarm for insurance companies that may be this case, show fraud.

Population density: usually, an accident occurs in a big city with a vast population.

Airbags: number of airbags can reduce the risk of death.

Model: some cars have more security, and some have less, so it's essential to which model of the car you insure.

2. Question 1

Chosen Numerical Variable : age_of_car: Normalized age of the car in years.

2.1 A

Plot a histogram with an appropriate bin size, then overlay that with the curve of a density, is the distribution approximately normal?

I use two different methods to estimate the number of bins.

1. Freedman Diaconis rule:

The general equation for the rule is: $2 \frac{IQR(x)}{\sqrt[3]{n}}$

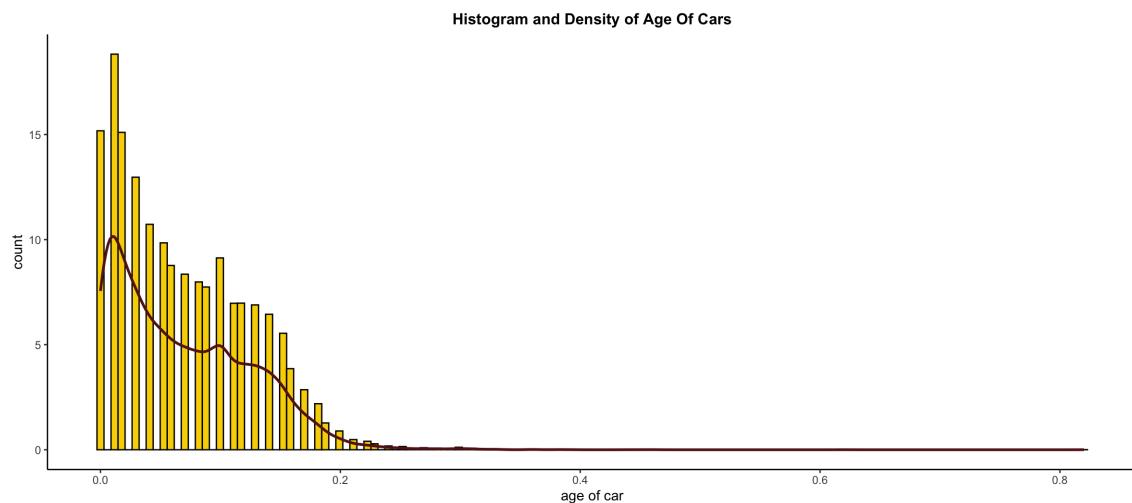


Figure 5: Histogram of Age Of Car with Freedman Diaconis rule bandwidth

2. $\sqrt[2]{n}$ rule

The number of bins corresponds to the square root of the number of observations

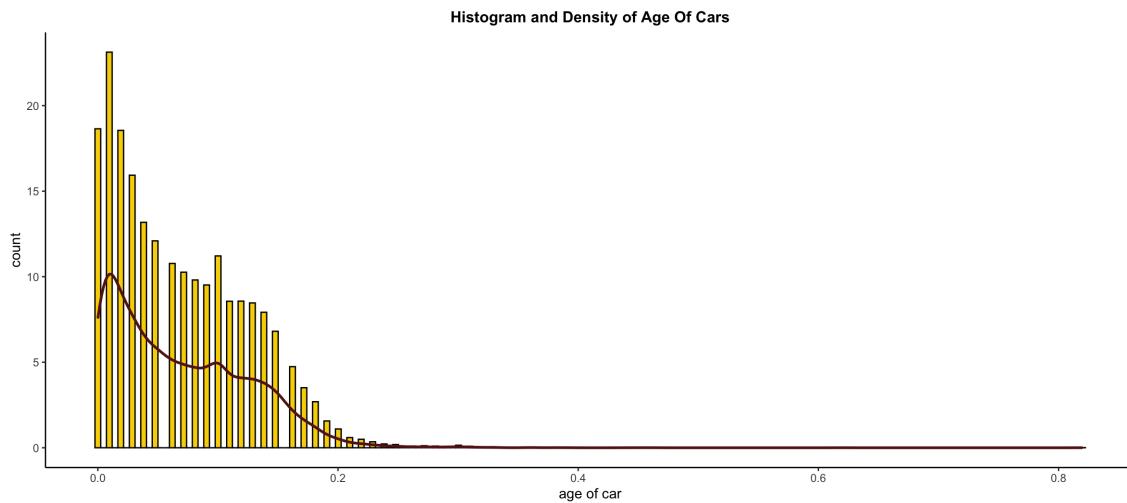


Figure 6: Histogram of Age Of Car with $\sqrt[2]{n}$ rule bandwidth

3. default in R

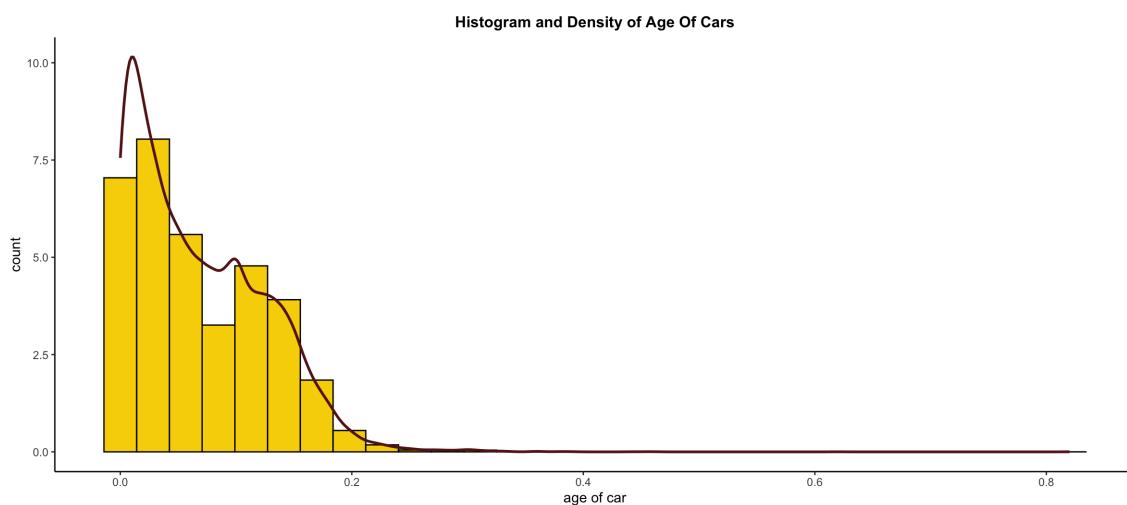


Figure 7: Histogram of Age Of Car with default in R bandwidth

It seems to be the default setting in r is a better choice for select bandwidth, and I think this because the age of the car variable is normalized.

We can see in the histogram that the distribution is highly right-skewed normal

To more accurately compare it to a normal distribution, we use QQ Plot.

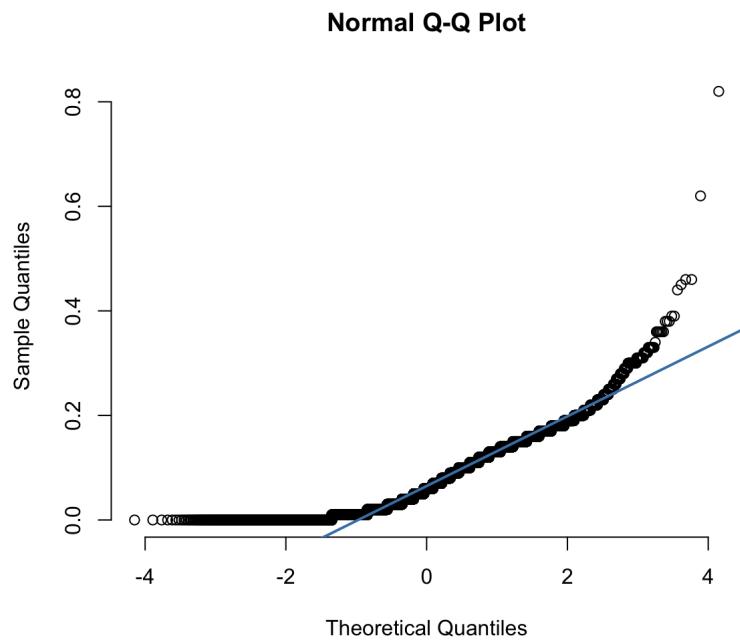


Figure 8: Normal Q-Q plot of Age Of Car

This plot confirms that the distribution is definitely right-skewed.

2.2 B

Based on the previous plot in section “A” talk about the modularity and skewness of this variable and describe it.

The plot above shows a right-skewed unimodal distribution which means that there are more recent cars in this dataset.

We rarely have worn-out cars either because, naturally, after a while, cars are replaced with other cars, which must have a different reason, like an accident or people like append changes in their life and so on or insurance companies just insure just new cars.

2.3 C

Determine the upper and lower quartiles, whiskers, and IQR by drawing a boxplot.

```
> quantile(car$age_of_car)
 0% 25% 50% 75% 100%
0.00 0.02 0.06 0.11 0.82
```

Figure 9: Important percentiles

$min = 0$, $lower\ extreme = 0$

$Q1 = 0.02$ so the lower whisker is of length 0.02, and we don’t have outlier here.

$Q3 = 0.11$, $max = 0.84$

$upper\ extreme = 0.24$, is not equal max which means $1.5 * IQR + Q3$ is 0.24 and the rest of upper part are outliers. the upper whisker is of length 0.13.

$IQR = Q3 - Q1 = 0.9$

```
> boxplot
$stats
[,1]
[1,] 0.00
[2,] 0.02
[3,] 0.06
[4,] 0.11
[5,] 0.24

$n
[1] 30000

$conf
[,1]
[1,] 0.05917901
[2,] 0.06082099
```

Figure 10: five-number information about age of car

BoxPlot for Age Of Car

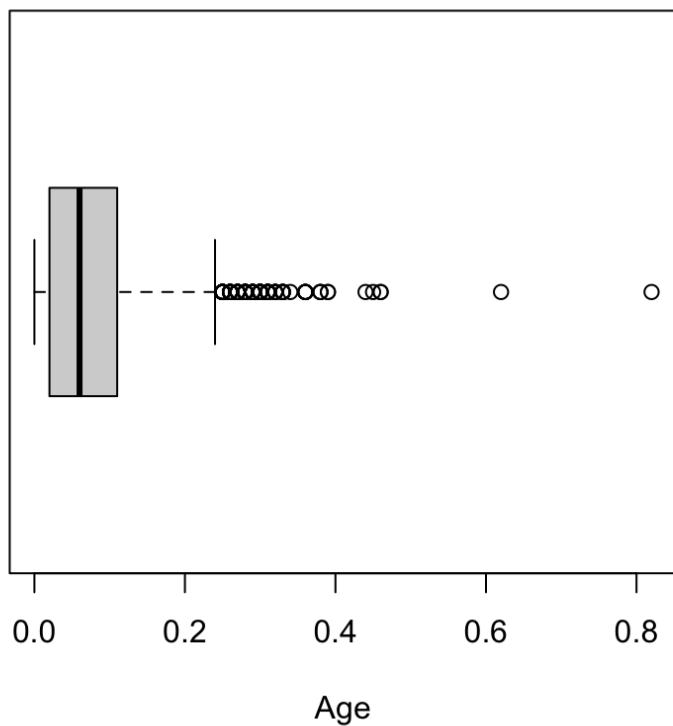


Figure 11: Important percentiles

2.4 D

What are the outliers in this variable? Determine the outliers and their quantity, then try to find what is the meaning of them.

Mention these in the last part. On the left, we haven't outliers, but on the right, we have 147 outliers which means 147 out of 30,000 cars are relatively worn out.

```
> aoc <- car$age_of_car  
> length(aoc[which(aoc < boxplot$stats[1] | aoc > boxplot$stats[5])])  
[1] 147
```

Figure 12: Important percentiles

Because the numbers of the age of cars are normalized, it is impossible to express the actual number, but if the age of your car is more than 0.24, it means that your car has a longer lifespan compared to all existing cars or If it doesn't work correctly, it is worn out.

2.5 E

Calculate the mean, median, variance, and standard deviation and describe each one.

Mean : mean is also known as average of all the numbers in the data set which is calculated by below equation.

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$$

So the sum of cars' age is divided by the population's size.

Median : median is mid value in this ordered data set.

So, we can see half of the cars are 0.06, which means that almost cars are new.

Mean is larger than median so the distribution is right-skewed.

Variance : variance is the numerical values that describe the variability of the observations from its arithmetic mean.

Variance measure how far individuals in the group are spread out, in the set of data from the mean.

$$\sigma^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$$

Standard Deviation : It is a measure of dispersion of observation within dataset relative to their mean. It is square root of the variance.

Standard deviation is expressed in the same unit as the values in the dataset so it measure how much observations of the data set differs from its mean, on average, data have 0.056, far from the mean.

$$\sigma = \sqrt{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2}$$

```
> mean(car$age_of_car)
[1] 0.06948467
> median(car$age_of_car)
[1] 0.06
> var(car$age_of_car)
[1] 0.003199781
> sd(car$age_of_car)
[1] 0.05656661
```

Figure 13: Some Statistics for Age Of Car

2.6 F

Categorize this variable into four intervals based on its mean and plot a pie chart that visualizes the frequency of these four categories. Your chart should be colorized, and the labels should contain each category with its percentage.

Our data are categorized according to a range of values.

$$\left\{ \begin{array}{ll} min \leq \leq \frac{\mu}{2} & First \\ \frac{\mu}{2} \leq < \mu & Second \\ \frac{\mu}{2} \leq < \frac{3\mu}{4} & Third \\ \frac{3\mu}{4} \leq \leq max & Fourth \end{array} \right.$$

```
1 aoc <- car$age_of_car
2 mean = mean(aoc)
3
4 Fo = aoc[aoc > .75*mean]
5 Th = aoc[aoc >= .5*mean & aoc <= .75*mean]
6 Se = aoc[aoc < .5*mean & aoc >= .25*mean]
7 Fi = aoc[aoc < .25*mean]
8
9 frequencies<-c(length(Fo),length(Th),length(Se),length(Fi))
10 percentage <- round(100*frequencies/sum(frequencies), 2)
11
12 aoc.categorized <- data.frame(names = c("Fourth", "Third", "Second", "First"), value = percentage)
13
14 ggplot(aoc.categorized, aes(x="", y = value, fill = names)) +
15   geom_bar(stat = "identity") + coord_polar("y") +
16   geom_text(aes(label = paste0(value, "%")),
17             position = position_stack(vjust = 0.5)) +
18   labs(title="Age Of Car Pie Chart", x = 'Frequency', y = 'Age Of Car')
```

Age Of Car Pie Chart

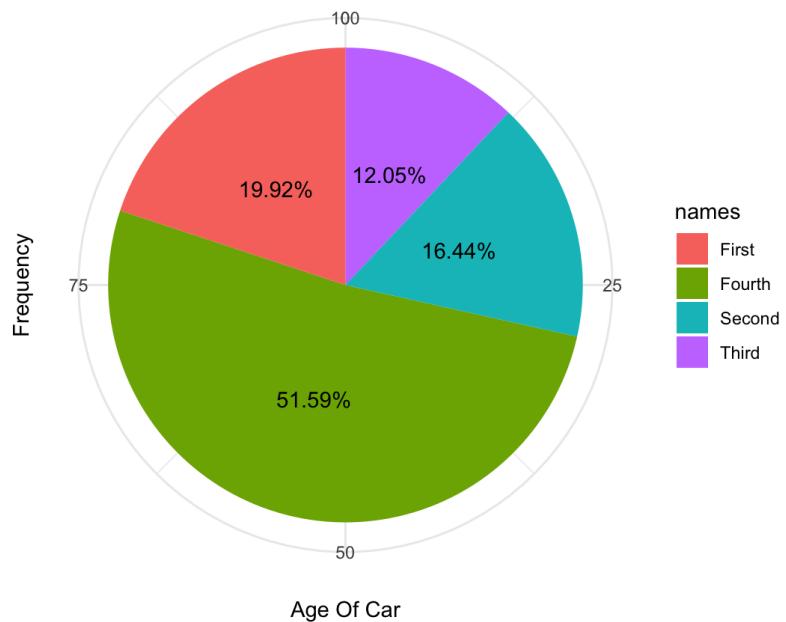


Figure 14: Age of Car Pie Chart

2.7 G

Draw a density plot of this variable and add lines for the mean and median to it. What is the relationship between the mean, median, and density of this variable?

```
1 ggplot(car, aes(x = age_of_car)) + geom_density( size = 1) +
2   geom_vline(aes(xintercept = median(age_of_car)), linetype = "dashed"
3   , size = .7, col= 'red') +
4   geom_vline(aes(xintercept = mean(age_of_car)), linetype = "dashed",
5   size = .7,col= 'green') +
6   theme_classic() +
7   theme(
8     plot.title = element_text(size = 12, face = "bold", hjust = 0.5),
9     plot.caption = element_text(face = "italic")
10   )
11 +
12   annotate("text", x = mean(car$age_of_car) + .02 , label = "mean", y
13   = 6, size = 3, angle = 90 , color = 'green') +
14   annotate("text", x = median(car$age_of_car) - 0.02 , label = "median
15   ", y = 0.6, size = 3, angle = 90, color = 'red')
16 +
17   labs(
18     title = "Age of Car Density",
19     caption = "green line is about mean and red is about median
20     variable on the diagram.",
21     x = "Age Of The Car",
22     y = "Count"
23   )
```

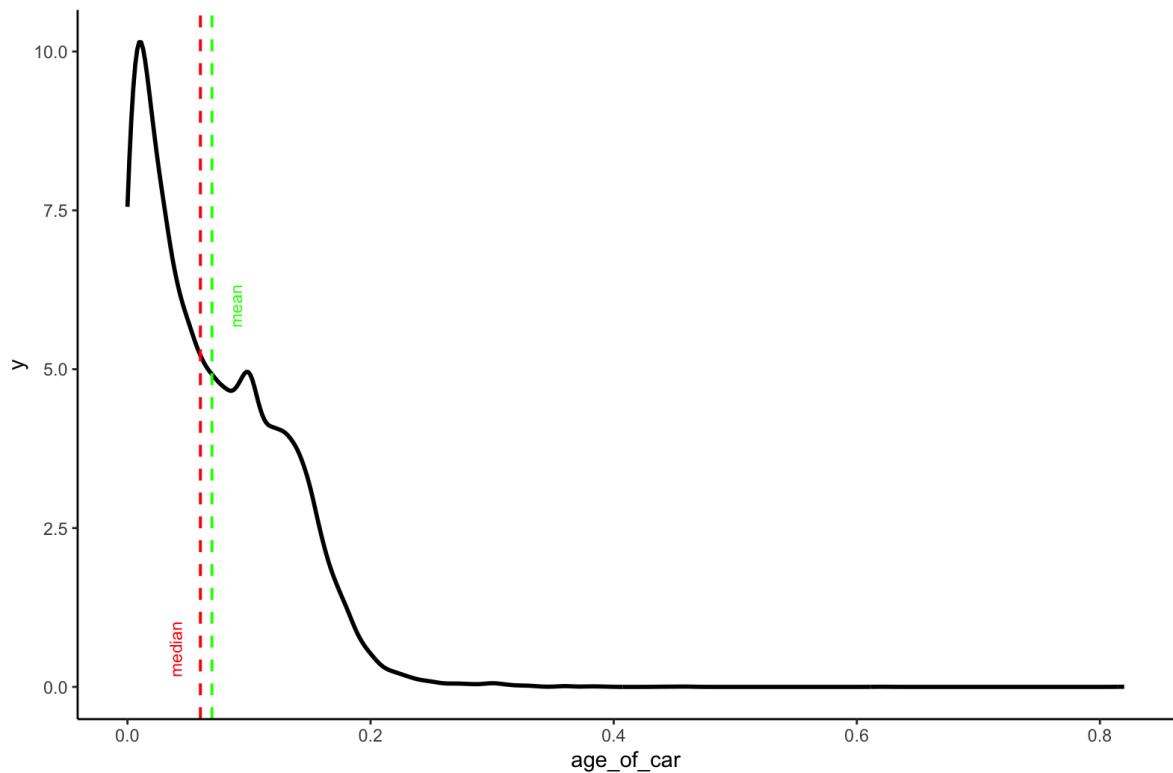


Figure 15: Age of Car Density diagram

Mean is larger than median so the distribution is right-skewed.
We can see this fact by calculating skewness:

```
> skewness(car$age_of_car)
[1] 0.9014617
```

Figure 16: skewness

It's closed to 1 so we have strong right-skewed.

3. Question 2

Chosen Categorical Variable : transmission_type: Transmission type of the car

3.1 A

Create a frequency table for this variable.

```
1 transmission <- car$transmission_type
2 gp <- unique(transmission)
3
4 Manual = transmission[transmission == 'Manual']
5 Automatic = transmission[transmission == 'Automatic']
6
7 frequencies<-c(length(Manual),length(Automatic))
8 percentage  <- round(100*frequencies/sum(frequencies), 2)
9
10 transmission.categorized <- data.frame(types = c("Manual", "Automatic"
    ),percentage = percentage, value = frequencies)
```

| | types | percentage | value |
|---|-----------|------------|-------|
| 1 | Manual | 64.99 | 19498 |
| 2 | Automatic | 35.01 | 10502 |

Figure 17: frequency table for transmission_type.

As you can see, most of the cars are still manual. It's about 64% of the population.

3.2 B

Plot a horizontal bar plot, sort by frequency, and use different colors for each category.

```
1 transmission.categorized <- transmission.categorized[order(percentage)
   ,]
2 ggplot(data=transmission.categorized, aes(x=types, y=frequencies, fill
   =types)) +
3   geom_bar(stat="identity", alpha = 0.7, width = 0.7) +
4   labs(title="Barplot of Transmission Type") +
5   coord_flip()
```

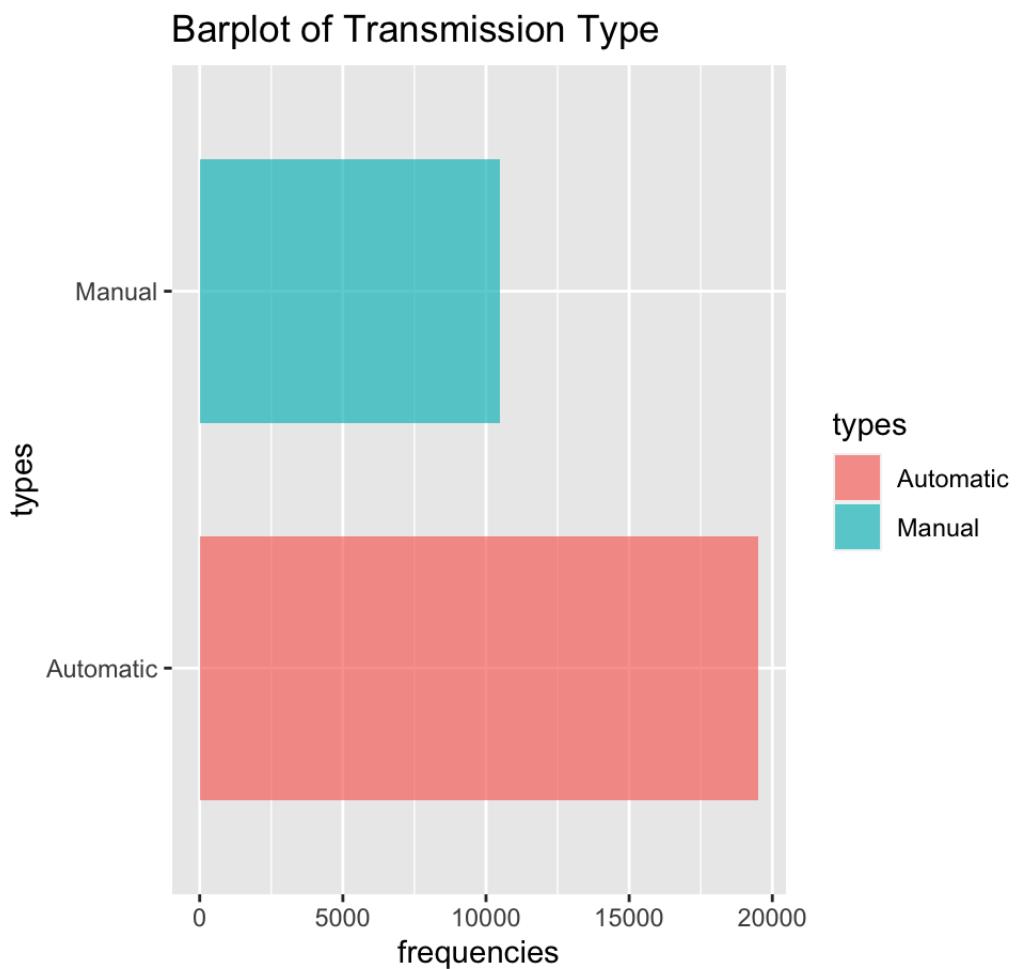


Figure 18: Horizontal Bar Plot of Transmission Type

3.3 C

Plot a bar plot for this variable and add percentage marks to it.

```
1 ggplot(data=transmission.categorized, aes(x=types, y=percentage, fill=
2   types)) +
3   geom_bar(stat="identity", alpha = 0.7, width = 0.6) +
4   labs(title="Barplot of Transmission Type") +
5   geom_text(aes(label=paste(percentage, "%")), vjust=-0.4, size=4)
```

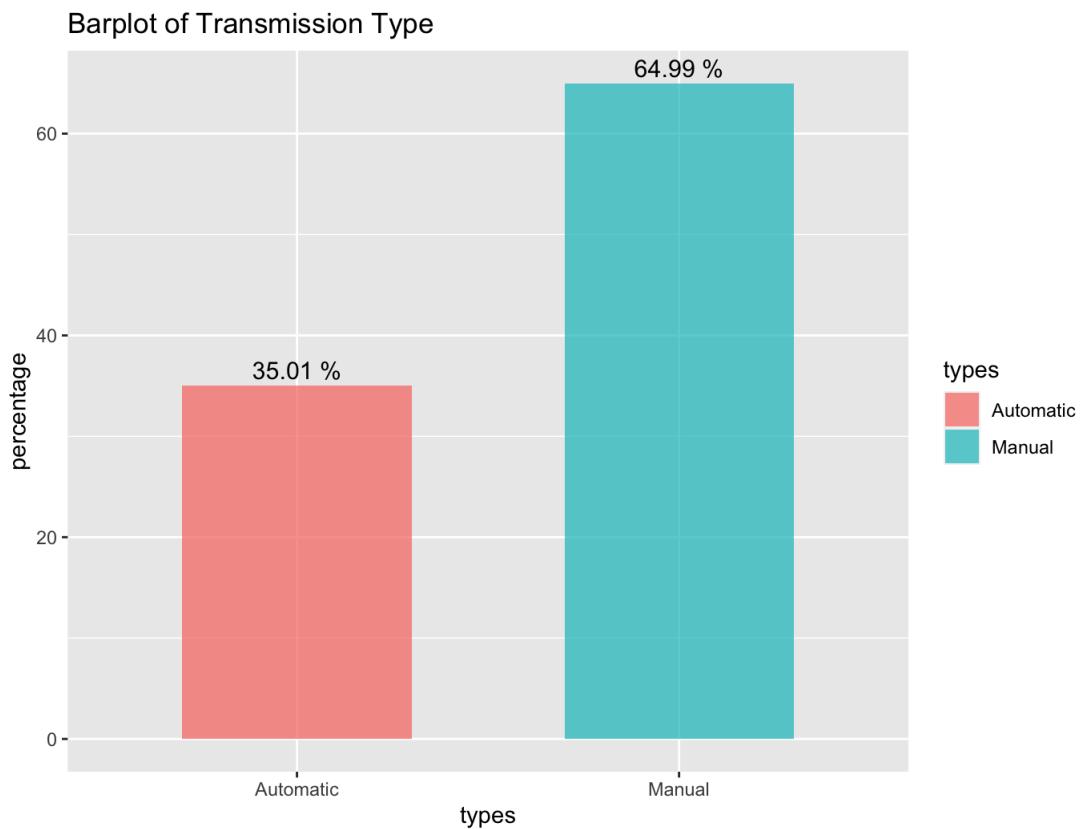


Figure 19: Bar Plot of Transmission Type

3.4 D

Plot a violin plot for this variable.

```
1 ggplot(data = car, aes(x=transmission_type, y=age_of_car, fill=
2   transmission_type))+  
2   geom_violin(trim=FALSE, alpha = 0.7)+  
3   labs(title="Violin Plot of age of car by transmission type")
```

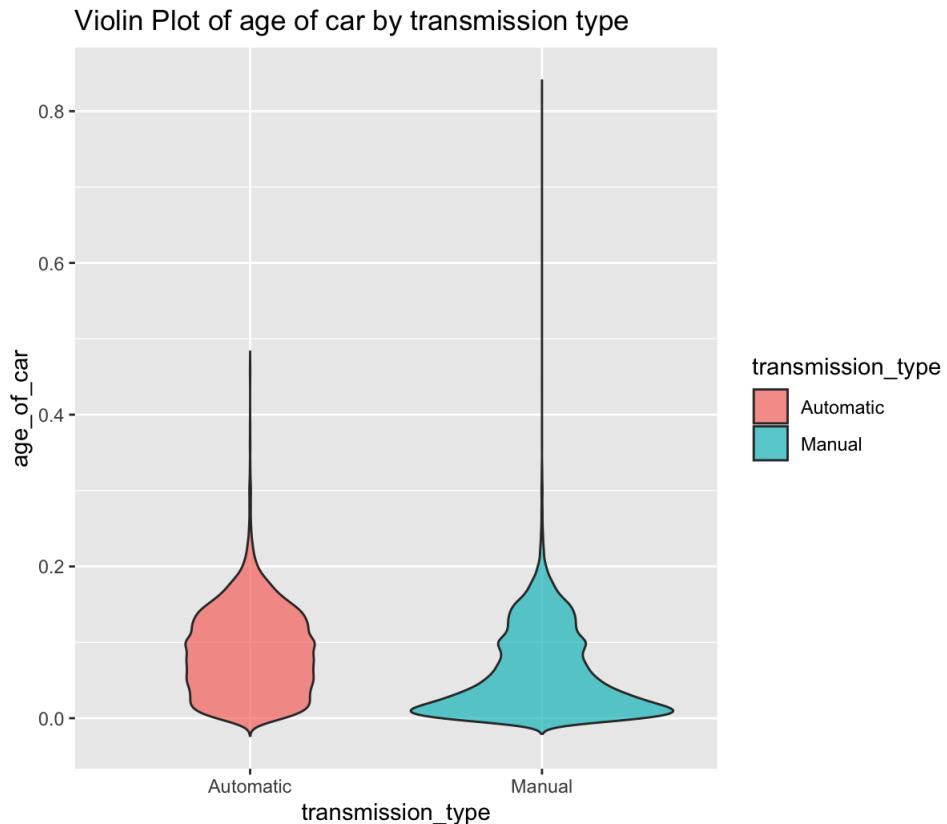


Figure 20: Violin Plot

Automatic cars have mostly normal distribution and show that these models produce with a fixed range during the time.

Manual cars recently produced or bought more.

4. Question 3

Chosen Numerical Variable : age_of_car and age_of_policyholder

4.1 A

Draw a scatter plot for two variables and describe the relationship between them based on the plot.

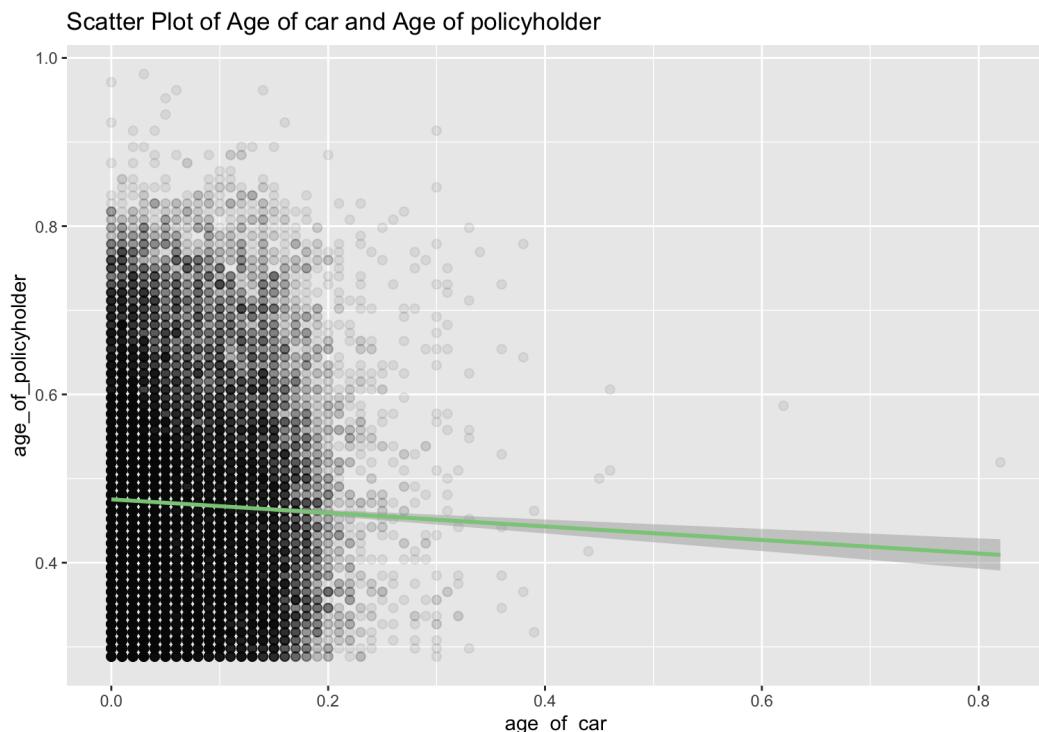


Figure 21: Scatter plot

As you can see, the relationship between these two variables is not easy. We can't describe a clear function between them. Some outliers are detected.

Because both of my variables represent age, the majority of points can be found in a specific area between 0 to 0.2 for cars' age and 0 to 0.8 for policyholders' age.

About outliers: people who are old can't drive, and old cars should be scrapped, so they occur in my dataset rarely.

4.2 B

Select a categorical variable, and determine the samples either by the symbol or by the color (or by both) in a scatter plot that has been drawn in section “A”. Does the relation still hold for different categories?

Select fuel type as categorical variable.

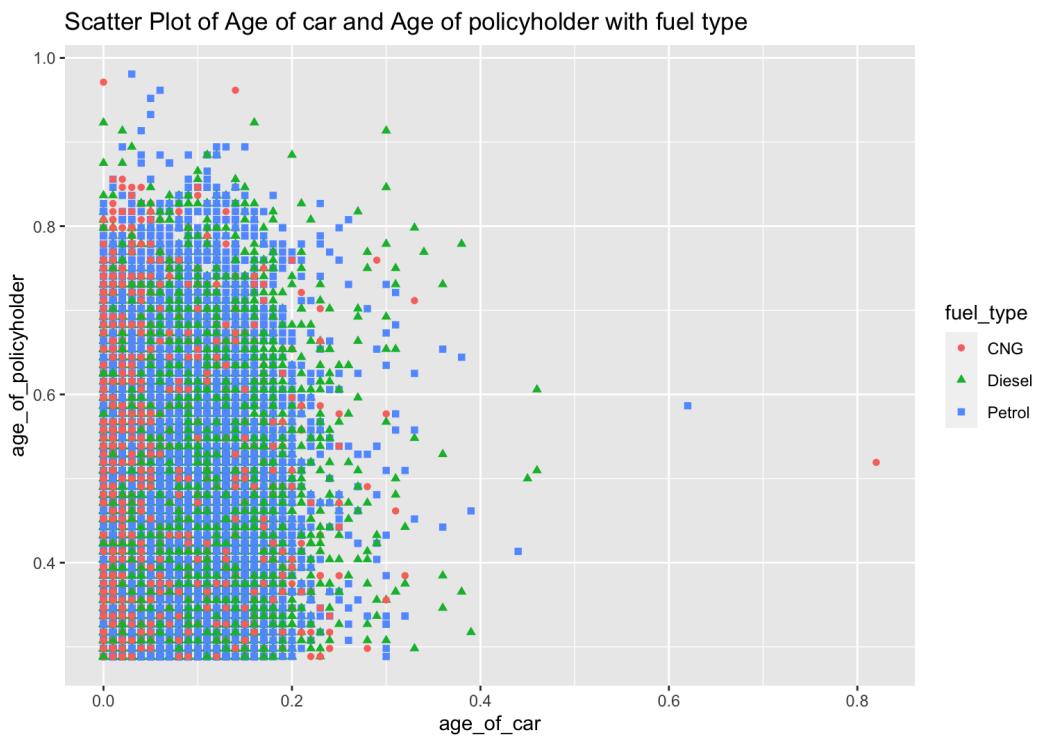


Figure 22: Scatter plot

Because there was no relationship in the previous part, this is not observed either. It can only be said that newer cars use CNG fuel.

4.3 C

Calculate the correlation coefficient for these two variables. Using the “cor.test” function, we can also test the significance of a correlation. Are the variables correlated? According to the test, what is shown by the p value, and what is the intuition of the p value?

```
> cor.test(car$age_of_car, car$age_of_policyholder)

Pearson's product-moment correlation

data: car$age_of_car and car$age_of_policyholder
t = -6.3982, df = 29998, p-value = 1.595e-10
alternative hypothesis: true correlation is not equal to 0
95 percent confidence interval:
-0.04821195 -0.02561091
sample estimates:
      cor
-0.03691615
```

Figure 23: Scatter plot

As expected, there is no correlation between these two variables, and it is close to zero. P value define on hypothesis test.

On average, is there a relationship between two means of these variables?

H₀ : means are close two each other.

H_a : means are not close two each other.

And we can see p value is close to 0, so we reject H₀.

These two expressions have the same Concept in this context, the population has no powerful linear correlation between age of car and age of policyholder.

4.4 D

A hexbin (Figure 1) plot with marginal distribution is like a two-dimensional histogram. The data is divided into bins, and the color strength of each bin represents the number of data points in that bin. Also, each dimension has its own distribution in front of its axis. Draw the hexbin plot with marginal distribution for chosen variables. What is your interpretation? Discuss the bin size and how it changes the result.

```
1 p12 <- ggplot(car, aes(age_of_car, age_of_policyholder)) +
2   geom_point() +
3   geom_hex(bins = 12) +
4   geom_smooth(color='#FFFF00', se=FALSE) +
5   ggtitle("Hexbin Plot - Age of car and policyholder with binsize 12")
6
7 ggMarginal(p12, type="histogram", size=3, fill='lightblue')
```

Figure 24: hexbin plot with bin size 7
 Hexbin Plot - Age of car and policyholder with binsize 7

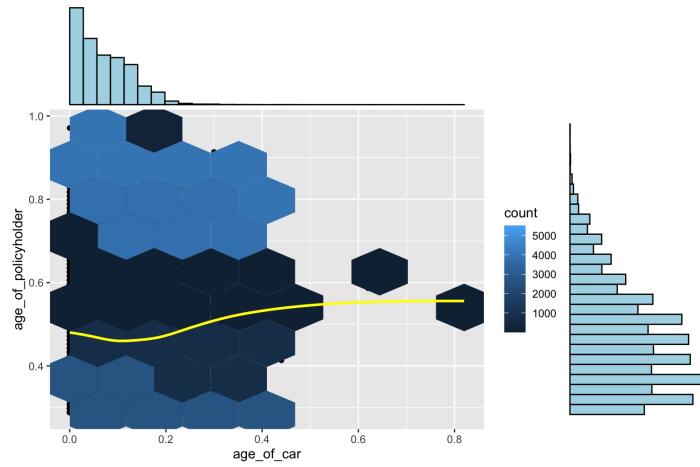


Figure 25: hexbin plot with bin size 12
 Hexbin Plot - Age of car and policyholder with binsize 12

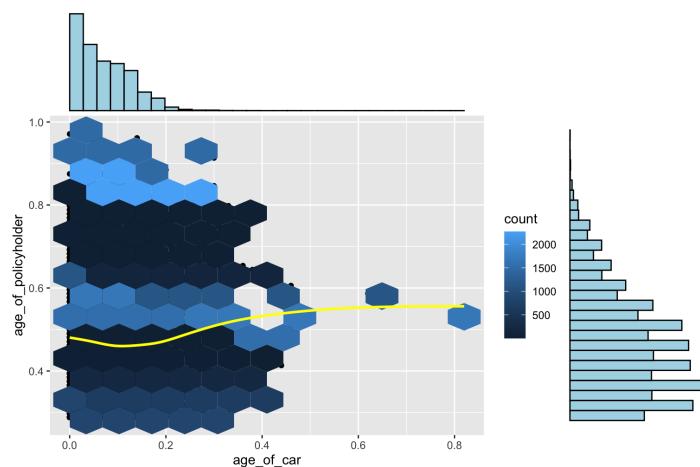


Figure 26: hexbin plot with bin size 30
 Hexbin Plot - Age of car and policyholder with binsize 30

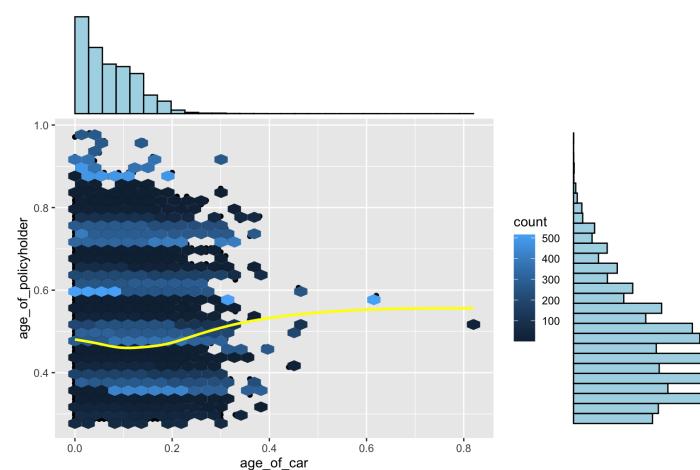


Figure 27: hexbin plot of age_of_car and age_of_policyholder with different bin size.

This graph is used when the number of samples is large, and the scatter plot cannot determine the distribution of the data locally.

In this diagram, the coordinate space is divided into bins, where the number of samples in each bin is indicated by color.

So we can see some interesting information, for example, older people who are policy-holders support more cars and, in other words, have more consumers, and this may be because of their experience in this market or people trust them more because of their age.

About bin size, when the number of bins increases, the hexbin size decreases, and the hexbin becomes sparser.

According to the figure, it seems bin size 12 is more informative, and when increasing bin size, this graph is close to a scatter plot.

4.5 E

Draw the 2D density plot for chosen variables. How do you interpret the resulting graph? Describe the advantages and disadvantages of the 2D density and hexbin graph.

Density can be represented in the form of 2D density graphs or density plots. A 2d density chart displays the relationship between 2 numeric variables, where one variable is represented on the X-axis, the other on the Y axis, like for a scatterplot.

```
1 ggplot(car, aes(x=age_of_car, y=age_of_policyholder) ) +  
2   stat_density_2d(aes(fill = ..level..), geom = "polygon") +  
3   scale_fill_continuous(type = "viridis") +  
4   labs(title = "2D density plot Age of car and Age of policyholder")  
5
```

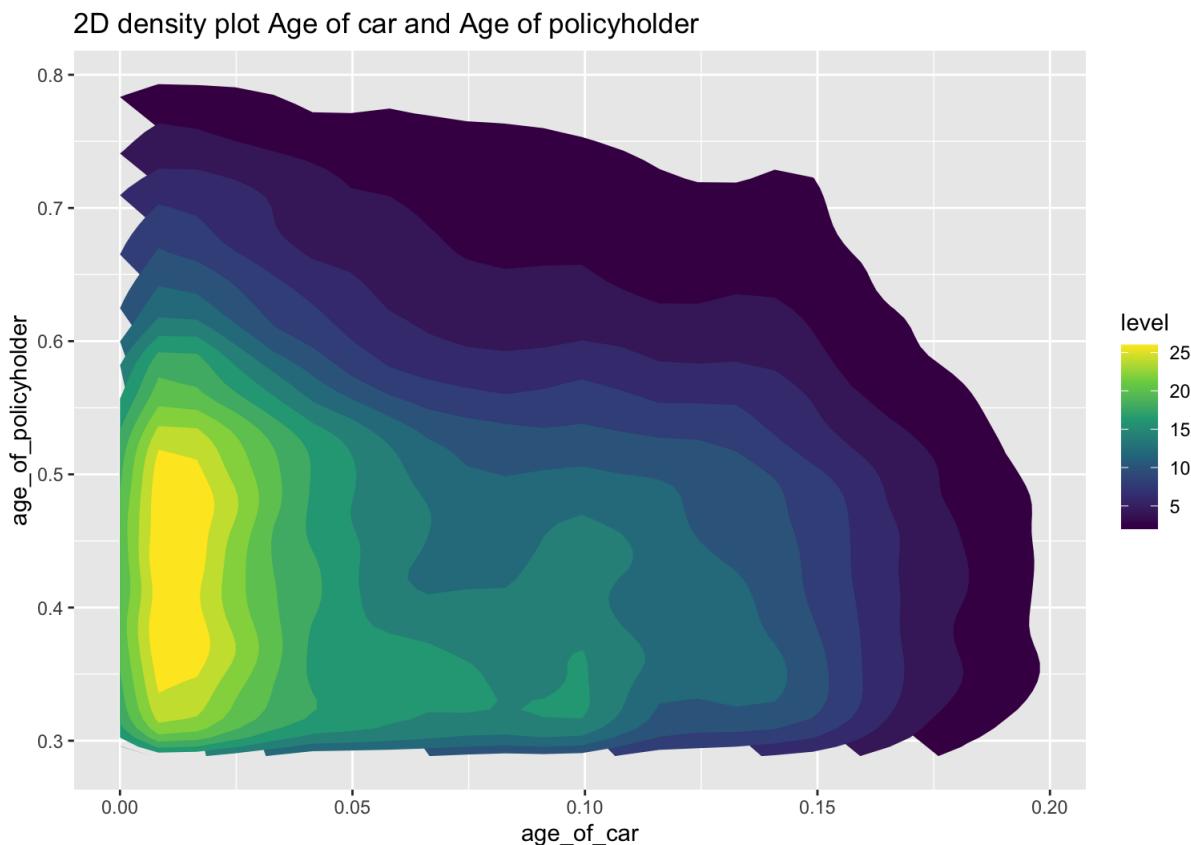


Figure 28: 2D density plot

As can be concluded from the graph, this data and relation is too depth, and learning this is so hard.

In hexbin, the number of bins can be controlled, which changes the shape of the distribution but not in the density plot.

In density and hexbin plots can see the pick of data but in hexbin show the local distribution in a better way by use color.

5. Question 4

5.1 A

Create a heatmap correlogram from your variables.

Annotate each cell with their corresponding Pearson's correlation coefficients and p-value as well.

Use red for the positive correlation and blue for the negative correlation.

Highlight significant correlations

```
1 #some preprocessing
2 car$is_claim <- as.logical(car$is_claim)
3 car <- subset(car, select = -c(X))
4
5 #separate numerical features
6 car.numeric <- Filter(is.numeric, car)
7
8 #correlation matrix
9 cor_mat <- cor(as.matrix(car.numeric), method="pearson")
10 cor_mat <- round(cor_mat, 2)
11
12 #calc p_values
13 p.mat <- cor_pmat(car.numeric)
14
15 ggcorrplot(
16   cor_mat, hc.order = TRUE, type = "upper",
17   lab = TRUE , p.mat = p.mat, sig.level = 0.05,
18   title = "heatmap correlogram",
19   colors = c("blue", "white", "red"),
20   hc.method = "complete", lab_size = 2,
21   as.is = FALSE)
22
```

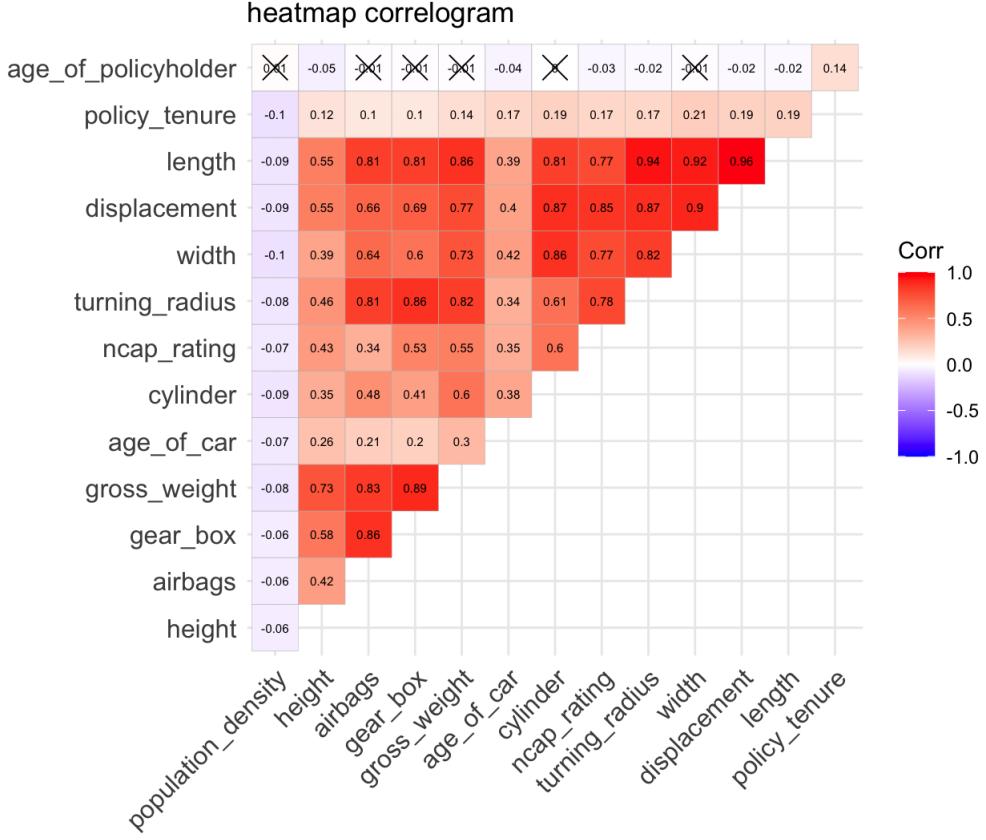


Figure 29: heatmap correlogram for numeric features

$p - value > 0.05$ show with red box.

| | policy_tenure | age_of_car | age_of_policyholder | population_density | airbags | displacement | cylinder | gear_box | turning_radius | length | width | height | gross_weight | ncap_rating |
|---------------------|---------------|---------------|---------------------|--------------------|---------------|---------------|---------------|---------------|----------------|---------------|---------------|--------------|---------------|---------------|
| policy_tenure | 0.000000E+00 | 8.161691E-187 | 4.944304E-138 | 7.665155E-65 | 1.356004E-70 | 1.098616E-251 | 3.290997E-238 | 2.883208E-61 | 4.609704E-187 | 2.386836E-242 | 2.315608E-299 | 6.097440E-94 | 6.375608E-127 | 6.443607E-203 |
| age_of_car | 8.161691E-187 | 0.000000E+00 | 1.595233E-10 | 8.414501E-34 | 7.760921E-291 | 0.000000E+00 | 0.000000E+00 | 2.962429E-270 | 0.000000E+00 | 0.000000E+00 | 0.000000E+00 | 0.000000E+00 | 0.000000E+00 | 0.000000E+00 |
| age_of_policyholder | 4.944304E-138 | 1.595233E-10 | 0.000000E+00 | 2.023697E-01 | 5.564222E-02 | 3.051618E-05 | 4.615045E-01 | 3.789340E-01 | 2.807660E-03 | 2.293577E-04 | 2.574386E-01 | 4.039960E-20 | 1.888151E-01 | 1.402864E-07 |
| population_density | 7.665155E-65 | 8.414501E-34 | 2.023697E-01 | 0.000000E+00 | 2.965172E-29 | 3.844415E-56 | 3.280059E-57 | 9.925498E-23 | 2.903633E-41 | 1.885420E-57 | 9.383007E-64 | 7.551083E-29 | 5.423320E-43 | 6.334509E-32 |
| airbags | 1.356004E-70 | 7.760921E-291 | 5.564222E-02 | 2.965172E-29 | 0.000000E+00 | 0.000000E+00 | 0.000000E+00 | 0.000000E+00 | 0.000000E+00 | 0.000000E+00 | 0.000000E+00 | 0.000000E+00 | 0.000000E+00 | 0.000000E+00 |
| displacement | 1.098616E-251 | 0.000000E+00 | 3.051618E-05 | 3.844415E-56 | 0.000000E+00 | 0.000000E+00 | 0.000000E+00 | 0.000000E+00 | 0.000000E+00 | 0.000000E+00 | 0.000000E+00 | 0.000000E+00 | 0.000000E+00 | 0.000000E+00 |
| cylinder | 3.290997E-238 | 0.000000E+00 | 4.615045E-01 | 3.280059E-57 | 0.000000E+00 | 0.000000E+00 | 0.000000E+00 | 0.000000E+00 | 0.000000E+00 | 0.000000E+00 | 0.000000E+00 | 0.000000E+00 | 0.000000E+00 | 0.000000E+00 |
| gear_box | 2.883208E-61 | 2.962429E-270 | 3.789340E-01 | 9.925498E-23 | 0.000000E+00 | 0.000000E+00 | 0.000000E+00 | 0.000000E+00 | 0.000000E+00 | 0.000000E+00 | 0.000000E+00 | 0.000000E+00 | 0.000000E+00 | 0.000000E+00 |
| turning_radius | 4.609704E-187 | 0.000000E+00 | 2.807660E-03 | 2.903633E-41 | 0.000000E+00 | 0.000000E+00 | 0.000000E+00 | 0.000000E+00 | 0.000000E+00 | 0.000000E+00 | 0.000000E+00 | 0.000000E+00 | 0.000000E+00 | 0.000000E+00 |
| length | 2.386836E-242 | 0.000000E+00 | 2.293577E-04 | 1.885420E-57 | 0.000000E+00 | 0.000000E+00 | 0.000000E+00 | 0.000000E+00 | 0.000000E+00 | 0.000000E+00 | 0.000000E+00 | 0.000000E+00 | 0.000000E+00 | 0.000000E+00 |
| width | 2.315608E-299 | 0.000000E+00 | 2.574386E-01 | 9.383007E-64 | 0.000000E+00 | 0.000000E+00 | 0.000000E+00 | 0.000000E+00 | 0.000000E+00 | 0.000000E+00 | 0.000000E+00 | 0.000000E+00 | 0.000000E+00 | 0.000000E+00 |
| height | 6.097440E-94 | 0.000000E+00 | 4.039960E-20 | 7.551083E-29 | 0.000000E+00 | 0.000000E+00 | 0.000000E+00 | 0.000000E+00 | 0.000000E+00 | 0.000000E+00 | 0.000000E+00 | 0.000000E+00 | 0.000000E+00 | 0.000000E+00 |
| gross_weight | 6.375608E-127 | 0.000000E+00 | 1.888151E-01 | 5.423320E-43 | 0.000000E+00 | 0.000000E+00 | 0.000000E+00 | 0.000000E+00 | 0.000000E+00 | 0.000000E+00 | 0.000000E+00 | 0.000000E+00 | 0.000000E+00 | 0.000000E+00 |
| ncap_rating | 6.443607E-203 | 0.000000E+00 | 1.402864E-07 | 6.334509E-32 | 0.000000E+00 | 0.000000E+00 | 0.000000E+00 | 0.000000E+00 | 0.000000E+00 | 0.000000E+00 | 0.000000E+00 | 0.000000E+00 | 0.000000E+00 | 0.000000E+00 |

Figure 30: report p values

Most of the features have little p value.

5.2 B

Display all the bivariate relations between the variables using a correlogram where each element is a scatter plot between two variables. Can you find any meaningful pattern between them?

Select some important features because the computation is high.

```

1 selected_var = data.frame("policy_tenure" = car$policy_tenure,
2   "car's age" = car$age_of_car,
3   "policyholder's age"=car$age_of_policyholder,
4   "population density"= car$population_density,
5   "displacement"=car$displacement,
6   "turning radius"=car$turning_radius,
7   "length"=car$length,
8   "width"=car$width,
9   "height"=car$height,
10  "gross_weight" = car$gross_weight)
11
12 ggpairs(selected_var, title = "Correlogram",
13   upper = list(continuous = wrap("density", alpha = 0.3)),
14   lower = list(continuous = wrap("points", colour="grey", alpha =
0.3), combo = "dot_no_facet"))

```

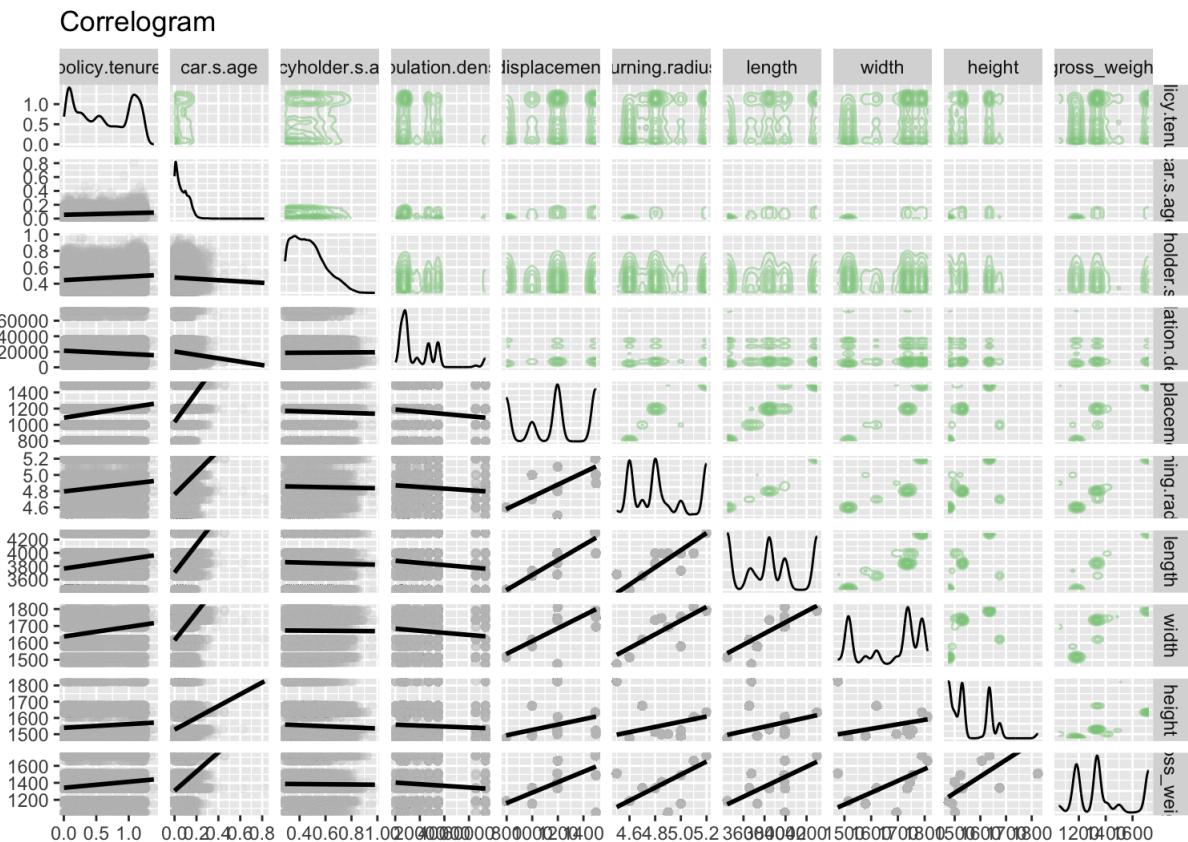


Figure 31: correlogram

Some conclusions:

Age has right-skew distribution.

As expected, width, length, and height have a positive correlation. the higher the width, the higher the length and height, and some categories, like driving cars and cranes.

Policy tenure has two picks because of short-term and long-term options for consumers. Population show we have one crowded and two medium population and others are Low population.

Gross weight has a positive correlation with the physical features of the cars.

The length has a high tail, and the width has a low tail.

5.3 C

Choose 3 numerical and 1 categorical variable from your dataset. Draw a 3D scatter plot for the numerical variables and use the categorical variable as the points' color. Describe the relation between them

Choose population-density, age-of-policyholder and displacement as numerical variable and fuel-type as categorical.

```
1 colors <- c("#999999", "#E69F00", "#56B4E9")
2 fuel_type <- colors[as.numeric(factor(car$fuel_type))]
3
4 scatterplot3d(
5   main="3D Scatter Plot",
6   car$population_density,
7   car$age_of_policyholder,
8   car$displacement,
9   color=fuel_type,
10  xlab = "population_density",
11  ylab = "age_of_policyholder",
12  zlab = "displacement",
13  pch = 16)
14
15 legend("right", legend = unique(car$fuel_type), pch = 16, col =
  colors)
```

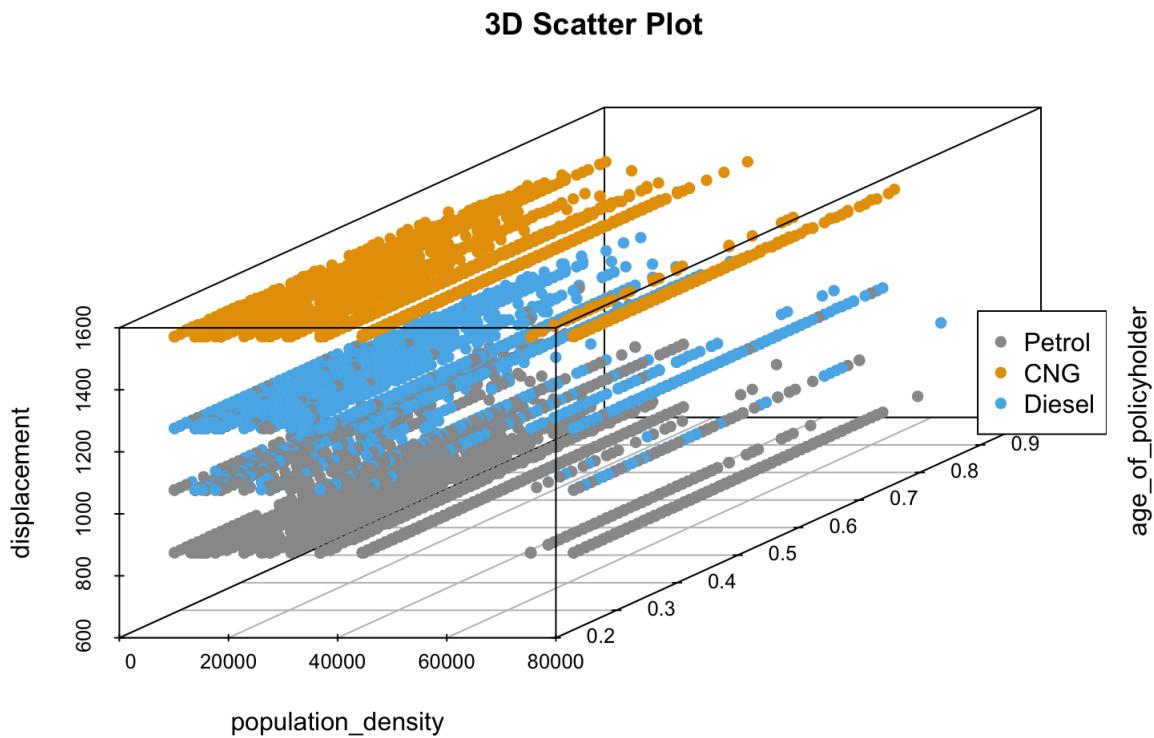


Figure 32: 3D scatter plot

As you can see, the type of fuel and displacement have a good relationship, cars with low displacement use diesel and cars with medium displacement use petrol, and cars with high displacement use CNG.

Cars like trucks usually use CNG, which drives between the city and travels a long distance, and they don't pollute the air.

Cars in the town typically use diesel fuel because their price is low, like vans or mini trucks.

And people use petrol fuel for daily transportation, which usually has low displacement.

6. Question 5

Chosen Categorical Variables : fuel_type and transmission_type

6.1 A

Contingency table

```
1 Total <- sum
2 t<-as.data.frame.matrix(addmargins(table(car[,c('fuel_type',
3 transmission_type')]), FUN = Total))
```

| | Automatic | Manual | Total |
|--------|-----------|--------|-------|
| CNG | 0 | 10370 | 10370 |
| Diesel | 7171 | 1931 | 9102 |
| Petrol | 3331 | 7197 | 10528 |
| Total | 10502 | 19498 | 30000 |

Figure 33: Contingency table of fuel_type and transmission_type

Automatic cars don't use CNG.

Diesel is used in the automatic car more than manual ones.

This relation is reversed for petrol fuel.

6.2 B

Grouped Bar Charts

```
1 ggplot(car, aes(x = fuel_type , fill = transmission_type)) +
2   geom_bar(position = "dodge", alpha = 0.6) +
3   geom_text(aes(label = ..count..), stat = "count", vjust = -.3, size
4             = 3, color = 'black', position = position_dodge(width = 1)) +
5   labs(title="Grouped Bar Chart", x="fuel_type")
```

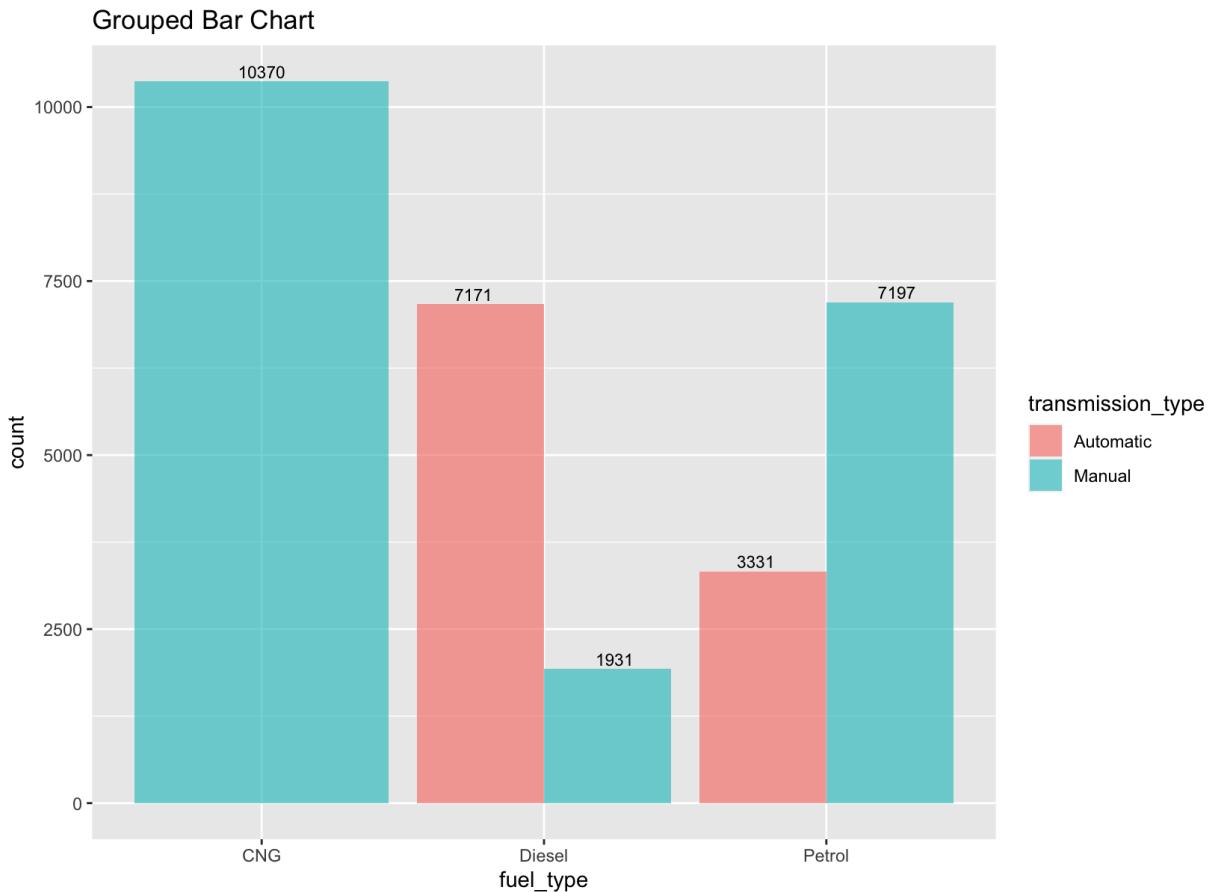


Figure 34: Grouped Barplot of fuel_type and transmission_type

6.3 C

Segmented Bar Chart

```
1 ggplot(car, aes(x = fuel_type, fill = transmission_type)) +
2   geom_bar(position = "stack", alpha = 0.6) +
3   geom_text(aes(label=..count..), stat='count', vjust = -.3, size =
4     3, color = 'black', position = position_stack(0.5)) +
5   labs(title=paste("Segmented barplot"))
```

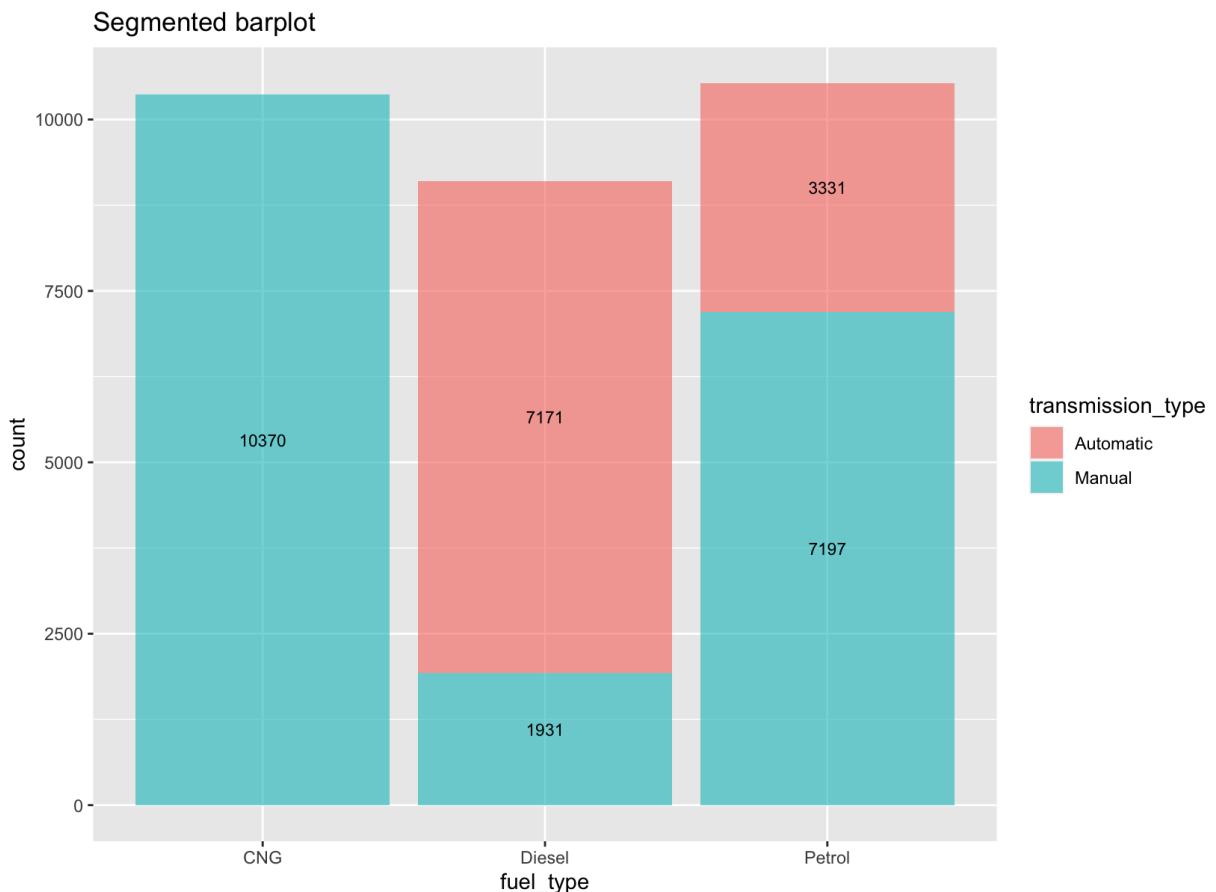


Figure 35: Segmented Barplot of fuel_type and transmission_type

6.4 D

Mosaic Plott

```
1 ggplot(data = car) +  
2   geom_mosaic(aes(x = product(transmission_type), fill = fuel_type ),  
3   alpha = 0.6) +  
4   annotate("text", x = 0.195, y = 0.35,label= paste(round(t[2,1]/t  
[4,1],2)*100,"%"),fontface = "bold",size=4,color="black") +  
5   annotate("text", x = 0.195, y = 0.855,label= paste(round(t[3,1]/t  
[4,1],2)*100,"%"), fontface = "bold",size = 4,color="black") +  
6   annotate("text", x = 0.7, y = 0.25,label= paste(round(t[1,2]/t  
[4,2],2)*100,"%"),fontface = "bold",size = 4,color="black") +  
7   annotate("text", x = 0.7, y = 0.57,label= paste(round(t[2,2]/t  
[4,2],2)*100,"%"), fontface = "bold",size = 4,color="black") +  
8   annotate("text", x = 0.7, y = 0.8,label= paste(round(t[3,2]/t[4,2],2)  
*100,"%"), fontface = "bold",size = 4,color="black") +  
9   labs(title= paste("Mosaic plot"))
```

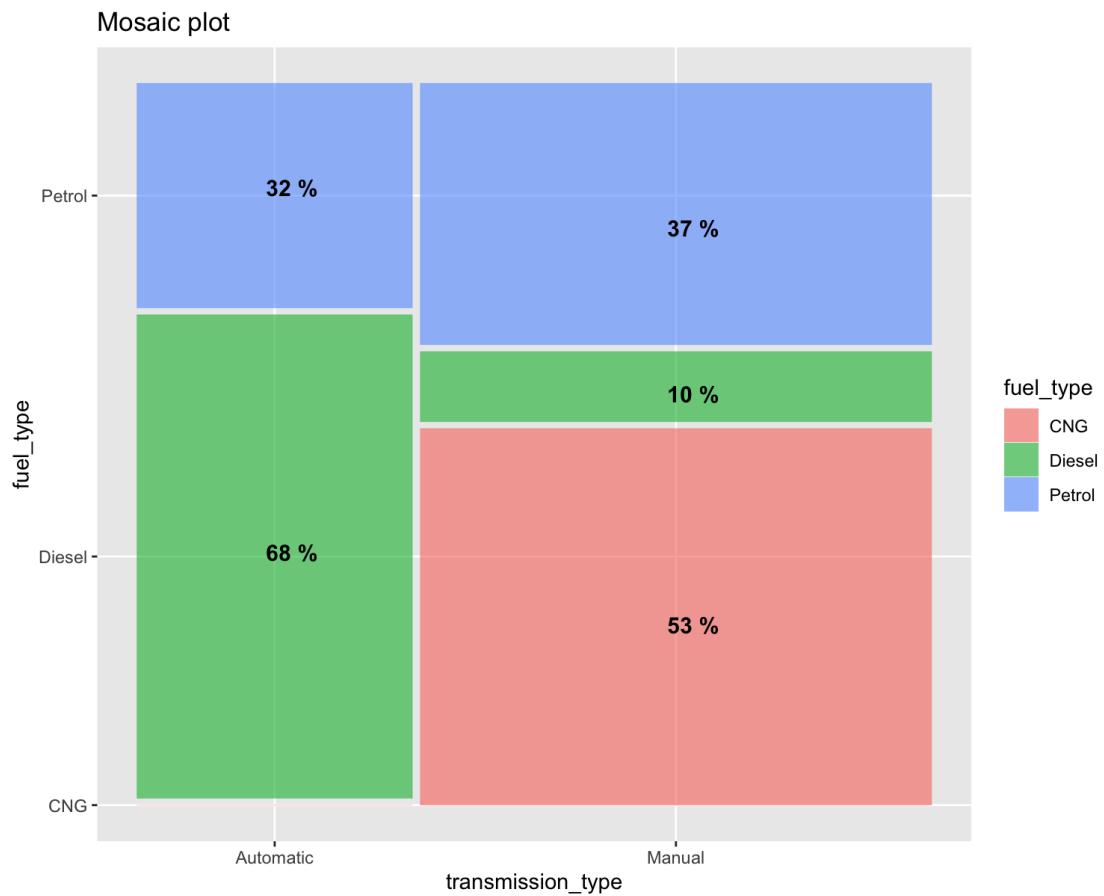


Figure 36: Mosaic plot of fuel_type and transmission_type

7. Question 6

In this question, you will conduct a hypothesis test for two numerical variables. Choose a random sample of 25 data points from the dataset and choose two numerical variables that are not of a corresponding quantity. This data can be used to compare the average quantity between the two variables.

Chosen Numerical Variables : displacement and policy_tenure

7.1 A

What is the best method for testing our hypothesis? Is it a t-test or a z-test? Explain it.

$n = 25$

Checking conditions:

1. Independence: random sample and $n < 10\%$ of all population (30000)
* Suppose the tenures of policies are independent because they could be dependent on one another. For example, a company has two types of insurance policies for insuring cars in a small town, but we assume, they are independence.
2. Sample size/skew: $n < 30$.

So, we can't assume that the sampling distribution from this dataset will be nearly normal.

But we can use Student's t-distribution for this problem.

7.2 B

Design a hypothesis test to see if these data provide convincing evidence of a difference between mean values. Does the result agree with the 95% confidence interval?

Conditions for inference for comparing two independent means

1. Independence:
 - 1.1 within groups: sampled observations must be independent
 - 1.1.1 random sample/assignment
 - 1.1.2 if sampling without replacement, $n < 10\%$ of population
 - 1.2 between groups: the two groups must be independent of each other (non-paired)
2. Sample size/skew: The more skew in the population distributions, the higher the sample size needed

All conditions are met.

Let's estimating the Mean

Reject $H_0 \rightarrow, P_{value} < 0.05$ so H_0 reject,
it means between mean of $\mu_{displacement}$ and $\mu_{policy\ tenure}$ are not equal.

Calculate Confidence interval 95%:

$$df = n - 1 = 24$$

$$SE_{displacement-policy\ tenure} = \sqrt{\frac{s_{displacement}^2}{n_{displacement}} + \frac{s_{policy\ tenure}^2}{n_{policy\ tenure}}} = 29249.43$$

```
> margin <- qt(0.975,df=n+n-1)*sqrt(sp/n + sp/n)
> margin
[1] 97.20933
>
> lowerinterval <- (xbar1-xbar2) - margin
> lowerinterval
[1] -1197.122
>
> upperinterval <- (xbar1-xbar2) + margin
> upperinterval
[1] -1002.703
>
> cat("95% confidence interval is= (", lowerinterval, ", ",upperinterval, ")")
95% confidence interval is= (-1197.122 , -1002.703 )
```

Figure 37: Calculate CI 95%

$$0 \notin (-1197.122, -1002.703)$$

T-test result and 95% confidence interval agree.

8. Question 7

Chosen Numerical Variable : age_of_policyholder

$$n = 100, \bar{x} = 0.4760577, s = 0.01150644$$

Checking conditions:

1. Independence: random sample and $n < 10\%$ of all population of people who has car in this dataset($100 < 3000$).
2. Sample size/skew: $n \geq 30$.

Can see histogram of age_of_policyholder below:

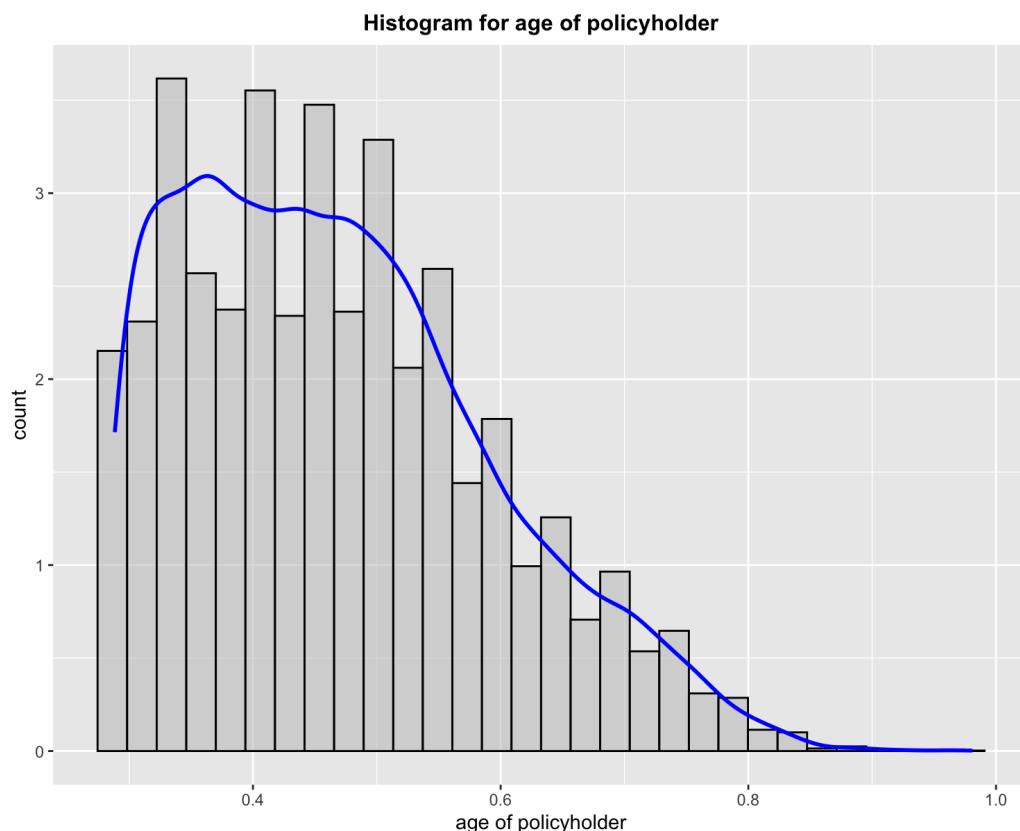


Figure 38: Histogram of age_of_policyholder

We have a right-skewed distribution, and it's far from normal, so for using CLT, we must increase the number of samples by choosing $n = 100$, and we can assume that the sampling distribution of size 100 will be nearly normal.

Central Limit Theorem (CLT)

$$\bar{x} \sim N \left(\text{mean} = \mu, SE = \frac{\sigma}{\sqrt{n}} \right) = N(0.476057, 0.01150644)$$

8.1 A

Calculate a 98% confidence interval for the mean of this variable

```
1 n = 100
2
3 #reads the dataset 'car' and take the 100 rows as sample
4 sdf<- sample(1:nrow(car), n)
5
6 #sample 10 rows
7 sub_car <- car[sdf,]
8
9 xbar = mean(sub_car$age_of_policyholder)
10
11 ci = 0.98
12 z = qnorm((1-ci)/2)
13
14 se <- sd(sub_car$age_of_policyholder) / sqrt(n) #SE = s/sqrt(n)
15
16 lower <- xbar + z * se
17 upper <- xbar - z * se
18
19 cat("98%\% confidence interval is = (", lower, ", ", upper, ")")
```

approximate 98% CI for μ : 0.4760577 ± 2.32 SE

$$(0.4492897, 0.5028257)$$

We are 98% confidence that the mean point of age of policyholder is between \$0.4492897 and \$0.5028257.

8.2 B

Interpret this confidence interval. In this context, what does a 98% confidence level mean?

We are 98% confidence that the mean point of age of policyholder is between \$0.4492897 and \$0.5028257.

```
> mean(car$age_of_policyholder)
[1] 0.4696782
```

Figure 39: Actual mean of age of policyholder.

As we can see, this claim is valid, and the actual mean is in this interval.

8.3 C

Plot the histogram of the variable and mark the mean of all the samples as a vertical line on top of the histogram plot. You must also mark the confidence intervals on the plot as two vertical lines

```
1 ggplot(sub_car, aes(x = age_of_policyholder)) +
2   geom_histogram(aes(y = ..density..), alpha=0.5, colour = "black") +
3   labs(
4     title = "Histogram for sample of age of policyholder",
5     x = "age of policyholder",
6     y = "count"
7   ) +
8   geom_vline(aes(xintercept = mean(age_of_policyholder)), linetype = "
9     dashed", size = .7, col= 'red') +
10  geom_vline(aes(xintercept = lower), linetype = "dashed", size = .7,
11    col= 'green') +
12  geom_vline(aes(xintercept = upper), linetype = "dashed", size = .7,
13    col= 'blue') +
14  geom_vline(aes(xintercept = mean(car$age_of_policyholder)), linetype =
15    "dashed", size = .7, col= 'orange') +
16  theme(
17    plot.title = element_text(size = 12, face = "bold", hjust = 0.5),
18    plot.caption = element_text(face = "italic") ) +
19  annotate("text", x = mean(sub_car$age_of_policyholder) + 0.007 ,
20    label = "sample mean", y = 6, size = 4, angle = 90 , color = 'red'
21  ) +
22  annotate("text", x = lower - 0.01 , label = "lower", y = 6, size =
23    4, angle = 90, color = 'green') +
24  annotate("text", x = upper - 0.01 , label = "upper", y = 6, size =
25    4, angle = 90, color = 'blue') +
26  annotate("text", x = mean(car$age_of_policyholder) - 0.01 , label =
27    "actual mean", y = 6, size = 4, angle = 90, color = 'orange')
```

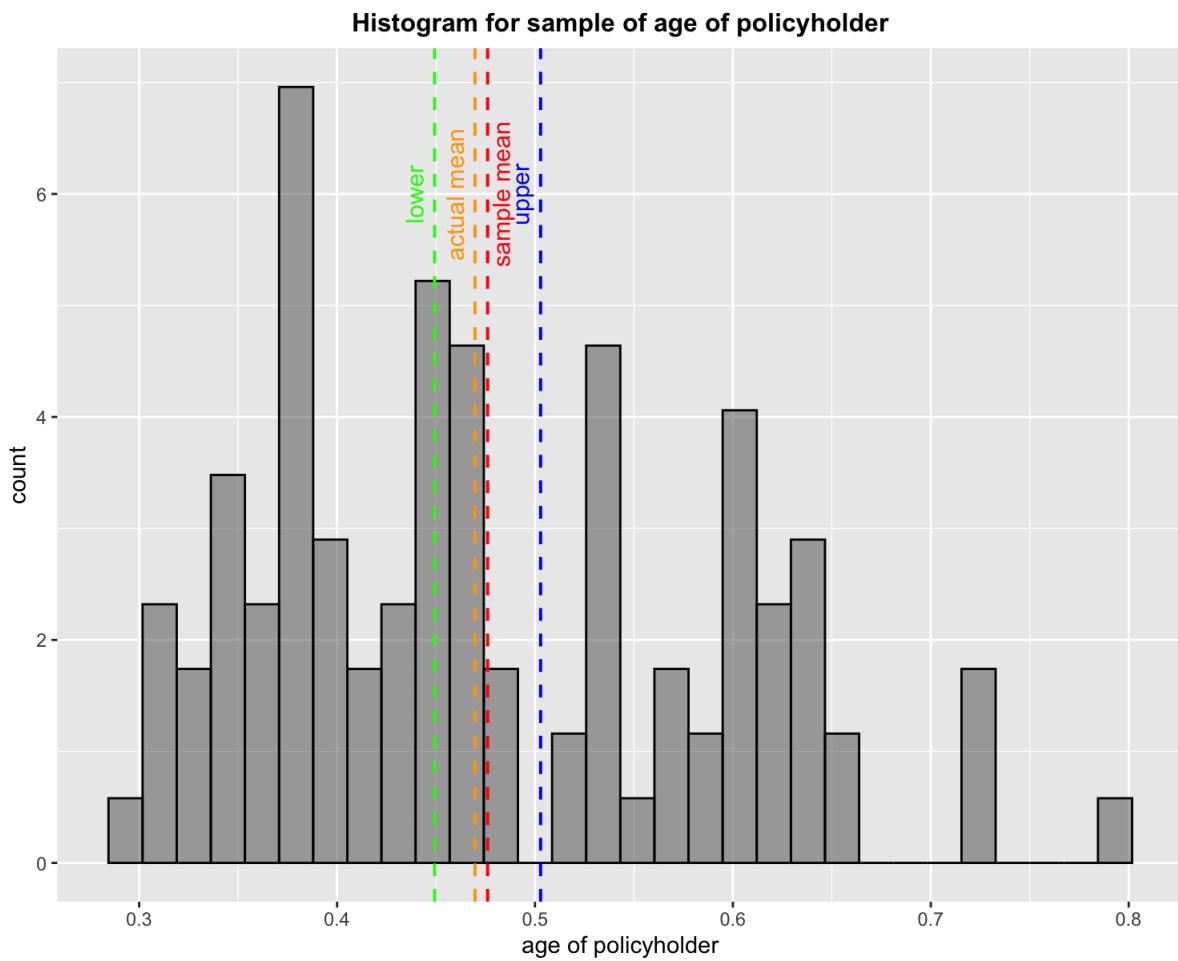


Figure 40: Histogram of age of policyholder marked with CI and sample mean

8.4 D

For the mean value of this numerical variable, design a hypothesis test and by finding the p-value, confirm or reject your assumption. What does this p-value signify?

The actual mean of the population is equal to 0.469. The mean of the sample is equal to 0.476. It is therefore necessary to design a test to measure it. Is it bigger than the actual one?

In the first part, we checked the central limit condition.

We have a random variable with a Normal Distribution. We should standardize it with Z score:

$$Z = \frac{\text{observation} - \mu}{SE} = \frac{0.469 - 0.476}{0.0115} = -0.608$$

Setting the hypothesis:

$$H_0 : \mu = 0.469 \quad \text{on average, the age of policyholder is 0.469.}$$

$$H_a : \mu > 0.469 \quad \text{on average, the age of policyholder is bigger than 0.469.}$$

and p-value like:

$$\begin{aligned} p\text{-value} &= P(\text{observed or more extreme outcome} \mid H_0 \text{ true}) \\ &= P(\bar{x} \geq 0.476 \mid H_0 : \mu = 0.469) \\ &= P(Z \geq -0.608) \\ &\simeq 0.27 \quad \alpha = 0.02 \end{aligned}$$

Null Hypothesis is not rejected in favor of H_a . we can believe that mean age of policyholder is not overestimate and the mean age of policyholder is less than 0.476.

```
1 null.value = 0.469
2 alpha = 0.02
3
4 x_bar <- mean(sub_car$age_of_policyholder)
5 s <- sd(sub_car$age_of_policyholder)
6 se <- s/sqrt(length(sub_car$age_of_policyholder))
7
8 z_score <- abs((x_bar - null.value)) / se
9
10 pvalue <- pnorm(z_score, lower.tail = FALSE)
11
12 if (pvalue < alpha) {
13   print("Reject null hypothesis.")
14 } else {
15   print("Fail to reject null hypothesis.")
16 }
```

Result: Fail to reject null hypothesis. There needs to be more evidence to say the mean of the samples is higher.

8.5 E

Based on the confidence interval you calculated in part “A”, does the data support the hypothesis that you have designed? Explain it.

Statistical results can be analyzed using P-values and confidence intervals, and their effects usually agree.

We cannot reject H_0 because the sample is in CI with a confidence level of 98%.

In addition, we cannot reject H_0 at a significance level of 0.02.

Since H_0 represents the population mean, it was expected to be the same.

8.6 F

Calculate type II error. What does this value mean?

We know that the actual mean number is 0.4696782 so:

$$\mu_a = 0.4696782 \rightarrow \bar{x} \sim N \left(\mu = 0.4696782, SE = \frac{s}{\sqrt{n}} = 0.79 \right)$$

$\alpha = 0.02$, one sided test $\rightarrow z_a = ?$

$$P(Z < z_a) = 0.02 \rightarrow z_a = qnorm(0.02) = -2.053$$

$$\begin{aligned} \text{Type II Error} &= \beta = P(\text{fail to reject } H_0 | \mu = \mu_a) \\ &= P\left(\frac{\bar{x} - 0.476}{0.79} < -2.053 | \bar{x} \sim N(\mu = 0.4696782, SE = 0.79)\right) \\ &= P(\bar{x} < -1.45 | \bar{x} \sim N(\mu = 0.4696782, SE = 0.79)) \\ &= P\left(Z < \frac{-1.45 - 0.4696782}{0.79}\right) \\ &= pnorm(-1.79) = 0.036 \end{aligned}$$

$Power = 1 - \beta = 0.96$

```

1 ci = 0.98 # alpha = 0.02
2 z = qnorm((1-ci)/2)
3
4 s = sd(sub_car$age_of_policyholder)
5 se <- s / sqrt(n) #SE = s/sqrt(n)
6
7 lower <- xbar + z * se
8 upper <- xbar - z * se
9
10 muactual = mean(car$age_of_policyholder)
11
12 Zleft <- (lower-muactual)/(se)
13 Zright <- (upper-muactual)/(se)
14 b <- pnorm(Zright)-pnorm(Zleft, lower.tail = FALSE)
15
16 power = 1-b
17 power

```

Result: Error Type II = 0.034 and Power is 96.3%.

Because the value we considered in the null hypothesis was very close to the actual mean value. As a result, the type 2 error is expected to be very small.

8.7 G

Calculate the power and Explain the relationship between the power and the effect size

It was calculated that 96% of the power was obtained from the previous section.

Power and effect size are closely related in statistical tests.

It is easy to reject the null hypothesis since there are significant differences between the two groups.

The null hypothesis becomes more likely to fail with increasing effect size.

Test power is increased as a result.

You can see relationship between effect size and power below, when we reach 0.3 effect size power increases to 100%.

```
1 # calc effectsize ~ power
2 d_seq <- seq(0, 2, by = 0.1)
3 pwr_list <- lapply(d_seq, function(d){
4   power.t.test(n = 100, delta = d, sd = sd(sub_car$age_of_policyholder),
5   type="one.sample")
6 })
7 pwr <- sapply(pwr_list, '[[, 'power')
8 dfpwr <- data.frame(power = pwr, effect.size = d_seq)
9
10 ggplot(dfpwr, aes(effect.size, power)) +
11   geom_point(size = 2, colour = "#8ACA88") +
12   geom_line(size = 0.5, colour = "black") +
13   scale_y_continuous(labels = scales::percent) +
14   xlab("effect size") +
15   ylab(expression("test power =" ~ 1 - beta))
```

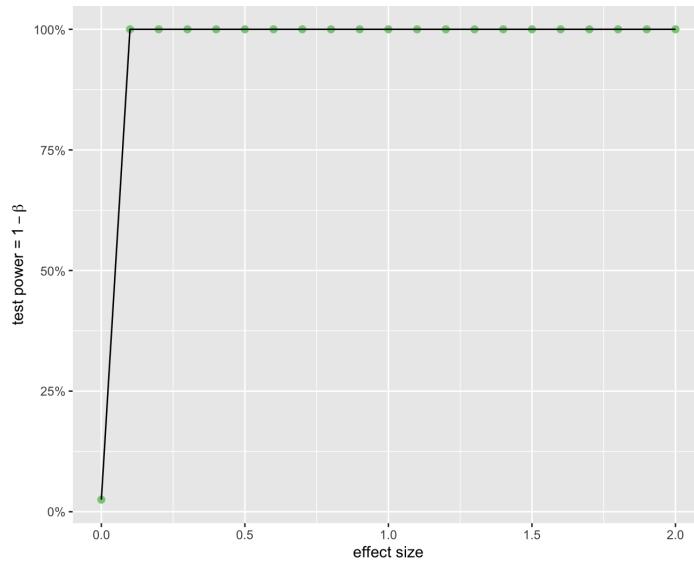


Figure 41: Effect size vs. Power

9. Question 8

Choose a numerical variable that has outliers, and we cannot apply CLT-based methods we have learned so far.

Chosen Numerical Variables : age_of_car

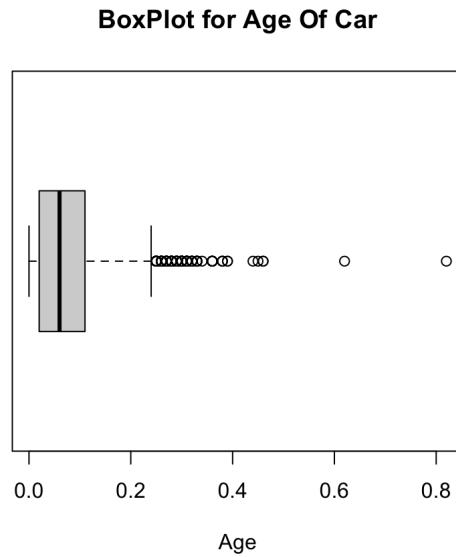


Figure 42: Boxplot of age of car

This variable has some outliers so we can't use CLT for this.

9.1 A

Calculate a 95% confidence interval for the mean of this variable using the percentile method and show the interval on the histogram

```
> quantile(car$age_of_car,c(.025,.975))
2.5% 97.5%
0.00 0.19
```

Figure 43: Percentile Method to find CI

Result: Confidence Interval for Original population!! : (0, 0.19)

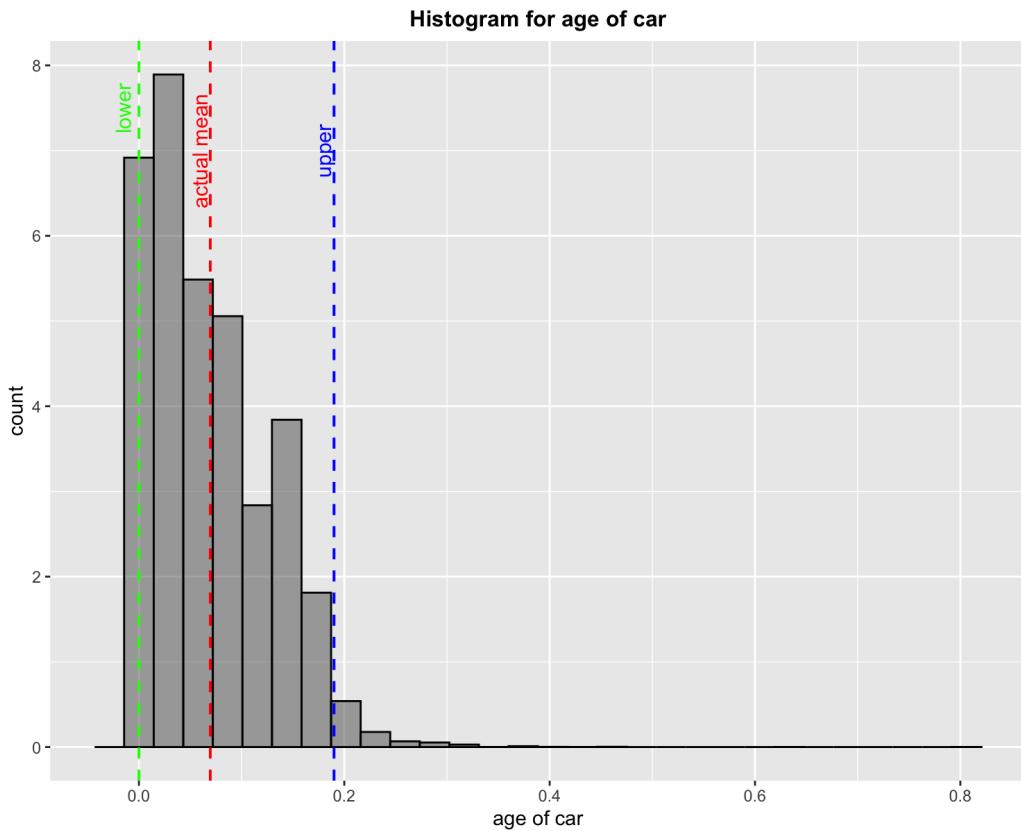


Figure 44: Histogram of age of car with CI 95% and actual mean

Using the percentile method, we can see that the confidence interval is between 0 and 0.19, so this interval covers the minimum and nearly maximum (because it has a right-skewed distribution), so it is ineffective and useless since it encompasses the entire range of possible values.

9.2 B

Pick a random sample of size 20. Then, using the bootstrapping method, calculate a 95% confidence interval for the mean of this variable using the standard error method. Also, plot the bootstrap distribution on the dot-plot.

```

1 car.sample = sample_n(car, 20)
2 boots.dist <- replicate(1000, mean(sample(car.sample$age_of_car, size
  = 20, replace=TRUE) ))
3
4 boots.mean <- mean(boots.dist)
5 boots.sd = sd(boots.dist)
6 boots.se = boots.sd / sqrt(1000)
7
8 upper = boots.mean + qt((1-0.95)/2, df=1000-1, lower.tail = FALSE) *
  boots.se
9 lower = boots.mean - qt((1-0.95)/2, df=1000-1, lower.tail = FALSE) *
  boots.se
10
11 cat('Confidence Interval for bootstrap distribution : (', lower,upper,
  ') ')

```

Result: Confidence Interval for bootstrap distribution : (0.06881241, 0.07061859).

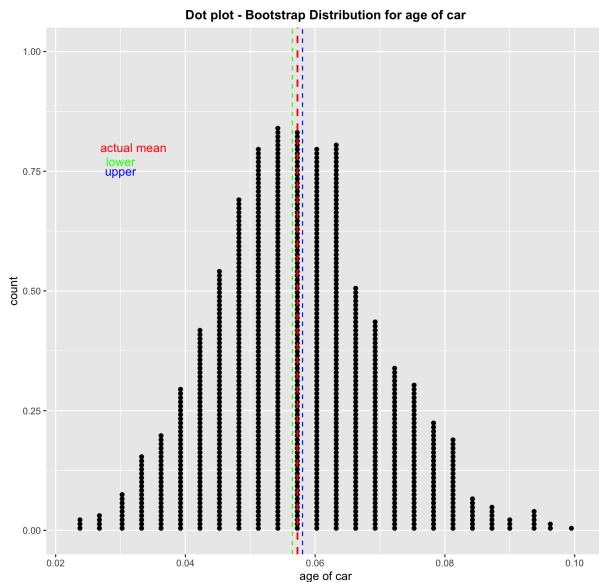


Figure 45: Bootstrap distribution for age of car - find CI with standard error method

Based on 20 samples of the age of a car, calculate the margin and bootstrap distribution using the standard error method.

This interval is much better and more informative than the previous part, as seen by the confidence interval between 0.06881241 and 0.07061859. As a result, we have a better and more accurate idea of what to expect.

9.3 C

Is there any noticeable difference between these two calculated confidence intervals?
Explain your reasoning

Yes, there is a noticeable difference between these two intervals. Can see qqplot from two distributions.

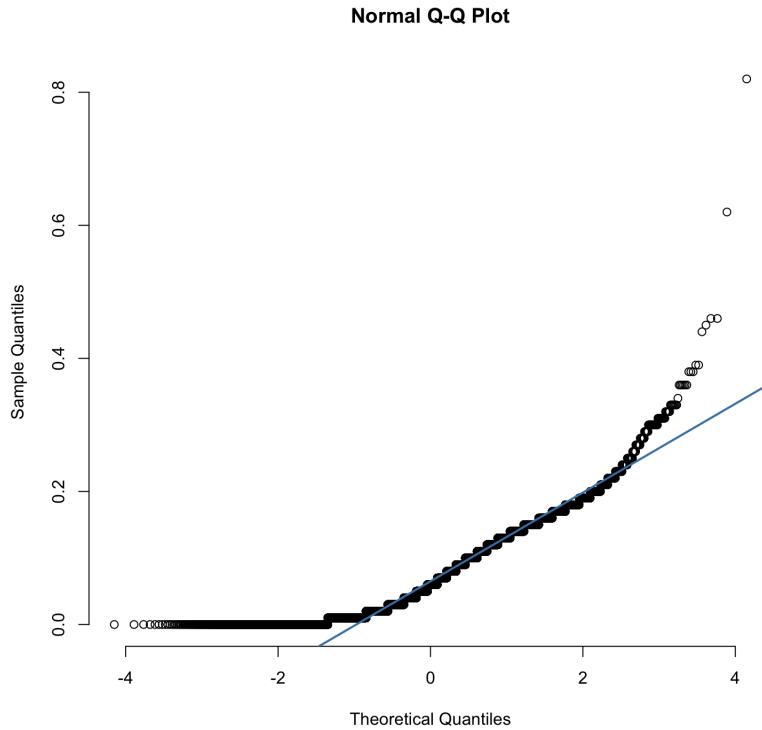


Figure 46: QQPlot - Original distribution - age of car

As you can see, this distribution has outliers. It is skewed, so when we use the percentile method, the result is not accurate and informative because we, only in normal distributions, can use this method.

Therefore use bootstrapping method to reach a normal distribution. Furthermore, use the standard error to approximate the margin of error, which is better than the percentage method.

Here is the Q-Q plot of bootstrap, which is a nearly normal distribution and has no skewness:

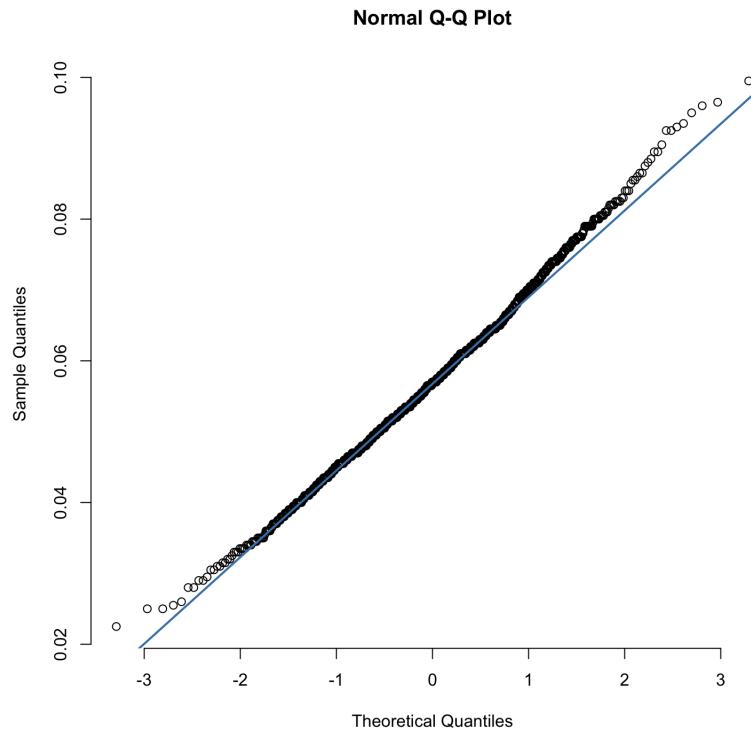


Figure 47: QQPlot - bootstrap distribution - age of car

10. Question 9

Choose a numerical and a categorical variable with more than two levels.
Divide observations of this dataset into different groups such that each group represents a level of the chosen categorical variable.

Chosen Numerical Variables : age_of_car Chosen Categorical Variables : fuel_type
A random sample without replacement was taken from the original sample, n = 200.
You can see the Box plot of three variables by fuel type below:

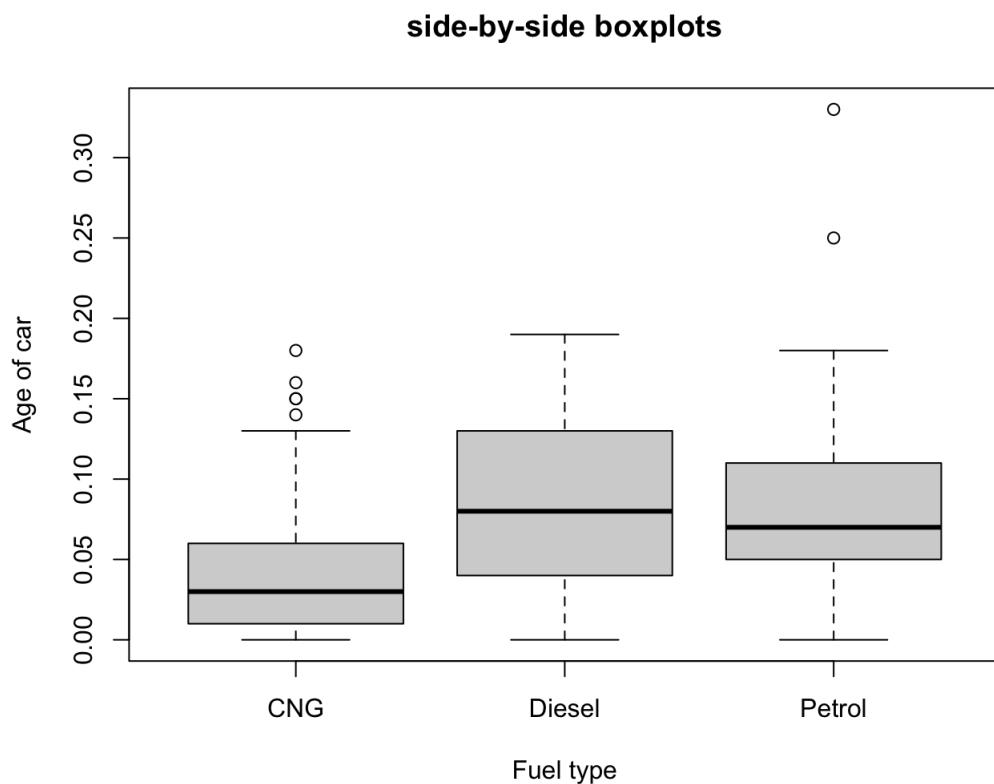


Figure 48: Three Types of fuel

Compared to petrol and diesel, CNG has a significant difference. In the next section, we will examine these cases in more detail.

10.1 A

Use the ANOVA test and compare the mean value of the numerical variable in the groups

In ANOVA, the null hypothesis is that there is no difference among group means. If any group differs significantly from the overall group mean, then the ANOVA will report a statistically significant result.

Significant differences among group means are calculated using the F statistic, which is the ratio of the mean sum of squares (the variance explained by the independent variable) to the mean square error (the variance left over).

If the F statistic is higher than the critical value (the value of F that corresponds with your alpha value, usually 0.05), then the difference among groups is deemed statistically significant.

Conditions for ANOVA:

1. Independence:

1.1 within groups: sampled observations is independent.

1.1.1 random sample / assignment.

1.1.2 each n_j less than 10% of population.

1.2 between groups: the groups is independent of each other (non-paired).

2. Approximate normality: distributions should be nearly normal within each group.

3. Equal variance: groups should have roughly equal variability.

Based on these conditions, we can assume that they are met.

Based on the probability distribution, establish the null hypothesis and alternative hypothesis.

- $$H_0 : \text{On average the mean of all group means are equal.}$$
- $$H_a : \text{At least one group mean is different from the rest.}$$

```
> car.sample = sample_n(car, 200)
> fisher <- aov(age_of_car ~ as.factor(fuel_type), data = car.sample)
> summary(fisher)
      Df Sum Sq Mean Sq F value    Pr(>F)
as.factor(fuel_type)   2 0.1586 0.07931   36.48 3.32e-14 ***
Residuals            197 0.4282 0.00217
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

Figure 49: ANOVA Test

Because p-value is smaller than significance level, we reject the null hypothesis.

Conclusion:

This means that there is a difference between at least two of the groups.

10.2 B

Choose two of the groups, perform a hypothesis test for the mean difference of the selected numerical variable in these groups and calculate the p-value.

Make a decision and explain the result using a significance level of 5%.

Based on the probability distribution, establish the null hypothesis and alternative hypothesis.

We have independence and approximate normality conditions.

$H_0 : \mu_A - \mu_B = 0$ On average the mean of these two groups are equal.

$H_a : \mu_A - \mu_B \neq 0$ On average the mean of these two groups are not equal.

First check Petrol and Diesel

```
1 gp <- unique(car.sample$fuel_type)
2
3 Petrol <- car.sample$age_of_car[which(car.sample$fuel_type==gp[1])]
4 Diesel <- car.sample$age_of_car[which(car.sample$fuel_type==gp[2])]
5
6 car.t.test <- t.test(Petrol , Diesel)
7 car.t.test
```

Welch Two Sample t-test

```
data: Petrol and Diesel
t = -0.2742, df = 101.22, p-value = 0.7845
alternative hypothesis: true difference in means is not equal to 0
95 percent confidence interval:
-0.02250769 0.01704102
sample estimates:
mean of x mean of y
0.08226667 0.08500000
```

Figure 50: T-test for Petrol and Diesel

$$\alpha^* = \frac{\alpha}{K} = \frac{0.05}{3} = 0.016, \alpha = 0.05, K = \binom{3}{2} = 3$$

$\alpha^* = 0.78 > 0.016 \rightarrow \text{Reject } H_0$ mean of type Petrol is different from the mean of type Diesel.

Based on t-test there is no differences between the means of them.

Second check Petrol and CNG

```
1 Petrol <- car.sample$age_of_car[which(car.sample$fuel_type==gp[1])]  
2 CNG <- car.sample$age_of_car[which(car.sample$fuel_type==gp[3])]  
3  
4 car.t.test <- t.test(Petrol , CNG)  
5 car.t.test
```

```
Welch Two Sample t-test  
  
data: Petrol and CNG  
t = 8.3535, df = 118.76, p-value = 1.425e-13  
alternative hypothesis: true difference in means is not equal to 0  
95 percent confidence interval:  
 0.04374431 0.07092601  
sample estimates:  
 mean of x mean of y  
0.08226667 0.02493151
```

Figure 51: T-test for Petrol and CNG

$\alpha^* = 0 < 0.016 \rightarrow \text{Reject } H_0$ mean of type Petrol is different from the mean of type CNG.

Based on t-test there is a differences between the means of them.

* that the t-test was used because when the df is increased, it goes toward normal, and there is no matter.

11. R codes

R Codes

```
1 # Import Libraries
2
3 require(qqplotr)
4 library(gridExtra)
5 library(moments)
6 library(magrittr)
7 library(ggfortify)
8 library(ggplot2)
9 library(Hmisc)
10 library(plyr)
11 library(hexbin)
12 library("plot3D")
13 library(plotly)
14 library(scatterplot3d)
15 library(RNHANES)
16 library(GGally)
17 library(dplyr)
18 require(ggpubr)
19 library(ggmosaic)
20 require(corrplot)
21 library(patchwork)
22 library(ggExtra)
23 library(ggcorrplot)
24 library(reshape)
25 theme_set(theme_minimal())
26
27
28 # read dataset
29 path = './Car_Insurance_Claim_Prediction.csv'
30 car<-read.csv(path)
31
32
33 #-----> Question 0 <-----
34 #Part A -----
35 #Part B -----
36 summary(car)
37
38 #Part C -----
39 any(is.na(car))
40 missingValue <- data.frame(colSums(is.na.data.frame(car)))
41
42 #Part D -----
43
44 #-----> Question 1 <-----
45 #Part A -----
46 # Freedman-Diaconis rule
47 bins_fd <- function(vec) {
48   diff(range(vec)) / (2 * IQR(vec) / length(vec)^(1/3))}
49 bin_fd = bins_fd(car$age_of_car)
50
51 # plot
```

```

52 ggplot(car, aes(x = age_of_car)) +
53   geom_histogram(aes(y = ..density..), colour = "black", fill = "#F7D302",
54                 bins = bin_fd) +
55   labs(
56     title = "Histogram and Density of Age Of Cars",
57     x = "age of car",
58     y = "count"
59   ) +
60   geom_density(color = "#631919", size = 1) +
61   theme_classic() +
62   theme(
63     plot.title = element_text(size = 12, face = "bold", hjust = 0.5),
64     plot.caption = element_text(face = "italic")
65   )
66
67 # sqrt(n) rule
68 bins_sqrt <- function(vec) {
69   sqrt(length(vec))}
70 bin_sq <- bins_sqrt(car$age_of_car)
71
72 # plot
73 ggplot(car, aes(x = age_of_car)) +
74   geom_histogram(aes(y = ..density..), colour = "black", fill = "#F7D302",
75                 bins = bin_sq) +
76   labs(
77     title = "Histogram and Density of Age Of Cars",
78     x = "age of car",
79     y = "count"
80   ) +
81   geom_density(color = "#631919", size = 1) +
82   theme_classic() +
83   theme(
84     plot.title = element_text(size = 12, face = "bold", hjust = 0.5),
85     plot.caption = element_text(face = "italic")
86   )
87
88 # default
89 ggplot(car, aes(x = age_of_car)) +
90   geom_histogram(aes(y = ..density..), colour = "black", fill = "#F7D302") +
91   labs(
92     title = "Histogram and Density of Age Of Cars",
93     x = "age of car",
94     y = "count"
95   ) +
96   geom_density(color = "#631919", size = 1) +
97   theme_classic() +
98   theme(
99     plot.title = element_text(size = 12, face = "bold", hjust = 0.5),
100    plot.caption = element_text(face = "italic")
101  )
102
103 # qqplot
104 qqnorm(car$age_of_car, pch = 1, frame = FALSE)
105 qqline(car$age_of_car, col = "steelblue", lwd = 2)

```

```

106 #Part A -----
107
108
109 #Part C -----
110 quantile(car$age_of_car)
111
112 boxplot(car$age_of_car,
113   main = "BoxPlot for Age Of Car",
114   names = c("age_of_car"),
115   xlab = "Age",
116   horizontal = TRUE)
117
118
119 #Part D -----
120 aoc <- car$age_of_car
121 length(aoc[which(aoc < boxplot$stats[1] | aoc > boxplot$stats[5])])
122
123 #Part E -----
124 mean(car$age_of_car)
125 median(car$age_of_car)
126 var(car$age_of_car)
127 sd(car$age_of_car)
128
129 #Part F -----
130 aoc <- car$age_of_car
131 mean = mean(aoc)
132
133
134 Fo = aoc[aoc > .75*mean]
135 Th = aoc[aoc >= .5*mean & aoc <= .75*mean]
136 Se = aoc[aoc < .5*mean & aoc >= .25*mean]
137 Fi = aoc[aoc < .25*mean]
138
139 frequencies<-c(length(Fo),length(Th),length(Se),length(Fi))
140 percentage <- round(100*frequencies/sum(frequencies), 2)
141
142 aoc.categorized <- data.frame(names = c("Fourth", "Third", "Second", "First"), value = percentage)
143
144 ggplot(aoc.categorized, aes(x="", y = value, fill = names)) +
145   geom_bar(stat = "identity") + coord_polar("y") +
146   geom_text(aes(label = paste0(value, "%")),
147             position = position_stack(vjust = 0.5)) +
148   labs(title="Age Of Car Pie Chart", x = 'Frequency', y = 'Age Of Car')
149
150 #Part G -----
151 ggplot(car, aes(x = age_of_car)) + geom_density( size = 1) +
152   geom_vline(aes(xintercept = median(age_of_car)), linetype = "dashed",
153   , size = .7,col= 'red') +
154   geom_vline(aes(xintercept = mean(age_of_car)), linetype = "dashed",
155   , size = .7,col= 'green') +
156   theme_classic() +
157   theme(
158     plot.title = element_text(size = 12, face = "bold", hjust = 0.5),
159     plot.caption = element_text(face = "italic")
160   ) +

```

```

159  annotate("text", x = mean(car$age_of_car) + .02 , label = "mean" , y
160   = 6, size = 3, angle = 90 , color = 'green') +
161  annotate("text", x = median(car$age_of_car) - 0.02 , label = "median"
162   , y = 0.6, size = 3, angle = 90, color = 'red') +
163  labs(
164    title = "Age of Car Density",
165    caption = "green line is about mean and red is about median
166    variable on the diagram.",
167    x = "Age Of The Car",
168    y = "Count"
169  )
170
171 #-----> Question 2 <-----#
172 #Part A -----
173 transmission <- car$transmission_type
174 gp <- unique(transmission)
175
176 Manual = transmission[transmission == 'Manual']
177 Automatic = transmission[transmission == 'Automatic']
178
179 frequencies<-c(length(Manual),length(Automatic))
180 percentage <- round(100*frequencies/sum(frequencies), 2)
181
182 transmission.categorized <- data.frame(types = c("Manual", "Automatic"
183   ),percentage = percentage, value = frequencies)
184
185 #Part B -----
186 transmission.categorized <- transmission.categorized[order(percentage)
187   ,]
188 ggplot(data=transmission.categorized, aes(x=types, y=frequencies, fill
189   =types)) +
190   geom_bar(stat="identity", alpha = 0.7, width = 0.7) +
191   labs(title="Barplot of Transmission Type") +
192   coord_flip()
193
194 #Part C -----
195 ggplot(data=transmission.categorized, aes(x=types, y=percentage, fill=
196   types)) +
197   geom_bar(stat="identity", alpha = 0.7, width = 0.6)++
198   labs(title="Barplot of Transmission Type")+
199   geom_text(aes(label=paste(percentage, "%")),vjust=-0.4,size=4)
200
201 #Part D -----
202 ggplot(data = car, aes(x=transmission_type, y=age_of_car, fill=
203   transmission_type))++
204   geom_violin(trim=FALSE, alpha = 0.7)++
205   labs(title="Violin Plot of age of car by transmission type")
206
207 #-----> Question 3 <-----#
208 #Part A -----
209 ggplot(car, aes(y=age_of_policyholder, x=age_of_car)) +
210   geom_point(size=2, alpha = .07) +
211   geom_smooth(method=lm, col = '#8ACA88') +
212   labs(title=paste("Scatter Plot of Age of car and Age of policyholder
213   "))

```

```

207 #Part B -----
208 ggplot(car, aes(y=age_of_policyholder, x=age_of_car, shape = fuel_type
209   , color = fuel_type)) +
210   geom_point(size=1.5, alpha = 1) +
211   labs(title=paste("Scatter Plot of Age of car and Age of policyholder
212     with fuel type"))
213
214 #Part C -----
215 cor.test(car$age_of_car, car$age_of_policyholder)
216
217 #Part D -----
218 # bin 7
219 p7 <- ggplot(car, aes(age_of_car, age_of_policyholder)) +
220   geom_point() +
221   geom_hex(bins = 7) +
222   geom_smooth(color='#FFFF00', se=FALSE) +
223   ggtitle("Hexbin Plot - Age of car and policyholder with binsize 7")
224
225 ggMarginal(p7, type="histogram", size=3, fill='lightblue')
226 #-----
227 # bin 12
228 p12 <- ggplot(car, aes(age_of_car, age_of_policyholder)) +
229   geom_point() +
230   geom_hex(bins = 12) +
231   geom_smooth(color='#FFFF00', se=FALSE) +
232   ggtitle("Hexbin Plot - Age of car and policyholder with binsize 12")
233
234 ggMarginal(p12, type="histogram", size=3, fill='lightblue')
235 #-----
236 # bin 30
237 p30 <- ggplot(car, aes(age_of_car, age_of_policyholder)) +
238   geom_point() +
239   geom_hex(bins = 30) +
240   geom_smooth(color='#FFFF00', se=FALSE) +
241   ggtitle("Hexbin Plot - Age of car and policyholder with binsize 30")
242
243 ggMarginal(p30, type="histogram", size=3, fill='lightblue')
244 #-----
245 # bin 100
246 p100 <- ggplot(car, aes(age_of_car, age_of_policyholder)) +
247   geom_point() +
248   geom_hex(bins = 100) +
249   geom_smooth(color='#FFFF00', se=FALSE) +
250   ggtitle("Hexbin Plot - Age of car and policyholder with binsize 100")
251
252 ggMarginal(p100, type="histogram", size=3, fill='lightblue')
253
254 #Part E -----
255 ggplot(car, aes(x=age_of_car, y=age_of_policyholder) ) +
256   stat_density_2d(aes(fill = ..level..), geom = "polygon") +
257   scale_fill_continuous(type = "viridis") +
258   labs(title = "2D density plot Age of car and Age of policyholder")
259
260

```

```

261 #-----> Question 4 <-----  

262 #Part A -----  

263  

264 # Some preprocessing  

265 car$is_claim <- as.logical(car$is_claim)  

266 car <- subset(car, select = -c(X))  

267  

268 # Separate Numerical features  

269 car.numeric <- Filter(is.numeric, car)  

270  

271 # correlation matrix  

272 cor_mat <- cor(as.matrix(car.numeric), method="pearson")  

273 cor_mat <- round(cor_mat, 2)  

274  

275 # calc p_values  

276 p.mat <- cor_pmat(car.numeric)  

277  

278 ggcorrplot(  

279   cor_mat, hc.order = TRUE, type = "upper",  

280   lab = TRUE , p.mat = p.mat, sig.level = 0.05,  

281   title = "heatmap correlogram",  

282   colors = c("blue", "white", "red"),  

283   hc.method = "complete", lab_size = 2,  

284   as.is = FALSE)  

285  

286 #Part B -----  

287 selected_var = data.frame("policy_tenure" = car$policy_tenure,  

288                           "car's age" = car$age_of_car,  

289                           "policyholder's age_"=car$age_of_  

290                           policyholder,  

291                           "population density"= car$population_density  

292                           ,  

293                           "displacement"=car$displacement,  

294                           "turning radius"=car$turning_radius,  

295                           "length"=car$length,  

296                           "width"=car$width,  

297                           "height"=car$height,  

298                           "gross_weight" = car$gross_weight  

299 )  

300  

301  

302 #ggpairs(dplyr::select_if(car, is.numeric), title = "Correlogram",  

303           options(expressions=10000))  

304  

305 ggpairs(selected_var, title = "Correlogram",  

306           lower = list(continuous = wrap("smooth", alpha = 0.1, colour="gray")),  

307           upper = list(continuous = wrap("density",colour="#8ACA88",  

308                         alpha = 0.5), combo = "dot_no_facet"))  

309  

310  

311 #Part C -----  

312 colors <- c("#999999", "#E69F00", "#56B4E9")  

313  

314 fuel_type <- colors[as.numeric(factor(car$fuel_type))]  

315  

316 scatterplot3d(

```

```

313 main="3D Scatter Plot",
314 car$population_density,
315 car$age_of_policyholder,
316 car$displacement,
317 color=fuel_type,
318 xlab = "population_density",
319 ylab = "age_of_policyholder",
320 zlab = "displacement",
321 pch = 16)
322
323 legend("right", legend = unique(car$fuel_type), pch = 16, col =
324 colors)
325 #-----> Question 5 <-----#
326 #Part A -----
327 Total <- sum
328 t<-as.data.frame.matrix(addmargins(table(car[,c('fuel_type',
329 transmission_type)])), FUN = Total))
330
331 #Part B -----
332 ggplot(car, aes(x = fuel_type , fill = transmission_type)) +
333 geom_bar(position = "dodge", alpha = 0.6) +
334 geom_text(aes(label = ..count..), stat = "count", vjust = -.3, size =
335 = 3, color = 'black', position = position_dodge(width = 1)) +
336 labs(title="Grouped Bar Chart", x="fuel_type")
337
338 #Part C -----
339 ggplot(car, aes(x = fuel_type, fill = transmission_type)) +
340 geom_bar(position = "stack", alpha = 0.6) +
341 geom_text(aes(label=..count..), stat='count', vjust = -.3, size =
342 = 3, color = 'black', position = position_stack(0.5)) +
343 labs(title=paste("Segmented barplot"))
344
345 #Part D -----
346 ggplot(data = car) +
347 geom_mosaic(aes(x = product(transmission_type), fill = fuel_type ),
348 alpha = 0.6) +
349 annotate("text", x = 0.195, y = 0.35,label=paste(round(t[2,1]/t
350 [4,1],2)*100,"%"),fontface = "bold" ,size=4,color="black") +
351 annotate("text", x = 0.195, y = 0.855,label=paste(round(t[3,1]/t
352 [4,1],2)*100,"%"), fontface = "bold",size = 4,color="black") +
353 annotate("text", x = 0.7, y = 0.25,label=paste(round(t[1,2]/t
354 [4,2],2)*100,"%"),fontface = "bold" ,size = 4,color="black") +
355 annotate("text", x = 0.7, y = 0.57,label=paste(round(t[2,2]/t
356 [4,2],2)*100,"%"), fontface = "bold",size = 4,color="black") +
357 annotate("text", x = 0.7, y = 0.8,label=paste(round(t[3,2]/t[4,2],2)
358 *100,"%"), fontface = "bold",size = 4,color="black") +
359 labs(title=paste("Mosaic plot"))
360
361 #-----> Question 6 <-----#
362 #Part A -----
363 car_sample <- sample_n(car, 25)
364 car_sample$diff <- car_sample$policy_tenure - car_sample$displacement
365
366 #Part B -----
367 t.test(car_sample$policy_tenure, car_sample$displacement, paired =
368 TRUE, var.equal = TRUE)

```

```

359 # calculate CI
360 xbar1 <- mean(car_sample$policy_tenure)
361 s1 <- sd(car_sample$policy_tenure)
362 xbar2 <- mean(car_sample$displacement)
363 s2 <- sd(car_sample$displacement)
364
365
366 sp = ((n-1)*s1^2+(n-1)*s2^2)/(n+n-2)
367 sp
368
369 margin <- qt(0.975, df=n+n-1)*sqrt(sp/n + sp/n)
370 margin
371
372 lowerinterval <- (xbar1-xbar2) - margin
373 lowerinterval
374
375 upperinterval <- (xbar1-xbar2) + margin
376 upperinterval
377
378 cat("95%\% confidence interval is= (", lowerinterval, ", ", upperinterval,
     ")")
379
380 #-----> Question 7 <-----
```

381

```

382 ggplot(car, aes(x = age_of_policyholder)) +
  geom_histogram(aes(y = ..density..), alpha=0.5, colour = "black",
    fill = "gray")+
  labs(
    title = "Histogram for age of policyholder",
    x = "age of policyholder",
    y = "count"
  ) +
  geom_density(color = "blue", size = 1) +
  theme(
    plot.title = element_text(size = 12, face = "bold", hjust = 0.5)
  )
393
394 #Part A -----
395 n = 100
396
397 # Reads the dataset 'car' and take the 100 rows as sample
398 sdf<- sample(1:nrow(car), n)
399
400 #sample 10 rows
401 sub_car <- car[sdf,]
402
403 xbar = mean(sub_car$age_of_policyholder)
404
405 ci = 0.98
406 z = qnorm((1-ci)/2)
407
408 se <- sd(sub_car$age_of_policyholder) / sqrt(n) #SE = s/sqrt(n)
409
410 lower <- xbar + z * se
411 upper <- xbar - z * se
412
413 cat("98%\% confidence interval is= (", lower, ", ", upper, ")")
```

```

414
415 #Part B -----
416 ggplot(sub_car, aes(x = age_of_policyholder)) +
417   geom_histogram(aes(y = ..density..), alpha=0.5, colour = "black")+
418   labs(
419     title = "Histogram for sample of age of policyholder",
420     x = "age of policyholder",
421     y = "count"
422   ) +
423   geom_vline(aes(xintercept = mean(age_of_policyholder)), linetype = "dashed", size = .7,col= 'red') +
424   geom_vline(aes(xintercept = lower), linetype = "dashed", size = .7,
425   col= 'green') +
426   geom_vline(aes(xintercept = upper), linetype = "dashed", size = .7,
427   col= 'blue') +
428   geom_vline(aes(xintercept = mean(car$age_of_policyholder)), linetype =
429   = "dashed", size = .7,col= 'orange') +
430   theme(
431     plot.title = element_text(size = 12, face = "bold", hjust = 0.5),
432     plot.caption = element_text(face = "italic")
433   ) +
434   annotate("text", x = mean(sub_car$age_of_policyholder) + 0.007 ,
435   label = "sample mean", y = 6, size = 4, angle = 90 , color = 'red'
436   ) +
437   annotate("text", x = lower - 0.01 , label = "lower", y = 6, size =
438   4, angle = 90, color = 'green') +
439   annotate("text", x = upper - 0.01 , label = "upper", y = 6, size =
440   4, angle = 90, color = 'blue') +
441   annotate("text", x = mean(car$age_of_policyholder) - 0.01 , label =
442   "actual mean", y = 6, size = 4, angle = 90, color = 'orange')
443
444 #Part D -----
445 null.value = 0.469
446 alpha = 0.02
447
448 x_bar <-mean(sub_car$age_of_policyholder)
449 s <- sd(sub_car$age_of_policyholder)
450 se <- s/sqrt(length(sub_car$age_of_policyholder))
451
452 z_score <- abs((x_bar - null.value)) / se
453
454 pvalue <- pnorm(z_score, lower.tail = FALSE)
455
456 if (pvalue < alpha) {
457   print("Reject null hypothesis.")
458 } else {
459   print("Fail to reject null hypothesis.")
460 }
461
462 #Part F -----
463 ci = 0.98 # alpha = 0.02
464 z = qnorm((1-ci)/2)
465
466 s = sd(sub_car$age_of_policyholder)
467 se <- s / sqrt(n) #SE = s/sqrt(n)
468
469 lower <- xbar + z * se

```

```

462 upper <- xbar - z * se
463
464 muactual = mean(car$age_of_policyholder)
465
466 Zleft <- (lower-muactual)/(se)
467 Zright <- (upper-muactual)/(se)
468 b <- pnorm(Zright)-pnorm(Zleft, lower.tail = FALSE)
469
470 power = 1-b
471 power
472 #-----
473
474 # calc effectsize ~ power
475 d_seq <- seq(0, 2, by = 0.1)
476 pwr_list <- lapply(d_seq, function(d){
477   power.t.test(n = 100, delta = d, sd = sd(sub_car$age_of_policyholder),
478   type="one.sample")
478 })
479 pwr <- sapply(pwr_list, '[[' , 'power')
480
481 dfpwr <- data.frame(power = pwr, effect.size = d_seq)
482
483 ggplot(dfpwr, aes(effect.size, power)) +
484   geom_point(size = 2, colour = "#8ACA88") +
485   geom_line(size = 0.5, colour = "black") +
486   scale_y_continuous(labels = scales::percent) +
487   xlab("effect size") +
488   ylab(expression("test power = " ~ 1 - beta))
489
490 #-----> Question 8 <-----+
491 boxplot(car$age_of_car,
492           main = "BoxPlot for Age Of Car",
493           names = c("age_of_car"),
494           xlab = "Age",
495           horizontal = TRUE)
496
497 #Part A -----
498 q <- quantile(car$age_of_car,.025,.975)
499 cat('Confidence Interval for Original population!! : (', q[1],q[2], ')')
500
501 ggplot(car, aes(x = age_of_car)) +
502   geom_histogram(aes(y = ..density..),alpha=0.5, colour = "black")+
503   labs(
504     title = "Histogram for age of car",
505     x = "age of car",
506     y = "count"
507   ) +
508   geom_vline(aes(xintercept = mean(age_of_car)), linetype = "dashed",
509   size = .7,col= 'red') +
510   geom_vline(aes(xintercept = q[1]), linetype = "dashed", size = .7,
511   col= 'green') +
512   geom_vline(aes(xintercept = q[2]), linetype = "dashed", size = .7,
513   col= 'blue') +
514   theme(
515     plot.title = element_text(size = 12, face = "bold", hjust = 0.5),
516     plot.caption = element_text(face = "italic")

```

```

514 ) +
515   annotate("text", x = mean(car$age_of_car) - 0.01 , label = "actual
      mean", y = 7, size = 4, angle = 90 , color = 'red') +
516   annotate("text", x = q[1] - 0.015 , label = "lower", y = 7.5, size
      = 4, angle = 90, color = 'green') +
517   annotate("text", x = q[2] - 0.01 , label = "upper", y = 7, size = 4,
      angle = 90, color = 'blue')
518
519 #Part B -----
520 car.sample = sample_n(car, 20)
521
522 boots.dist <- replicate(1000, mean(sample(car.sample$age_of_car, size
      = 20, replace=TRUE) ))
523
524 boots.mean <- mean(boots.dist)
525 boots.sd = sd(boots.dist)
526 boots.se = boots.sd / sqrt(1000)
527
528 upper = boots.mean + qt((1-0.95)/2, df=1000-1, lower.tail = FALSE) *
      boots.se
529 lower = boots.mean - qt((1-0.95)/2, df=1000-1, lower.tail = FALSE) *
      boots.se
530
531 cat('Confidence Interval for bootstrap distribution : (', lower,upper,
      ')')
532 df.boots.dist = data.frame(age_of_car = matrix(data = boots.dist, ncol
      = 1, byrow = TRUE))
533
534 ggplot(df.boots.dist, aes(x = age_of_car)) +
535   geom_dotplot( stackratio = 1.01,dotsize = .25)+ 
536   labs(
537     title = "Dot plot - Bootstrap Distribution for age of car",
538     x = "age of car",
539     y = "count"
540   ) +
541   geom_vline(aes(xintercept = mean(age_of_car)), linetype = "dashed",
      size = .8,col= 'red') +
542   geom_vline(aes(xintercept = lower), linetype = "dashed", size = .5,
      col= 'green') +
543   geom_vline(aes(xintercept = upper), linetype = "dashed", size = .5,
      col= 'blue') +
544   theme(
545     plot.title = element_text(size = 12, face = "bold", hjust = 0.5),
546     plot.caption = element_text(face = "italic")
547   ) +
548   annotate("text", x = .032 , label = "actual mean", y = .8, size = 4,
      angle = 0 , color = 'red') +
549   annotate("text", x = .03 , label = "lower", y = .77, size = 4, angle
      = 0, color = 'green') +
550   annotate("text", x = .03 , label = "upper", y = .75, size = 4, angle
      = 0, color = 'blue')
551
552 #Part C -----
553 qqnorm(car$age_of_car, pch = 1, frame = FALSE)
554 qqline(car$age_of_car, col = "steelblue", lwd = 2)
555
556

```

```

557 qqnorm(df.boots.dist$age_of_car, pch = 1, frame = FALSE)
558 qqline(df.boots.dist$age_of_car, col = "steelblue", lwd = 2)
559
560
561 #-----> Question 9 <-----
```

562 car.sample = sample_n(car, 200)

563

564 #draw box plot

565 boxplot(car.sample\$age_of_car ~ car.sample\$fuel_type, ylab = "Age of car", xlab = "Fuel type",
 main = "side-by-side boxplots")

566

567 #Part A -----

568 fisher <- aov(age_of_car ~ as.factor(fuel_type), data = car.sample)

569 summary(fisher)

570

571 #Part B -----

572 gp <- unique(car.sample\$fuel_type)

573 gp

574 # petrol vs, diesel

575 Petrol <- car.sample\$age_of_car[which(car.sample\$fuel_type==gp[1])]

576 Diesel <- car.sample\$age_of_car[which(car.sample\$fuel_type==gp[2])]

577

578 car.t.test <- t.test(Petrol , Diesel)

579 car.t.test

580

581

582 # petrol vs, CNG

583 Petrol <- car.sample\$age_of_car[which(car.sample\$fuel_type==gp[1])]

584 CNG <- car.sample\$age_of_car[which(car.sample\$fuel_type==gp[3])]

585

586 car.t.test <- t.test(Petrol , CNG)

587 car.t.test