# STA2005S - Regression Assignment

true          true

2024-10-17

## 0.1 Part One : Analysis

# 1 Section 1: Introduction

Air pollution, particularly high levels of particulate matter (PM), is a major environmental and public health issue in South Africa's urban centers. Exposure to elevated PM levels is linked to respiratory diseases and other serious health conditions. Understanding the factors influencing PM concentrations is crucial for developing policies that improve air quality and protect public health. This analysis seeks to identify the key drivers of air pollution in South Africa's cities, focusing on how various urban, environmental, and socioeconomic factors affect particulate matter levels.

Unknown Factors to Investigate:

Traffic Density: How do varying levels of vehicle traffic contribute to PM levels in different areas?

Industrial Activity: What is the impact of industrial activity near monitoring stations on air quality?

Temperature & Humidity: How do changes in weather conditions, like temperature and humidity, influence PM concentrations?

Wind Speed: How does wind speed affect the dispersion or accumulation of particulate matter in urban areas?

Day of the Week & Public Holidays: Do patterns of human activity on weekdays, weekends, and holidays significantly influence pollution levels?

Urban Greenery: How effective are green spaces in reducing air pollution in densely populated areas?

# 2 Objective

The goal of this analysis is to explore the relationships between PM levels and these explanatory variables. By identifying the most influential factors, we aim to inform urban planning and public health strategies that address air pollution and improve the quality of life in South African cities.
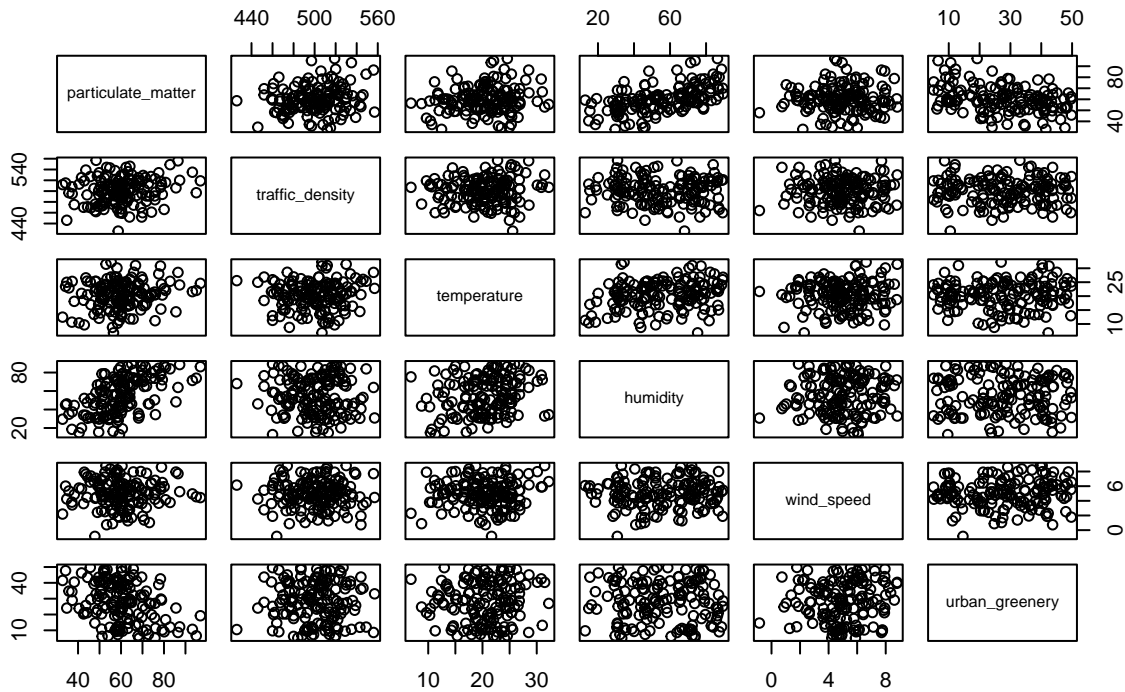
## 2.1 Section 2 : Data Exploration

density plot

pairwsie plots

```
continuous_vars <- data_tidy_air_quality[, sapply(data_tidy_air_quality, is.numeric)]
pairs(continuous_vars, main = "Pairwise Scatterplots of Continuous Variables")
```

## Pairwise Scatterplots of Continuous Variables



categorial variable plots

```r
data_tidy_air_quality$industrial_activity <- factor(data_tidy_air_quality$industrial_activity,
                              levels = c("None","Low", "Moderate", "High"))  # Adjust the levels a

data_tidy_air_quality$day_of_week <- factor(data_tidy_air_quality$day_of_week,
                        levels = c("Monday", "Tuesday", "Wednesday",
                                   "Thursday", "Friday", "Saturday", "Sunday"))

data_tidy_air_quality$holiday <- factor(data_tidy_air_quality$holiday,
                        levels = c("Yes", "No"))

categorical_vars <- names(data_tidy_air_quality)[sapply(data_tidy_air_quality, is.factor)]


for (var in categorical_vars) {
  plt<- ggplot(data_tidy_air_quality, aes_string(x = var, y = "particulate_matter")) +
    geom_boxplot() +
    labs(title = paste("Particulate Matter vs", var),
         x = var,
         y = "Particulate Matter") +
    theme_minimal()


  print(plt)  # Print the plot
}
```
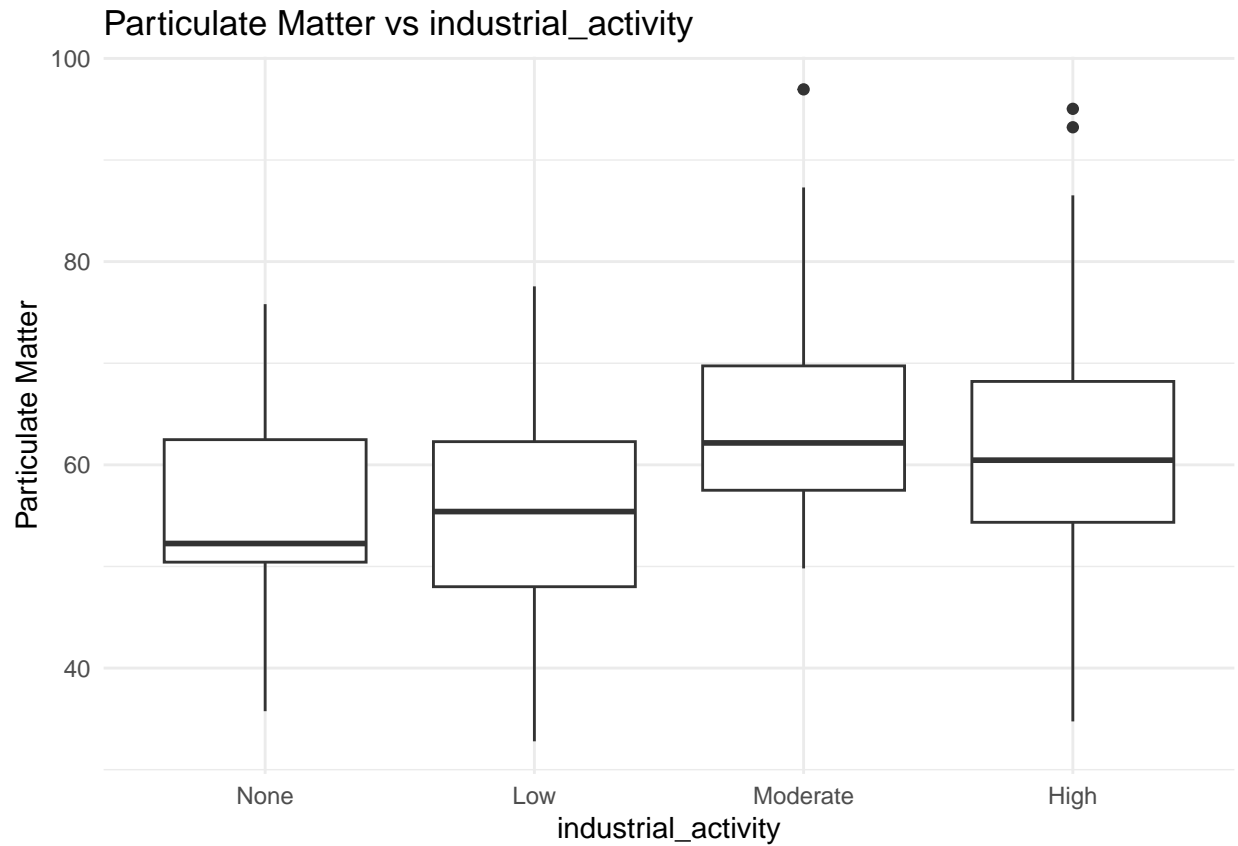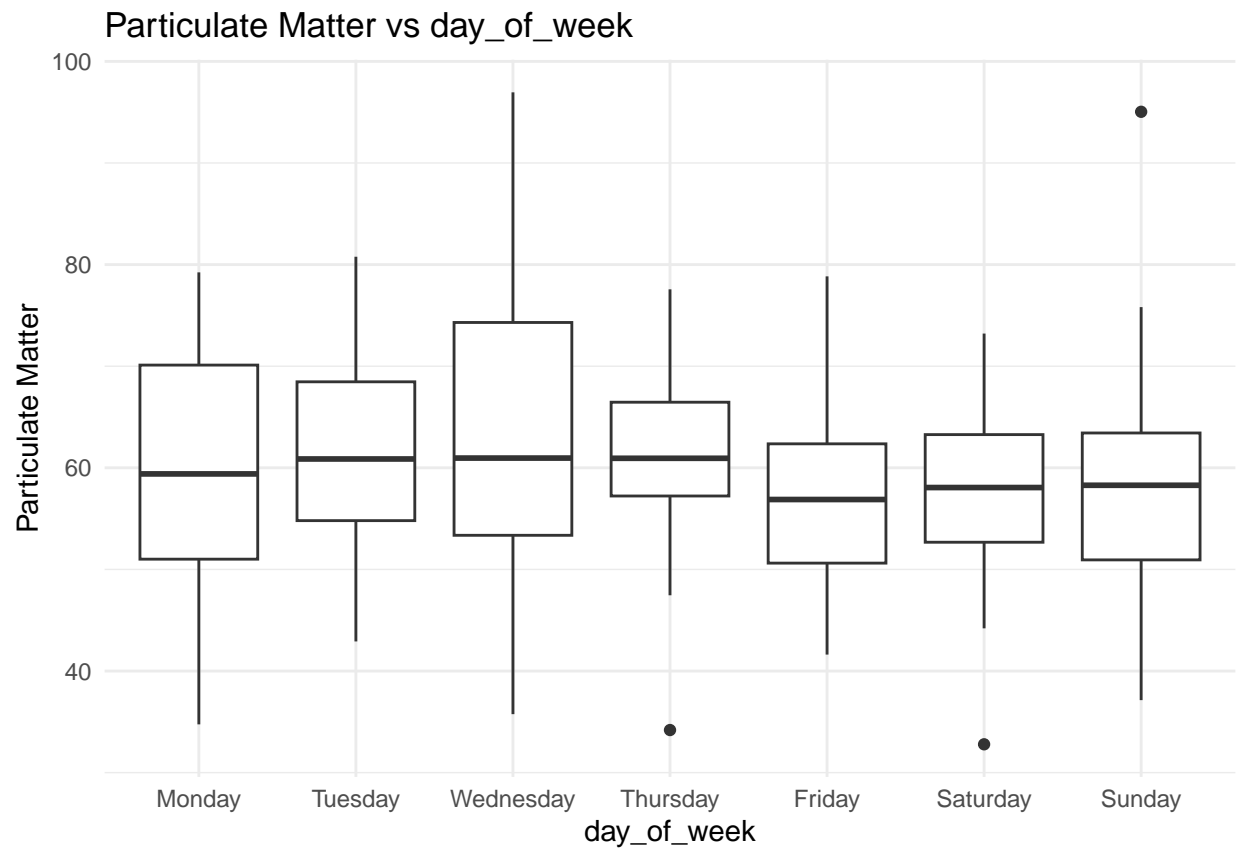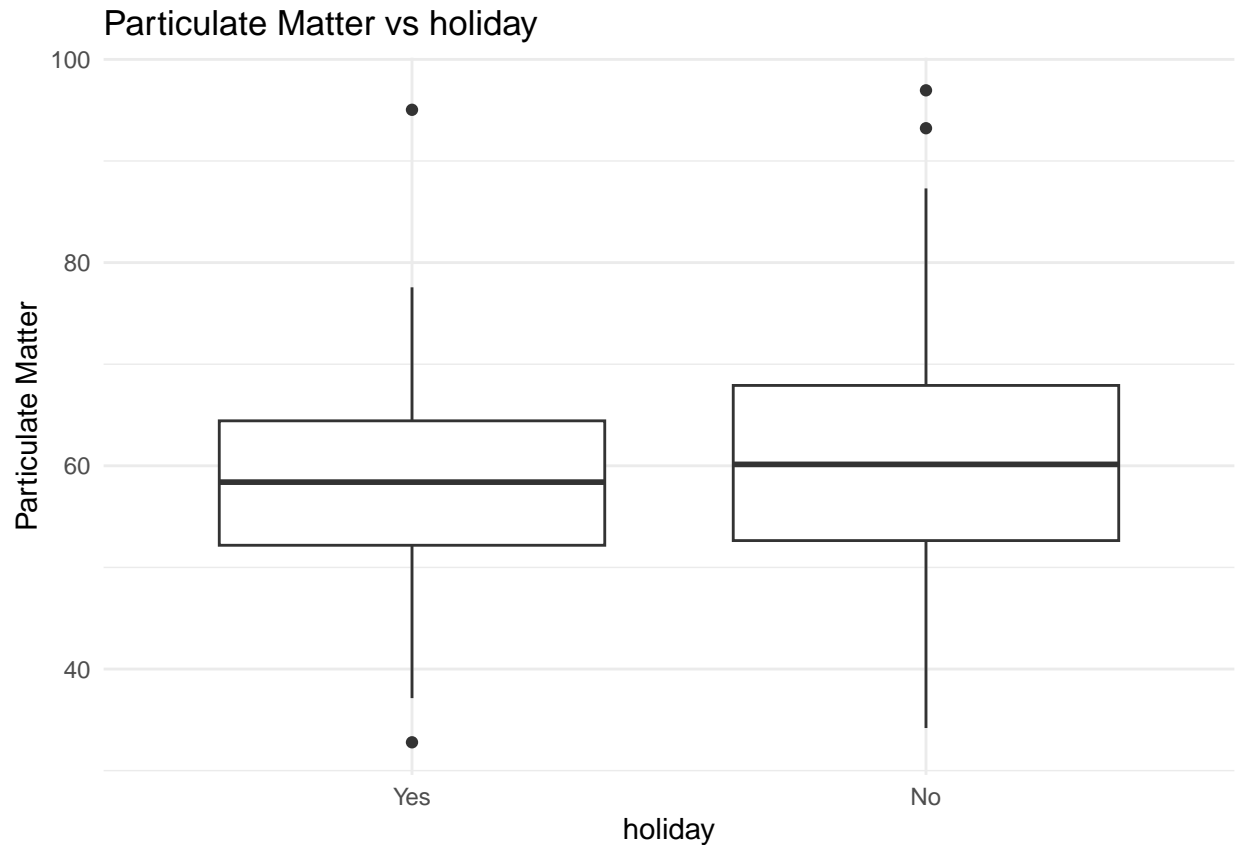
```
## Warning: 'aes_string()' was deprecated in ggplot2 3.0.0.
## i Please use tidy evaluation idioms with 'aes()'.
## i See also 'vignette("ggplot2-in-packages")' for more information.
## This warning is displayed once every 8 hours.
## Call 'lifecycle::last_lifecycle_warnings()' to see where this warning was
## generated.
```



Particulate Matter vs industrial_activity

Particulate Matter vs day_of_week

## Particulate Matter vs holiday



tabular representation of relationship between categorial variables

```r
for (i in 1:(length(categorical_vars)-1)) {
  for (j in (i+1):length(categorical_vars)) {
    cat("Contingency Table for", categorical_vars[i], "and", categorical_vars[j], "\n")
    print(table(data_tidy_air_quality[[categorical_vars[i]]], data_tidy_air_quality[[categorical_vars[j]
    cat("\n")
  }
}
```

```
## Contingency Table for industrial_activity and day_of_week
##
##           Monday Tuesday Wednesday Thursday Friday Saturday Sunday
##   None         2       0         3        3      2        0      4
##   Low          5       6         4        7      6        9      4
##   Moderate     4       4        10        8      6        4      3
##   High        11       7         9        5      8       10      6
##
## Contingency Table for industrial_activity and holiday
##
##           Yes No
##   None       5  9
##   Low       17 24
##   Moderate   9 30
##   High      21 35
##
## Contingency Table for day_of_week and holiday
```
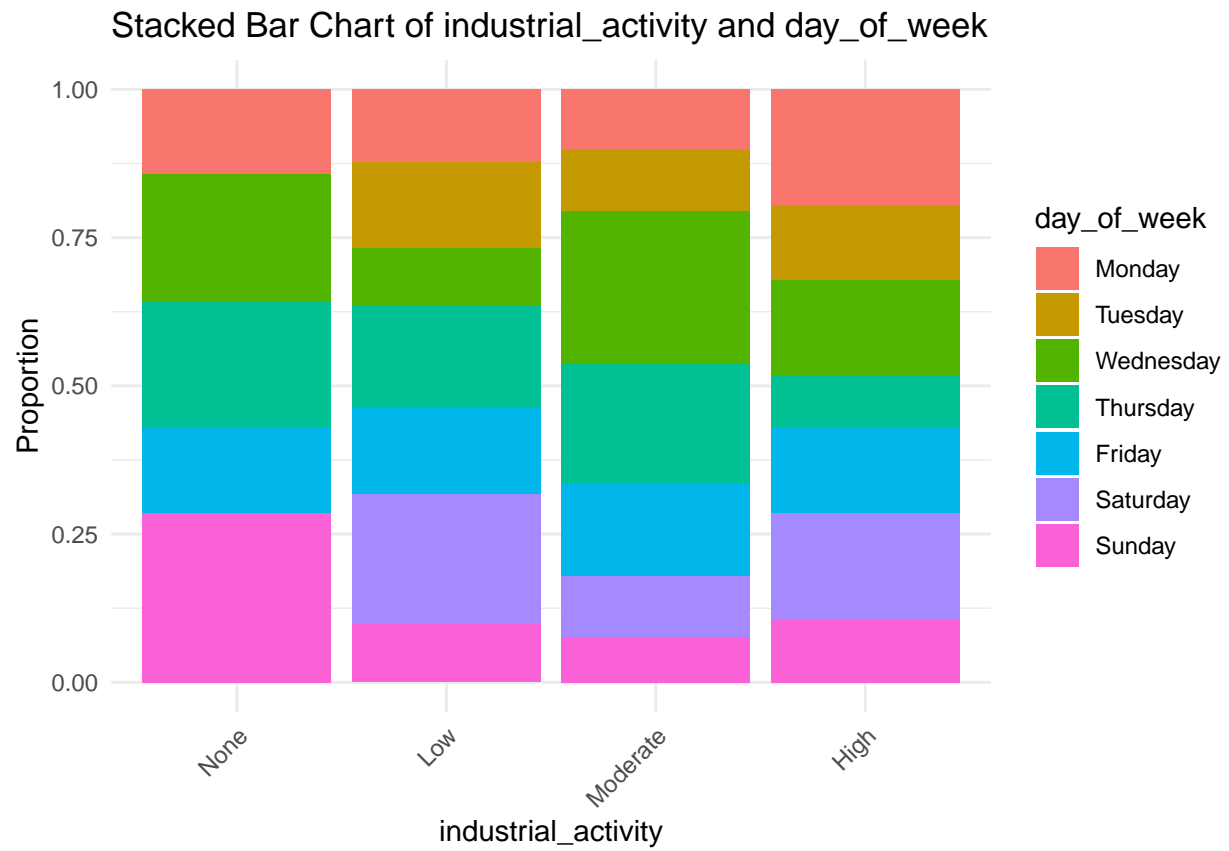
```
##
##              Yes No
##   Monday       1 21
##   Tuesday      1 16
##   Wednesday    3 23
##   Thursday     4 19
##   Friday       3 19
##   Saturday    23  0
##   Sunday      17  0
```
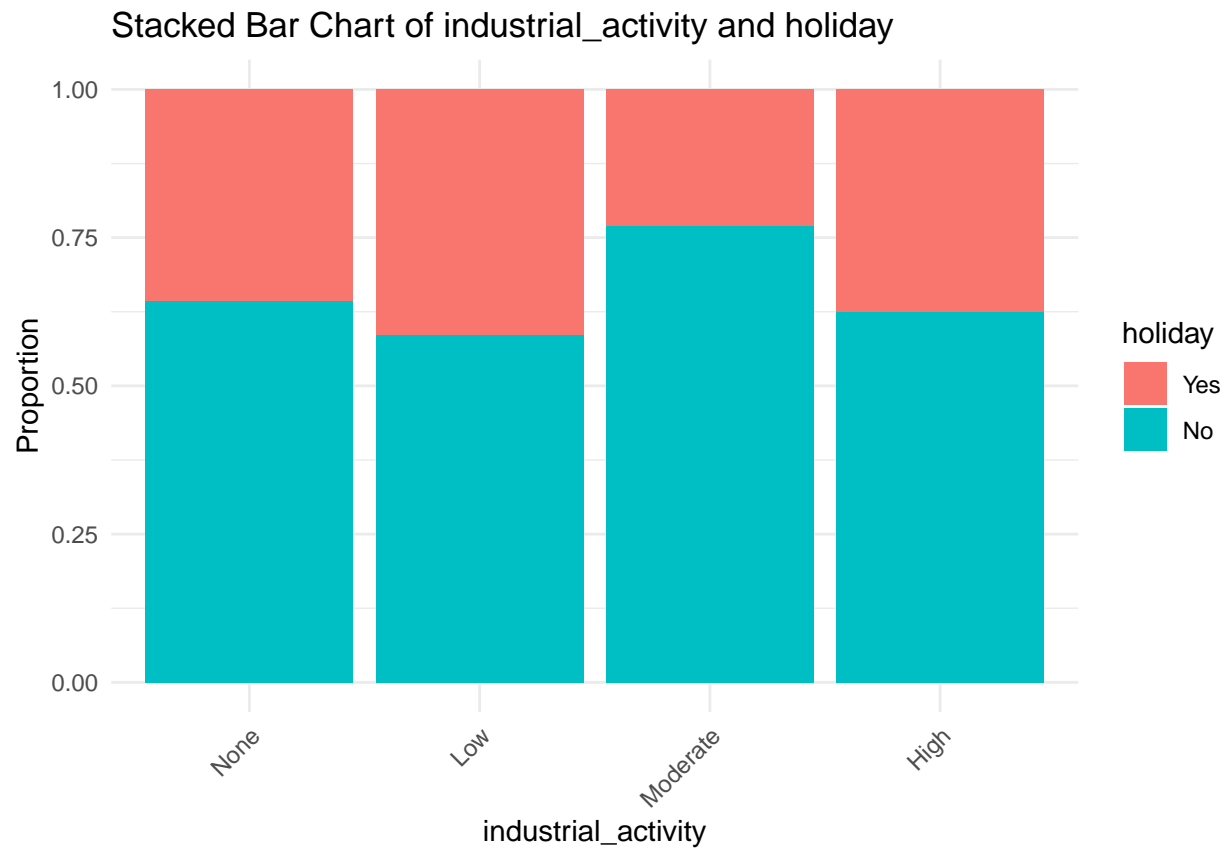
visual representation of relationship between categorial variables

```r
for (i in 1:(length(categorical_vars) - 1)) {
  for (j in (i + 1):length(categorical_vars)) {
    # Create the plot
    p <- ggplot(data_tidy_air_quality, aes_string(x = categorical_vars[i], fill = categorical_vars[j]))
      geom_bar(position = "fill") +  # Use "fill" to make it a stacked bar chart (proportions)
      labs(title = paste("Stacked Bar Chart of", categorical_vars[i], "and", categorical_vars[j]),
           x = categorical_vars[i],
           y = "Proportion") +
      theme_minimal() +
      theme(axis.text.x = element_text(angle = 45, hjust = 1))

    # Print the plot
    print(p)
  }
}
```

Stacked Bar Chart of industrial_activity and day_of_week

# Stacked Bar Chart of industrial_activity and holiday

## Stacked Bar Chart of day_of_week and holiday



comments

distribution characterisitcs

The distribution of particulate matter levels is generally right-skewed, indicating that a small number of observations have significantly high levels of particulate matter while most observations are clustered at lower levels. The presence of outliers suggests variations in local conditions affecting air quality.

Observed Relationships

1. Traffic Density: A positive correlation exists between particulate matter levels and traffic density, suggesting that areas with higher vehicle traffic tend to experience elevated levels of particulate matter.

2. Urban Greenery: A negative trend is observed, where higher urban greenery correlates with lower particulate matter, indicating that vegetation may help mitigate air pollution.

3. Temperature and Wind Speed: No strong relationship was identified between particulate matter and temperature. However, there is a slight negative correlation with wind speed, indicating that higher wind speeds may help disperse particulate matter.

Potential Collinearity

Some potential collinearity is observed among the explanatory variables, particularly between traffic density and urban greenery. High traffic areas often have less vegetation, leading to a relationship that may confound the analysis. Additionally, temperature and wind speed may also exhibit collinearity, as changes in one could affect the other.

# 3    Section 3

simple linear regression

```r
X <- cbind(1,data_tidy_air_quality$traffic_density)

Y <-data_tidy_air_quality$particulate_matter
bhat <- solve(t(X) %*% X) %*% t(X) %*% Y

Cmat <- solve(t(X) %*% X)

k <- ncol(X)
rss <- t(Y - X %*% bhat) %*% (Y - X %*% bhat)
# Calculate s2 = RSS/(n-k)
s2 <- as.numeric((rss)/148)
s2
```

```
## [1] 143.5745
```

```r
c_ii <- diag(Cmat)

std.error <- sqrt(s2 * c_ii)
std.error
```

```
## [1] 20.37801682  0.04065266
```

```r
mod1<-lm(data_tidy_air_quality$particulate_matter ~ data_tidy_air_quality$traffic_density, data = data_

summary(mod1)
```

```
##
## Call:
## lm(formula = data_tidy_air_quality$particulate_matter ~ data_tidy_air_quality$traffic_density,
##     data = data_tidy_air_quality)
##
## Residuals:
##     Min      1Q  Median      3Q     Max
## -28.332  -7.561  -1.050   6.110  35.243
##
## Coefficients:
##                                       Estimate Std. Error t value Pr(>|t|)
## (Intercept)                            18.11537   20.37802   0.889   0.3755
## data_tidy_air_quality$traffic_density   0.08400    0.04065   2.066   0.0406 *
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 11.98 on 148 degrees of freedom
## Multiple R-squared:  0.02804,    Adjusted R-squared:  0.02147
## F-statistic: 4.269 on 1 and 148 DF,  p-value: 0.04055
```

hypthesis test

```r
# Summary of ANOVA results
summary(aov(particulate_matter ~ industrial_activity, data = data_tidy_air_quality))
```

11

```
##                    Df Sum Sq Mean Sq F value Pr(>F)
## industrial_activity   3   2182   727.3   5.396 0.0015 **
## Residuals           146  19680   134.8
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```r
# Calculate F-statistic and p-value manually
group_means <- tapply(data_tidy_air_quality$particulate_matter, data_tidy_air_quality$industrial_activi
overall_mean <- mean(data_tidy_air_quality$particulate_matter)

# Calculate SST
SST <- sum((data_tidy_air_quality$particulate_matter - overall_mean)^2)

# Calculate SStreatment
n <- table(data_tidy_air_quality$industrial_activity)
SStreatment <- sum(n * (group_means - overall_mean)^2)

# Calculate SSerror
group_means_vector <- unlist(tapply(data_tidy_air_quality$particulate_matter, data_tidy_air_quality$ind
SSerror <- sum((data_tidy_air_quality$particulate_matter - group_means_vector)^2)

# Calculate degrees of freedom
k <- length(unique(data_tidy_air_quality$industrial_activity))
N <- nrow(data)
DFtreatment <- k - 1
DFerror <- 150 - k

# Calculate Mean Squares
MStreatment <- SStreatment / DFtreatment
MSerror <- SSerror / DFerror


# Calculate F-statistic
F_statistic <- MStreatment/MSerror

# Output F-statistic
F_statistic
```

```
## [1] 5.395959
```

```r
# Calculate p-value
p_value <- pf(F_statistic, DFtreatment, DFerror, lower.tail = FALSE)
p_value
```

```
## [1] 0.001502236
```

# 4 Question 4

Table 1: Confidence Interval for each Coefficients

|  | 2.5 % | Estimate | 97.5 % |
|---|---|---|---|
| **Intercept** | | | |
| (Intercept) | -21.0568 | 13.7937 | 48.6442 |
| **Traffic Density** | | | |
| traffic_density | 0.0155 | 0.0799 | 0.1444 |
| **Industrial Activity** | | | |
| industrial_activityLow | -3.1721 | 2.6589 | 8.4900 |
| industrial_activityModerate | 0.6047 | 6.4545 | 12.3043 |
| industrial_activityHigh | -0.2503 | 5.3652 | 10.9806 |
| **Natural Factors** | | | |
| temperature | -1.1521 | -0.2815 | 0.5891 |
| humidity | -0.1111 | 0.1926 | 0.4962 |
| wind_speed | -0.8040 | 0.0193 | 0.8426 |
| temperature:humidity | -0.0088 | 0.0061 | 0.0209 |
| **Day of Week** | | | |
| day_of_weekTuesday | -5.9877 | 0.0133 | 6.0142 |
| day_of_weekWednesday | -5.3501 | 0.1565 | 5.6630 |
| day_of_weekThursday | -5.5367 | 0.1662 | 5.8690 |
| day_of_weekFriday | -8.0602 | -2.4221 | 3.2161 |
| day_of_weekSaturday | -12.3605 | -4.4832 | 3.3940 |
| day_of_weekSunday | -10.2167 | -2.0885 | 6.0396 |
| **Holiday** | | | |
| holidayNo | -6.7151 | -0.9961 | 4.7228 |
| **Urban Greenery** | | | |
| urban_greenery | -0.4142 | -0.2954 | -0.1766 |

### 4.0.1 Hypothesis Testing

We'd like to perform hypothesis tests on the following variables: Temperature, Humidity, Industrial Levels, and Day of Week.

We'll start by examining whether Temperature has an effect on the concentration of Particulate Matter. We'll use the following We use the following set of hypothesis:

```
#\begin{align}
#$$H_0: \beta_{temp} = \beta_{hum:temp} = 0$$
#$$H_A: \beta_{temp} \neq 0 \text{ and } \beta_{hum:temp} \neq 0$$
#\end{align}
#\
#This can be done by comparing the restricted and un restricted model:
#$$Y_R = \beta_0 + \beta_{traffic}X + $$
```

```
model_unrestricted <- lm(particulate_matter ~ . +
                         temperature:humidity,
                         data=data_tidy_air_quality)
model_restricted <- update(model_unrestricted, .~.
                         - temperature
                         - temperature:humidity)
anova(model_unrestricted, model_restricted)
```

```
## Analysis of Variance Table
##
## Model 1: particulate_matter ~ traffic_density + industrial_activity +
##     temperature + humidity + wind_speed + day_of_week + holiday +
##     urban_greenery + temperature:humidity
## Model 2: particulate_matter ~ traffic_density + industrial_activity +
##     humidity + wind_speed + day_of_week + holiday + urban_greenery
##   Res.Df   RSS Df Sum of Sq      F Pr(>F)
## 1    133 11032
## 2    135 11096 -2   -63.801 0.3846 0.6815
```

Using the anova function in R, we compare the two models with F test, yielding a P value 0.6815, suggesting that temperature doesn't have a significant effect on the concentration of particular matter.

# 5   part 2

# 6   Scenario A:

## 6.1   Methodology and discussion of results for Scenario A

Simulation under the null hypothesis ($\beta_1 = 0$):

We simulate the data assuming $\beta_0 = 30$, $\beta_1 = 0$, and errors are uniformly distributed. The errors will be sampled from a uniform distribution where $e \sim U(a, b)$ with the constraint that $Var(e) = 100$

For a uniform distribution $e \sim U(a, b)$ with a variance $\sigma^2 = 100$, $Var(e) = \frac{(b-a)^2}{12}$. Solving for $a$ and $b$ we get $a = -17.32$ and $b = 17.32$

To simulate our data, we ran a loop 1000 times, and in each simulation, we used runif(length(temperature), min = a, max = b) function in R to generate random errors. Next calculated $Y_i = 30 + e$. We then used the lm function to fit our regression model. We extacted p values from our model, and ran ifelse statement to check whether our null hypothesis was rejected or not. Finally, we counted the number of null hypothesises rejected and we observed that our type 1 error rate for this scenario is 0.043

Type I error is the probability of incorrectly rejecting the null hypothesis when it is true (false positive). Under the null hypothesis, the expected Type I error rate should be equal to the chosen significance level, typically 0.05.

Under a uniform distribution, the error terms tend to be more tightly spread compared to the tails of a normal distribution (which has more extreme values due to its longer tails). This can result in:

Underestimated variability in the model, leading to more frequent rejections of the null hypothesis when it should not be rejected. Inflated Type I error rate: The test may incorrectly reject the null hypothesis more often than expected under the nominal significance level.

# 7 Scenario B:

## 7.1 Methodology and discussion of results for Scenario B

Again we run a loop of 1000 simulations. In each simulation, We simulate the error variances using the normal distribution with mean of 100 and variance of 50. We then ensure that the variances are positive by taking the absolute value of the simulated variances. Then we simulated error terms using a normal distribution with a mean of 0 and a variance of our error_variance calculated before. This ensures that our errors have a non constant variance. We then repeated the same steps as above to obtain our type 1 error rate. Our observated type 1 error rate is 0.047.

Impact of Heteroscedasticity: Effect on Type I Error Rate: Heteroscedasticity violates the assumption of constant variance in regression models. This leads to incorrect estimates of standard errors and, consequently, incorrect p-values. As a result, the Type I error rate will likely increase beyond the nominal 0.05 level.

Our error rate is lower than 0,05. There could be a few reasons for this:

Conservative Hypothesis Test: Our test may be too conservative. This means the test is less likely to reject the null hypothesis even when it might be warranted. It could happen if: The variability of the errors isn't large enough to make the standard errors inaccurate in a way that increases the Type I error.

The model has adjusted in such a way that it becomes harder to reject the null hypothesis, potentially overcorrecting for the heteroscedasticity.

# 8 Scenario C:

## 8.1 Methodology and discussion of results for Scenario C:

Using our loop again, we run 1000 simulations,but this time we test the dependence of errors in our model. To generate Correlated errrors in our model we, use the provided function given to us. As above we repeated the same steps in order to obtain the proportion of the number of times our null hypothesis was incorrectly rejected. OUr type 1 errot rate was 0.053.

Autocorrelation violates the assumption of independent errors in linear regression. This can lead to incorrect standard errors for the regression coefficients, making hypothesis tests less reliable.

Typically, autocorrelated errors inflate Type I errors if unaccounted for because the model underestimates the true error variability. However, depending on the specific pattern of autocorrelation and sample size, this effect can vary.