

STAT0028 ICA

Group 1 - 24185360 & 20082011

19 Nov 2024

Question 1

Data Description

For this question, we are investigating a dataset which has 2022 measurements of **NOx** (nitrogen oxide) pollution content in the ambient air. There are three response variables are collected:

- noxem** : Sum of NOx emission of cars on this motorway.
- ws** : Wind speed in m/s.
- humidity** : Absolute humidity in the air in g/kg air.

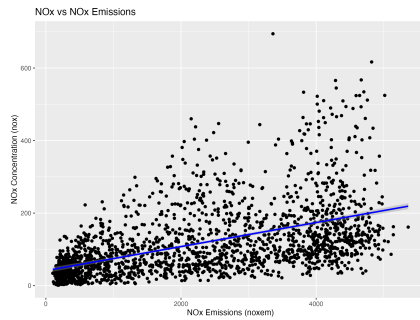


Figure 1: NOx vs noxem

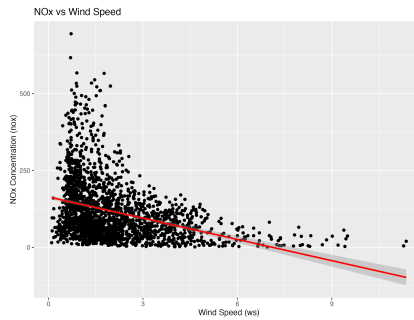


Figure 2: NOx vs windspeed

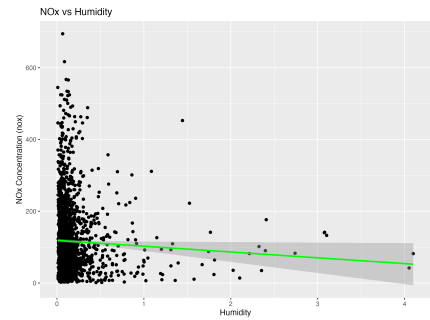


Figure 3: NOx vs humidity

Figure 4: Relationship between target variables and response variables. (Code: Appendix 1)

Figure 4 illustrates the relationship between the target variable (NOx) and the three predictor variables. While the first plot in Figure 5 (NOx vs noxem) suggests a positive linear relationship, the other two plots in Figure 6 and Figure 7 (NOx vs ws and NOx vs humidity) indicate significant non-linear associations as well as some potential outliers that might affect the model's performance. As a result, appropriate transformations will be applied to the data to improve model fitting and capture these non-linear effects more effectively.

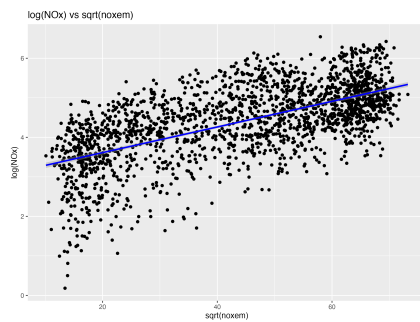


Figure 5: $\log(\text{NOx})$ vs $\sqrt{\text{noxem}}$

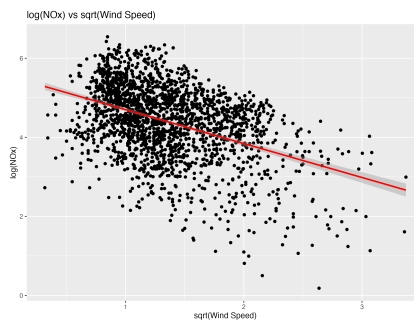


Figure 6: $\log(\text{NOx})$ vs $\sqrt{\text{ws}}$

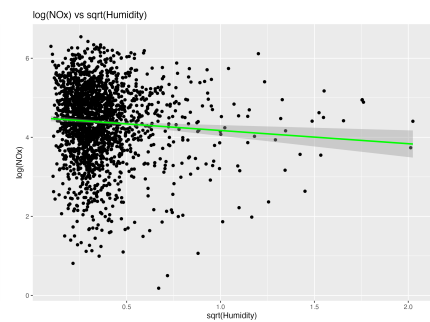


Figure 7: $\log(\text{NOx})$ vs $\sqrt{\text{hum}}$

Figure 8: Relationship between target variables and response variables after transformation. (Code: Appendix 2)

Statistical Models

Linear Model (Baseline Model)

The **linear model**, used as a baseline, fits the response variable nox_i as a linear combination of the original predictors humidity $_i$, noxem $_i$, and ws $_i$ without any transformations. It assumes linear relationships, normally distributed residuals, and constant variance, but may struggle with non-linear patterns or skewed data. The mathematical equation of linear model performs as followed.

$$nox_i = \beta_0 + \beta_1 \cdot \text{humidity}_i + \beta_2 \cdot \text{noxem}_i + \beta_3 \cdot \text{ws}_i + \epsilon_i$$

Where:

- nox_i : The target variable (NOx concentration for observation i).
- β_0 : The intercept term.
- β_1 : Effect of humidity $_i$ (absolute humidity).
- β_2 : Effect of noxem $_i$ (NOx emissions).
- β_3 : Effect of ws $_i$ (wind speed).
- ϵ_i : The residual error term for observation i .

Generalised Linear Models

The **generalized linear model (GLM)** applies transformations to address these limitations. The response variable nox_i is log-transformed, and the predictors are square-root transformed to capture non-linear relationships. Two distributions will be considered: the **Gamma distribution** and the **Gaussian distribution**. The mathematical equation of GLM is presended below.

$$\log(nox_i) = \beta_0 + \beta_1 \cdot \sqrt{\text{noxem}_i} + \beta_2 \cdot \sqrt{\text{ws}_i} + \beta_3 \cdot \sqrt{\text{humidity}_i} + \epsilon_i$$

Where:

- nox_i : The target variable (NOx concentration for observation i).
- $\sqrt{\text{noxem}_i}$: The square root transformation of NOx emissions for observation i .
- $\sqrt{\text{ws}_i}$: The square root transformation of wind speed for observation i .
- $\sqrt{\text{humidity}_i}$: The square root transformation of absolute humidity for observation i .
- β_0 : The intercept term.
- $\beta_1, \beta_2, \beta_3$: The coefficients of the respective predictors.
- ϵ_i : The residual error term for observation i .

Results of Fitted Models

Before fitting the models, the original dataset was split into an 80% training set and a 20% testing set to enable robust model evaluation. Model comparison was conducted using MSE, MAE, R-squared and AIC as evaluation metrics. Since data transformations were applied to the GLMs, the predicted values are on a different scale compared to the baseline linear model. To ensure a fair comparison, the MSE and MAE values were recalculated on the original scale of the target variable. Table 1 below summarizes the performance of the three models on both the training and testing datasets. The code for this result in the table is given in Appendix 3, 4, 5)

Model	Dataset	MSE	MAE	R-squared	AIC
Linear Model (Baseline)	Training	5765.019	52.53818	0.439158	–
	Testing	4215.312	47.67469	0.410928	–
GLM with Gamma Distribution	Training	7038.046	53.15608	0.616754	–
	Testing	5440.196	47.15623	0.548968	3442.646
GLM with Gaussian Distribution	Training	5499.205	46.83815	0.6578908	–
	Testing	3803.385	40.61301	0.5974634	2735.202

Table 1: Models Performance Results including MSE, MAE, R-squared, and AIC.

Non-recommended Models

Based on the model performance results, we do not recommend the linear model fitted with the original variables (baseline model) or the GLM with a Gamma distribution. The linear model is not recommended due to its relatively poor performance across all evaluation metrics compared to the GLM with a Gaussian distribution. Specifically, its R-square values of 0.439 for the training set and 0.411 for the testing set indicate that it explains only a modest proportion of the variance in the target variable.

Similarly, while the GLM with a Gamma distribution achieves higher R-square values for both training and testing sets compared to the linear model, its MSE and MAE values are significantly higher than those of the baseline model. Therefore, this model is also not recommended.

From the data description section, it is evident that the relationships between the target variable and the predictors do not strictly follow a linear pattern. This observation is further supported by the low R-square values of the linear model, which shows the necessity for data transformations such as the logarithmic and square root transformations applied in the GLMs. Moreover, the high MSE and MAE values of the GLM with a Gamma distribution emphasizes the importance of selecting an appropriate distribution for the target variable. This result ultimately led to the adoption of the GLM with a Gaussian distribution, which is discussed further in the next section.

Recommended Models

The GLM with a Gaussian distribution is the recommended model based on the best performance metrics compared with the previous two models. As shown in Table 1, this model achieved the lowest MSE and MAE values across both the training and testing datasets, which indicates higher predictive accuracy. On the other hand, the R-square values of 0.658 for training and 0.597 for testing show that this model explains a substantial proportion of the variance in the target variable, which is better than the alternative models. The low AIC value of 2735.202 compared to GLM with Gamma distribution further supports the model's efficiency in balancing goodness-of-fit and complexity.

The Figure 9 demonstrates the diagnostic plot for Gaussian model. The **Residuals vs Fitted** plot shows no clear patterns and suggests that the model reasonably satisfies the assumptions of linearity and homoscedasticity. The **Normal Q-Q** plot indicates that residuals follow a normal distribution, supporting the appropriateness of the Gaussian assumption. The **Scale-Location** plot demonstrates consistent variance across fitted values, with minor heteroscedasticity that does not significantly affect model reliability. The **Residuals vs Leverage** plot shows that most residuals are centered around zero with no clear pattern, which demonstrates unbiased predictions and a good model fit.

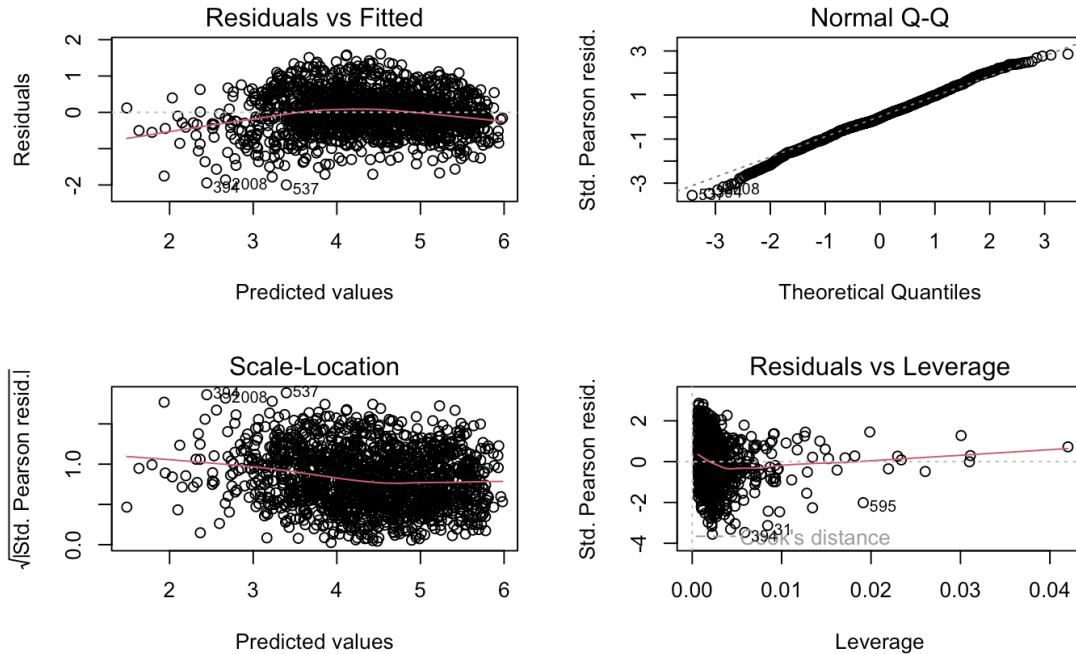


Figure 9: Diagnostic Plot for GLM Gaussian Model. (Code: Appendix 6)

Hypothesis Testing (Code: Appendix 7 and 8)

In the GLM with a Gaussian distribution, we aim to test the null hypothesis that the effects of $\sqrt{\text{noxem}}$ and $\sqrt{\text{ws}}$ as well as $\sqrt{\text{humidity}}$ and $\sqrt{\text{ws}}$ are identical. First, we conduct the test for the first pair. From the model summary, we observe that the estimated coefficient for $\sqrt{\text{noxem}}$ is 0.036, while for $\sqrt{\text{ws}}$, it is -1.041. To determine whether the difference in these coefficients is statistically significant, we formulate the hypotheses as follows:

$$H_0 : \beta_{\text{noxem}} = \beta_{\text{ws}} \quad (\text{The effects of } \sqrt{\text{noxem}} \text{ and } \sqrt{\text{ws}} \text{ are the same})$$

$$H_a : \beta_{\text{noxem}} \neq \beta_{\text{ws}} \quad (\text{The effects of } \sqrt{\text{noxem}} \text{ and } \sqrt{\text{ws}} \text{ are different})$$

For the hypothesis test, we reformulate the null hypothesis to assess the difference between the coefficients:

$$H_0 : \beta_{\text{noxem}} - \beta_{\text{ws}} = 0$$

$$H_a : \beta_{\text{noxem}} - \beta_{\text{ws}} \neq 0$$

Given the estimated coefficients for $\sqrt{\text{noxem}}$ and $\sqrt{\text{ws}}$ from the GLM with Gaussian distribution, we aim to calculate the 95% confidence interval for the difference between their effects. The difference in the estimated coefficients is denoted as:

$$\Delta = \hat{\beta}_{\sqrt{\text{noxem}}} - \hat{\beta}_{\sqrt{\text{ws}}}$$

The standard error of this difference is given as:

$$\text{SE}_{\Delta} = 0.0299854$$

The 95% confidence interval is then computed as:

$$CI = \Delta \pm t_{\alpha/2} \cdot \text{SE}_{\Delta}$$

Plugging in the values, we get:

$$CI = 1.076882 \pm 1.96 \cdot 0.0299854$$

$$CI = 1.076882 \pm 0.058811$$

$$CI = [1.018281, 1.135483]$$

This interval provides the range of plausible values for the true difference between the effects of $\sqrt{\text{noxem}}$ and $\sqrt{\text{ws}}$. Since the interval does not include 0, we reject the null hypothesis and conclude that there is a statistically significant difference between the effects of $\sqrt{\text{noxem}}$ and $\sqrt{\text{ws}}$ at the 95% confidence level.

Using the same method, the null hypothesis for $\sqrt{\text{humidity}}$ and $\sqrt{\text{ws}}$ is:

$$H_0 : \beta_{\text{humidity}} - \beta_{\text{ws}} = 0$$

After computation, the 95% confidence interval is given as:

$$CI = [0.7629484, 1.063325]$$

The result also does not include 0 in the confidence interval, therefore, we can conclude that the effect of $\sqrt{\text{humidity}}$ and $\sqrt{\text{ws}}$ is distinct.

NOx Pollution Report

Introduction

The purpose of this report is to analyze the influence of the three variables **noxem** NOx emissions (amount of NOx produced by cars), wind speed, and absolute humidity on the response variable **NOx** concentrations (a key measure of air pollution) using three different typed of models. The models evaluated include a simple linear model and two generalized linear models (GLMs) using different approaches. These models aim to find the best fit to the data with the results explained in accessible terms.

Overview of the Models

The linear model directly relates NOx concentrations to the regressors without transforming the data. This straightforward approach assumes that relationships between the variables are linear and provides a baseline model for furthre comparison.

To improve model's reliability, the GLMs apply transformations to the data. NOx concentrations are log-transformed to address variability and skewness, while the regressors, such as NOx emissions, wind speed, and humidity, are square-root transformed to better capture non-linear relationships. Two GLMs with different distribution assumptions were evaluated.

Analysis of Model Results

For the results, the linear model struggled to capture the complexity of the relationships between NOx concentrations and the regressors. It did not account for non-linear patterns in the data and showed limited accuracy and reliability.

The GLM with Gamma distribution showed improvement over the linear model but still had limitations. Its assumptions about the data distribution were not fully consistent with the observed characteristics, resulting in higher errors.

The GLM with Gaussian distribution performed the best overall. By applying log and square-root transformations, it effectively captured the non-linear relationships in the data and aligned well with the distribution of NOx concentrations. This made it the most reliable model in this study.

Interpreting the model's results, the sum of NOx emissions from cars on this motorway shows the most significant positive correlation with NOx pollution in the air. On the other hand, wind speed has a negative correlation with pollution, indicating that stronger wind speeds result in lower levels of NOx pollution in the air. While humidity does not have a significant impact on pollution, it still demonstrates a negative relationship with pollution levels.

To further investigate the effects of the regression coefficients, we conducted hypothesis tests to compare the effects between **noxem** and **ws**, as well as between **humidity** and **ws**. After statistical testing, we found that the effects of both pairs are significantly distinct.

Limitations and Improvements

Considering the potential limitations, the log transformation of NOx concentrations and square-root transformations of predictors improved model performance but may complicate direct interpretation of results. Moreover, the current model assumes that the patterns observed in the data remain consistent. Any significant changes in data characteristics may require adjustments to the model or additional transformations. Lastly, though the GLM with Gaussian distribution achieves high performance, it still does not capture a substantial proportion of the variance. More complex models, such as Generalized Additive Models (GAMs), could be applied to this study to better capture intricate non-linear relationships and improve model performance.

Conclusion

This study analyzed NOx concentrations using a baseline linear model and two GLMs with Gaussian and Gamma distributions. The baseline model performed poorly due to its inability to capture non-linear relationships. The GLM with Gaussian distribution, using log-transformed responses and square-root-transformed predictors, proved most effective. Furthermore, the analysis highlights the dominant role of NOx emissions in contributing to pollution, the mitigating effect of wind speed, and the minor influence of humidity, with hypothesis tests confirming the distinct impacts of these factors.

Question 2

We dispose of the following equation for the expression of D_i :

$$D_i = \frac{1}{p\hat{\sigma}^2} \left(\hat{\beta} - \hat{\beta}_{(i)} \right)^T \mathbf{X}^T \mathbf{X} \left(\hat{\beta} - \hat{\beta}_{(i)} \right) \quad (1)$$

We are given the following:

$$\mathbf{C}^{-1} = (\mathbf{X}^T \mathbf{X})^{-1} \quad (2)$$

$$e_i = y_i - \hat{y}_i \quad (3)$$

$$\hat{y}_i = \mathbf{x}_i^T \hat{\beta} \quad (4)$$

And:

$$\mathbf{H} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X} \quad (5)$$

where \mathbf{H} satisfies $\mathbf{H}\mathbf{Y} = \hat{\mathbf{y}}$ and $h_{i,i}$ is the i^{th} diagonal term of the matrix, which can be written as:

$$h_{i,i} = \mathbf{x}_i^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{x}_i = \mathbf{x}_i^T \mathbf{C}^{-1} \mathbf{x}_i \quad (6)$$

We start by writing the equation of $\hat{\beta}_{(i)}$, representing the least squares estimator of β without the i^{th} observation. We introduce $\mathbf{d}_{(i)} = \mathbf{d} - \mathbf{x}_i y_i = \mathbf{X}^T \mathbf{Y} - \mathbf{x}_i y_i$, where $\mathbf{d} = \mathbf{X}^T \mathbf{Y}$. Then, the least square estimate β is:

$$\hat{\beta} = \mathbf{C}^{-1} \mathbf{X}^T \mathbf{Y} = \mathbf{C}^{-1} \mathbf{d}$$

Then, removing the i^{th} observation:

$$\hat{\beta}_{(i)} = \mathbf{C}_{(i)}^{-1} \mathbf{d}_{(i)}$$

By definition:

$$\mathbf{C}_{(i)}^{-1} = \mathbf{C}^{-1} - \mathbf{x}_i \mathbf{x}_i^T$$

Which can be rewritten as:

$$\mathbf{C}_{(i)}^{-1} = \mathbf{C}^{-1} - \frac{\mathbf{C}^{-1} \mathbf{x}_i \mathbf{x}_i^T \mathbf{C}^{-1}}{1 - \mathbf{x}_i^T \mathbf{C}^{-1} \mathbf{x}_i}$$

Replacing in the difference expression:

$$\begin{aligned} \hat{\beta} - \hat{\beta}_{(i)} &= \mathbf{C}^{-1} \mathbf{d} - \mathbf{C}_{(i)}^{-1} \mathbf{d}_{(i)} \\ &= \mathbf{C}^{-1} \mathbf{d} - \left(\mathbf{C}^{-1} - \frac{\mathbf{C}^{-1} \mathbf{x}_i \mathbf{x}_i^T \mathbf{C}^{-1}}{1 - \mathbf{x}_i^T \mathbf{C}^{-1} \mathbf{x}_i} \right) \mathbf{d}_{(i)} \\ &= \mathbf{C}^{-1} \mathbf{d} - \left(\mathbf{C}^{-1} - \frac{\mathbf{C}^{-1} \mathbf{x}_i \mathbf{x}_i^T \mathbf{C}^{-1}}{1 - \mathbf{x}_i^T \mathbf{C}^{-1} \mathbf{x}_i} \right) (\mathbf{d} - \mathbf{x}_i y_i) \\ &= \mathbf{C}^{-1} \mathbf{d} - \mathbf{C}^{-1} \mathbf{d} + \mathbf{C}^{-1} \mathbf{x}_i y_i - \frac{\mathbf{C}^{-1} \mathbf{x}_i \mathbf{x}_i^T \mathbf{C}^{-1} \mathbf{d}}{1 - \mathbf{x}_i^T \mathbf{C}^{-1} \mathbf{x}_i} + \frac{\mathbf{C}^{-1} \mathbf{x}_i \mathbf{x}_i^T \mathbf{C}^{-1} \mathbf{x}_i y_i}{1 - \mathbf{x}_i^T \mathbf{C}^{-1} \mathbf{x}_i} \\ &= \mathbf{C}^{-1} \mathbf{x}_i y_i - \frac{\mathbf{C}^{-1} \mathbf{x}_i \mathbf{x}_i^T \mathbf{C}^{-1} \mathbf{d}}{1 - \mathbf{x}_i^T \mathbf{C}^{-1} \mathbf{x}_i} + \frac{\mathbf{C}^{-1} \mathbf{x}_i \mathbf{x}_i^T \mathbf{C}^{-1} \mathbf{x}_i y_i}{1 - \mathbf{x}_i^T \mathbf{C}^{-1} \mathbf{x}_i} \end{aligned} \quad (7)$$

Identifying the expression of $h_{i,i}$ and of $\hat{\beta}$:

$$\hat{\beta} - \hat{\beta}_{(i)} = \mathbf{C}^{-1} \mathbf{x}_i y_i + \frac{\mathbf{C}^{-1} \mathbf{x}_i (\mathbf{x}_i^T \mathbf{C}^{-1} \mathbf{x}_i) y_i}{1 - \mathbf{x}_i^T \mathbf{C}^{-1} \mathbf{x}_i} - \frac{\mathbf{C}^{-1} \mathbf{x}_i \mathbf{x}_i^T (\mathbf{C}^{-1} \mathbf{X}^T \mathbf{Y})}{1 - \mathbf{x}_i^T \mathbf{C}^{-1} \mathbf{x}_i}$$

$$\hat{\beta} - \hat{\beta}_{(i)} = \mathbf{C}^{-1} \mathbf{x}_i \mathbf{y}_i + \frac{\mathbf{C}^{-1} \mathbf{x}_i h_{i,i} \mathbf{y}_i}{1 - h_{i,i}} - \frac{\mathbf{C}^{-1} \mathbf{x}_i \mathbf{x}_i^T \hat{\beta}}{1 - h_{i,i}} \quad (9)$$

And therefore, remembering $\hat{y}_i = \mathbf{x}_i^T \hat{\beta}$, we deduce:

$$\hat{\beta} - \hat{\beta}_{(i)} = \mathbf{C}^{-1} \mathbf{x}_i y_i + \frac{\mathbf{C}^{-1} \mathbf{x}_i h_{i,i} y_i}{1 - h_{i,i}} - \frac{\mathbf{C}^{-1} \mathbf{x}_i (\mathbf{x}_i^T \hat{\beta})}{1 - h_{i,i}} \quad (10)$$

$$\begin{aligned} &= \mathbf{C}^{-1} \mathbf{x}_i \left(y_i + \frac{h_{i,i} y_i}{1 - h_{i,i}} \right) - \frac{\mathbf{C}^{-1} \mathbf{x}_i \hat{y}_i}{1 - h_{i,i}} \\ &= \mathbf{C}^{-1} \mathbf{x}_i \left(1 + \frac{h_{i,i}}{1 - h_{i,i}} \right) y_i - \frac{\mathbf{C}^{-1} \mathbf{x}_i \hat{y}_i}{1 - h_{i,i}} \\ &= \mathbf{C}^{-1} \mathbf{x}_i \left(\frac{1 - h_{i,i}}{1 - h_{i,i}} + \frac{h_{i,i}}{1 - h_{i,i}} \right) y_i - \frac{\mathbf{C}^{-1} \mathbf{x}_i \hat{y}_i}{1 - h_{i,i}} \\ &= \mathbf{C}^{-1} \mathbf{x}_i \left(\frac{1 - h_{i,i}}{1 - h_{i,i}} + \frac{h_{i,i}}{1 - h_{i,i}} \right) y_i - \frac{\mathbf{C}^{-1} \mathbf{x}_i \hat{y}_i}{1 - h_{i,i}} \\ &= \mathbf{C}^{-1} \mathbf{x}_i \left(\frac{1 - h_{i,i} + h_{i,i}}{1 - h_{i,i}} \right) y_i - \frac{\mathbf{C}^{-1} \mathbf{x}_i \hat{y}_i}{1 - h_{i,i}} \\ &\hat{\beta} - \hat{\beta}_{(i)} = \mathbf{C}^{-1} \mathbf{x}_i \left(\frac{1}{1 - h_{i,i}} \right) y_i - \frac{\mathbf{C}^{-1} \mathbf{x}_i \hat{y}_i}{1 - h_{i,i}} \end{aligned} \quad (11)$$

Then, factorizing and identifying the expression of the error e_i , we get:

$$\begin{aligned} \hat{\beta} - \hat{\beta}_{(i)} &= \mathbf{C}^{-1} \mathbf{x}_i \left(\frac{y_i}{1 - h_{i,i}} - \frac{\hat{y}_i}{1 - h_{i,i}} \right) \\ &= \mathbf{C}^{-1} \mathbf{x}_i \left(\frac{y_i - \hat{y}_i}{1 - h_{i,i}} \right) \\ \hat{\beta} - \hat{\beta}_{(i)} &= \mathbf{C}^{-1} \mathbf{x}_i \frac{e_i}{1 - h_{i,i}} \end{aligned} \quad (12)$$

which is the result we were expecting to find.

We can rewrite the expression of D_i as:

$$D_i = \frac{1}{p\hat{\sigma}^2} \left(\mathbf{C}^{-1} \mathbf{x}_i \frac{e_i}{1 - h_{i,i}} \right) \mathbf{X}^T \mathbf{X} \left(\mathbf{C}^{-1} \mathbf{x}_i \frac{e_i}{1 - h_{i,i}} \right) \quad (13)$$

Appendix: Q1 R code

Listing 1: Correlations between variables

```
1 library(ggplot2)
2
3 ggplot(emissionssw, aes(x = noxem, y = nox)) +
4   geom_point() +
5   geom_smooth(method = "lm", col = "blue") +
6   labs(title = "NOx vs NOx Emissions",
7         x = "NOx Emissions (noxem)",
8         y = "NOx Concentration (nox)")
9 ggsave("NOx_vs_noxem.png", plot = plot1, width = 8, height = 6, dpi = 300)
10
11 ggplot(emissionssw, aes(x = ws, y = nox)) +
12   geom_point() +
13   geom_smooth(method = "lm", col = "red") +
14   labs(title = "NOx vs Wind Speed",
15         x = "Wind Speed (ws)",
16         y = "NOx Concentration (nox)")
17 ggsave("NOx_vs_WindSpeed.png", plot = plot2, width = 8, height = 6, dpi = 300)
18
19
20 ggplot(emissionssw, aes(x = humidity, y = nox)) +
21   geom_point() +
22   geom_smooth(method = "lm", col = "green") +
23   labs(title = "NOx vs Humidity",
24         x = "Humidity",
25         y = "NOx Concentration (nox)")
```

Listing 2: Correlations between variables after transformation

```

1 ggplot(emissionssw, aes(x = sqrt(noxem), y = log(nox))) +
2   geom_point() +
3   geom_smooth(method = "lm", col = "blue") +
4   labs(title = "log(NOx) vs sqrt(noxem)",
5         x = "sqrt(noxem)",
6         y = "log(NOx)")
7 ggsave("NOx_vs_sqrtnoxem.png", plot = plot4, width = 8, height = 6, dpi = 300)
8
9 ggplot(emissionssw, aes(x = sqrt(ws), y = log(nox))) +
10  geom_point() +
11  geom_smooth(method = "lm", col = "red") +
12  labs(title = "log(NOx) vs sqrt(Wind Speed)",
13        x = "sqrt(Wind Speed)",
14        y = "log(NOx)")
15 ggsave("NOx_vs_sqrtws.png", plot = plot5, width = 8, height = 6, dpi = 300)
16
17
18 ggplot(emissionssw, aes(x = sqrt(humidity), y = log(nox))) +
19  geom_point() +
20  geom_smooth(method = "lm", col = "green") +
21  labs(title = "log(NOx) vs sqrt(Humidity)",
22        x = "sqrt(Humidity)",
23        y = "log(NOx)")
24 ggsave("NOx_vs_sqrtthum.png", plot = plot6, width = 8, height = 6, dpi = 300)

```

Listing 3: Linear Model (Baseline Model)

```

1 # Model fitting using lm function
2 model_linear <- lm(nox ~ humidity + noxem + ws, data = train_data)
3 summary(model_linear)
4
5 # Calculate and print the model performance metrics
6 train_predictions <- predict(model_linear, newdata = train_data)
7 train_actual <- train_data$nox
8
9 ss_total_train <- sum((train_actual - mean(train_actual))^2)
10 ss_residual_train <- sum((train_actual - train_predictions)^2)
11
12 r_squared_train <- 1 - (ss_residual_train / ss_total_train)
13 mse_train <- mean((train_actual - train_predictions)^2)
14
15 mae_train <- mean(abs(train_actual - train_predictions))
16
17 test_predictions <- predict(model_linear, newdata = test_data)
18 test_actual <- test_data$nox
19
20 mse_test <- mean((test_actual - test_predictions)^2)
21 ss_total_test <- sum((test_actual - mean(test_actual))^2)
22 ss_residual_test <- sum((test_actual - test_predictions)^2)
23 r_squared_test <- 1 - (ss_residual_test / ss_total_test)
24
25 mae_test <- mean(abs(test_actual - test_predictions))
26
27 cat("Linear Model (Baseline Model) Performance:\n")
28 cat("Training Data:\n")
29 cat("Mean Squared Error (Training):", mse_train, "\n")
30 cat("Mean Absolute Error (Training):", mae_train, "\n")
31 cat("R-squared (Training):", r_squared_train, "\n")
32 cat("\nTesting Data:\n")
33 cat("Mean Squared Error (Testing):", mse_test, "\n")
34 cat("Mean Absolute Error (Testing):", mae_test, "\n")
35 cat("R-squared (Testing):", r_squared_test, "\n")
36 cat("AIC:", AIC(model_linear), "\n")

```

Listing 4: GLM with Gamma Distribution

```

1 # Model fitting using glm function
2 model_gamma <- glm(log(nox) ~ sqrt(noxem) + sqrt(ws) + sqrt(humidity),
3   family = Gamma(link = "log"), data = train_data)
4
5 summary(model_gamma)
6
7 # Calculate and print the model performance metrics
8 # MSE and MAE for original scale
9 train_predictions <- predict(model_gamma, newdata = train_data, type = "response")
10 train_predictions_original <- exp(train_predictions)
11 train_actual_original <- train_data$nox
12
13 mse_train_original <- mean((train_actual_original - train_predictions_original)^2)
14 mae_train <- mean(abs(train_actual_original - train_predictions_original))
15
16 test_predictions <- predict(model_gamma, newdata = test_data, type = "response")
17 test_predictions_original <- exp(test_predictions)
18 test_actual_original <- test_data$nox
19
20 mse_test_original <- mean((test_actual_original - test_predictions_original)^2)
21 mae_test <- mean(abs(test_actual_original - test_predictions_original))
22
23 cat("GLM with Gamma Distribution Model Performance:\n")
24 cat("Mean Squared Error (Original Scale):\n")
25 cat("Training Data:", mse_train_original, "\n")
26 cat("Testing Data:", mse_test_original, "\n")
27 cat("Mean Absolute Error (Training):", mae_train, "\n")
28 cat("Mean Absolute Error (Testing):", mae_test, "\n")
29
30 # Model performance for transformed scale
31 train_actual <- log(train_data$nox)
32
33 ss_total_train <- sum((train_actual - mean(train_actual))^2)
34 ss_residual_train <- sum((train_actual - train_predictions)^2)
35
36 r_squared_train <- 1 - (ss_residual_train / ss_total_train)
37
38 mse_train <- mean((train_actual - train_predictions)^2)
39 mae_train <- mean(abs(train_actual - train_predictions))
40
41 test_actual <- log(test_data$nox)
42
43 mse_test <- mean((test_actual - test_predictions)^2)
44 mae_test <- mean(abs(test_actual - test_predictions))
45
46 ss_total_test <- sum((test_actual - mean(test_actual))^2)
47 ss_residual_test <- sum((test_actual - test_predictions)^2)
48 r_squared_test <- 1 - (ss_residual_test / ss_total_test)
49
50 cat("GLM with Gamma Distribution Model Performance:\n")
51 cat("Training Data:\n")
52 cat("Mean Squared Error (Training):", mse_train, "\n")
53 cat("Mean Absolute Error (Training):", mae_train, "\n")
54 cat("R-squared (Training):", r_squared_train, "\n")
55 cat("\nTesting Data:\n")
56 cat("Mean Squared Error (Testing):", mse_test, "\n")
57 cat("Mean Absolute Error (Testing):", mae_test, "\n")
58 cat("R-squared (Testing):", r_squared_test, "\n")
59 cat("AIC:", AIC(model_gamma), "\n")

```

Listing 5: GLM with Gaussian Distribution

```

1 # Fit the model
2 model_gaussian <- glm(
3   log(nox) ~ sqrt(noxem) + sqrt(ws) + sqrt(humidity),
4   family = gaussian(link = "identity"),
5   data = train_data
6 )
7
8 summary(model_gaussian)
9
10 # Calculate and print the model performance metrics
11 # MSE and MAE for original scale
12
13 train_predictions <- predict(model_gaussian, newdata = train_data, type = "response")
14 train_predictions_original <- exp(train_predictions)
15 train_actual_original <- train_data$nox
16
17 mse_train_original <- mean((train_actual_original - train_predictions_original)^2)
18 mae_train <- mean(abs(train_actual_original - train_predictions_original))
19
20 test_predictions <- predict(model_gaussian, newdata = test_data, type = "response")
21 test_predictions_original <- exp(test_predictions)
22 test_actual_original <- test_data$nox
23
24 mse_test_original <- mean((test_actual_original - test_predictions_original)^2)
25 mae_test <- mean(abs(test_actual_original - test_predictions_original))
26
27 cat("GLM with Gaussian Distribution Model Performance:\n")
28 cat("Mean Squared Error (Original Scale):\n")
29 cat("Training Data:", mse_train_original, "\n")
30 cat("Testing Data:", mse_test_original, "\n")
31 cat("Mean Absolute Error (Training):", mae_train, "\n")
32 cat("Mean Absolute Error (Testing):", mae_test, "\n")
33
34 # For transformed scale
35 train_actual <- log(train_data$nox)
36
37 ss_total_train <- sum((train_actual - mean(train_actual))^2)
38 ss_residual_train <- sum((train_actual - train_predictions)^2)
39
40 r_squared_train <- 1 - (ss_residual_train / ss_total_train)
41
42 mse_train <- mean((train_actual - train_predictions)^2)
43
44 test_actual <- log(test_data$nox)
45
46 mse_test <- mean((test_actual - test_predictions)^2)
47
48 ss_total_test <- sum((test_actual - mean(test_actual))^2)
49 ss_residual_test <- sum((test_actual - test_predictions)^2)
50 r_squared_test <- 1 - (ss_residual_test / ss_total_test)
51
52 cat("GLM with Gaussian Distribution Model Performance:\n")
53 cat("Training Data:\n")
54 cat("Mean Squared Error (Training):", mse_train, "\n")
55 cat("R-squared (Training):", r_squared_train, "\n")
56 cat("\nTesting Data:\n")
57 cat("Mean Squared Error (MSE):", mse_test, "\n")
58 cat("R-squared (Testing):", r_squared_test, "\n")
59 cat("AIC:", AIC(model_gaussian), "\n")

```

Listing 6: Diagnostic plot for GLM with Gaussian Distribution

```
1 par(mfrow = c(2, 2))
2 plot(model_gaussian)
```

Listing 7: Hypothesis Testing 1

```
1 coefs <- coef(model_gaussian)
2 cov_matrix <- vcov(model_gaussian)
3
4 coef_diff <- coefs["sqrt(noxem)"] - coefs["sqrt(ws)"]
5
6 se_diff <- sqrt(
7   cov_matrix["sqrt(noxem)", "sqrt(noxem)"] +
8   cov_matrix["sqrt(ws)", "sqrt(ws)"] -
9   2 * cov_matrix["sqrt(noxem)", "sqrt(ws)"]
10 )
11
12 ci_lower <- coef_diff - 1.96 * se_diff
13 ci_upper <- coef_diff + 1.96 * se_diff
14
15 cat("Difference in coefficients (sqrt(noxem) - sqrt(ws)):", coef_diff, "\n")
16 cat("95% Confidence Interval:", ci_lower, ", ", ci_upper, "]\n")
17 cat("se:", se_diff, "\n")
18
19 if (ci_lower <= 0 && ci_upper >= 0) {
20   cat("Fail to reject the null hypothesis:
21   The effects of noxem and ws are not significantly different.\n")
22 } else {
23   cat("Reject the null hypothesis:
24   The effects of noxem and ws are significantly different.\n")
25 }
```

Listing 8: Hypothesis Testing 2

```
1 coefs <- coef(model_gaussian)
2 cov_matrix <- vcov(model_gaussian)
3
4 coef_diff <- coefs["sqrt(humidity)"] - coefs["sqrt(ws)"]
5
6 se_diff <- sqrt(
7   cov_matrix["sqrt(humidity)", "sqrt(humidity)"] +
8   cov_matrix["sqrt(ws)", "sqrt(ws)"] -
9   2 * cov_matrix["sqrt(humidity)", "sqrt(ws)"]
10 )
11
12 ci_lower <- coef_diff - 1.96 * se_diff
13 ci_upper <- coef_diff + 1.96 * se_diff
14
15 cat("Difference in coefficients (sqrt(humidity) - sqrt(ws)):", coef_diff, "\n")
16 cat("95% Confidence Interval:", ci_lower, ", ", ci_upper, "]\n")
17 cat("se:", se_diff, "\n")
18
19 if (ci_lower <= 0 && ci_upper >= 0) {
20   cat("Fail to reject the null hypothesis:
21   The effects of noxem and ws are not significantly different.\n")
22 } else {
23   cat("Reject the null hypothesis: The effects of noxem and ws are significantly different.\n")
24 }
```