



Analysis of Trends in Temperatures

Authors:

Thomas Beurksens - i6296221

Timothy Hanson - i6286194

Alexandra Holikova - i6309418

Mathematical Statistics

Faculty of Business and Economics

March 2025

Contents

1	Introduction	1
2	Problem Statement	1
2.1	Research Question	1
3	Annual Temperature Data Analysis	1
4	Subsample Comparison Using t -Tests	3
4.1	Hypotheses	3
4.2	Assumptions	3
4.3	One-Sample t -Test	3
4.4	Two-Sample t -Test	4
4.5	Confidence Intervals	4
4.6	Comparison between one and two-sample t -Tests	5
5	Linear Regression Model	5
6	Bootstrapping Implementation	9
6.1	Relevance to Temperature Trend Analysis	9
6.2	Bootstrap Methods: Theory and Application	9
7	Extension to Monthly Data	11
7.1	Data Grouping and Seasonal Averages	12
7.2	Linear Regression Analysis	12
7.3	Hypothesis Testing	13
8	Conclusion	14
9	Appendix	15

1 Introduction

In recent decades, climate change has become one of the most pressing global challenges. A key indicator of this phenomenon is the rise in average temperatures across the world. Understanding local temperature trends over time is important for identifying how global climate patterns affect specific regions.

This project analyzes long-term temperature data from the Netherlands to investigate whether there is statistical evidence of a significant upward trend. We focus primarily on annual average temperatures from De Bilt, a centrally located and well-documented weather station Visser [1], and compare the findings with data from Eelde and Maastricht to assess the consistency of patterns across regions.

Our approach starts with visual analysis to identify potential trends. We then apply statistical tools including confidence intervals, hypothesis testing, and linear regression to evaluate whether the observed changes are statistically significant. We use bootstrap methods to improve the reliability of our results, especially in cases where traditional assumptions may not hold.

By combining these techniques, we aim to draw meaningful conclusions about temperature trends in the Netherlands and contribute to the broader understanding of how climate change may be reflected in long-term local data.

2 Problem Statement

The aims of this paper is to determine whether there is statistical evidence of a long-term increase in average temperatures in the Netherlands. Using historical annual temperature data from De Bilt, we investigate whether observed changes over time can be confirmed through, hypothesis testing, linear regression, and bootstrap methods.

2.1 Research Question

Is there statistically significant evidence of an upward trend in average annual temperatures in the Netherlands over the past century, based on data from De Bilt ?

3 Annual Temperature Data Analysis

In this section we analyze the annual temperatures, and split our timespan into decades. The goal is to identify trends across decades and estimate the reliability of these averages using confidence intervals.

To better understand the reliability of each 10-year average temperature, we calculated a 95% confidence interval for every period. This interval gives a range within which we expect the true average temperature to fall, based on the variation in the data and the size of the sample. Since we only have a limited number of years per decade (usually 10), we used the t-distribution:

$$T(X) = \frac{\bar{x} - \mu}{S/\sqrt{n}} \sim t_{n-1}$$

which is more appropriate for smaller sample sizes. The confidence interval takes into account both how much the temperatures vary within a decade and how many data points are available. In general, more variation or fewer data points will lead to wider intervals, indicating more uncertainty.

These intervals were then visualized as error bars in the final graph to show not just the estimated averages, but also the level of confidence we have in each one. Our confidence interval is given by:

$$C(X) = \left[\bar{x} - t_{0.025, n-1} \times \frac{S}{\sqrt{n}}, \bar{x} + t_{0.025, n-1} \times \frac{S}{\sqrt{n}} \right]$$

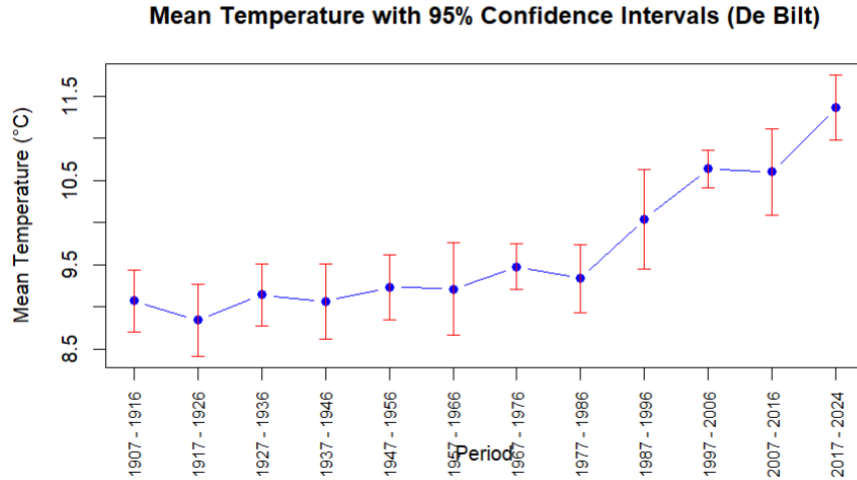


Figure 1: Mean temperature per decade in De Bilt with 95% confidence intervals.

Figure 1 displays the average temperature in De Bilt for each 10-year period from 1907 to 2024. The points represent the mean temperature for each decade, with the connecting line showing overall trends.

The plot illustrates a general upward trend in temperatures over the past century, with some variability across decades. The use of 10-year periods and confidence intervals provides a clearer picture of long-term warming, while also highlighting the degree of statistical uncertainty. Based on this, we proceed with formal inference methods to validate our initial observations.

4 Subsample Comparison Using t -Tests

In this part, we investigate whether there is a statistically significant increase in average annual temperatures in De Bilt after 1970. 1970 was chosen as a cutoff point since it is often used as a climate benchmark IPCC [2]. To assess this, we perform two types of t -tests: a one-sample t -test using the pre-1970 mean as a benchmark, and a two-sample t -test assuming equal variances. We then compare the results and discuss which test is more appropriate for our context.

4.1 Hypotheses

To assess changes in mean temperature over time, we divided the dataset into two groups:

- Pre-1970: years 1907–1970
- Post-1970: years 1971–2022

We test the following hypotheses for both tests:

$$H_0 : \mu_{\text{post}} \leq \mu_{\text{pre}}$$

$$H_A : \mu_{\text{post}} > \mu_{\text{pre}}$$

4.2 Assumptions

Both tests rely on the following assumptions:

- **Independence:** Temperature observations from different years are independent.
- **Identically distributed:** Each group is assumed to come from a distribution with constant variance.
- **Normality:** The data within each group are approximately normally distributed. Due to large sample sizes, the Central Limit Theorem supports this assumption.

4.3 One-Sample t -Test

We begin by performing a one-sample t -test using the pre-1970 mean as the benchmark to test whether the mean temperature after 1970 is significantly higher.

The test statistic is calculated as:

$$t = \frac{\bar{X}_{\text{post}} - \mu_0}{s/\sqrt{n}}, \quad \text{where } \mu_0 = \bar{X}_{\text{pre}} \quad (1)$$

Results

- $t = 9.618$, $df = 53$, critical value $t_{0.05} = 1.674$
- $p \approx 1.61 \times 10^{-13}$

Since $t > t_{0.05}$ and $p < 0.05$, we reject H_0 and conclude that post-1970 temperatures are significantly higher.

4.4 Two-Sample t -Test

We also perform a two-sample t -test, assuming equal variances between the two groups. The test statistic is defined as:

$$t = \frac{\bar{X}_{\text{post}} - \bar{X}_{\text{pre}}}{s_p \sqrt{\frac{1}{n_{\text{post}}} + \frac{1}{n_{\text{pre}}}}} \quad (2)$$

where s_p is the pooled standard deviation:

$$s_p = \sqrt{\frac{(n_{\text{post}} - 1)s_{\text{post}}^2 + (n_{\text{pre}} - 1)s_{\text{pre}}^2}{n_{\text{post}} + n_{\text{pre}} - 2}} \quad (3)$$

The degrees of freedom are computed as:

$$df = n_{\text{post}} + n_{\text{pre}} - 2 \quad (4)$$

Results

- $t = 8.519$, $df = 116$, critical value $t_{0.05} = 1.658$
- $p \approx 3.386 \times 10^{-14}$

Since $t > t_{0.05}$ and $p < 0.05$, we reject H_0 and conclude that there is a significant increase in temperature after 1970.

4.5 Confidence Intervals

To account for uncertainty in our estimates, we compute 95% confidence intervals for both group means:

$$\bar{X} \pm t_{1-\alpha/2, df} \cdot \frac{s}{\sqrt{n}} \quad (5)$$

The resulting intervals are:

- Pre-1970: [8.967, 9.253]
- Post-1970: [10.026, 10.501]

As these intervals do not overlap, they further support the conclusion of a significant temperature increase.

4.6 Comparison between one and two-sample t -Tests

Both the one-sample and two-sample t -tests provide strong evidence that temperatures in De Bilt have increased significantly since 1970. However, the two-sample t -test is more appropriate for our context because it accounts for estimation uncertainty in both groups. The one-sample test treats the pre-1970 mean as fixed, which may underestimate variability and lead to overconfident conclusions.

Therefore, we base our main conclusion on the two-sample t -test, which shows strong statistical evidence of a significant and meaningful increase in average annual temperatures in De Bilt after 1970.

5 Linear Regression Model

In this section, we investigate the presence of a linear upward trend in average annual temperatures using regression analysis. Compared to simple averages, regression offers a more precise measure of the trend over time. We estimate the following linear model:

$$Y_i = \alpha + \beta x_i + \varepsilon_i \quad \forall i = 1, \dots, n$$

In this formulation:

- Y_i is the dependent variable representing the average annual temperature in year i ,
- x_i is the explanatory variable denoting the year (e.g., 1907, 1908, ...),
- α is the intercept, representing the estimated average temperature when $x_i = 0$,
- β is the slope coefficient, indicating the average change in temperature per year,
- ε_i is the error term, capturing unobserved random variation and measurement noise.

Since a linear regression model is being used the following assumptions must be made. We assume that the linear model is correct $\forall i = 1, \dots, n$ and we assume that x_i is a constant and not random.

Furthermore, it is also assumed that the mean of ε_i is equal to zero and has some variance σ^2 . Furthermore, $\varepsilon_1, \dots, \varepsilon_n$ are all uncorrelated and its distribution is left unspecified. The model is estimated as an ordinary least squared estimation which aims to minimize the sum of the squared distances between Y_i and $\alpha + \beta x_i$, by doing so we get least squares estimators for α and β which are α_{LS} and β_{LS} . This was implemented and the visual representation of the regression function can be found in Figure 2.

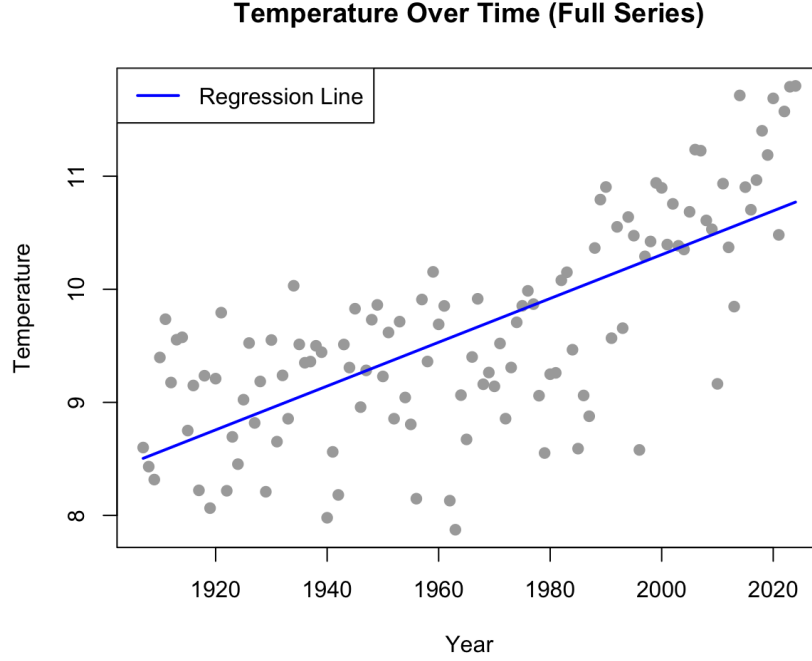


Figure 2: Temperature observations in De Bilt from 1907 to 2022 with a fitted linear regression line.

Once estimated the following poluation model was found

$$Y_i = -28.4346 + 0.0193x_i \quad (6)$$

Although the plot suggests an upward-sloping trend, this visual observation alone is not sufficient to draw definitive conclusions. Estimation uncertainty remains, especially given that the sample size, while moderate (117 observations), may not be large enough to guarantee precise inference. Additionally, time series data can violate classical assumptions such as independence and homoscedasticity. To address this, we rely on asymptotic analysis, which assumes that the sample size is sufficiently large for statistical approximations to hold.

The asymptotic distribution of the OLS estimator $\hat{\beta}$ is derived using the central limit theorem (CLM), which ensures that under suitable conditions, the standardized estimator converges in

distribution to the normal distribution as the sample size increases.

$$\frac{\hat{\beta}_{n,LS} - \beta}{\sqrt{S_n^2/S_{xx}}} \xrightarrow{d} \mathcal{N}(0, 1)$$

Since the true parameter β is unobserved, we estimate it using a 95% confidence interval based on the asymptotic normality of the OLS estimator. The confidence interval is given by:

$$C(X) = \left\{ \hat{\beta}_{n,LS} \pm z_{0.025} \sqrt{\frac{S_n^2}{S_{xx}}} \right\}$$

where $\hat{\beta}_{n,LS}$ is the least squares estimator, S_n^2 is the estimator of the error variance, $S_{xx} = \sum_{i=1}^n (x_i - \bar{x})^2$, and $z_{0.025} \approx 1.96$ is the 97.5th percentile of the standard normal distribution.

By substituting the values obtained from the regression estimation, the following confidence interval is found:

$$C(X) = \{0.01589; 0.02284\}$$

It can be said that the true value of β lies between these two values with 95% certainty. The estimator $\hat{\beta}_{n,LS}$ lies in this interval as well. Furthermore, from this confidence interval it can be concluded that β is greater than 0 providing the evidence needed for an upward trend. In order to solidify this claim a hypothesis test is conducted.

The null and alternative hypotheses are defined as:

$$H_0 : \beta \leq 0 \quad \text{vs.} \quad H_1 : \beta > 0$$

We reject the null hypothesis if the test statistic exceeds the critical value from the standard normal distribution:

$$\frac{\hat{\beta}_{n,LS} - \beta_0}{\sqrt{S_n^2/S_{xx}}} > z_{0.025}$$

Using our computed values, we find that:

$$\frac{\hat{\beta}_{n,LS} - 0}{\sqrt{S_n^2/S_{xx}}} = 10.93183 > z_{0.025} = 1.959$$

Since the test statistic exceeds the critical value, we reject the null hypothesis at the 5% significance level. This provides strong statistical evidence in favor of the alternative hypothesis, supporting the conclusion that there is a significant upward trend in temperatures over time.

So far, the complete time series has been used it spans from 1907 to 2024. However, there is reason to believe that this temperature change is more pronounce after the 1970s. Hartmann et al. [3], says that there was a brief cooling period between 1970 and 1979 where certain parts of the northern hemisphere including the Netherlands had inconsistent temperature. So while there wasn't a dramatic cooling plateau between 1970-1979, this period is seen as a transition period. Given this the analysis will be for the years of 1979 until the 2024.

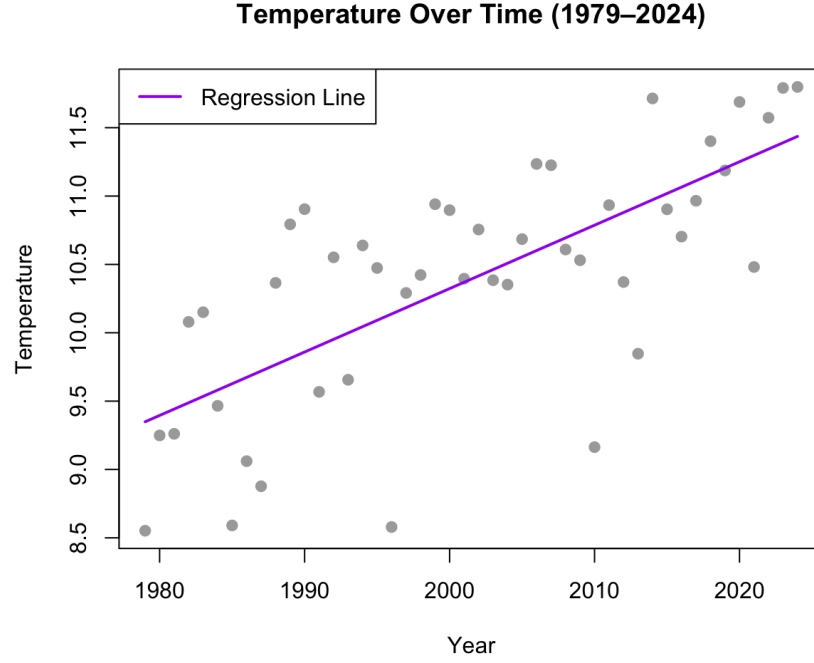


Figure 3: Average annual temperatures in De Bilt from 1979 to 2024. The plot is restricted to this period to investigate whether the upward trend becomes more prominent in recent decades.

In Figure 3, a linear regression line is fitted using the same procedure described earlier.

This new regression line has a slope of 0.04638959 and an intercept of -82.4561 .Furthermore, in order to prove the validity of the results, an asymptotic analysis is conducted in the same manner as before,where we assume n to be large enough. We find the following confidence interval:

$$C(X) = \{0.0327257, 0.06005349\}$$

Performing a one-sided z -test, we obtain a test statistic of 6.654, which is greater than the critical value $z_{0.025} = 1.959$. Therefore, we reject the null hypothesis H_0 . This result provides strong evidence suggesting an upward trend in average annual temperatures.

When comparing the two plots, the regression line in Figure 3 is notably steeper than that in Figure 2. Specifically, the estimated slope for the restricted post-1979 sample is 0.0463, compared to

0.0194 for the full time series. Given that the estimated slope in the restricted period is substantially higher, there is compelling evidence that the rate of temperature increase has accelerated since the 1970s.

That said, this does not imply that there was no increase in temperature before 1970—rather, the data suggest that the warming trend has intensified in more recent decades.

6 Bootstrapping Implementation

6.1 Relevance to Temperature Trend Analysis

The bootstrap method (Efron et al. [4]) provides a robust approach to statistical inference that is particularly valuable for our analysis of temperature trends for several reasons:

- **Assumption robustness:** Traditional regression inference assumes normally distributed, homoscedastic errors. Climate data may violate these assumptions due to potential autocorrelation and changing variability over time.
- **Small sample concerns:** While our full dataset spans 1907–2022, subsamples (e.g., post-1979 data) have limited observations where asymptotic approximations may be less reliable.
- **Complex inference:** The bootstrap provides straightforward confidence intervals for our trend estimates without relying on parametric assumptions.

6.2 Bootstrap Methods: Theory and Application

The bootstrap is a non-parametric resampling technique that allows us to estimate the sampling distribution of a statistic by drawing repeated samples from the data. Unlike classical inference, bootstrap methods do not require assumptions about normality or known variances.

General Bootstrap Algorithm

Let $X = \{x_1, x_2, \dots, x_n\}$ be a sample and let $\hat{\theta} = s(X)$ be a statistic of interest (e.g., the sample mean, a regression coefficient, etc.).

The bootstrap proceeds as follows:

1. Generate B bootstrap samples:

$$X^{*(b)} = \{x_1^{*(b)}, x_2^{*(b)}, \dots, x_n^{*(b)}\}, \quad \text{sampled with replacement from } X, \quad b = 1, \dots, B$$

2. For each bootstrap sample, compute the statistic:

$$\hat{\theta}^{*(b)} = s(X^{*(b)})$$

3. Use the empirical distribution $\{\hat{\theta}^{*(1)}, \dots, \hat{\theta}^{*(B)}\}$ to:

- Approximate the sampling distribution of $\hat{\theta}$
- Estimate the standard error:

$$\widehat{\text{SE}}_{\text{boot}}(\hat{\theta}) = \sqrt{\frac{1}{B-1} \sum_{b=1}^B \left(\hat{\theta}^{*(b)} - \bar{\theta}^* \right)^2}$$

- Construct confidence intervals using the percentile method:

$$\text{CI}_{1-\alpha}^{\text{boot}} = \left[\hat{\theta}^{*(\alpha/2)}, \hat{\theta}^{*(1-\alpha/2)} \right]$$

Bootstrap Hypothesis Testing

To test a null hypothesis $H_0 : \theta = \theta_0$, we compute the test statistic under the observed data:

$$T = \hat{\theta} - \theta_0$$

Then we generate bootstrap replicates under H_0 (e.g., by resampling from a centered distribution or pooled data) and calculate:

$$T^{*(b)} = \hat{\theta}^{*(b)} - \hat{\theta}, \quad b = 1, \dots, B$$

The p-value is then estimated as:

$$p_{\text{boot}} = \frac{1}{B} \sum_{b=1}^B 1(T^{*(b)} \geq T)$$

Application to Temperature Means

Using $B = 1000$ replications, we applied the bootstrap to estimate confidence intervals for average temperatures before and after 1970 in De Bilt:

Pre-1979 (bootstrap): [8.959, 9.249]

Post-1979 (bootstrap): [10.007, 10.476]

These closely match the classical confidence intervals and show no overlap, supporting a statistically significant increase.

Bootstrap Test for Difference in Means

We bootstrapped the difference in means $\Delta = \bar{x}_{\text{post}} - \bar{x}_{\text{pre}}$, obtaining:

$$\hat{\Delta} = 1.153^{\circ}\text{C}$$

The one-sided bootstrap p-value (from 1000 resamples) was:

$$p_{\text{boot}} = 0.000$$

indicating a highly significant increase in post-1970 temperatures.

Bootstrap for Linear Trend Estimation

We applied the bootstrap to estimate the slope $\hat{\beta}_1$ from the linear regression model:

$$Y_t = \beta_0 + \beta_1 t + \varepsilon_t$$

Bootstrap confidence intervals for the slope were:

- **Full sample (1907–2024):** [0.01602, 0.02299]
- **Post-1979:** [0.03028, 0.0566]

This suggests that not only is there a warming trend, but that the rate of warming has increased in more recent decades.

Conclusion

The bootstrap results confirm and reinforce our findings from classical inference. The consistent patterns across both methodologies enhance the robustness of our conclusions. The bootstrap's flexibility and minimal assumptions make it particularly valuable in climatological studies, where standard assumptions (e.g., homoscedasticity, normality) may not hold.

7 Extension to Monthly Data

In this section, we extend our analysis from annual to monthly temperature data. This allows us to examine seasonal variations more closely.

7.1 Data Grouping and Seasonal Averages

To study seasonal effects, we split the monthly data into two key seasons:

- **Winter:** January–March
- **Summer:** June–August

We then grouped observations by decade and calculated the mean temperature for each season and decade. For each average, we construct **95% confidence intervals**.

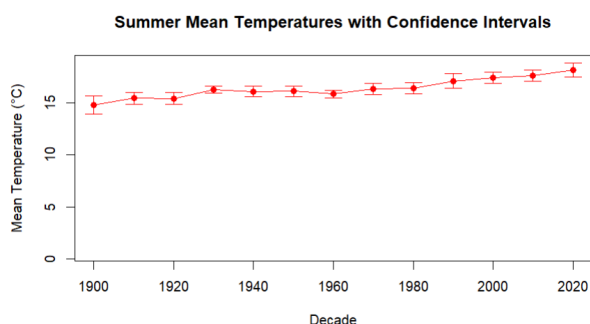


Figure 4: Summer mean temperature with 95% confidence intervals by decade

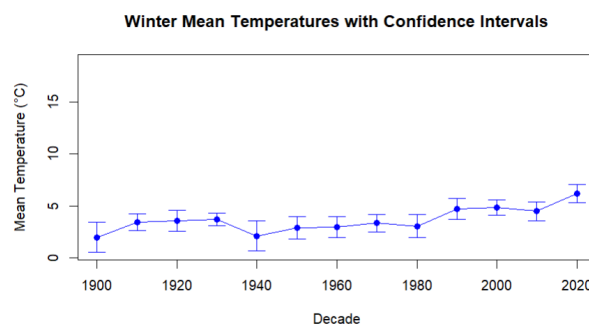


Figure 5: Winter mean temperature with 95% confidence intervals by decade

Figures 4 and 5 display the seasonal means and their associated confidence intervals by decade.

Observed Trends

By comparing the first and last observations for each season, we find:

$$\begin{aligned} \text{Winter 1900: } 1.99^{\circ}\text{C} \quad \text{vs} \quad \text{Winter 2020: } 6.21^{\circ}\text{C} \\ \text{Summer 1900: } 14.76^{\circ}\text{C} \quad \text{vs} \quad \text{Summer 2020: } 18.13^{\circ}\text{C} \end{aligned}$$

While the overall increase in temperature is similar for both seasons, winter exhibits greater decade-to-decade variability, resulting in narrower confidence intervals.

7.2 Linear Regression Analysis

We perform linear regressions using the monthly temperature Y_i as the dependent variable, and a dummy variable x_i for the year:

$$\begin{aligned}\text{Winter: } Y_i &= -33.62 + 0.02 \cdot x_i \\ \text{Summer: } Y_i &= -27.03 + 0.02 \cdot x_i\end{aligned}$$

In both cases, the positive slope confirms a long-term upward trend in seasonal temperatures.

7.3 Hypothesis Testing

We formally tested whether average temperatures after 1970 were significantly higher than those before 1970, using the following hypotheses:

$$\begin{aligned}H_0 : \mu_{\text{post}} &\leq \mu_{\text{pre}} \\ H_a : \mu_{\text{post}} &> \mu_{\text{pre}}\end{aligned}$$

Given our large sample size, we apply an asymptotic z-test. By the CLT, the sampling distribution of the sample mean approximates a normal distribution. The test statistic is given by:

$$Z = \frac{\bar{X}_{\text{post}} - \mu_0}{s/\sqrt{n}} \quad \text{where } \mu_0 = \bar{X}_{\text{pre}}$$

Calculated z-statistics:

$$\begin{aligned}\text{Winter: } Z &= 5.197 \\ \text{Summer: } Z &= 8.95\end{aligned}$$

In both cases, the null hypothesis is rejected at typical significance levels, indicating a statistically significant increase in mean temperature after 1970.

Conclusion

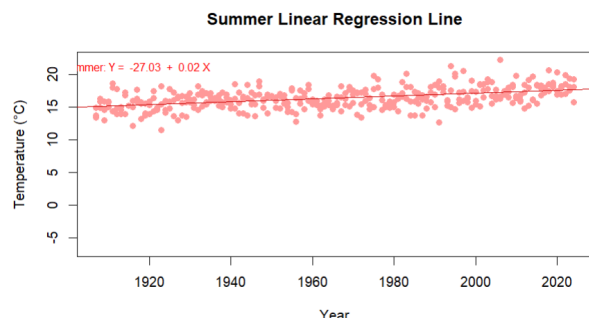


Figure 6: Summer temperature regression line

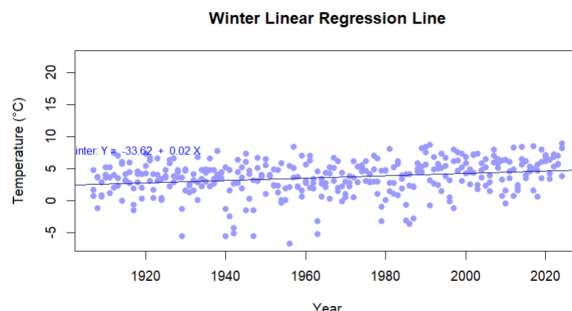


Figure 7: Winter temperature regression line

Both winter and summer months exhibit an upward trend in temperature. Summer temperatures are more consistent across decades, as evidenced by tighter confidence intervals, while winter temperatures show more variability.

8 Conclusion

Our analysis provides strong statistical evidence of a long-term increase in average temperatures in the Netherlands. Both one-sample and two-sample t-tests revealed significantly higher temperatures after 1970, with non-overlapping confidence intervals confirming this trend.

Linear regression further supported these findings, showing a positive and statistically significant slope over the full sample. When restricting the sample to post-1979, the slope nearly doubled, indicating an acceleration in warming in recent decades.

Bootstrap methods confirmed the robustness of these results, yielding similar confidence intervals and significance levels without relying on strict distributional assumptions. The extension to monthly data showed consistent increases in both summer and winter temperatures, with clear upward trends across decades. Together, the evidence points to a clear and statistically significant warming trend in Dutch temperatures over the past century, with a marked increase since the 1970s. These results have further implications within policy making, in regard to greenhouse gas emissions. And provide an eerie warning for the future, but with the implementation of the correct type of policy, things could slow down or even revert.

References

- [1] H. Visser. *The significance of climate change in the Netherlands: An analysis of historical and future trends (1901–2020) in weather conditions, weather extremes and temperature-related impacts*. Tech. rep. 550002007/2005. Project S/550002/01/TO, Uncertainties, Transparency and Communication: Tools for Uncertainty Analysis. Bilthoven, Netherlands: RIVM - National Institute for Public Health and the Environment, 2005.
- [2] IPCC. *Climate Change 2013: The Physical Science Basis*. Working Group I Contribution to the Fifth Assessment Report. 2013. URL: <https://www.ipcc.ch/report/ar5/wg1/>.
- [3] D. L. Hartmann et al. *Climate Change 2013: The Physical Science Basis*. Technical Report. Contribution of Working Group I to the Fifth Assessment Report of the IPCC. Intergovernmental Panel on Climate Change, 2013.
- [4] Bradley Efron et al. *An Introduction to the Bootstrap*. New York: Chapman & Hall/CRC, 1993. ISBN: 9780412042317.

9 Appendix

```
1 ---
2 title: "Final_assignment_Mathematical_statistics"
3 format: html
4 Authors: Timothy, Alexandra, Thomas: (TAT)
5 ---
6
7 ## Code final project Mathematical statistics
8
9 Load data
10
11 ```{r}
12 # --- Load Annual Data ---
13 data <- read.csv2("/Users/timothyhanson/Documents/Temperature_Data/AnnualTemp.csv"
14 )
15 colnames(data)[1:2] <- c("Year", "DeBilt")
16 data$Year <- as.numeric(data$Year)
17 data$DeBilt <- as.numeric(gsub(",", ".", gsub("_", "", data$DeBilt)))
18
19
20 Sampling + averages
21
22 ```{r}
23 # Extract years
24 years <- data[, 1] # First column assumed to be years
25
```

```

26 # Define time range
27 start_year <- 1907
28 end_year <- 2024
29
30 # Initialize result table
31 results_df <- data.frame(Period = character(), De_Bilt = numeric(),
32   stringsAsFactors = FALSE)
33
34 # Compute 10-year means
35 for (start in seq(start_year, end_year, by = 10)) {
36   end <- min(start + 9, end_year) # Cap at 2024
37   period_indices <- which(years >= start & years <= end) # Indices for this range
38   subset_DB <- data[period_indices, 2] # De Bilt temps
39
40   # Store mean
41   results_df <- rbind(results_df, data.frame(
42     Period = paste(start, "-", end),
43     De_Bilt = mean(subset_DB, na.rm = TRUE)
44   ))
45 }
46
47 # Show means
48 print(results_df)
49
50 # CI function for De Bilt temps
51 calc_CI <- function(data, confidence_level = 0.95) {
52   data <- na.omit(data)
53   n <- length(data)
54   if (n == 0) return(c(NA, NA)) # No data case
55
56   # Compute CI
57   mean_val <- mean(data)
58   sd_val <- sd(data)
59   t_value <- qt(1 - (1 - confidence_level) / 2, df = n - 1)
60   margin_error <- t_value * (sd_val / sqrt(n))
61
62   return(c(mean_val - margin_error, mean_val + margin_error))
63 }
64
65 # Initialize CI table
66 ci_df <- data.frame(Period = character(), De_Bilt_CI = character(),
67   stringsAsFactors = FALSE)
68
69 # Compute 10-year CIs
70 for (start in seq(start_year, end_year, by = 10)) {
71   end <- min(start + 9, end_year)
72   period_indices <- which(years >= start & years <= end)
73   subset_DB <- data[period_indices, 2]

```

```

73 # Store CI
74 ci_DB <- calc_CI(subset_DB)
75 ci_df <- rbind(ci_df, data.frame(
76   Period = paste(start, "-", end),
77   De_Bilt_CI = paste("[", round(ci_DB[1], 2), ",", round(ci_DB[2], 2), "]")
78 ))
79 }
80
81 # Show CIs
82 print(ci_df)
83
84 # Merge mean and CI data
85 final_df <- merge(results_df, ci_df, by = "Period")
86
87 # Extract CI bounds
88 ci_bounds <- do.call(rbind, strsplit(gsub("\\[|\\]", "", final_df$De_Bilt_CI), ","))
89
90 final_df$CI_Lower <- as.numeric(ci_bounds[, 1])
91 final_df$CI_Upper <- as.numeric(ci_bounds[, 2])
92
93 # Define x-axis positions
94 x_pos <- 1:nrow(final_df)
95
96 # Plot mean temps with CI
97 plot(
98   x_pos, final_df$De_Bilt,
99   type = "b",
100   ylim = range(final_df$CI_Lower, final_df$CI_Upper),
101   xaxt = "n",
102   xlab = "Period",
103   ylab = "Mean Temperature ( C )",
104   main = "Mean Temperature with 95% Confidence Intervals (De_Bilt)",
105   col = "blue",
106   pch = 16
107 )
108
109 axis(1, at = x_pos, labels = final_df$Period, las = 2, cex.axis = 0.8)
110
111 arrows(
112   x0 = x_pos, y0 = final_df$CI_Lower,
113   x1 = x_pos, y1 = final_df$CI_Upper,
114   angle = 90, code = 3, length = 0.05, col = "red"
115 )
116
117 ***
118 **Hypothesis Testing**
119
120 {r}

```

```

121 # Split groups
122 pre <- data$DeBilt[data$Year <= 1970]
123 post <- data$DeBilt[data$Year > 1970]
124
125 # Sample sizes & stats
126 n_pre <- length(pre)
127 n_post <- length(post)
128 mean_pre <- mean(pre)
129 mean_post <- mean(post)
130 var_pre <- var(pre)
131 var_post <- var(post)
132
133 # One-sample t-test (post vs. pre mean)
134 t1 <- (mean_post - mean_pre) / (sqrt(var_post / n_post))
135 df1 <- n_post - 1
136 p1 <- 1 - pt(t1, df1)
137
138 cat("One-sample t-test:\n")
139 cat("t=", round(t1, 4), "|df=", df1, "|p=", format(p1, scientific = TRUE), "\n\n")
140
141 # Two-sample t-test (equal variance)
142 df2 <- n_pre + n_post - 2
143 pooled_sd <- sqrt(((n_pre - 1)*var_pre + (n_post - 1)*var_post) / df2)
144 t2 <- (mean_post - mean_pre) / (pooled_sd * sqrt(1/n_pre + 1/n_post))
145 p2 <- 1 - pt(t2, df2)
146
147 cat("Two-sample t-test:\n")
148 cat("t=", round(t2, 4), "|df=", df2, "|p=", format(p2, scientific = TRUE), "\n\n")
149
150 # CIs
151 alpha <- 0.05
152 t_crit <- qt(1 - alpha / 2, df2)
153 ci_pre <- mean_pre + c(-1, 1) * t_crit * sqrt(var_pre / n_pre)
154 ci_post <- mean_post + c(-1, 1) * t_crit * sqrt(var_post / n_post)
155
156 cat("95% CI for pre-1970 mean:", round(ci_pre, 2), "\n")
157 cat("95% CI for post-1970 mean:", round(ci_post, 2), "\n")
158
159 '''
160
161 **Linear Regression**
162
163 '''{r}
164
165
166 # --- Full dataset regression ---
167 x <- data$year

```

```

168 y <- data$temp
169 x_bar <- mean(x)
170 y_bar <- mean(y)
171
172 b1 <- sum((x - x_bar) * (y - y_bar)) / sum((x - x_bar)^2)
173 b0 <- y_bar - b1 * x_bar
174 y_pred <- b0 + b1 * x
175 residuals <- y - y_pred
176
177 # --- Compute confidence interval and stats ---
178 SST <- sum((y - y_bar)^2)
179 SSE <- sum(residuals^2)
180 R_squared <- 1 - (SSE / SST)
181
182 n <- length(x)
183 s_squared <- SSE / (n - 2)
184 SE_b1 <- sqrt(s_squared / sum((x - x_bar)^2))
185 SE_b0 <- sqrt(s_squared * (1/n + x_bar^2 / sum((x - x_bar)^2)))
186
187 t_b1 <- b1 / SE_b1
188 t_b0 <- b0 / SE_b0
189
190 # 95% Z-confidence interval (asymptotic)
191 alpha <- 0.05
192 z_alpha <- qnorm(1 - alpha / 2)
193 CI_lower <- b1 - z_alpha * SE_b1
194 CI_upper <- b1 + z_alpha * SE_b1
195
196 # --- Print full regression output ---
197 cat("FULL_DATASET_REGRESSION(all_years)\n")
198 cat("Intercept(b0):", b0, "\n")
199 cat("Slope(b1):", b1, "\n")
200 cat("R-squared:", R_squared, "\n")
201 cat("Standard_Error(slope):", SE_b1, "\n")
202 cat("Standard_Error(intercept):", SE_b0, "\n")
203 cat("Z-test_statistic(slope_0):", (1 - 0) / SE_k1, "\n")
204 cat("Z-based_95%CI_for_slope:", CI_lower, ",", CI_upper, "]\n\n")
205
206 # --- Plot full data and regression line ---
207 plot(x, y, main = "Temperature Over Time(Full Series)", xlab = "Year", ylab = "
    Temperature", pch = 19, col = "darkgray")
208 lines(x, y_pred, col = "blue", lwd = 2)
209 legend("topleft", legend = "Regression Line", col = "blue", lwd = 2)
210
211 # --- Subset data: 1979 to 2024 ---
212 filtered_data <- subset(data, year >= 1979 & year <= 2024)
213 c <- filtered_data$year
214 z <- filtered_data$temp
215 c_bar <- mean(c)

```

```

216 z_bar <- mean(z)
217
218 # --- Regression on subset ---
219 k1 <- sum((c - c_bar) * (z - z_bar)) / sum((c - c_bar)^2)
220 k0 <- z_bar - k1 * c_bar
221 z_pred <- k0 + k1 * c
222 residuals_subset <- z - z_pred
223
224 SST_subset <- sum((z - z_bar)^2)
225 SSE_subset <- sum(residuals_subset^2)
226 R_squared_subset <- 1 - (SSE_subset / SST_subset)
227
228 n_subset <- length(c)
229 s_squared_subset <- SSE_subset / (n_subset - 2)
230 SE_k1 <- sqrt(s_squared_subset / sum((c - c_bar)^2))
231 SE_k0 <- sqrt(s_squared_subset * (1/n_subset + c_bar^2 / sum((c - c_bar)^2)))
232 t_k1 <- k1 / SE_k1
233 t_k0 <- k0 / SE_k0
234
235
236 # --- Print subset regression output ---
237 cat("SUBSET REGRESSION (1979 2024 )\n")
238 cat("Intercept (k0):", k0, "\n")
239 cat("Slope (k1):", k1, "\n")
240 cat("R-squared:", R_squared_subset, "\n")
241 cat("Standard Error (slope):", SE_k1, "\n")
242 cat("Standard Error (intercept):", SE_k0, "\n")
243 cat("Z-value (slope):", t_k1, "\n\n")
244
245 # --- Plot subset regression ---
246 plot(c, z, main = "Temperature Over Time (1979 2024 )", xlab = "Year", ylab = "
    Temperature", pch = 19, col = "darkgray")
247 lines(c, z_pred, col = "purple", lwd = 2)
248 legend("topleft", legend = "Regression Line", col = "purple", lwd = 2)
249 # --- 95% Z-confidence interval (asymptotic) for slope from 1979 2024 ---
250 CI_lower_subset <- k1 - z_alpha * SE_k1
251 CI_upper_subset <- k1 + z_alpha * SE_k1
252
253 # Print confidence interval and test statistic
254 cat("Z-based 95% CI for slope (1979 2024 ):", CI_lower_subset, ",", CI_upper_
    subset, "]\n")
255 cat("Z-test statistic (slope 0):", (k1 - 0) / SE_k1, "\n")
256
257 ' '
258
259 **Bootstrap**
260
261 ' '{r}
262

```

```

263 # Split data into pre-1970 and post-1970
264 pre_data <- subset(data, Year < 1970)
265 post_data <- subset(data, Year >= 1970)
266
267 # Define temperature vectors
268 pre <- pre_data$DeBilt
269 post <- post_data$DeBilt
270
271 # Define sample sizes
272 n_pre <- length(pre)
273 n_post <- length(post)
274
275 # Means
276 mean_pre <- mean(pre)
277 mean_post <- mean(post)
278
279 # --- BOOTSTRAP SECTION ---
280 set.seed(123) # for reproducibility
281 B <- 1000 # number of bootstrap replications
282
283 # --- Bootstrap Confidence Intervals for the Mean ---
284 boot_mean <- function(data, B) {
285   n <- length(data)
286   replicate(B, mean(sample(data, n, replace = TRUE)))
287 }
288
289 boot_pre <- boot_mean(pre, B)
290 boot_post <- boot_mean(post, B)
291
292 # Confidence intervals
293 ci_boot_pre <- quantile(boot_pre, c(0.025, 0.975))
294 ci_boot_post <- quantile(boot_post, c(0.025, 0.975))
295
296 cat("Bootstrap 95% CI for pre-1970 mean:", round(ci_boot_pre, 3), "\n")
297 cat("Bootstrap 95% CI for post-1970 mean:", round(ci_boot_post, 3), "\n\n")
298
299 # --- Bootstrap Hypothesis Test (mean difference post - pre) ---
300 obs_diff <- mean_post - mean_pre
301 combined <- c(pre, post)
302
303 boot_diff <- replicate(B, {
304   sample_combined <- sample(combined, length(combined), replace = TRUE)
305   sample_post <- sample_combined[1:n_post]
306   sample_pre <- sample_combined[(n_post+1):(n_post + n_pre)]
307   mean(sample_post) - mean(sample_pre)
308 })
309
310 p_boot <- mean(boot_diff >= obs_diff)
311

```

```

312 cat("Bootstrap hypothesis test:\n")
313 cat("Observed mean difference (post-pre):", round(obs_diff, 3), "\n")
314 cat("Bootstrap p-value (one-sided):", format(p_boot, scientific = TRUE), "\n\n")
315
316 # --- Bootstrap for Regression Slope ---
317 x <- data$Year
318 y <- data$DeBilt
319 n_total <- length(y)
320
321 boot_regression <- function(x, y, B) {
322   replicate(B, {
323     idx <- sample(1:n_total, n_total, replace = TRUE)
324     coef(lm(y[idx] ~ x[idx]))[2]
325   })
326 }
327
328 boot_slope_all <- boot_regression(x, y, B)
329 ci_slope_all <- quantile(boot_slope_all, c(0.025, 0.975))
330
331 # Post-1979 regression
332 post1979_data <- subset(data, Year > 1979)
333 x_post <- post1979_data$Year
334 y_post <- post1979_data$DeBilt
335
336 boot_reg_post <- function(x, y, B) {
337   replicate(B, {
338     idx <- sample(1:length(y), length(y), replace = TRUE)
339     coef(lm(y[idx] ~ x[idx]))[2]
340   })
341 }
342
343 boot_slope_post <- boot_reg_post(x_post, y_post, B)
344 ci_slope_post <- quantile(boot_slope_post, c(0.025, 0.975))
345
346 cat("Bootstrap 95% CI for slope (1907-2024):", round(ci_slope_all, 5), "\n")
347 cat("Bootstrap 95% CI for slope (post-1979):", round(ci_slope_post, 5), "\n")
348
349 '''
350
351 **Monthly Data**
352
353 '''{r}
354 # Load and clean data
355 data <- read.csv2("/Users/User/OneDrive/E&OR_year2/Mathematical_Statistics/
    MonthlyTemp.csv")
356 colnames(data) <- trimws(colnames(data)) # Trim column names
357 data$Date <- as.character(data$Date) # Convert Date to character
358 data$Year <- as.integer(substr(data$Date, 1, 4)) # Get year
359 data$Month <- as.integer(substr(data$Date, 5, 6)) # Get month

```



```

360
361 # Label seasons and filter for Winter and Summer
362 data$Season <- ifelse(data$Month %in% c(1, 2, 3), "Winter", ifelse(data$Month %in%
363   c(6, 7, 8), "Summer", "Other"))
364 winter_summer_data <- data[data$Season %in% c("Winter", "Summer"), ]
365 winter_summer_data$Decade <- floor(winter_summer_data$Year / 10) * 10
366
367 # Prepare results table
368 results <- data.frame(Decade = integer(), Season = character(), Mean_Temperature =
369   numeric(),
370   CI_Lower = numeric(), CI_Upper = numeric(), Slope = numeric
371   (), Intercept = numeric(), RSS = numeric())
372
373 # Loop over seasons
374 for (season in c("Winter", "Summer")) {
375   season_data <- winter_summer_data[winter_summer_data$Season == season, ] #
376   Filter by season
377
378   # Loop over decades
379   for (decade in unique(season_data$Decade)) {
380     subset_data <- season_data[season_data$Decade == decade, ]
381     if (nrow(subset_data) < 2) next # Skip if too few observations
382
383     # Mean temp and 95% CI
384     mean_temp <- mean(subset_data$De.Bilt, na.rm = TRUE)
385     sd_temp <- sd(subset_data$De.Bilt, na.rm = TRUE)
386     n <- length(subset_data$De.Bilt)
387     se_temp <- sd_temp / sqrt(n)
388     t_value <- qt(0.975, df = n - 1)
389     ci_lower <- mean_temp - t_value * se_temp
390     ci_upper <- mean_temp + t_value * se_temp
391
392     results <- rbind(results, data.frame(Decade = decade, Season = season,
393       Mean_Temperature = mean_temp, CI_Lower =
394       ci_lower, CI_Upper = ci_upper,
395       Slope = NA, Intercept = NA, RSS = NA))
396
397     # Linear regression for current decade
398     season_data$Decade_Dummy <- ifelse(season_data$Decade == decade, 1, 0)
399     X <- season_data$Decade_Dummy
400     Y <- season_data$De.Bilt
401     X_bar <- mean(X, na.rm = TRUE)
402     Y_bar <- mean(Y, na.rm = TRUE)
403     numerator <- sum((X - X_bar) * (Y - Y_bar), na.rm = TRUE)
404     denominator <- sum((X - X_bar)^2, na.rm = TRUE)
405     beta1 <- numerator / denominator
406     beta0 <- Y_bar - beta1 * X_bar
407
408     # Save regression results

```

```

404 results$Slope[results$Decade == decade & results$Season == season] <- beta1
405 results$Intercept[results$Decade == decade & results$Season == season] <-
    beta0
406
407 # Residual sum of squares
408 Y_hat <- beta0 + beta1 * X
409 residuals <- Y - Y_hat
410 rss <- sum(residuals^2, na.rm = TRUE)
411 results$RSS[results$Decade == decade & results$Season == season] <- rss
412
413 # Output RSS
414 cat("Season:", season, "Decade:", decade, "RSS:", rss, "\n")
415 }
416 }
417
418 # Display results
419 print(results)
420
421 # Plot Winter mean temps with confidence intervals
422 plot(results$Decade[results$Season == "Winter"], results$Mean_Temperature[results$
    Season == "Winter"],
423       type = "o", col = "blue", xlab = "Decade", ylab = "Mean_Temperature_( C )",
424       main = "Winter_Mean_Temperatures_with_Confidence_Intervals",
425       ylim = range(results$CI_Lower, results$CI_Upper, na.rm = TRUE), pch = 16)
426 arrows(results$Decade[results$Season == "Winter"], results$CI_Lower[results$Season
    == "Winter"],
427         results$Decade[results$Season == "Winter"], results$CI_Upper[results$Season
    == "Winter"],
428         angle = 90, code = 3, length = 0.1, col = "blue")
429
430 # Plot Summer mean temps with confidence intervals
431 plot(results$Decade[results$Season == "Summer"], results$Mean_Temperature[results$
    Season == "Summer"],
432       type = "o", col = "red", xlab = "Decade", ylab = "Mean_Temperature_( C )",
433       main = "Summer_Mean_Temperatures_with_Confidence_Intervals",
434       ylim = range(results$CI_Lower, results$CI_Upper, na.rm = TRUE), pch = 16)
435 arrows(results$Decade[results$Season == "Summer"], results$CI_Lower[results$Season
    == "Summer"],
436         results$Decade[results$Season == "Summer"], results$CI_Upper[results$Season
    == "Summer"],
437         angle = 90, code = 3, length = 0.1, col = "red")
438
439 # Plot Winter data and regression line
440 plot(winter_summer_data$Year[winter_summer_data$Season == "Winter"],
441       winter_summer_data$De.Bilt[winter_summer_data$Season == "Winter"],
442       type = "p", col = rgb(0.6, 0.6, 1), pch = 16, xlab = "Year", ylab = "
    Temperature_( C )",
443       main = "Winter_Linear_Regression_Line",
444       ylim = range(winter_summer_data$De.Bilt, na.rm = TRUE))

```

```

445
446 # Add regression line and formula for Winter
447 winter_model <- lm(De.Bilt ~ Year, data = winter_summer_data[winter_summer_data$
    Season == "Winter", ])
448 abline(winter_model, col = rgb(0, 0, 0.5), lty = 1)
449 text(x = min(winter_summer_data$Year[winter_summer_data$Season == "Winter"]) + 10,
450      y = max(winter_summer_data$De.Bilt[winter_summer_data$Season == "Winter"]) -
    1,
451      labels = paste("Winter: Y=", round(coef(winter_model)[1], 2), "+", round(
    coef(winter_model)[2], 2), "X"),
452      col = "blue", cex = 0.8)
453
454 # Plot Summer data and regression line
455 plot(winter_summer_data$Year[winter_summer_data$Season == "Summer"],
456      winter_summer_data$De.Bilt[winter_summer_data$Season == "Summer"],
457      type = "p", col = rgb(1, 0.6, 0.6), pch = 16, xlab = "Year", ylab = "
    Temperature ( C )",
458      main = "Summer Linear Regression Line",
459      ylim = range(winter_summer_data$De.Bilt, na.rm = TRUE))
460
461 # Add regression line and formula for Summer
462 summer_model <- lm(De.Bilt ~ Year, data = winter_summer_data[winter_summer_data$
    Season == "Summer", ])
463 abline(summer_model, col = rgb(0.8, 0, 0), lty = 1)
464 text(x = min(winter_summer_data$Year[winter_summer_data$Season == "Summer"]) + 10,
465      y = max(winter_summer_data$De.Bilt[winter_summer_data$Season == "Summer"]) -
    1,
466      labels = paste("Summer: Y=", round(coef(summer_model)[1], 2), "+", round(
    coef(summer_model)[2], 2), "X"),
467      col = "red", cex = 0.8)
468
469 # Split data before and after 1970
470 pre_1970_winter <- winter_summer_data[winter_summer_data$Year >= 1900 & winter_
    summer_data$Year < 1970 & winter_summer_data$Season == "Winter", ]
471 post_1970_winter <- winter_summer_data[winter_summer_data$Year >= 1970 & winter_
    summer_data$Year <= 2020 & winter_summer_data$Season == "Winter", ]
472 pre_1970_summer <- winter_summer_data[winter_summer_data$Year >= 1900 & winter_
    summer_data$Year < 1970 & winter_summer_data$Season == "Summer", ]
473 post_1970_summer <- winter_summer_data[winter_summer_data$Year >= 1970 & winter_
    summer_data$Year <= 2020 & winter_summer_data$Season == "Summer", ]
474
475 # Z-test function
476 z_test <- function(pre_temp, post_temp) {
477   mean_post <- mean(post_temp)
478   sd_post <- sd(post_temp)
479   n_post <- length(post_temp)
480   mean_pre <- mean(pre_temp)
481   se_post <- sd_post / sqrt(n_post)
482   z_stat <- (mean_post - mean_pre) / se_post

```

```

483   p_value <- 1 - pnorm(z_stat)
484   return(list(z_stat = z_stat, p_value = p_value))
485 }
486
487 # Run z-tests and print results
488 winter_test <- z_test(pre_1970_winter$De.Bilt, post_1970_winter$De.Bilt)
489 cat("Winter- Z-statistic:", round(winter_test$z_stat, 3), "\n")
490 cat("Winter- P-value:", round(winter_test$p_value, 3), "\n")
491
492 summer_test <- z_test(pre_1970_summer$De.Bilt, post_1970_summer$De.Bilt)
493 cat("Summer- Z-statistic:", round(summer_test$z_stat, 3), "\n")
494 cat("Summer- P-value:", round(summer_test$p_value, 3), "\n")
495
496 # Interpret Winter result
497 if (winter_test$p_value < 0.05) {
498   cat("Winter- Significant increase after 1970.\n")
499 } else {
500   cat("Winter- No significant change after 1970.\n")
501 }
502
503 # Interpret Summer result
504 if (summer_test$p_value < 0.05) {
505   cat("Summer- Significant increase after 1970.\n")
506 } else {
507   cat("Summer- No significant change after 1970.\n")
508 }
509
510 ' ' '

```