

MIDDLE EAST TECHNICAL UNIVERSITY

FACULTY OF SCIENCE AND LITERATURE

STATISTICS

FACTORS AFFECTING HAPPINESS

ZEYNEP FENERCİOĞLU, BAŞAK UĞURLU, DAMLA BAŞARMIŞ

INTRODUCTION

In this project, Factors Affecting Happiness have been researched. Observed the variables that affect happiness in the dataset are the Healthcare Index, Air Quality Index, Green Space Area, Cost of Living, and Traffic Density. Firstly, some research questions were created about the relationship between the Happiness Score and variables. After that, some statistical tests were used. These are comparisons of means, one-way or two-way ANOVA and multiple comparisons, inferences about mean, inferences about proportion, comparisons of proportion, and simple and multiple linear regression. Performing Data analysis was done by visualizing these values in the graph. Appropriate tests were used to reach the conclusion of our research questions. However, there are some assumptions to perform tests. Checked whether these assumptions are appropriate. Therefore, the result has been reached.

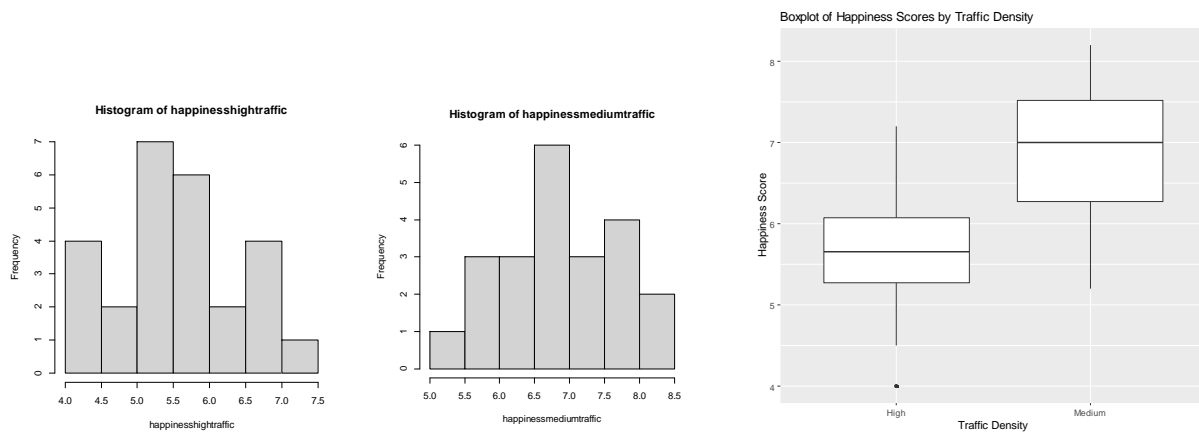
COMPARISONS OF MEANS (Two-sample hypothesis testing)

Is there a significant difference between means of happiness scores of cities with high traffic density and medium traffic density?

1) Performing exploratory data analysis:

There are 26 observations for happiness scores of high traffic density cities. The median is 5.65.

There are 22 observations for happiness scores of medium traffic density cities. The median is 7.



2) Assumptions:

Samples are independent there is no relationship between cities with medium traffic density and high traffic density.

We don't know the population standard deviations, we do not know if the happiness score population is normally distributed, and the sample sizes are smaller than 30. Therefore, we use the Mann-Whitney- Wilcoxon test since the requirement of normal distribution for the t-test is not met. Mean ranks approximate the median, so we are testing the median differences. We check whether the distributions are identical so that we can compare medians. As seen from the histograms, they have similar shapes so it is appropriate to use this test.

3) Performing the Mann-Whitney-Wilcoxon Test:

- Null hypothesis: The two groups (happiness scores of cities with medium and high traffic density) are sampled from populations with identical distributions.
- Alternative hypothesis: The two groups are sampled from populations with different distributions.

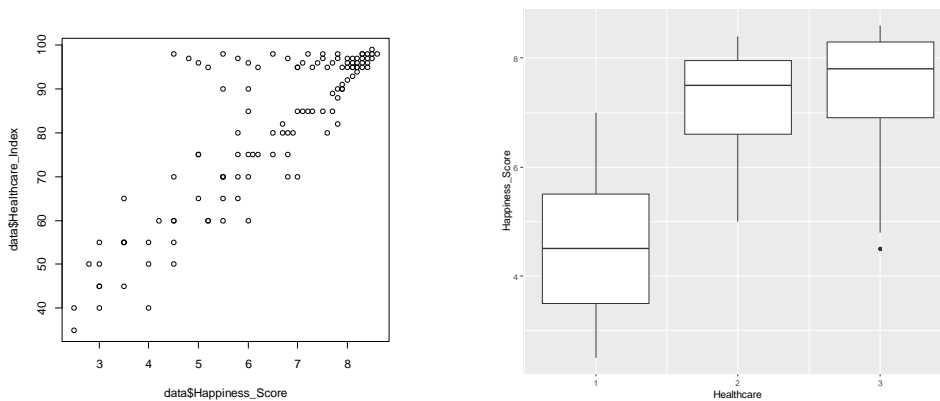
Calculated p value is 0.0000255, thus we reject the null hypothesis and we can say that happiness scores from cities with high and medium traffic density have different population distributions. Therefore we can say that their medians are significantly different.

1.1. ONE-WAY or TWO-WAY ANOVA AND MULTIPLE COMPARISONS

Do means of happiness scores differ depending on the healthcare index?

1) Performing exploratory data analysis:

There are 117 observations. The mean happiness score is 6.37265. The mean of the healthcare index is 80.59829. Since both are quantitative, we will group the healthcare index into 3 groups where each group has the same number of observations, 39. Below 75, 75-95, Above 95. In the scatter plot, we see that points are partially spread. Thus, we wanted to check if the differences are significant.



2) Assumptions:

Populations and cases within each sample are independent. They do not affect each other.

We do not know that populations are normally distributed. Each sample size is bigger than 30; we can use the central limit theorem and approximate to normal distribution. By Bartlett's test of homogeneity of variances, variances are equal, with the p-value being 0.1787 we fail to reject the null hypothesis that the samples have equal variances. We can use One-Way ANOVA.

3) Performing One-Way ANOVA:

- Null hypothesis: Population means of happiness scores from different healthcare index groups are equal.
- Alternative hypothesis: Population means of happiness scores from different healthcare index groups are not equal.

Calculated p value is $2e-16$. We reject null hypothesis and conclude that populations means are significantly different.

To see which ones are different, we used Tukey's procedure of multiple comparison:

For $H_0: 2 - 1 = 0$ p-value is $1e-04$, population means of the 1st and 2nd groups are significantly different. 2nd group has a higher population mean.

For $H_0: 3 - 1 = 0$ p-value is $1e-04$, population means of the 1st and 3rd groups are significantly different. 3rd group has a higher population mean.

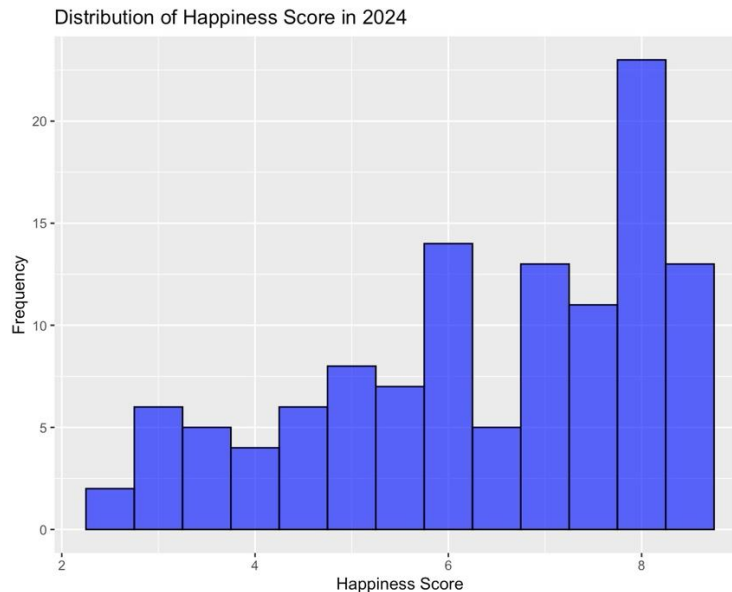
For $H_0: 3 - 2 = 0$ p-value is 0.759, and the population means of the 2nd and 3rd groups are not significantly different. 2nd and 3rd groups can have equal population means.

1.2. INFERENCES ABOUT MEAN (One sample hypothesis testing)

Is the average happiness score significantly different from the hypothesized value 5?

1) Performing Exploratory Data Analysis:

The mean happiness score is 6.37265. It is ranging from 2.5 to 8.6. It is quantitative.



2) Assumptions:

First, we need to check normality of happiness score. To check the normality, we use Shapiro test. Since its value is less than 0.05, normality assumption is not met for happiness score. We should use a non-parametric test instead of the one-sample t-test.

3) Performing Wilcoxon Signed-Rank Test:

- Null hypothesis: The median happiness score is equal to 5.
- Alternative hypothesis: The median happiness score is not equal to 5.

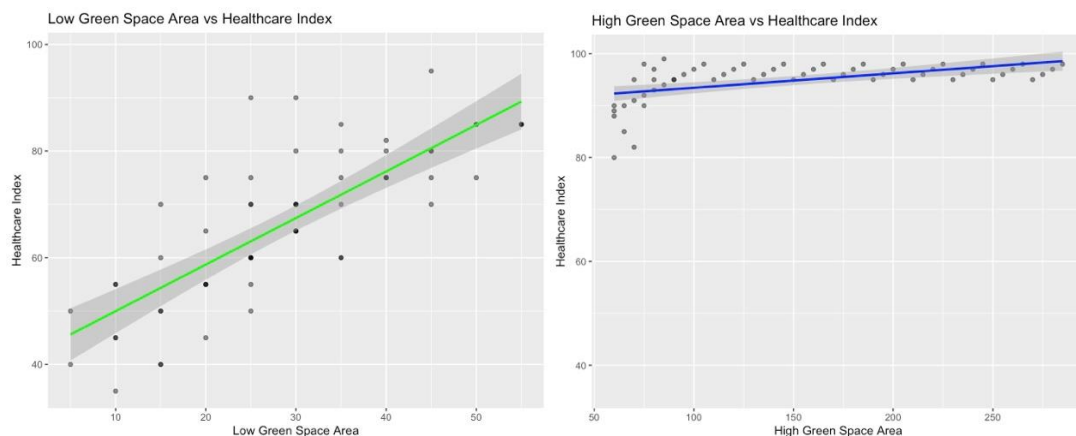
Calculated p value is 1.068e-11, thus we reject the null hypothesis and we can say that the median happiness score is significantly different from 5.

1.3. COMPARISONS OF PROPORTIONS (*Two-sample hypothesis testing*)

Is there a significant difference in the mean healthcare index between cities with high green space area and low green space area?

1) Performing Exploratory Data Analysis:

The mean of the healthcare index is 80.59829. We will split the data based on the median green space area. Median green space area is 55. Low green space area is below 55 and high green area is above 55. In regions with relatively lower green space, an increase in green space area correlates with an improvement in the healthcare index. This might suggest that even small increases in green space in areas with initially low green space can significantly impact healthcare outcomes. In regions with high green space, increases in green space area correlate with slight improvements in the healthcare index. This may suggest that while green space continues to positively impact healthcare outcomes, the marginal benefit may decrease as green space area becomes very large.



2) Assumptions:

We need to check normality of healthcare index for each group and equal variances between the two groups. To check normality for low and high green space with healthcare index we use Shapiro test. Since p value is smaller than 0.05, we can say that the normality assumptions are not met. For equal variances assumption, we use Levine test. Since p value is smaller than 0.05, we conclude that the variances are not equal. We should use a non-parametric test instead of the two-sample t-test. The Mann-Whitney U test (also known as the Wilcoxon rank-sum test) is appropriate for comparing two independent groups when the assumptions for the t-test are not met.

3) Performing Mann-Whitney U Test:

- Null hypothesis: The distribution of healthcare index is the same for high and low green space area groups.
- Alternative hypothesis: The distribution of healthcare index is different for high and low green space area groups.

Calculated p value is $2.2e-16$, thus we reject the null hypothesis. There is a significant difference in the healthcare index between high and low green space area groups.

INFERENCES ABOUT PROPORTION(One-Sample Hypothesis -Test)

Do countries with air quality above %60 degrees have higher happiness rates?

1)Performing exploratory data analysis:

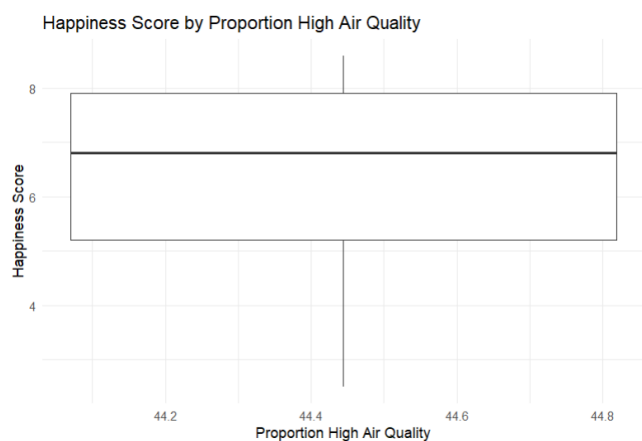
Proportion of cities with an air quality index below %60 is %55.55 in dataset.

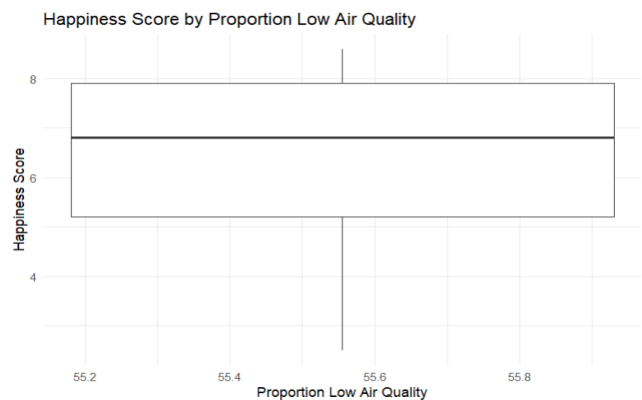
Proportion of cities with an air quality index above %60 is %44.44 in dataset.

Average happiness score of cities with an air quality index below %60 is 7.4969.

Average happiness score of cities with an air quality index above %60 is 4.967.

There are 52 observations that Air Quality Index is more than or equal to 60. There are 65 observations that Air Quality Index is less than 60





2) Assumptions:

Since each observation or sample is independent of the other, there is no relationship between cities with low air quality and cities with high air quality. Additionally, the data comes from a random sample. Since $np \geq 10$ and $n(1-p) \geq 10$, a normal approximation method can be used. But; We don't know how the happiness score is normally distributed. Also, Central Limit theorem is used if the sample size is more than 30 and our sample size is more than 30. Therefore, central limit theorem can be used. The Shapiro-Wilk Test is suitable for our data because it is quite sensitive for small and medium-sized samples.

Since standard deviations of the population is not known, the t-test can be used.

- Shapiro Wilk Test:

Null Hypothesis: The sample data follow a normal distribution

Alternative Hypothesis: The sample data does not follow a normal distribution.

The result of the Shapiro-Wilk Test indicates that W is approximately 0.84205. Also, the p -value is associated with 0.0000008124. Because of the p -value is less than 0.05, we reject the null hypothesis. We conclude that data is not normally distributed so it can be approached to normal distribution.

- t-test:

Null Hypothesis: Countries with air quality above %60 degrees have higher happiness rates

Alternative Hypothesis: Countries with air quality less %60 degrees have higher happiness rates

The result of the t test, p value is 1 so it is greater than significance level. We can say that fail to reject null hypothesis.

Result:

This dataset shows that more cities have an air quality index below 60 compared to cities above 60. Shapiro-Wilk test is used and the result of the Shapiro-Wilk Test W is approximately 0.84205 and

the p-value is associated with 0.0000008124. Because the p-value is less than 0.05, the null hypothesis is rejected. As a result, data is not normally distributed but we assume it is normally distributed. After that t-test is used and the result of the t-test, the p-value is 1 and fails to reject the null hypothesis. As a result, Countries with air quality above %60 degrees have higher happiness rates.

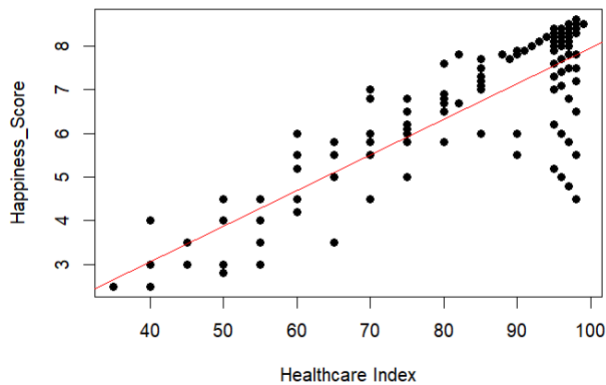
Conclusion:

Research question variables visualized by box plot. It was checked whether the conditions were met and the necessary tests were performed. As a result, Countries with air quality above 60% degrees have higher happiness rates.

Simple Linear Regression:

How does Healthcare Index affect Happiness Score?

Scatter Plot of Healthcare Index vs Happiness_Score



Call:
`lm(formula = Happiness_Score ~ Healthcare_Index, data = df)`

Coefficients:
 (Intercept) Healthcare_Index
 97.322 -1.527

$\beta_0 = 97.322$

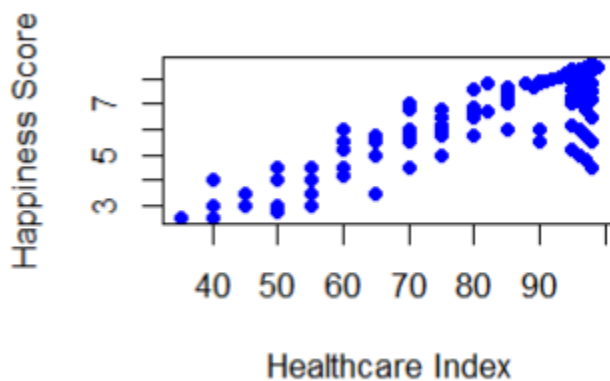
$\beta_1 = -1.527$

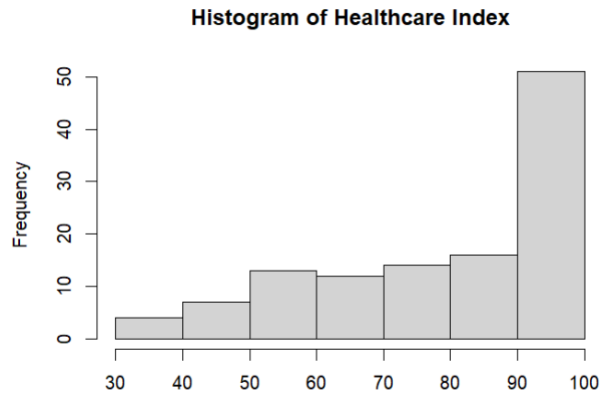
$\text{Happiness_Score} = 97.322 - 1.527 \cdot \text{Healthcare_Index}$

1) Perform exploratory data analysis

There are 117 variables for Healthcare index

There are 117 variables for Happiness Score





Histogram of Healthcare Index and Happiness score have left skewed distribution. Therefore, density is higher on the left side.

2) Assumptions:

-Linearity:

The relationship between the predictors and the response variable should be linear. We should create scatter plot and add regression line to visualize the linear relationship. If it close to 1, -1 we can see that it is strong relationship. If it close to 0, there is a weak relationship. Our correlation coefficient is 0.8470444, so there is a strong positive relationship.

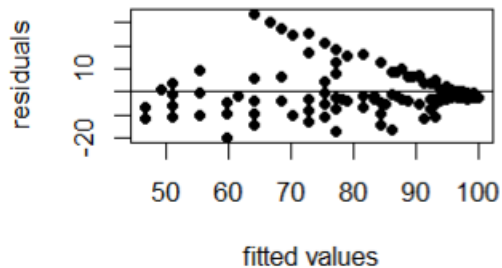
-Independence:

Tests such as the Durbin-Watson test can be used to evaluate this assumption.

Durbin-Watson Test

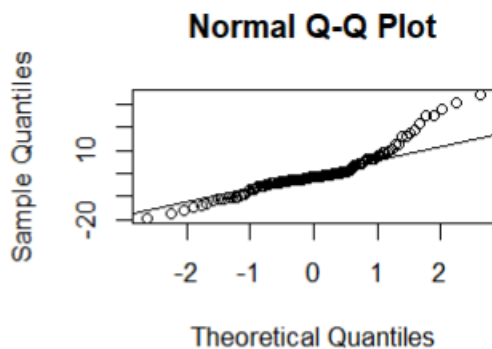
The statistics of Durbin Watson test is 0.6766644 and p value is 0.533. Because of the ($p \text{ value} > 0.5$), It is assumed that there is no autocorrelation between the terms. As a result, it can be said that there is no significant autocorrelation between the error terms in simple linear regression model and the independence assumption is met.

-Homoscedasticity:



The graph of homoscedasticity with a smoothly running graph may indicate that our model satisfies the homoscedasticity assumption. Representing homoscedasticity with a smoothly running graph may indicate that our model satisfies the homoscedasticity assumption. This situation shows that the error terms have a constant variance at the levels of the independent variables and the distribution of the errors does not change.

-Normality of Residuals:



Q-Q Plot

A straight line forms in the Q-Q graph, which indicates that the error periods comply with the normal distribution. A solid line indicates that the observed error terms fit well within a normal distribution. Thus, our regression model has a normal distribution.

Result:

Healthcare Index and Happiness Score is visualized by scatter plot and there is a positive linear relationship between them. Also, the equation is $\text{Happiness_Score} = 97.322 - 1.527 \cdot \text{Healthcare_Index}$. The intercept is 97.322 and $\beta_1 = -1.527$

Also, Healthcare index and Happiness Score is visualized by histogram. They have a left skewed distribution. They need to meet certain assumptions in order to examine their contents. These assumptions are linearity, independence, homoscedasticity, normality of residuals. The result of the linearity, correlation coefficient is 0.8470444 so there is a positive strong relationship. To check independence, Durbin Watson Test is used and the independence assumption is met. The graph of homoscedasticity is a smoothly running graph. It indicates that model satisfies the homoscedasticity assumption. Q-Q plot has a straight line forms so model has a normal distribution.

Conclusion:

We proved that it satisfies the assumptions for applying linear regression. As a result, equation is $\text{Happiness_Score} = 97.322 - 1.527 \cdot \text{Healthcare_Index}$. We made tests according to these assumptions. As a result, there is a strong positive relationship between Healthcare Index and Happiness Score.

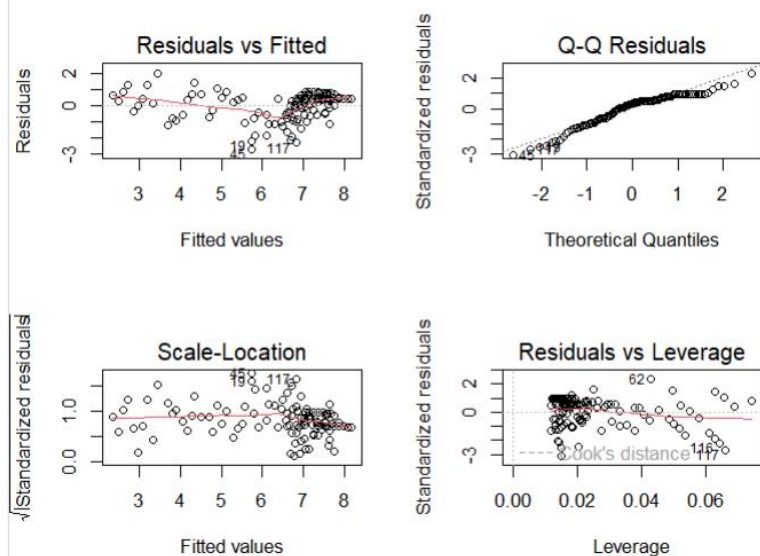
Multiple Linear Regression:

What Air Quality Index and Healthcare Index affect Happiness Score?

1) Perform exploratory data analysis:

The mean of the Air Quality Index is 76.58, median is 45, the max value is 245, and min value is 5.

The mean of the Green Space Area is 89.87, the median is 55, the max value is 285, and min value is 5.



```
call:
lm(formula = Happiness_Score ~ Air_Quality_Index + Green_Space_Area,
    data = df)
```

```
Residuals:
    Min       1Q   Median       3Q      Max
-2.7282 -0.5587  0.3051  0.6553  2.0518
```

```
Coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept)   8.883840   0.241746  36.749  < 2e-16 ***
Air_Quality_Index -0.026509   0.001685 -15.731  < 2e-16 ***
Green_Space_Area -0.005353   0.001351  -3.963  0.000129 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
Residual standard error: 0.8934 on 114 degrees of freedom
Multiple R-squared:  0.7332, Adjusted R-squared:  0.7285
F-statistic: 156.7 on 2 and 114 DF, p-value: < 2.2e-16
```

$\beta_0 = 8.883840$

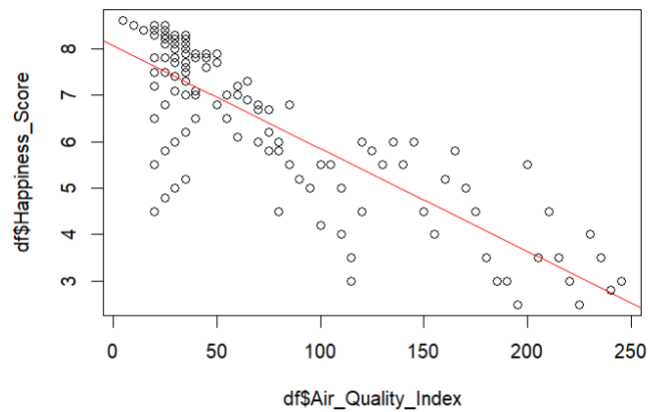
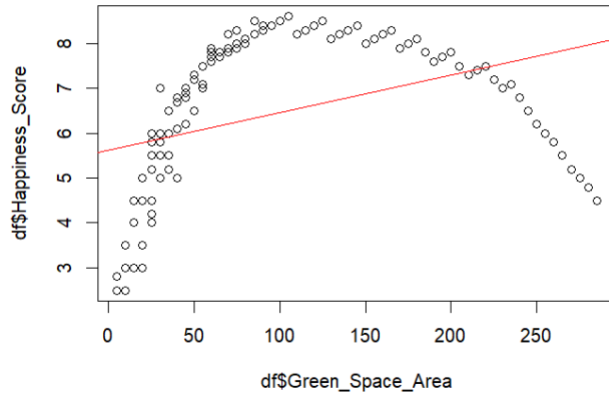
$\beta_1 = -0.026509$

$\beta_2 = -0.005353$

Equation is Happiness Score = $8.883840 - 0.026509(\text{Air Quality Index}) - 0.005353(\text{Green Space Area})$

2) Assumptions:

-Linearity:



The linearity assumption indicates that the relationship between the independent variables and the dependent variables is linear. There is a positive weak relationship between Green_Space Area and Happiness Score because the correlation coefficient is 0.3925. There is a strong negative relationship between the Air Quality Index and Happiness Score because the correlation coefficient is -0.8345

-Normality of Residuals:

The Q-Q plot is often used to visualize how well a normal distribution of error terms fits another distribution. The graph has a straight line and it shows that there is a normal distribution between the two residuals.

Shapiro-Wilk Test

As a result of the Shapiro-Wilk test, $W = 0.93209$ and $p\text{-value} = 1.647e-05$. W is close to 1 so the closer it is to a normal distribution and $p\text{-value}$ is smaller than 0.05 so null hypothesis is rejected. Thus, residuals are not

normality distributed

- Homoscedasticity:

The distribution graph is straight and a constant line forms. It shows that homoscedasticity is met. That is, predictions of this model have a similar variance for different independent variable values. Also, Having homoscedasticity increases the accuracy and robustness of the regression model.

-Multicollinearity

The variation inflation factor is 1.715366 . The VIF values are smaller than 5 so there is no multicollinearity problem.

Independence of Errors:

Durbin-Watson Test:

The Durbin-Watson test is often used to check for autocorrelation of errors in the regression model. As a result of the Durbin-Watson test, $DW = 0.8633241$. This indicates that there is a strong positive autocorrelation. The p-value is 0 and smaller than the general alpha value 0.05. It indicates that the null hypothesis is rejected.

Result:

We researched the effect of the Air Quality index and Green Space Area on happiness scores. The intercept of the coefficients is 8.883840 and the standard error is 0.241746. As a result the regression equation is $\text{Happiness Score} = 8.883840 - 0.026509(\text{Air Quality Index}) - 0.005353(\text{Green Space Area})$. Also, multiple R-squared is 0.7332 and adjusted R-squared 0.7285 so High values indicate that the model explains the variance of the dependent variable very well. There are some assumptions for multiple models. These are linearity, normality of residuals, homoscedasticity, multicollinearity, and independence of errors. Checked whether they were met or not. Firstly, The relationship between variables and happiness score was examined. There is a positive weak linear relationship between Green Space Area and Happiness Score because the correlation coefficient is 0.3925. There is a negative strong relationship between the air quality index because the correlation coefficient is -0.8345. After that Q-Q plot is visualized and the graph has a straight line so there is a normal distribution between the two residuals. Homoscedasticity is visualized and the distribution graph is straight and a constant line forms. It indicates that homoscedasticity is met. Multicollinearity is checked and the variation inflation factor is 1.715366 so there is moderate multicollinearity and there is no multicollinearity problem. To check the independence of errors Durbin Watson Test is used. The result of the test, the null hypothesis is rejected. Therefore, this situation violates the independence of errors assumption. However, we assume that the independence of errors assumption is satisfied.

Conclusion:

The necessary assumptions for multiple linear regression were examined and the necessary tests were applied. These tests were visualized with a Q-Q plot and homoscedasticity scatter plot. Green space area affects positively but air quality index affects negatively. Then, other assumptions are tested. Also, regression line equation is Happiness Score= $8.883840 - 0.026509(\text{Air Quality Index}) - 0.005353(\text{Green Space Area})$.