

Statistics Software Lab Report - 9

Name of the Student: Shatansh Patnaik
Roll No: 20MA20067

IIT Kharagpur
Statistics Software Lab

Linear Regression

Linear regression is a statistical method used to model the relationship between a dependent variable and one or more independent variables. It assumes a linear relationship between the variables, aiming to find the best-fitting line that predicts the dependent variable based on the independent variables.

Algorithm to Compute Coefficients

The following algorithm outlines the steps to compute the coefficients of a linear regression model using Ordinary Least Squares (OLS) method:

Algorithm 1 Compute Coefficients of Linear Regression Model

```
1: procedure OLS( $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n), n$ )
2:   Calculate the means:  $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$ ,  $\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$ 
3:   Calculate the slope (coefficient):  $\beta_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}$ 
4:   Calculate the intercept (coefficient):  $\beta_0 = \bar{y} - \beta_1 \bar{x}$ 
5:   return  $\beta_0, \beta_1$ 
6: end procedure
```

Analysis of the Algorithm

The algorithm employs the Ordinary Least Squares (OLS) method to estimate the coefficients of the linear regression model. It starts by calculating the means of the independent and dependent variables. Then, it computes the slope (coefficient) of the regression line by finding the ratio of the covariance of the variables to the variance of the independent variable. Finally, it calculates the intercept (coefficient) using the means and the slope. The resulting coefficients represent the parameters of the linear regression model.

Computation of Confidence Intervals of the Parameters of the Model

Confidence intervals for the parameters of a linear regression model provide a range of values within which we can be reasonably confident that the true value of the parameter lies. These intervals are essential for assessing the uncertainty associated with the estimated coefficients of the model. The computation involves estimating the standard error of the coefficients and then using critical values from the t-distribution to determine the bounds of the interval.

Algorithm for the Computation of Confidence Intervals of the Parameters of the Model

The following algorithm outlines the steps for computing confidence intervals for the parameters of a linear regression model:

Algorithm 2 Compute Confidence Intervals

```
1: procedure CONFIDENCEINTERVALS( $\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_n, SE(\hat{\beta}_0), SE(\hat{\beta}_1), \dots, SE(\hat{\beta}_n), df$ )
2:   Set desired confidence level (e.g., 95%)
3:   Find critical value  $t^*$  from t-distribution with  $df$  degrees of freedom
4:   for  $i \leftarrow 0$  to  $n$  do
5:     Calculate margin of error  $ME_i = t^* \times SE(\hat{\beta}_i)$ 
6:     Compute confidence interval  $[\hat{\beta}_i - ME_i, \hat{\beta}_i + ME_i]$ 
7:   end for
8:   return Confidence intervals for each parameter
9: end procedure
```

This algorithm provides a systematic approach to calculate confidence intervals for the parameters of a linear regression model, allowing researchers to make inferences about the true values of the coefficients with a specified level of confidence.

Significance Testing

Significance testing for the parameters of a linear regression model assesses whether the estimated coefficients are significantly different from zero. This helps determine the importance of each predictor variable in explaining the variation in the dependent variable. The computation involves comparing the t-statistic for each coefficient to a critical value from the t-distribution at a specified significance level (e.g., $\alpha = 0.05$).

Algorithm for performing significance testing

The following algorithm outlines the steps for conducting significance testing for the parameters of a linear regression model:

Algorithm 3 Significance Testing

```
1: procedure SIGNIFICANCETESTING( $\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_n, SE(\hat{\beta}_0), SE(\hat{\beta}_1), \dots, SE(\hat{\beta}_n), df, \alpha$ )
2:   for  $i \leftarrow 0$  to  $n$  do
3:     Calculate the t-statistic:  $t_i = \frac{\hat{\beta}_i}{SE(\hat{\beta}_i)}$ 
4:     Find critical value  $t^*$  from t-distribution with  $df$  degrees of freedom at significance level
       $\alpha$ 
5:     if  $|t_i| > t^*$  then
6:        $\hat{\beta}_i$  is statistically significant
7:     else
8:        $\hat{\beta}_i$  is not statistically significant
9:     end if
10:  end for
11: end procedure
```

This algorithm provides a systematic approach to conduct significance testing for the parameters of a linear regression model. It helps identify which predictors have a significant impact on the

dependent variable, allowing researchers to make informed decisions about the model's explanatory power.

Lack of Fitness Test

The lack of fitness test, also known as lack-of-fit test, assesses whether a simple linear regression model adequately fits the observed data. It is used to determine whether there is significant evidence that the relationship between the dependent and independent variables is not adequately described by a linear model. The test compares the variation explained by the regression model to the residual variation not explained by the model.

Algorithm for performing lack of fitness test

The following algorithm outlines the steps for conducting the lack of fitness test for a simple linear regression model:

Algorithm 4 Lack of Fitness Test

```

1: procedure LACKOFFITNESSTEST( $\hat{y}_1, \hat{y}_2, \dots, \hat{y}_n, y_1, y_2, \dots, y_n$ )
2:   Calculate the total sum of squares (TSS):  $TSS = \sum_{i=1}^n (y_i - \bar{y})^2$ 
3:   Calculate the regression sum of squares (RSS):  $RSS = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2$ 
4:   Calculate the residual sum of squares (ESS):  $ESS = \sum_{i=1}^n (y_i - \hat{y}_i)^2$ 
5:   Calculate the lack of fit sum of squares (LFSS):  $LFSS = ESS - RSS$ 
6:   Calculate the degrees of freedom for the lack of fit:  $df_{LF} = n - 2$ 
7:   Calculate the mean square for lack of fit:  $MS_{LF} = \frac{LFSS}{df_{LF}}$ 
8:   Calculate the mean square error:  $MSE = \frac{ESS}{n-1}$ 
9:   Perform an F-test using  $F = \frac{MS_{LF}}{MSE}$  and  $df_{LF}$  and  $n - 1$  degrees of freedom
10:  if F is significant then
11:    Lack of fit is significant; the linear model does not adequately fit the data
12:  else
13:    Lack of fit is not significant; the linear model adequately fits the data
14:  end if
15: end procedure

```

Conclusion

The lack of fitness test is an important diagnostic tool for evaluating the adequacy of a simple linear regression model. By comparing the variation explained by the model to the residual variation, researchers can determine whether the linear model adequately captures the relationship between the variables. If the lack of fit is significant, it indicates that the model does not adequately explain the data, and alternative modeling approaches may be necessary.

Correlation analysis

Covariance analysis for a linear regression model assesses the relationships between variables and their contributions to predicting the dependent variable. It involves analyzing the covariance matrix

to understand the relationships between predictors and to identify multicollinearity, which can affect the model's stability and interpretation.

The following algorithm outlines the steps for covariance analysis for a linear regression model:

Algorithm 5 Covariance Analysis for Linear Regression Model

```
1: procedure COVARIANCEANALYSIS( $X_1, X_2, \dots, X_n$ )
2:   Compute covariance matrix  $Cov(X)$  from the predictor variables
3:   Compute correlation matrix  $Corr(X)$  from the covariance matrix
4:   for each pair of predictors  $(X_i, X_j)$  do
5:     if  $|Corr(X_i, X_j)| > \text{threshold}$  then
6:       Identify multicollinearity
7:     end if
8:   end for
9:   Analyze eigenvalues of  $Corr(X)$  to determine magnitude of multicollinearity
10:  Perform diagnostic tests and address multicollinearity if present
11: end procedure
```

Covariance analysis is an essential step in assessing the quality of a linear regression model. By examining the relationships between predictors and identifying multicollinearity, researchers can ensure the stability and reliability of their regression analysis.

Plotting of Scatter Plots for the Given Datasets

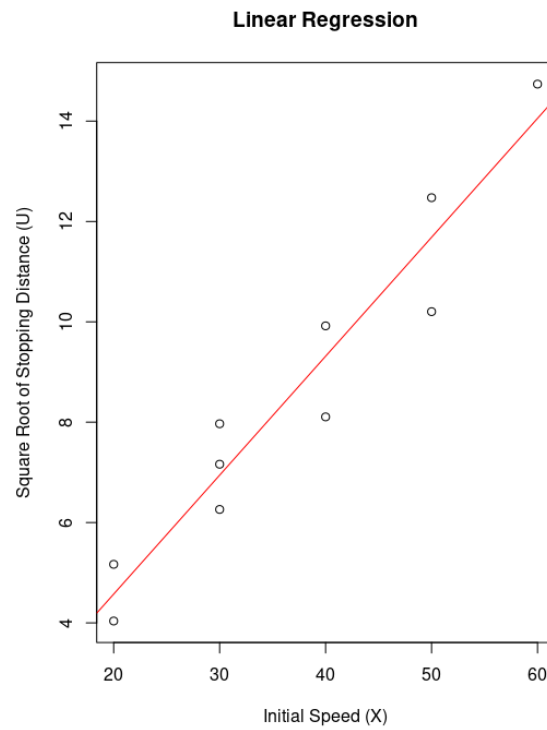


Figure 1: Scatter Plot for Y vs X

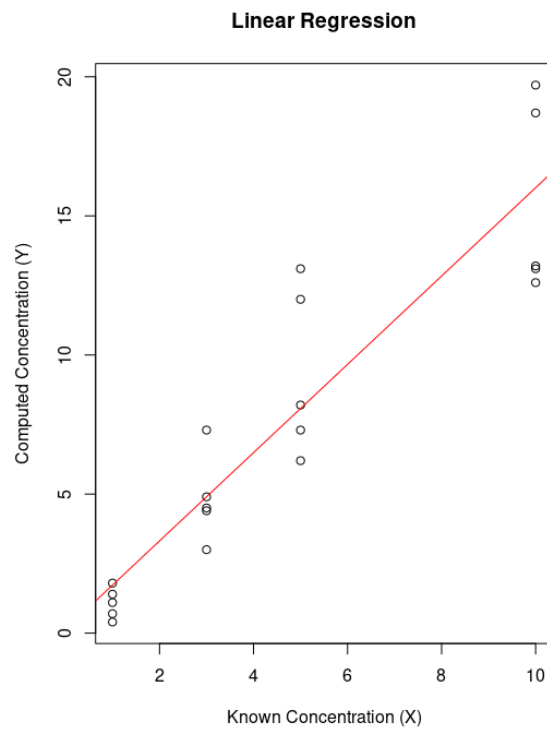
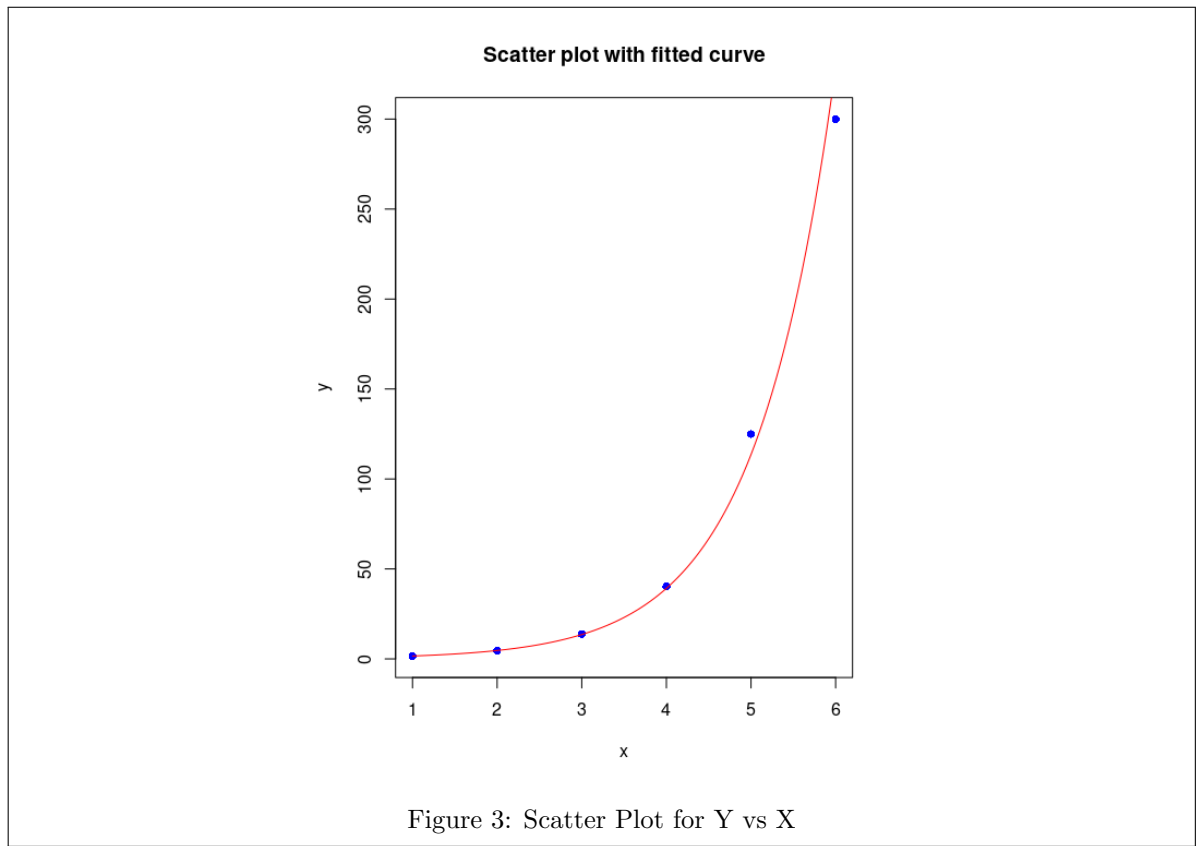


Figure 2: Scatter Plot for Y vs X



Implementation of the above algorithms using R

The following code demonstrates the implementation of computation of parameters of the model, computation of the respective confidence and prediction interval, significance testing and confidence interval testing.

Code for Exercise-1

```

1  # Exercise - 1
2  X <- c(20, 20, 30, 30, 30, 40, 40, 50, 50, 60)
3  y <- c(16.3, 26.7, 39.2, 63.5, 51.3, 98.4, 65.7, 104.1, 155.6, 217.2)
4
5  # a) Scatter plot
6  plot(X, y, xlab = "X", ylab = "y", main = "Scatter plot")
7
8  # b) Fitting a linear regression model
9  U <- sqrt(y)
10 meanX <- sum(X) / length(X)
11 meanU <- sum(U) / length(U)
12

```



```

13 beta1 <- sum((X - meanX) * (U - meanU)) / sum((X - meanX)^2)
14 beta0 <- meanU - beta1 * meanX
15 UPredicted <- beta0 + beta1 * X
16 cat("Intercept (beta0):", beta0, "\n")
17 cat("Slope (beta1):", beta1, "\n")
18
19 plot(X, U, xlab = "Initial Speed (X)", ylab = "Square Root of Stopping
    Distance (U)", main = "Linear Regression")
20 abline(beta0, beta1, col = "red")
21
22 # c) Confidence Intervals for beta0, beta1 and sigma squared
23 n <- length(X)
24 df <- n - 2
25 RSS <- sum((U - UPredicted)^2)
26 sigma2 <- RSS / df
27 SESlope <- sqrt(sigma2 / sum((X - meanX)^2))
28 SEIntercept <- sqrt(sigma2 * (1 / n + meanX^2 / sum((X - meanX)^2)))
29
30 t <- qt(0.975, df)
31 CIForSlope <- c(beta1 - t * SESlope, beta1 + t * SESlope)
32 CIForIntercept <- c(beta0 - t * SEIntercept, beta0 + t * SEIntercept)
33
34 chiLower <- qchisq(0.025, df)
35 chiUpper <- qchisq(0.975, df)
36
37 CISigma2 <- c((df * sigma2) / chiUpper, (df * sigma2) / chiLower)
38
39 cat("95% Confidence Interval for Slope (beta1):", CIForSlope, "\n")
40 cat("95% Confidence Interval for Intercept (beta0):", CIForIntercept, "\n")
41 cat("95% Confidence Interval for Sigma^2:", CISigma2, "\n")
42
43 # d) Test of Significance for beta1 and beta0
44 # For Beta0:
45 alpha <- 0.05
46 SSE <- sum((U - UPredicted)^2)
47 StandardErrorBeta0 <- sqrt((SSE / df) * ((1 / n) + (meanX^2 / sum((X - meanX
    )^2))))
48 tObserved <- beta0 / StandardErrorBeta0
49 pValue <- 2 * pt(abs(tObserved), n-2, lower.tail = FALSE)
50
51 if (pValue < alpha) {
52   cat("Reject null hypothesis: Beta0 is significant.\n")
53 } else {
54   cat("Accept the null hypothesis: Beta0 is not significant.\n")
55 }
56
57 # For Beta1
58 alpha <- 0.05
59 df <- n-2
60 SSE <- sum((U - UPredicted)^2)
61 StandardErrorBeta1 <- sqrt((SSE / df) / sum((X - meanX)^2))
62 tObserved <- beta1 / StandardErrorBeta1

```

```

63 pValue <- 2 * pt(abs(tObserved), n-2, lower.tail = FALSE)
64
65 if (pValue < alpha) {
66   cat("Reject null hypothesis: Beta1 is significant.\n")
67 } else {
68   cat("Accept the null hypothesis: Beta1 is not significant.\n")
69 }
70
71 # e) Confidence Interval for mean of y given X
72 X0 <- 35
73 UPredictedForX0 <- beta0 + beta1 * X0
74 SEPred <- sqrt(sigma2 * ( 1 / n + (X0 - meanX)^2 / sum((X - meanX)^2)))
75 df <- length(X) - 2
76 t <- qt(0.975, df)
77
78 CIForU <- c(UPredictedForX0 - t * SEPred, UPredictedForX0 + t * SEPred)
79 CIForY <- CIForU^2
80 cat("95% Confidence Interval for the expected stopping distance when the
      initial speed is 35:", CIForY, "\n")
81
82 # f) Prediction Interval for a new observation of y given X
83 X0 <- 35
84 UPredictedForX0 <- beta0 + beta1 * X0
85 SEPred <- sqrt(sigma2 * (1 + 1 / n + (X0 - meanX)^2 / sum((X - meanX)^2)))
86 df <- length(X) - 2
87 t <- qt(0.975, df)
88 PIForU <- c(UPredictedForX0 - t * SEPred, UPredictedForX0 + t * SEPred)
89
90 PIForY <- PIForU^2
91 cat("95% Confidence Interval for the expected stopping distance when the
      initial speed is 35:", PIForY, "\n")
92
93 # g) Lack of fit analysis
94 n <- length(X)
95 m <- length(unique(X))
96 SSE <- sum((U - beta0 - beta1*X)^2)
97
98 SSPE <- 0
99 for(i in 1:m){
100   SSPE <- SSPE + sum((sqrt(y[X==unique(X)[i]]) - mean(sqrt(y[X==unique(X)[i]
      ]))))^2)
101 }
102
103 SSLOF <- SSE - SSPE
104
105 fStat <- (SSLOF/(m-2))/(SSPE/(n-m))
106 fCritical <- qf(0.95,m-2,n-m)
107
108 if(fStat > fCritical){
109   cat("The model caught lacking","\n")
110 }else{
111   cat("The model doesnt not lack in fitting the data","\n")

```

Code for Exercise-2

```

1  # Exercise - 2
2  n <- 20
3  X <- c(1, 1, 1, 1, 1, 3, 3, 3, 3, 3, 5, 5, 5, 5, 5, 10, 10, 10, 10, 10)
4  y <- c(1.1, 0.7, 1.8, 0.4, 1.4, 3.0, 4.5, 4.9, 4.4, 7.3, 7.3, 8.2, 6.2,
5      13.1, 12.0, 12.6, 13.2, 13.1, 18.7, 19.7)
6  plot(X, y, xlab = "X", ylab = "y", main = "Scatter plot")
7
8  # Fitting a linear regression model
9  meanX <- sum(X) / length(X)
10 meanY <- sum(y) / length(y)
11
12 beta1 <- sum((X - meanX) * (y - meanY)) / sum((X - meanX)^2)
13 beta0 <- meanY - beta1 * meanX
14 YPredicted <- beta0 + beta1 * X
15 cat("Intercept (beta0):", beta0, "\n")
16 cat("Slope (beta1):", beta1, "\n")
17
18 plot(X, y, xlab = "Known Concentration (X)", ylab = "Computed
19      Concentration (Y)", main = "Linear Regression")
20 abline(beta0, beta1, col = "red")
21
22 # Lack of fit test
23 m <- length(unique(X))
24 SSE <- sum((y - beta0 - beta1*X)^2)
25
26 SSPE <- 0
27 for(i in 1:m){
28   SSPE <- SSPE + sum(((y[X==unique(X)[i]]) - mean((y[X==unique(X)[i]])))
29     ^2)
30 }
31
32 SSLOF <- SSE - SSPE
33
34 fStat <- (SSLOF/(m-2))/(SSPE/(n-m))
35 fCritical <- qf(0.95,m-2,n-m)
36
37 if(fStat > fCritical){
38   cat("The model caught lacking","\n")
39 }else{
40   cat("The model doesnt not lack in fitting the data","\n")
41 }
42
43 # Correlation Coefficient
44 correlation <- (n * sum(X * y) - sum(X) * sum(y)) / sqrt((n * sum(X^2) -
45   sum(X)^2) * (n * sum(y^2) - sum(y)^2))
46 cat("Correlation Coefficient:", correlation, "\n")

```

```

44 # Test of significance for beta1 and beta0
45 # For Beta0:
46 alpha <- 0.05
47 df <- n-2
48 SSE <- sum((y - YPredicted)^2)
49 StandardErrorBeta0 <- sqrt((SSE / df) * ((1 / n) + (meanX^2 / sum((X -
50   meanX)^2))))
51 tObserved <- beta0 / StandardErrorBeta0
52 pValue <- 2 * pt(abs(tObserved), n-2, lower.tail = FALSE)
53
54 if (pValue < alpha) {
55   cat("Reject null hypothesis: Beta0 is significant.\n")
56 } else {
57   cat("Accept the null hypothesis: Beta0 is not significant.\n")
58 }
59
60 # For Beta1
61 alpha <- 0.05
62 SSE <- sum((y - YPredicted)^2)
63 StandardErrorBeta1 <- sqrt((SSE / df) / sum((X - meanX)^2))
64 tObserved <- beta1 / StandardErrorBeta1
65 pValue <- 2 * pt(abs(tObserved), n-2, lower.tail = FALSE)
66
67 if (pValue < alpha) {
68   cat("Reject null hypothesis: Beta1 is significant.\n")
69 } else {
70   cat("Accept the null hypothesis: Beta1 is not significant.\n")
71 }
72
73 # Confidence Intervals for all the parameters
74 sigma2 <- SSE / df
75 SESlope <- sqrt(sigma2 / sum((X - meanX)^2))
76 SEIntercept <- sqrt(sigma2 * (1 / n + meanX^2 / sum((X - meanX)^2)))
77
78 t <- qt(0.975, df)
79 CIForSlope <- c(beta1 - t * SESlope, beta1 + t * SESlope)
80 CIForIntercept <- c(beta0 - t * SEIntercept, beta0 + t * SEIntercept)
81
82 chiLower <- qchisq(0.025, df)
83 chiUpper <- qchisq(0.975, df)
84
85 CISigma2 <- c((df * sigma2) / chiUpper, (df * sigma2) / chiLower)
86
87 cat("95% Confidence Interval for Slope (beta1):", CIForSlope, "\n")
88 cat("95% Confidence Interval for Intercept (beta0):", CIForIntercept, "\n")
89 cat("95% Confidence Interval for Sigma^2:", CISigma2, "\n")

```

Code for Exercise-4

```

1 # Exercise - 4

```

```

2 X <- c(1, 2, 3, 4, 5, 6)
3 y <- c(1.60, 4.50, 13.80, 40.20, 125.00, 300.00)
4 plot(X, y, pch = 16, col = "blue", xlab = "x", ylab = "y", main = "Scatter
   plot with fitted curve")
5 n <- length(X)
6
7 Y <- log(y)
8 # We note that beta0 is log a and b is beta1
9 meanX <- sum(X) / length(X)
10 meanY <- sum(Y) / length(Y)
11
12 beta1 <- sum((X - meanX) * (Y - meanY)) / sum((X - meanX)^2)
13 beta0 <- meanY - beta1 * meanX
14 a <- exp(beta0)
15
16 yPredicted <- beta0 + beta1 * X
17
18 # Lack of Fitness Test
19 SSLF <- sum((Y - yPredicted)^2)
20 SSE <- sum((Y - meanY)^2)
21
22 df1 <- n - 2
23 df2 <- n - 3
24 F_statistic <- (SSLF / df1) / (SSE / df2)
25
26 alpha <- 0.05
27 critical_value <- qf(1 - alpha, df1, df2)
28
29 cat("F-statistic:", F_statistic, "\n")
30 cat("Critical value:", critical_value, "\n")
31
32 if (F_statistic > critical_value) {
33   cat("Reject null hypothesis: Lack of fit is significant.\n")
34 } else {
35   cat("Accept the null hypothesis: Lack of fit is not significant.\n")
36 }
37
38 # Plotting
39 curve(a * exp(beta1 * x), from = min(X), to = max(X), col = "red", add =
   TRUE)
40
41 # Significance Testing for Parameters
42 # For Beta0:
43 alpha <- 0.05
44 SSE <- sum((Y - yPredicted)^2)
45 StandardErrorBeta0 <- sqrt((SSE / df2) * ((1 / n) + (meanX^2 / sum((X -
   meanX)^2))))
46 tObserved <- beta0 / StandardErrorBeta0
47 pValue <- 2 * pt(abs(tObserved), n-2, lower.tail = FALSE)
48
49 if (pValue < alpha) {
50   cat("Reject null hypothesis: Beta0 is significant.\n")

```

```

51 } else {
52   cat("Accept the null hypothesis: Beta0 is not significant.\n")
53 }
54
55 # For Beta1
56 alpha <- 0.05
57 df <- n-2
58 SSE <- sum((Y - yPredicted)^2)
59 StandardErrorBeta1 <- sqrt((SSE / df) / sum((X - meanX)^2))
60 tObserved <- beta1 / StandardErrorBeta1
61 pValue <- 2 * pt(abs(tObserved), n-2, lower.tail = FALSE)
62
63 if (pValue < alpha) {
64   cat("Reject null hypothesis: Beta1 is significant.\n")
65 } else {
66   cat("Accept the null hypothesis: Beta1 is not significant.\n")
67 }

```

Code for Exercise-5

```

1  # Exercise - 5
2  X <- c(2, 3, 4, 5, 6)
3  y <- c(144.0, 172.80, 207.40, 248.50, 298.50)
4  plot(X, y, pch = 16, col = "blue", xlab = "X", ylab = "Y", main = "Scatter
   plot with fitted curve")
5
6  Y <- log(y)
7  n <- length(X)
8
9
10 # We note that beta0 is log a and b is exp(beta1)
11 meanX <- sum(X) / length(X)
12 meanY <- sum(Y) / length(Y)
13
14 beta1 <- sum((X - meanX) * (Y - meanY)) / sum((X - meanX)^2)
15 beta0 <- meanY - beta1 * meanX
16 a <- exp(beta0)
17 b <- exp(beta1)
18
19 yPredicted <- beta0 + beta1 * X
20
21 # Lack of Fitness Test
22 SSLF <- sum((Y - yPredicted)^2)
23 SSE <- sum((Y - meanY)^2)
24
25 df1 <- n - 2
26 df2 <- n - 3
27 F_statistic <- (SSLF / df1) / (SSE / df2)
28
29 alpha <- 0.05
30 critical_value <- qf(1 - alpha, df1, df2)

```

```

31
32 cat("F-statistic:", F_statistic, "\n")
33 cat("Critical value:", critical_value, "\n")
34
35 if (F_statistic > critical_value) {
36   cat("Reject null hypothesis: Lack of fit is significant.\n")
37 } else {
38   cat("Accept the null hypothesis: Lack of fit is not significant.\n")
39 }
40
41 # Significance Testing for Parameters
42 # For Beta0:
43 alpha <- 0.05
44 SSE <- sum((Y - yPredicted)^2)
45 StandardErrorBeta0 <- sqrt((SSE / df1) * ((1 / n) + (meanX^2 / sum((X -
46   meanX)^2))))
47 tObserved <- beta0 / StandardErrorBeta0
48 pValue <- 2 * pt(abs(tObserved), n-2, lower.tail = FALSE)
49
50 if (pValue < alpha) {
51   cat("Reject null hypothesis: Beta0 is significant.\n")
52 } else {
53   cat("Accept the null hypothesis: Beta0 is not significant.\n")
54 }
55
56 # For Beta1
57 alpha <- 0.05
58 SSE <- sum((Y - yPredicted)^2)
59 StandardErrorBeta1 <- sqrt((SSE / df1) / sum((X - meanX)^2))
60 tObserved <- beta1 / StandardErrorBeta1
61 pValue <- 2 * pt(abs(tObserved), n-2, lower.tail = FALSE)
62
63 if (pValue < alpha) {
64   cat("Reject null hypothesis: Beta1 is significant.\n")
65 } else {
66   cat("Accept the null hypothesis: Beta1 is not significant.\n")
67 }

```

Polynomial Regression

Polynomial regression is a type of regression analysis in which the relationship between the independent variable x and the dependent variable y is modeled as an n th degree polynomial. It is used when the relationship between the variables is nonlinear, and a straight line cannot adequately capture the relationship.

Algorithm 6 Polynomial Regression Algorithm

```
1: Input: Data points  $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$  and degree of polynomial  $d$ 
2: Output: Coefficients  $a_0, a_1, \dots, a_d$  of the polynomial regression model  $y = a_0 + a_1x + a_2x^2 + \dots + a_dx^d$ 
3: Initialize matrix  $X$  with dimensions  $n \times (d + 1)$ 
4: for  $i \leftarrow 1$  to  $n$  do
5:   for  $j \leftarrow 0$  to  $d$  do
6:      $X[i, j] \leftarrow x_i^j$ 
7:   end for
8: end for
9: Compute the transpose of  $X$ :  $X^T$ 
10: Compute the product of  $X^T$  and  $X$ :  $X^T X$ 
11: Compute the inverse of  $X^T X$ :  $(X^T X)^{-1}$ 
12: Compute the product of  $(X^T X)^{-1}$  and  $X^T$ :  $(X^T X)^{-1} X^T$ 
13: Compute the product of  $(X^T X)^{-1} X^T$  and  $y$ :  $(X^T X)^{-1} X^T y$ 
14: return Coefficients  $a_0, a_1, \dots, a_d$ 
```

Algorithm for the Computation of Polynomial Regression Coefficients

Analysis of the Algorithm

The algorithm for computing the coefficients of polynomial regression follows the least squares method. It involves constructing a design matrix X where each row corresponds to a data point and each column corresponds to a power of the independent variable x . Then, the coefficients are computed by performing matrix operations involving the transpose and inverse of the design matrix.

Algorithm for the Computing Confidence Intervals for Polynomial Regression Parameters

Algorithm 7 Confidence Intervals for Polynomial Regression Parameters

```
1: Input: Data points  $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$ , polynomial regression coefficients  $a_0, a_1, \dots, a_d$ , significance level  $\alpha$ 
2: Output: Confidence intervals for the polynomial regression coefficients  $a_0, a_1, \dots, a_d$ 
3: Compute the residuals:  $e_i = y_i - (a_0 + a_1x_i + a_2x_i^2 + \dots + a_dx_i^d)$ 
4: Compute the mean squared error:  $MSE = \frac{1}{n-d-1} \sum_{i=1}^n e_i^2$ 
5: Compute the standard error of the coefficients:  $SE(a_j) = \sqrt{\frac{MSE}{\sum_{i=1}^n (x_i - \bar{x})^2(j+1)}}$ , where  $\bar{x}$  is the mean of the  $x$  values
6: Compute the t-statistic for the given significance level:  $t_{\alpha/2, n-d-1}$ 
7: Compute the confidence interval for each coefficient:  $(a_j - t_{\alpha/2, n-d-1} \times SE(a_j), a_j + t_{\alpha/2, n-d-1} \times SE(a_j))$ 
8: return Confidence intervals for the polynomial regression coefficients  $a_0, a_1, \dots, a_d$ 
```

Analysis of the Algorithm

The algorithm for computing confidence intervals for parameters of a polynomial regression model is based on the standard error of the coefficients and the t-statistic. It involves computing the residuals, mean squared error, standard error of the coefficients, and the t-statistic for the given significance level. The confidence intervals are then calculated based on the standard error and t-statistic.

Algorithm: Lack of Fitness Test for Polynomial Regression

Algorithm 8 Lack of Fitness Test

```
1: Input: Residuals  $e_i$  from the polynomial regression model, degrees of freedom  $df_1$  and  $df_2$ 
2: Output: Result of the lack of fitness test
3: Compute the sum of squared residuals:  $SSR = \sum_{i=1}^n e_i^2$ 
4: Compute the mean squared error:  $MSE = \frac{SSR}{df_2}$ 
5: Compute the lack of fitness statistic:  $F = \frac{MSE}{\text{Mean Squared Error of the Regression}}$ 
6: Compute the critical value from the F-distribution:  $F_{\alpha, df_1, df_2}$ 
7: if  $F > F_{\alpha, df_1, df_2}$  then
8:   return Reject the null hypothesis: Lack of fit is significant
9: else
10:  return Accept the null hypothesis: Lack of fit is not significant
11: end if
```

Analysis of Lack of Fitness Test

The lack of fitness test in polynomial regression is used to determine whether the polynomial regression model adequately fits the data. It involves computing the lack of fitness statistic, which compares the mean squared error obtained from the model to the mean squared error of the regression. The test then compares this statistic to a critical value from the F-distribution at a given significance level. If the lack of fitness statistic exceeds the critical value, the null hypothesis is rejected, indicating that the lack of fit is significant. Otherwise, if the lack of fitness statistic does not exceed the critical value, the null hypothesis is accepted, indicating that the lack of fit is not significant.

Algorithm for Significance Testing for Polynomial Regression

Algorithm 9 Significance Testing

```
1: Input: Coefficients  $a_0, a_1, \dots, a_d$  from the polynomial regression model, standard errors  
    $SE(a_0), SE(a_1), \dots, SE(a_d)$ , degrees of freedom  $df$   
2: Output: Results of significance testing for the polynomial regression coefficients  
3: for  $j \leftarrow 0$  to  $d$  do  
4:   Compute the t-statistic:  $t_j = \frac{a_j}{SE(a_j)}$   
5:   Compute the critical value from the t-distribution:  $t_{\alpha/2, df}$   
6:   if  $|t_j| > t_{\alpha/2, df}$  then  
7:     return Reject the null hypothesis:  $a_j$  is significant  
8:   else  
9:     return Accept the null hypothesis:  $a_j$  is not significant  
10:  end if  
11: end for
```

Analysis of Significance Testing

The significance testing for polynomial regression coefficients is used to determine whether each coefficient in the polynomial regression model is statistically significant. It involves computing the t-statistic for each coefficient by dividing the coefficient by its standard error. The test then compares each t-statistic to a critical value from the t-distribution at a given significance level. If the absolute value of the t-statistic exceeds the critical value, the null hypothesis is rejected, indicating that the coefficient is significant. Otherwise, if the absolute value of the t-statistic does not exceed the critical value, the null hypothesis is accepted, indicating that the coefficient is not significant.

Scatter Plot for the Given Dataset

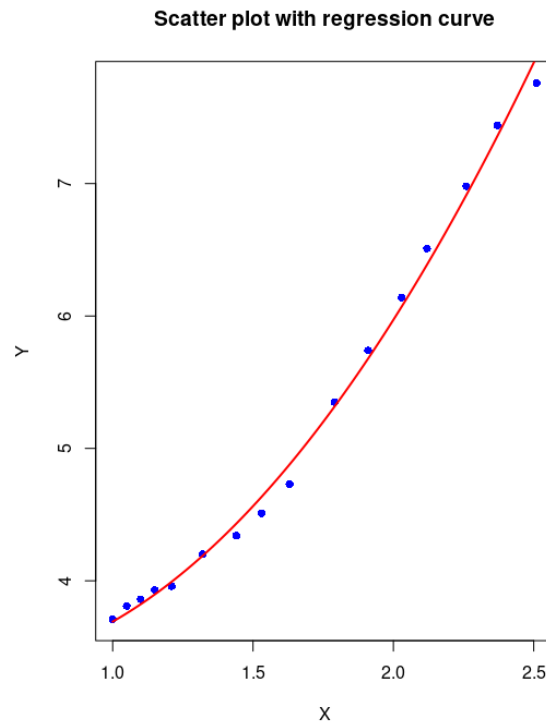


Figure 4: Scatter Plot for Y vs X

Implementation of the above algorithms using R

The following code demonstrates the implementation of computation of parameters of the model, computation of the respective confidence and prediction interval, significance testing and confidence interval testing.

Code for Exercise-3

```
1 # Exercise - 3
2 X <- c(1.00, 1.05, 1.10, 1.15, 1.21, 1.32, 1.44, 1.53, 1.63, 1.79, 1.91,
3       2.03, 2.12, 2.26, 2.37, 2.51)
4 y <- c(3.71, 3.81, 3.86, 3.93, 3.96, 4.20, 4.34, 4.51, 4.73, 5.35, 5.74,
5       6.14, 6.51, 6.98, 7.44, 7.76)
6
7 # a) Fitting a second degree polynomial regression model
8 XSquared <- X^2
9 designMatrix <- cbind(rep(1, length(X)), X, XSquared)
10 productMatrix <- (t(designMatrix)) %*% designMatrix
```

```

10 productY <- (t(designMatrix)) %*% y
11 coefficients <- solve(productMatrix, productY)
12 yPredicted <- coefficients[1] + coefficients[2] * X + coefficients[3] * X^2
13 cat("Coefficients:\n")
14 cat("beta_0: ", coefficients[1], "\nbeta_1: ", coefficients[2], "\nbeta_2: "
    , coefficients[3], "\n")
15
16
17 # Plotting of the curve
18 plot(X, y, pch=16, col="blue", xlab="X", ylab="Y", main="Scatter plot with
    regression curve")
19 curve(coefficients[1] * 1 + coefficients[2] * x + coefficients[3] * x^2, add
    = TRUE, col = "red", lwd = 2)
20
21 residuals <- y - (designMatrix %*% coefficients)
22 df <- length(y) - ncol(designMatrix)
23 standardErr <- sqrt(diag(solve(productMatrix)) * sum(residuals^2)/df)
24
25 alpha <- 0.05
26 tObserved <- coefficients / standardErr
27 t <- qt(1 - alpha/2, df)
28 lowerBound <- coefficients - t * standardErr
29 upperBound <- coefficients + t * standardErr
30
31
32 cat("\n95% Confidence Intervals:\n")
33 for (i in 1:length(coefficients)) {
34   cat("beta_", i-1, ": [", lowerBound[i], ", ", upperBound[i], "]\n")
35 }
36
37
38 # c) Test for significance of coefficients
39 # Test for significance of each coefficient (using if-else)
40 for (i in 1:length(coefficients)) {
41   if (abs(tObserved[i]) > t) {
42     cat("beta_", i-1, "is Significant (p < 0.05)\n")
43   } else {
44     cat("beta_", i-1, "is Not Significant (p >= 0.05)\n")
45   }
46 }
47
48
49 # d) Computation of R squared
50 SSRes <- sum((y - yPredicted)^2)
51
52 meanY <- mean(y)
53 SSTotal <- sum((y - meanY)^2)
54 R2 <- 1 - (SSRes / SSTotal)
55
56 cat("Coefficient of determination (R^2):", R2, "\n")

```