

# An electroencephalogram (EEG) statistical analysis

Adu Boahene Quarshie  
Faculty of Engineering, Environment and Computing  
MSc. Data Science and Computational Intelligence  
Coventry University, Coventry, United Kingdom  
[quarshiea@uni.coventry.ac.uk](mailto:quarshiea@uni.coventry.ac.uk)

**Abstract:** In this project we analyse an electroencephalogram dataset which we assume to be measured from two different regions of the brain. We then perform exploratory data analysis on the dataset and then move on to identify the non-linear regression model structure to estimate its parameters.

**Keywords:** EEG; Regression; R Code; Correlation;

## I. INTRODUCTION

Here, we would like to model and understand the degree (or order) of possible nonlinear interactions between those two EEG channels. The basic model can be considered as a generic nonlinear polynomial regression model with the following (exemplar) structure:

$$y = w_0 + w_1x + w_2x^2 + w_3x^3 + \dots + w_nx^n + \epsilon$$

We then go ahead to perform exploratory data analysis on the dataset by analysing the time series plots of the signals, distribution of each signal, the signals correlation and fit a linear model to the data to estimate its performance. We will then identify the non-linear regression model structure and estimate its parameters by splitting the input and output signals into train and test sets. We then apply forward subset selection to get the best model structure. We also estimate parameters using the least mean squared error. When the best model is determined, we estimate the model's covariance matrix and compute the model's output or prediction. We finally validate the model and apply Approximate Bayesian Computation to compute the posterior distribution.

## II. THE DATA SET

The dataset is an EEG time-series dataset with "X" and "Y" labels. It is assumed the signals are measured from two different regions of the brain. The output, Y, contains a corresponding additive noise with unknown variance. The dataset has 250 data points or entries.

## III. EXPLORATORY DATA ANALYSIS

For us to be able to work with the data we have, we have to make sure that we achieve all the information needed to understand the nature of the data.

### Renaming columns

When we load the data into R, we notice that the dataset is not labelled so R automatically column names V1 and V2. So, in order to make our dataset easy to work with, we do this with this R code:

```
1. #reading the csv file into R
2. data <- read.csv("x_y.csv", header = FALSE)
3. #renaming columns to input variable x and the output variable y
4. colnames(data)[colnames(data)=="V1"] <- "x"
5. colnames(data)[colnames(data)=="V2"] <- "y"
6. head(data)
```

### Plotting time series for both input and output signals

A time series shows a progression of fluctuations of a particular variable in a time space, in this coursework we plot time series input against output to analyse how the pair moves together in that time space.

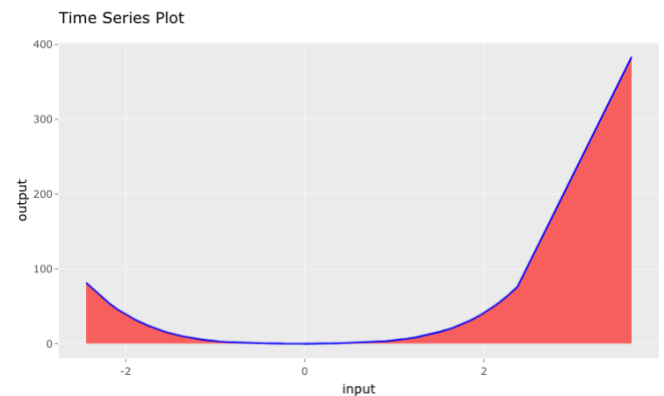


Fig. 1

The above plot shows the timeseries of both the input signal and its corresponding output signal. It is seen that when the input is  $-2.44104609$ , the output signal is  $8.153268e^{+01}$ .

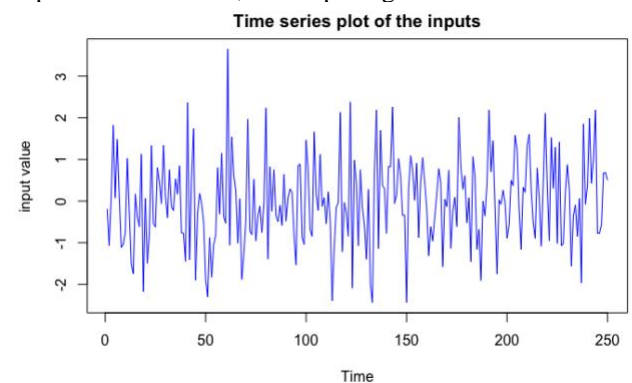


Fig. 2

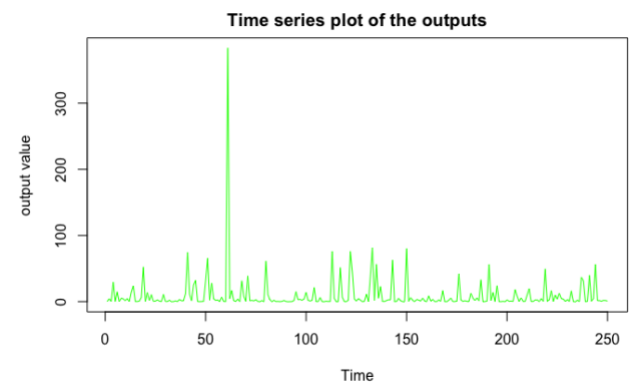


Fig 3

Plotting the individual time-series for both variables, the input signals has greater fluctuations as compared to the output signals. At about time,  $t = 60$ , the output signal is at its highest peak likewise the input signal.

### Distribution for input and output signal

The plots below is the distribution of the Output signal variable. It can be seen that the output variable is highly skewed to the left.

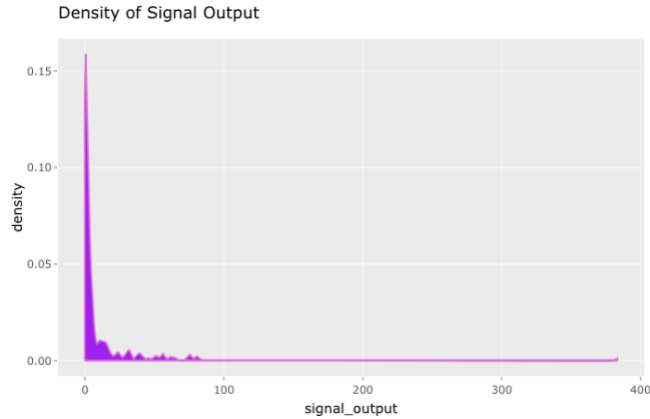


Fig 4.

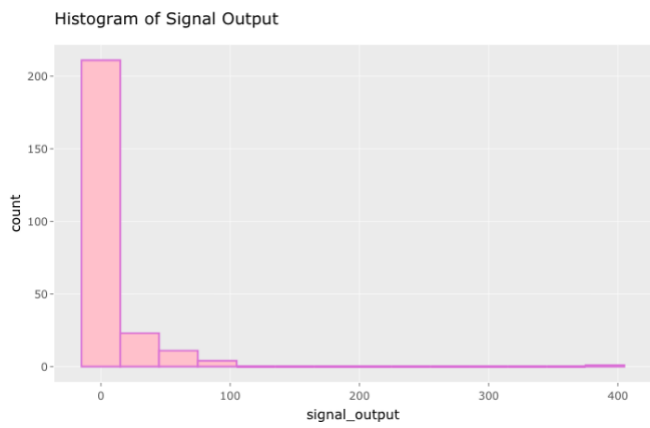


Fig 5.

As compared to the distribution of the output signal, the input signal is close to a normal distribution.

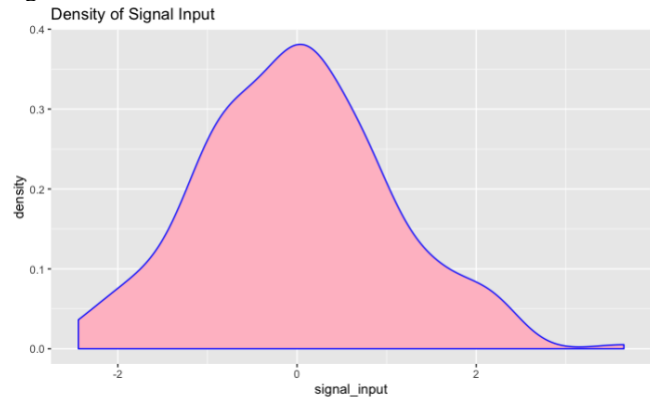


Fig 6.

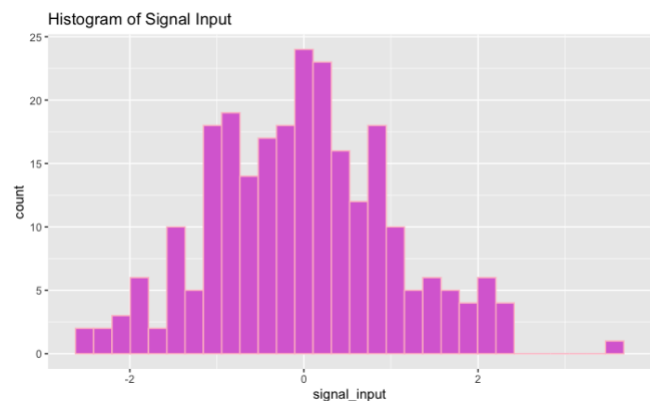


Fig 7.

### Input and output correlation and scatter plot

It can be seen from the plot below that there exists an extreme outlier. This plot also shows the correlation between the signals.

For the correlation to be calculated, we use the Pearson's correlation formula.

Pearson correlation coefficient for population is given as;

$$\rho_{x,y} = \frac{E[XY] - E[X]E[Y]}{\sqrt{E[X^2] - [E[X]]^2} \sqrt{E[Y^2] - [E[Y]]^2}}$$

Pearson correlation coefficient for sample is given as;

$$r_{xy} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}}$$

Using the 'cor()' function in R, it can be seen that there exists some sort of correlation between the input signal and the output signal from the diagram above. Performing the Pearson's correlation test, it is known that there exists a positive correlation of **0.2193661** between the input and output signals. R code below:

```
1. #correlation test
2. cor(data$X, data$y, method = "pearson", use = "complete.obs")
```

### Model fitting and its performance

Fitting a static linear regression model to the data, thetaHat (the estimate of theta) was calculated to be **5.844558**. And mean squared error (MSE) was **886.2422**. This MSE shows the performance of the linear model fitted and since it is high, there is also a high residual value, therefore this model is not good. For the simple linear regression this is the formula;

$$y = \theta_0 + \theta_1 x + \epsilon$$

With  $\theta_0$  being the intercept and being zero computing the estimate for  $\theta_1$

### IV. MODEL SELECTION

With the model selection, we create an intercept column, this creates a matrix of ones with 250 rows in one column. Then we create a matrix of 'X' that consist of the intercept and dependent variables of 'X'. We then split the dataset into test and train sets with 80% being for training and 20% being for testing.

#### Model selection (Forward model subset selection)

Forward selection is a kind of stepwise regression that starts with a null model and then adds variables in a succession. In each succession, a variable is added which makes our model better. For this part, we will create a data frame for Mean Squared Error, create a 'for' loop in R to find the minimum mean squared error values by using this code. We will then find the errors and append all the MSE's for each term in the data frame with the code below;

```
1. #Creating a data frame for MSE values
2. MSE_table = data.frame(terms = colnames(X), MSE = rep(0, length(6)))
3.
4. #Create a for loop to find the minimum MSE Values
```

```

5. for (i in 1:6){
6.   thetaHat = solve(t(X_train[,i]) %*% X_train[,i]) %*% t(X_train[,i]) %*% y_train
7.   y_hat = X_test[,i] %*% thetaHat
8.
9.   #Finding the errors
10.  error = y_test - y_hat
11.  MSE = mean(error^2)
12.  MSE_table[i, 2] = MSE
13. }

```

### Selecting minimum MSEs from data frame above

The X matrix has 250 values and I will be splitting it into a test and train set and choosing the ones with the minimum MSE from the MSE table computed below;

	Term	MSE
1	<i>Intercept</i>	177.950779
2	$X$	174.184985
3	$X^2$	93.988810
4	$X^3$	128.226234
5	$X^4$	5.232054
6	$X^5$	137.497411

Table 1.

Selecting the variables for the model which has the least MSE, from the first MSE table, it can be seen that  $X^4$  had the least MSE as compared to the rest, therefore  $X^4$  was selected and MSE was performed again on the rest of the variables, this was done repeatedly so as to select the variables with the least MSE, in order to create the best model.

### Final Model

We then compute for the final model with the code below after getting all the least MSE's.

```

1. #Final model
2. xfinal_train = cbind(new_train3, new_train2, new_train1)
3. xfinal_test = cbind(new_test3, new_test2, new_test1)
4. colnames(xfinal_train) = c("X", "X^2", "X^4")
5. colnames(xfinal_test) = c("X", "X^2", "X^4")
6.
7. final_thetaHat = solve(t(xfinal_train) %*% xfinal_train) %*% t(xfinal_train) %*% y_train
8.
9. y_hat_final = xfinal_train %*% final_thetaHat #Refer to coursework guide about computing model prediction based on the train set
10.
11. #Finding the errors
12. error = y_train - y_hat_final
13. MSE_final = mean(error^2)

```

The final model consists of  $X$ ,  $X^2$  and  $X^4$ . These variables/regressors give the least MSE.

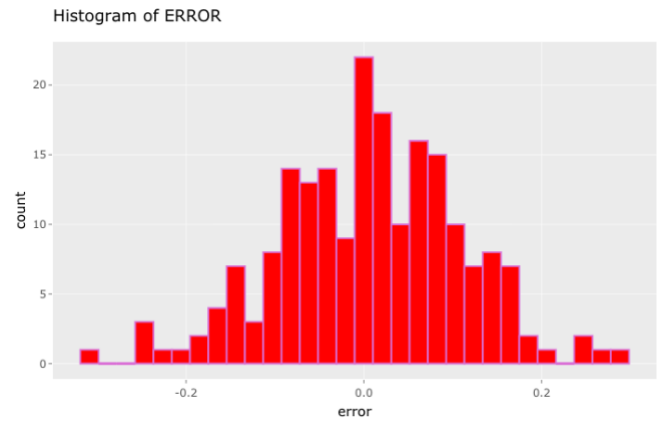


Fig 8.

### Covariance

Covariance is the matrix that shows the variability between the variables. We calculate this in R with the code below;

```

1. error = y_train - y_hat_final
2. residuals_sum_of_squares = sum((error)^2)
3. sigma_square = residuals_sum_of_squares/n - 1
4. #covariance matrix
5. cov_thetahat = sigma_square * (solve(t(xfinal_train) %*% xfinal_train)) #t(X) is t ranspose of X
6. colnames(cov_thetahat) = c("one_theta", "two_theta", "three_theta")
7. rownames(cov_thetahat) = c("one_theta", "two_theta", "three_theta")

```

### Contour plotting and confidence interval

We plot contours for our combinations and then go ahead to calculate the confidence intervals. The code for the contour plotting and the 3D plots will be found in the appendix and the code for calculating the confidence interval is below;

```

1. n = 200
2. var_y_hat = matrix(0 , n , 1)
3.
4. for( i in 1:n){
5.   X_i = matrix( xfinal_train[i,] , 1 , n_m_param ) # X[i,] creates a vector. Convert it to matrix
6.   var_y_hat[i,1] = X_i %*% cov_thetahat %*% t(X_i) # same as sigma_2 * ( X_i %*% ( solve(t(X) %*% X ) ) %*% t(X_i) )
7. }
8.
9. CI = 2 * sqrt(var_y_hat) # Confidence interval
10.
11. plot(data$X[1:200], y_hat_final, type = "p")
12. segments(data$X[1:200], y_hat_final - CI, data$X[1:200], y_hat_final + CI, col = "red") # Adds error bars to the individual data points

```

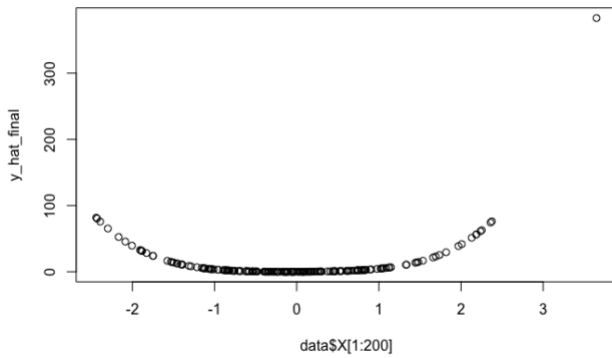


Fig 9.

#### IV. VALIDATION AND CONCLUSION

For validating our model, we will be splitting our dataset into 70% and 30% training and testing set respectively. We will then estimate the parameters for both test and train set and find the errors for both sets. The mean squared error is calculated to check the accuracy of the model.

From our computing with the code below, we find out that that the  $MSE_{final} = 0.0109631630515089$  and the error rate of our test\_MSE =  $0.00834867111029402$ . This proves with the very less error that our model is the best fit for prediction for our dataset.

#### V. REFERENCES

- [1] Brant, Rollin. Forward Selection. MDSC 643.02 Lecture Materials. Retrieved from <https://www.stat.ubc.ca/~rollin/teach/643w04/lec/node41.html> on July 7, 2018
- [2] Cook, Perry. Stepwise Selection. Human-Computer Interface Technology (CS436) Class Notes. Retrieved from <https://www.cs.princeton.edu/courses/archive/fall08/cos436/Duda/FS/stepwise.htm> on July 8, 2018.
- [3] Shalizi, Cosma. Lecture 26: Variable Selection. Modern Regression for Undergraduates Class Notes. <http://www.stat.cmu.edu/~cshalizi/mreg/15/lectures/26/lecture-26.pdf>
- [4] SAS Support. Forward Selection. The GLMSELECT Procedure. Retrieved from [http://support.sas.com/documentation/cdl/en/statug/66859/HTML/default/viewer.htm#statug\\_glmselect\\_details03.htm](http://support.sas.com/documentation/cdl/en/statug/66859/HTML/default/viewer.htm#statug_glmselect_details03.htm) on July 8, 2018

# APPENDIX

Full Code and Plots for the coursework

## EXPLORATORY DATA ANALYSIS

```
1. #reading the csv file into R
2. data <- read.csv("x_y.csv", header = FALSE)
3. #renaming columns to input variable x and the output variable y
4. colnames(data)[colnames(data)=="V1"] <- "X"
5. colnames(data)[colnames(data)=="V2"] <- "y"
6. head(data)
```

## PLOTTING BOTH SIGNALS TIME-SERIES

```
1. library(plotly)
2. library(dplyr)
3. time_series <-
4.   ggplot(data = data, aes(x=data$X, y=data$y)) +
5.     geom_area(fill="red", alpha = 0.6) +
6.     geom_line(color="blue") +
7.     labs(x = "input", y = "output", title = "Time Series Plot")
8.
9. # Turn it interactive with ggplotly
10. ggplotly(time_series)
11.
12. plot.ts(data$X, col = "blue", main = "Time series plot of the inputs", ylab = "input value")
13. plot.ts(data$y, col = "green", main = "Time series plot of the outputs", ylab = "output value")
```

## DISTRIBUTION FOR SIGNAL OUTPUT

```
1. library(ggplot2)
2. signal_output <- data$y
3. #plotting the density of the target variable
4. density <- ggplot(data = data, aes(signal_output)) +
5.   geom_density(kernel = "gaussian", col = "orchid", fill = "purple") +
6.   labs(title = "Density of Signal Output")
7. histogram <- ggplot(data = data, aes(signal_output)) +
8.   geom_histogram(col = "orchid", fill = "pink", binwidth = 30) +
9.   labs(title = "Histogram of Signal Output")
10. ggplotly(density)
11. ggplotly(histogram)
```

## DISTRIBUTION FOR SIGNAL INPUT

```
1. signal_input <- data$X
2. #plotting the density of the input variable
3. ggplot(data = data, aes(signal_input)) +
4.   geom_density(kernel = "gaussian", col = "blue", fill = "pink") +
5.   labs(title = "Density of Signal Input")
6. ggplot(data = data, aes(signal_input)) +
7.   geom_histogram(col = "pink", fill = "orchid") +
8.   labs(title = "Histogram of Signal Input")
```

## INPUT-OUTPUT CORRELATION AND TEST

```
1. ggplot(data, aes(x = data$X, y = data$y)) +
2.   geom_point(aes(color = "red")) +
3.   labs(x = "inputs", y = "outputs")
4.
5. #correlation test
6. cor(data$X, data$y, method = "pearson", use = "complete.obs")
```

## MODEL FITTING AND ITS PERFORMANCE

```
1. thetaHat = solve(t(data$X)%*% data$X) %*% t(data$X) %*% data$y
2. print(thetaHat)
3. y_hat = data$X %*% thetaHat
4. Y_hat = data.frame(y_hat)
5. y_hat = Y_hat$y_hat
6. histogram_pred <- ggplot(data = Y_hat, aes(y_hat)) +
7.   geom_histogram(col = "orchid", fill = "pink") +
8.   labs(title = "Histogram of the predictions")
9. ggplotly(histogram_pred)
10.
11. MSE = mean((data$y - y_hat)^2)
12. print (MSE)
13. error = data$y - y_hat
14. Error = data.frame(error)
15. error = Error$error
16. histogram_error <- ggplot(data = Error, aes(error)) +
17.   geom_histogram(col = "blue", fill = "red") +
18.   labs(title = "Histogram of the Error")
19. ggplotly(histogram_error)
```

## MODEL SELECTION

```
1. #Creating an intercept column
2. #This creates a matrix of ones with 250 rows in one column
3. library(optimbase)
4. intercept = ones(nx = 250, ny = 1)
5.
6. #Creating a matrix of X that consist of the intercept and dependent variables of X
7. X = data$X
8. X = cbind(intercept, X, X^2, X^3, X^4, X^5)
9. colnames(X) = c("intercept", "X", "X^2", "X^3", "X^4", "X^5")
10.
11. #This changes y into a matrix
12. y <- data$y
13. y <- matrix(y)
14.
15. #Splitting into train set and test set where 80% is for training and 20% is for testing
16. n = dim(X)[1]
17. #80% trainset
18. n_train = 0.8 * n
19. #splitting in train and test sets
20. X_train = X[1:n_train, ]
21. y_train = as.matrix(y[1:n_train, ])
22. X_test = X[(n_train+1):n, ]
23. y_test = as.matrix(y[(n_train+1):n, ])
```

## FORWARD MODEL SUBSET SELECTION

```
1. #Creating a dataframe for MSE values
2. MSE_table = data.frame(terms = colnames(X), MSE = rep(0, length(6)))
3.
4. #Create a for loop to find the minimum MSE Values
5. for (i in 1:6){
6.   thetaHat = solve(t(X_train[,i]) %*% X_train[,i]) %*% t(X_train[,i]) %*% y_train
7.   y_hat = X_test[,i] %*% thetaHat
8.
9.   #Finding the errors
10.  error = y_test - y_hati
11.  MSE = mean(error^2)
12.  MSE_table[i, 2] = MSE #Adding all the MSEs for each term in the dataframe
13. }
```

## SELECTING THE MINIMUM MSE VALUE FROM THE DATA FRAME ABOVE

```
1. new_term1 = MSE_table$terms[which.min(MSE_table$MSE)]
```

```

2. #The X matrix has 250 values and I will be splitting it into a test and train set, and choosing the ones with the minimum MSE from the MSE table computed above.
3. new_train1 = as.matrix(X_train[, which.min(MSE_table$MSE)])
4. new_test1 = as.matrix(X_test[, which.min(MSE_table$MSE)])
5.
6. #Removing X^4 from the model and computing the MSE for the other variables in the model
7. X_train1 = X_train[, - which.min(MSE_table$MSE)]
8. X_test1 = X_test[, - which.min(MSE_table$MSE)]
9. MSE_table1 = data.frame(terms = colnames(X_train1), MSE = rep(0, length(5))) #Creating a new data frame
10.
11. for (i in 1:5){
12.   train_x1 = cbind(new_train1, X_train1[, i])
13.   test_x1 = cbind(new_test1, X_test1[, i])
14.   thetaHat = solve(t(train_x1) %*% train_x1) %*% t(train_x1) %*% y_train
15.
16.   y_hat = test_x1 %*% thetaHat
17.
18.   #Finding the errors
19.   error = y_test - y_hat
20.   MSE1 = mean(error^2)
21.   MSE_table1[i, 2] = MSE1 #Adding all the MSEs for each term in the dataframe
22. }
23.
24. new_term2 = MSE_table1$terms[which.min(MSE_table1$MSE)]
25. #The X^2 has 250 values and I will be splitting it into a test and train set
26. new_train2 = as.matrix(X_train1[, which.min(MSE_table1$MSE)])
27. new_test2 = as.matrix(X_test1[, which.min(MSE_table1$MSE)])
28.
29. #Creating a new X matrix without X^2
30. X_train2 = X_train1[, - which.min(MSE_table1$MSE)]
31. X_test2 = X_test1[, - which.min(MSE_table1$MSE)]
32. MSE_table2 = data.frame(terms = colnames(X_train2), MSE = rep(0, length(4))) #Creating a new data frame
33.
34. for (i in 1:4){
35.   train_x2 = cbind(new_train2, new_train1, X_train2[, i])
36.   test_x2 = cbind(new_test2, new_test1, X_test2[, i])
37.   thetaHat = solve(t(train_x2) %*% train_x2) %*% t(train_x2) %*% y_train
38.
39.   y_hat = test_x2 %*% thetaHat
40.
41.   #Finding the errors
42.   error = y_test - y_hat
43.   MSE2 = mean(error^2)
44.   MSE_table2[i, 2] = MSE2 #Adding all the MSEs for each term in the dataframe
45. }
46.
47. #New term 3
48. new_term3 = MSE_table2$terms[which.min(MSE_table2$MSE)]
49. #The X has 250 values and I will be splitting it into a test and train set
50. new_train3 = as.matrix(X_train2[, which.min(MSE_table2$MSE)])
51. new_test3 = as.matrix(X_test2[, which.min(MSE_table2$MSE)])

```

## FINAL MODEL

```

1. #Final model
2. xfinal_train = cbind(new_train3, new_train2, new_train1)
3. xfinal_test = cbind(new_test3, new_test2, new_test1)
4. colnames(xfinal_train) = c("X", "X^2", "X^4")
5. colnames(xfinal_test) = c("X", "X^2", "X^4")
6.
7. final_thetaHat = solve(t(xfinal_train) %*% xfinal_train) %*% t(xfinal_train) %*% y_train
8.
9. y_hat_final = xfinal_train %*% final_thetaHat #Refer to coursework guide about computing model prediction based on the train set
10.
11. #Finding the errors
12. error = y_train - y_hat_final
13. MSE_final = mean(error^2)
14.

```



```

15. Error = data.frame(error)
16. error = Error$error
17. histogram_error <- ggplot(data = Error, aes(error)) +
18.   geom_histogram(col = "orchid", fill = "red") +
19.   labs(title = "Histogram of ERROR ")
20. ggplotly(histogram_error)

```

## COVARIANCE

```

1. error = y_train - y_hat_final
2. residuals_sum_of_squares = sum((error)^2)
3. sigma_square = residuals_sum_of_squares/n-1
4. #covariance matrix
5. cov_thetahat = sigma_square * (solve(t(xfinal_train) %*% xfinal_train)) #t(X) is transpose of X
6. colnames(cov_thetahat) = c("one_theta", "two_theta", "three_theta")
7. rownames(cov_thetahat) = c("one_theta", "two_theta", "three_theta")

```

## CONTOUR PLOTTING

```

1. num_param = 3
2. num_points = 50 # no point on the plot
3. one_theta = seq(0.489 , 0.509 , length=num_points)
4. two_theta = seq(1.993 , 2.013 , length=num_points)
5. three_theta = seq(1.990 , 2.010 , length=num_points)
6.
7. prob_den_func = matrix(0 , num_points , num_points)
8.
9. cov_thetahat_inv = (t(X) %*% X) * (1/sigma_square) # inverse of cov_thetaHat
10. det_cov_thetahat = det(cov_thetahat) # determinant of cov_thetaHat
11.
12. #theta_1 and theta_2 combination
13. for(r in 1:50){
14.   for(c in 1:50){
15.
16.     one_two_theta = matrix( c(one_theta[r] , two_theta[c], final_thetaHat[3,] ) , num_param , 1)
17.     thetahat_theta = one_two_theta - final_thetaHat
18.
19.     prob_den_func[r,c] = ( 1/sqrt( ( (2*pi)^num_param ) * det_cov_thetahat ) ) *
20.       exp( -0.5 * t(-thetahat_theta) %*% solve(cov_thetahat) %*% -thetahat_theta )
21.
22.   }
23. }
24.
25. contour(one_theta, two_theta, prob_den_func)
26. persp(one_theta, two_theta, prob_den_func, theta = 50 , phi = 50)
27.
28. #one_theta and three_theta combination
29.
30. for(r in 1:50){
31.   for(c in 1:50){
32.
33.     one_three_theta = matrix( c(one_theta[r], final_thetaHat[2,] , three_theta[c] ) , num_param, 1
34.     )
35.     thetahat_theta = one_three_theta - final_thetaHat
36.
37.     prob_den_func[r,c] = ( 1/sqrt( ( (2*pi)^num_param ) * det_cov_thetahat ) ) *
38.       exp( -0.5 * t(-thetahat_theta) %*% solve(cov_thetahat) %*% -thetahat_theta )
39.   }
40. }
41.
42. contour(one_theta, three_theta, prob_den_func)
43. persp(one_theta, three_theta, prob_den_func, theta = 50 , phi = 50)
44.
45.
46. #two_theta and three_theta combination
47.
48. for(r in 1:50){
49.   for(c in 1:50){

```



```

50.
51.     two_three_theta= matrix( c(final_thetaHat[1,], two_theta[r] , three_theta[c] ) , num_param , 1
52. )
53.     thetahat_theta = two_three_theta - final_thetaHat
54.
55.     prob_den_func[r,c] = ( 1/sqrt( ( (2*pi)^num_param ) * det_cov_thetahat ) ) *
56.     exp( -0.5 * t(-thetahat_theta) %>% solve(cov_thetahat) %>% -thetahat_theta )
57. }
58. }
59.
60. contour(two_theta, three_theta, prob_den_func)
61. persp(two_theta, three_theta, prob_den_func, theta = 50 , phi = 50)
62. ```
63.
64. #Confidence Interval
65. ```{r}
66. n = 200
67. var_y_hat = matrix(0 , n , 1)
68.
69. for( i in 1:n){
70.   X_i = matrix( xfinal_train[i,] , 1 , num_param ) # X[i,] creates a vector. Convert it to matrix
71.
72.   var_y_hat[i,1] = X_i %>% cov_thetahat %>% t(X_i) # same as sigma_2 * ( X_i %>% ( solve(t(X) %>%
73.   X ) ) %>% t(X_i) )
74. }
75.
76. CI = 2 * sqrt(var_y_hat) # Confidance interval
77.
78. plot(data$X[1:200], y_hat_final, type = "p")
79. segments(data$X[1:200], y_hat_final-
80.   CI, data$X[1:200], y_hat_final+CI, col = "red") # Adds error bars to the indivigual data points

```

## CONFIDENCE INTERVAL

```

1. n = 200
2. var_y_hat = matrix(0 , n , 1)
3.
4. for( i in 1:n){
5.   X_i = matrix( xfinal_train[i,] , 1 , num_param ) # X[i,] creates a vector. Convert it to matrix
6.
7.   var_y_hat[i,1] = X_i %>% cov_thetahat %>% t(X_i) # same as sigma_2 * ( X_i %>% ( solve(t(X) %>%
8.   X ) ) %>% t(X_i) )
9. }
10.
11. CI = 2 * sqrt(var_y_hat) # Confidance interval
12.
13. plot(data$X[1:200], y_hat_final, type = "p")
14. segments(data$X[1:200], y_hat_final-
15.   CI, data$X[1:200], y_hat_final+CI, col = "red") # Adds error bars to the indivigual data points

```

## VALIDATION

```

1. X_data_matrix = X[, c(2, 3, 5)]
2. colnames(X_data_matrix) = c("X", "X^2", "X^4")
3.
4. #Splitting the data into 70/30 sets
5. x_train_val = X_data_matrix[1:175,]
6. y_train_val = as.matrix(data$y[1:175])
7. x_test_val = X_data_matrix[176:250,]
8. y_test_val = as.matrix(data$y[176:250])
9.
10.
11. #Estimation of Parameters
12. val_thetaHat = solve(t(x_train_val) %>% x_train_val) %>% t(x_train_val) %>% y_train_val
13.
14. y_hat_test = x_test_val %>% val_thetaHat #for testing
15. y_hat_train = x_train_val %>% val_thetaHat #for training
16.

```

```

17. #Finding error
18. test_error = y_test_val - y_hat_test #Testing error
19. train_error = y_train_val - y_hat_train #Training error
20.
21. #MSE
22. test_MSE = mean((test_error)^2)
23. train_MSE = mean((train_error)^2)
24.
25. #Comment on the above to prove the validity of the model.
26. ```
27. Validation is checking if the chosen model is the best fit or accurate for the dataset
28. Compared the MSE
29.
30. #ABC
31. ```{r include=FALSE}
32. pacman::p_load(ggplot2, ggthemes, tidyr, patchwork, plotly, viridis)

```

## PRIORS

```

1. sampling_n = 4e7 #draw 30M from each prior
2.
3. prior_x1 = runif(sampling_n, -5, 5)
4. prior_x2 = runif(sampling_n, -5, 5)
5. prior_x4 = runif(sampling_n, -5, 5)
6.
7. #merge priors to one matrix so we can simply take rows as theta candidates
8. priors = cbind(prior_x1, prior_x2, prior_x4)
9. priors = as.data.frame(priors)
10. colnames(priors) = c("x1", "x2", "x4")
11.
12. rm(prior_x1)
13. rm(prior_x2)
14. rm(prior_x4)

```

## PLOTTING THE PRIORS

```

1. priors_long = gather(priors, prior, sample, x1:x4, factor_key = T)
2.
3. (priors_plot = ggplot(priors_long, aes(x = sample, color = prior))+
4.   geom_density()+
5.   facet_wrap(~prior, scales = "free")+
6.   theme_few()+
7.   scale_color_few()+
8.   guides(color = F)+
9.   xlab(""))
10.
11. ggsave("figures/10_priors.png", width = 7, height = 4)
12. rm(priors_long)
13. rm(priors_plot)
14.
15. X = cbind(data$x,
16.           data$x^2,
17.           data$x^4
18.           )

```

## REJECTION ABC

```

1. posterior =data.frame(x1 = numeric(),
2.                       x2 = numeric(),
3.                       x4 = numeric())
4.
5. #cross validation testing MSE +3*sd
6. #tolerance = 0.028
7.
8. tolerance = 0.2
9.
10.

```

```

11. for (i in 1:nrow(priors)) {
12.   candidate = t(as.matrix(priors[i,]))
13.   simulated = X %%% candidate
14.
15.   MSE = mean((data$y - simulated)^2)
16.
17.   if (MSE < tolerance) {
18.     append = t(candidate)
19.     colnames(append) = c("x1", "x2", "x4")
20.     posterior = rbind(posterior, append)
21.   }
22. }
23.
24. cat(nrow(posterior), "samples were accepted.")
25.
26. write.csv(posterior, "data/posterior.csv")
27.
28.
29.
30. h1 = ggplot(posterior, aes(x = x1))+
31.   geom_density()+
32.   labs(x = "", y="")+
33.   theme_few()
34. h2 = ggplot(posterior, aes(x = x2))+
35.   geom_density()+
36.   labs(x = "", y="")+
37.   theme_few()
38. h3 = ggplot(posterior, aes(x = x4))+
39.   geom_density()+
40.   labs(x = "", y="")+
41.   scale_x_continuous(expand = c(0, 0),
42.                       breaks = round(seq(min(posterior$x4), max(posterior$x4), by = 0.02),2))+
43.   theme_few()
44.
45. p1 = ggplot(posterior, aes(x=x1, y=x2) ) +
46.   stat_density_2d(aes(fill = ..density..), geom = "raster", contour = FALSE) +
47.   scale_fill_viridis() +
48.   scale_x_continuous(expand = c(0, 0),
49.                       breaks = round(seq(min(posterior$x1), max(posterior$x1), by = 0.2),1)) +
50.   scale_y_continuous(expand = c(0, 0),
51.                       breaks = round(seq(min(posterior$x2), max(posterior$x2), by = 0.2),1)) +
52.   theme(
53.     legend.position='none'
54.   )+
55.   guides(fill = F)+
56.   labs(x = "", y="")+
57.   theme_few()
58.
59. p2 = ggplot(posterior, aes(x=x1, y=x4) ) +
60.   stat_density_2d(aes(fill = ..density..), geom = "raster", contour = FALSE) +
61.   scale_fill_viridis() +
62.   scale_x_continuous(expand = c(0, 0),
63.                       breaks = round(seq(min(posterior$x1), max(posterior$x1), by = 0.2),1)) +
64.   scale_y_continuous(expand = c(0, 0),
65.                       breaks = round(seq(min(posterior$x4), max(posterior$x4), by = 0.05),2)) +
66.   theme(
67.     legend.position='none'
68.   )+
69.   guides(fill = F)+
70.   labs(x = "", y="")+
71.   theme_few()
72.
73. p3 = ggplot(posterior, aes(x=x2, y=x4) ) +
74.   stat_density_2d(aes(fill = ..density..), geom = "raster", contour = FALSE) +
75.   scale_fill_viridis() +
76.   scale_x_continuous(expand = c(0, 0),breaks = round(seq(min(posterior$x2), max(posterior$x2), by
77.   = 0.2),1)) +
78.   scale_y_continuous(expand = c(0, 0),
79.                       breaks = round(seq(min(posterior$x4), max(posterior$x4), by = 0.05),2)) +
80.   theme(
81.     legend.position='none'
82.   )+
82.   guides(fill = F)+

```

```

83. labs(x = "", y="")+
84. theme_few()
85.
86. final_plot =
87. (h1+(p1+coord_flip()+(p2+coord_flip()))/
88. (p1+h2+(p3+coord_flip()))/
89. (p2+p3+h3)
90. ggsave("figures/11_ABC_posterior.png", width = 7, height = 4)

```

## GETTING THE MAP

```

1. getmode <- function(v) {
2.   uniqv <- unique(v)
3.   uniqv[which.max(tabulate(match(v, uniqv)))]
4. }
5.
6. cat("x: ", getmode(posterior$x1))
7. cat("\nx2:", getmode(posterior$x2))
8. cat("\nx4:", getmode(posterior$x4))

```

## PREDICTION GENERATION

```

1. posterior = read.csv("data/posterior.csv", row.names = 1)
2. predictions = data.frame()
3.
4. for (i in 1:nrow(posterior)) {
5.   candidate = t(as.matrix(posterior[i,]))
6.   simulated = X %*% candidate
7.   append = data.frame(y = simulated, which = rep(i, length(simulated)))
8.   colnames(append)[1] = "y"
9.   predictions = rbind(predictions, append)
10. }
11.
12. predictions$which = as.factor(predictions$which)
13. predictions$x = rep(data$x, nrow(posterior))
14. predictions$truth = rep(data$y, nrow(posterior))
15.
16. (plot_predictions = ggplot(predictions, aes(x = x, y = y, color=which))+
17.   geom_line(alpha = 0.2)+
18.   geom_line(aes(y = truth), color="black")+
19.   guides(color = F)+
20.   theme_few())
21.
22. ggsave("figures/12_bayesian_prediction.png", plot_predictions, width = 7, height = 4)

```

## CONTOUR PLOTINGS

