

Analysis of the Correlation between Diet and COVID-19 Cases by Country

Zeng Fung Liew (lzfliew@ucdavis.edu)

December 12, 2020

Abstract

The goal of this analysis is to find out how a country's diet correlates with its number COVID-19 cases. With different food cultures across the world, it would be interesting to see what are the food categories that can best predict a country's rate of confirmed cases, deaths, recovery cases, and active cases respectively. It turns out that there are some food categories which, if the intake is high (or too low), highly correlates with the high percentage of COVID-19 cases per capita.

The next thing that was investigated was the fitting of a suitable linear model for the purpose of COVID-19 case predictions. Due to the large amount of predictor variables in this data, there were a few preliminary steps taken to eliminate some of them prior to the fitting of models, some of which include eliminating predictor variables of the same food type that have high correlations with each other, removing predictor variables with (almost) no distribution, and deleting predictor variables with large variance inflation factors (VIF), ie. predictor variables that are highly inter-correlated with the rest of the predictor variables. In the model selection process, the models were started off without any predictor variables and the forward stepwise regression method was applied with respect to both the AIC and BIC criteria. The resulting models were then analyzed based on their C_p , $Press_p$, adjusted R^2 , and mean square prediction error (MSPE) criteria to determine the model that generalizes better. Then, possible outliers were analyzed using the Cook's distance plot and subsequently removed before the selection of the best prediction model.

1 Introduction

The analysis of this data set was motivated by the rising COVID-19 cases across the world. Due to the different food and diet cultures across countries, it is interesting to find if any correlation exist between a country's typical diet and its number of COVID-19 confirmed, death, recovered, and active cases. Obviously, the correlations shown in this data analysis (if any) are not indications causation effect, ie. having certain food in large amount does not lead to a higher chance of contacting COVID-19, but it could be an indicator of strength of the general population's immune system, which could have an effect on the spread of COVID-19 in the country.

The data set used in this analysis was found on Kaggle, created and maintained by Marie Ren. It comes in four data frames containing 32 columns each, where each data frame takes a different measurement of the food categories, ie. fat quantity, food quantity, food energy, and protein quantity. The quantitative variables are represented in terms of percentages, eg. the percentage of fat content in eggs in a diet by country, or the percentage of the population of confirmed cases of COVID-19. However, not all of the columns will be used in this data analysis. As the goal of this analysis is to find the correlations between diet and COVID-19 cases, variables like obesity and undernourished rate, and country populations were ignored. Only the columns that describe predictor variables (food categories), response variables (confirmed, death, recovered, active cases), and countries were retained. This data analysis was started after combining the four data frames into one large data frame containing 170 rows and 100 columns.

The aim of this data analysis is to figure out the food categories that can best predict the number of confirmed, death, recovered, and active cases of COVID-19 in the world. It would also be interesting to see if

these food categories are the same across each of those response variables, ie. if the food categories that can predict the rate of recovered can also predict the rate of death cases.

This data analysis was done on the most recently updated data set at the time of writing, which was on December 2, 2020. Changes in the COVID-19 situations across countries could result in different analysis results.

2 Methods and Results

2.1 Exploratory Data Analysis

As the main goal of this data analysis is to fit the best possible Ordinary Least Squares linear models to predict the rate of confirmed, death, recovered, and active cases, it is important to make sure that its assumptions are satisfied prior to the model fitting step. Some of the important assumptions that have to be satisfied are:

- Normal distribution of the response variables and the constant variance among residuals. If this isn't the case in the given data, a transformation of variables has to be carried out.
- Any two responses are uncorrelated. This is assumed to be true on the given data set.
- Linearity between each predictor variable with each response variable.

Additionally, due to the large amount of predictor variables in this data set (92 food categories and measurements!), multicollinearity is bound to exist, which could severely hurt the model fitting process. To counter this issue, a preliminary variable elimination step was taken to remove "redundant" food categories/measurements.

Prior to the fitting of any model, it is also important to check if the data contains any missing values. In this data set, there were no missing values among the predictor variables, but there were some missing values in the response variables, ie. there were no reports on the confirmed, death, recovered, and active case rates on some of the countries in this data set. After the removal of those rows, only 163 rows of data remained in the Confirmed, Death, and Recovered cases data frames, while only 161 rows of data remained in the Active cases data frame.

2.1.1 Preliminary Variable Elimination

There were three criteria used to determine which predictor variables should be removed prior to model fitting. The goal here is to avoid the issue of multicollinearity in the later stages.

The first thing that was done was to select "representative" variables among the same food categories with different measurements. This is not only because of the multicollinearity issue caused by large correlation values, but also because if one of the variables is a good predictor of COVID-19 cases, the other highly correlated variables do not help much in further explaining the model. An example of the selection of "representative" variables is shown in Figure 1. In this step alone, the total number of predictor variables have been reduced from 92 to 51.

Another thing that was done to reduce the number of predictor variables prior to the model fitting step is to remove predictor variables with (close to) no distribution. Predictor variables like the fat content in alcoholic beverages can be removed as shown in Figure 2 since it wouldn't explain anything about the variability of COVID-19 cases across the world. In this step, another 6 predictor variables have been removed from the pool of predictor variables.

Lastly, predictor variables with high variance inflation factors (VIF) are also subsequently removed. While these variables might not be pairwise correlated with each predictor variable, they are highly intercorrelated with the rest of the predictor variables, which can also lead to multicollinearity among predictor variables. Some of the predictor variables with the largest VIF values that were removed are shown in Table 1. With those predictor variables removed, the number of predictor variables left prior to the start of the model fitting step is 41.

2.1.2 Variable Transformation

To determine if the response variables require transformation, their histograms are plotted as shown in Figure 3. If those histograms do not resemble normal distributions, it means that the first assumption stated in the previous section are not satisfied. To counter this issue, the Box-Cox procedure (Figure 4) is carried out to determine the best transformation of the response variables. In this case, the best transformation for all the response variables are all log transformations. This is subsequently carried out prior to the start of the model fitting process and the distribution of the transformed response variables are shown in Figure 5.

To determine if the predictor variables require transformation, each of the predictor variables were plotted against each response variables. A transformation of a predictor variable is needed if there is a significant non-linear relationship between the predictor variable and a transformed response variable. Since this is not the case in this data set, no transformation of the predictor variables is necessary.

2.1.3 Top Correlated Food Categories with the Number of Cases

Table 2 shows each response variable's top 5 correlated predictor variables. The goal of learning about the correlations between the response and predictor variables here is to observe how closely related they are. Note that the top pairwise correlated variables have absolute correlation values between 0.4 to 0.5, indicating that the number of COVID-19 cases has some correlation with some food categories, ie. the protein quantities of miscellaneous food and tree nuts, the calorie content in eggs, and animal fat.

2.2 Model Fitting and Validation

After gaining a slight understanding of the summary of the distribution of the data and removing redundant variables, the next step was to fit models that are able to do great jobs in predicting the rate of confirmed, death, recovered, and active cases respectively. There are two types of procedures to search for good models: the best subset selection and the stepwise regression methods. The best subset selection is guaranteed to work better since it parses through all the possible models with the given set of predictor variables, as compared to the stepwise regression method which is a greedy search. However, the best subset selection method is not a viable option in this case due to the large number of predictor variables. If the aim is to find the best model containing only linear terms and possible interaction terms, there would be a total $41 + \binom{41}{2} = 861$ variables to choose from, and a total of 2^{861} models to choose from. Not only is the process going to take infinitely long, the process is also not feasible since the number of variables could end up being more than the number of rows in our data.

Instead, the approach taken in this data analysis was to find the best model using the stepwise regression starting from the null model, ie. $\hat{Y}_i = \hat{\beta}_0$. Two different models were then built using the AIC and BIC criterion respectively. The AIC criterion, known to be less conservative as compared to the BIC criterion, produced larger models when fitting models for all the response variables, ie. confirmed, death, recovered, and active cases.

Before the model fitting procedure was carried out, a (70%-30%) random splitting of data was carried out into training and validation sets for all four data frames. Only the training sets will be used for the model fitting process, whereas the validation sets will be used to validate the prediction abilities of the "best" model acquired from the training sets.

2.2.1 Internal Validation

To determine the best model for each of the response variables, four tables were created as shown in Table 3 comparing the predictive abilities and the model bias and variability of the AIC and BIC selected models. The problem faced in calculating one of the metrics, Mallows's C_p , was the inability to fit a full model, ie. a model with all predictor variables and their two-way interactions, due to the large number of predictor variables here. This issue was handled by assuming that the model that contains all the variables of both the AIC and BIC models as the full model. In the case of model fittings for confirmed and recovered cases,

the AIC models were considered as the full models since they contain all the predictor variables present in their BIC counterparts.

For the confirmed cases models (Table 3(a)), the SSE and MSE values are found to be quite small in both the AIC and BIC models, which suggests a very good model fit in this case. This is backed up by their respective R_a^2 values. With the R_a^2 being higher than 0.6, it suggests that a huge variability of the number confirmed cases can be explained by the selected predictor variables. The $Press_p$ values for models are also very close to their respective SSE values, which suggests strong predictive abilities of both models. What allows for the AIC model to edge out in this case is its C_p value being significantly smaller than its BIC counterpart, suggesting very little model bias.

A similar conclusion can be made based on the AIC and BIC models for active cases, shown in Table 3(d). However, the difference between the models for active cases and the models for confirmed cases here is that they have higher R_a^2 values while also having significantly higher SSE, MSE, and $Press_p$ values. This suggests that while the predictor variables can explain the number of active cases better, its predictive ability is however significantly lower than that of the confirmed cases models.

On the other hand, the models explaining the rate of recovered cases (Table 3(c)) did not perform well. They all have extremely large SSE and MSE values, and very small R_a^2 values, suggesting that diet does not really explain the number of recovered cases well.

The AIC and BIC models for the number of death cases (Table 3(b)) is somewhat interesting. Unlike the aforementioned models, the difference in sizes between these two models is huge (a difference of 30 predictor variables/interaction terms). Here, the SSE, MSE, R_a^2 , and C_p values are very significantly different. This can be caused by the overfitting of the AIC model and the underfitting of the BIC model. However, the $Press_p$ values are not much different between the two models. As the $Press_p$ value for the AIC model is very far away from the SSE, while the opposite is true for the BIC model, it seems that the BIC model has a better predictive ability.

Based on the results obtained so far, it seems that all the AIC models are the better choices. However, a secondary step of model validation on the validation data is required to ensure that the models are not overfitting onto the training data.

2.2.2 External Validation

The goal of this procedure is to ensure that the models do not only work well on the training data, ie. the models selected should be able to generalize the trend of cases across all the countries in the world instead of the ones in the training data. The metrics used in the external validation process is simple, which is comparing the mean squared prediction error (MSPE) of the validation data with that of the training data. Mathematically, MSPE is defined as

$$MSPE = \frac{\sum_{j=1}^m (Y_j - \hat{Y}_j)^2}{m}$$

where m is the number of data either in the validation set or the training set. Essentially the MSPE of the training data can be thought of as SSE/n . The results of the MSPE comparisons between the AIC and BIC models are shown in Table 4.

From the table, it is obvious that in the confirmed cases models and active cases models, the $MSPE_{valid}$ produced by the BIC models are far closer to their respective $MSPE_{train}$ than their AIC counter parts. This is due to the result of overfitting in the training data set. In other words, using the AIC models on the training data set produces more accurate predictions only on the training data, not in general. In selecting the best model, however, requires the model to generalize well, hence the BIC models are selected in both of these cases.

In fact, this scenario also happens at an even bigger margin in the death cases models. Here, the $MSPE_{valid}$ is significantly larger than the $MSPE_{train}$ for the AIC model, whereas the opposite is true for the BIC model. This is due to the severe overfitting of the AIC model and the severe underfitting of the BIC model, which are both not ideal. However, the BIC model is selected as it is able to generalize better than

the AIC model.

As for the recovered cases models, since the $MSPE_{train}$ values are similar for both the AIC and BIC models, the model with the smaller $MSPE_{valid}$ value is selected since it is able to generalize better. Therefore, the BIC model is once again selected.

In summary, all the BIC models seem to generalize better than their AIC counterparts. What remains right now is that these BIC models satisfy the necessary model assumptions before fully concluding the best model.

2.3 Model Diagnostics

The goal of this procedure is to ensure the following:

- The BIC model residuals have equal variance.
- The BIC model residuals are normally distributed.
- Analyze possible outliers and decide if it is necessary to have them removed.

The first two assumptions above can be analyzed using the residuals vs fitted values plot and the QQ-plots respectively, as shown in Figure 6. From the set of plots, we find that the confirmed cases and active cases models are satisfy those assumptions reasonably well. Bar a few possible outliers, we can simply assume that the residuals of those models have equal variances and are normally distributed.

On the other hand, the same cannot be said about the recovered cases and death cases plots. Interestingly, the residuals vs fitted values plot for death cases seems to come in a pattern of two straight lines. Its QQ plot also showed a significant deviation from the normal line at theoretical quantiles less than -1. A similar situation also arises from the recovered cases model. This could be the reason why the model fit for the confirmed cases and active cases models are significantly better than the death cases and recovered cases models.

Moving on, to find the possible outliers of each model, the residuals, leverage values and Cook's distance are analyzed. One way to do so is by looking at the leverage vs residuals plots, shown in Figure 7. Essentially, points with large leverage values indicate outlying predictor variables, and points with large absolute values of standardized residuals indicated outlying response variables. If any of these points lie outside of the red dotted lines (aka. the line of Cook's distance), it becomes a strong indication that the points are outliers. Based on the figure, it can be seen that case 165 (Vanuatu) is a significant outlier for both the recovered cases and active cases models, whereas case 15 (Sri Lanka) is also a significant outlier for the recovered cases models.

Another thing that can be checked is the Cook's distance plot, as shown in Figures 8 to 11. Those plots make it easier to determine if a data point has a significantly huge Cook's distance value. Typically, a Cook's distance $D_i > \frac{4}{n-p}$ is used to indicate whether a data point i is a potential influential case. In more conservative cases, $D_i > 1$ is used as the cutoff instead. However, due to the small amount of data here, setting the cutoff to be $\frac{4}{n-p}$ means too many data points will be thrown away, whereas a cutoff of 1 would be too conservative. Instead, the "cutoff" used in this analysis is based off the Cook's distance plots, where data points that look to have Cook's distance significantly larger than the rest of the data points will be considered as potential outliers.

With these potential outliers in mind, the next thing to do is test if they really do heavily influence their respective models. This is done by calculating the percentage change between the fitted values based on the models with and without the potential outliers, ie.

$$\alpha_i = \frac{Y_i - Y_{i,(-J)}}{Y_i} \times 100\%, \text{ where } J \text{ are potential outliers}$$

Using the above equation, it can be found that the percentage changes between fitted values based on models with and without potential outliers are very significant (up to 30% for confirmed cases model, up to 164.9% for death cases model, and up to 260.9% for recovery cases model). Therefore, the potential outliers

were subsequently removed. On the other hand, this was not the case for the active cases model (close to 0% change) and hence the potential outliers were retained.

After eliminating the outliers from the models, it is time to move into the last step of this data analysis, which is fitting the final models.

2.4 Final Models

The final step of this data analysis is simply combining the training data (without outliers) and the validation data, and refit the data into the same models. The final models are shown as follows:

Confirmed cases model:

$$Y = \exp(-2.8 + 0.2X_1 + 0.17X_2 - 0.22X_3 - 0.26X_4 - 11.71X_5 - 0.22X_6 + 0.46X_7 + 2.67X_8 - 1.46X_6X_8)$$

where Y is rate of confirmed cases of COVID-19, X_1 is Miscellaneous (Protein), X_2 is Sweeteners (Calories), X_3 is Meat (Calories), X_4 is Pulses (Calories), X_5 is Stimulants (Protein), X_6 is Oilcrops (Calories), X_7 is Fruits (Protein), and X_8 is Eggs (Calories). The R_a^2 value for this model is 0.5631, which suggests that this model is a relatively good predictor of COVID-19 case rates of each country.

Death cases model:

$$Y = \exp(-14.86 + 0.53X_1 + 0.24X_2 + 3.12X_3)$$

where Y is rate of death cases of COVID-19, X_1 is Miscellaneous (Protein), X_2 is Vegetable Oil (Fat), and X_3 is Treenuts (Calories). The low R_a^2 value of 0.18 in this model indicates that diet only explains about 20% of the variability of COVID-19 deaths.

Recovery cases model:

$$Y = \exp(0.68 - 2.14X_1 - 90.13X_2 + 81.44X_1X_2)$$

where Y is rate of recovered cases of COVID-19, X_1 is Oilcrops (Calories), and X_2 is Vegetal Products (Protein). The low R_a^2 value of 0.21 in this model once again indicates that diet is not a good predictor of COVID-19 recovery.

Active cases model:

$$Y = \exp(-4.33 - 1.67X_1 + 0.21X_2 + 0.19X_3 + 0.21X_1X_2)$$

where Y is rate of active cases of COVID-19, X_1 is Oilcrops (Calories), X_2 is Miscellaneous (Protein), and X_3 is Alcoholic Beverages (Quantity). The low R_a^2 value of 0.58 in this model suggests that diet can explain about 58% of the variability of COVID-19 active cases.

3 Conclusion

In summary, a country's COVID-19 confirmed and active cases can somehow be explained relatively well by food categories such as the calorie contents of oilcrops, and the protein content in infant food and miscellaneous food. On the other hand, the same cannot be said about the death and recovered cases. This could be due to the fact that these models do not satisfy the necessary model assumptions of having equal variance and normally distributed residuals. However, it is also important to note that active cases are simply cases that have not had an outcome, and hence the model on active cases should only be taken as a grain of salt, or rather a model that follows closely with the model of confirmed cases. In short, the model predicting the rate of confirmed cases is the only reliable model obtained from this analysis.

However, recall that this model only talks about the correlation between food categories and the rate of confirmed cases. There is no evidence to suggest that a country's diet has an effect on the spread of COVID-19. Additionally, there are also many other factors causing the spread of COVID-19 that are totally uncorrelated with diet, eg. how active the general public are, the preventive measures implemented by the countries, density of population etc.

4 Appendices

4.1 Appendix 1: Figures and tables

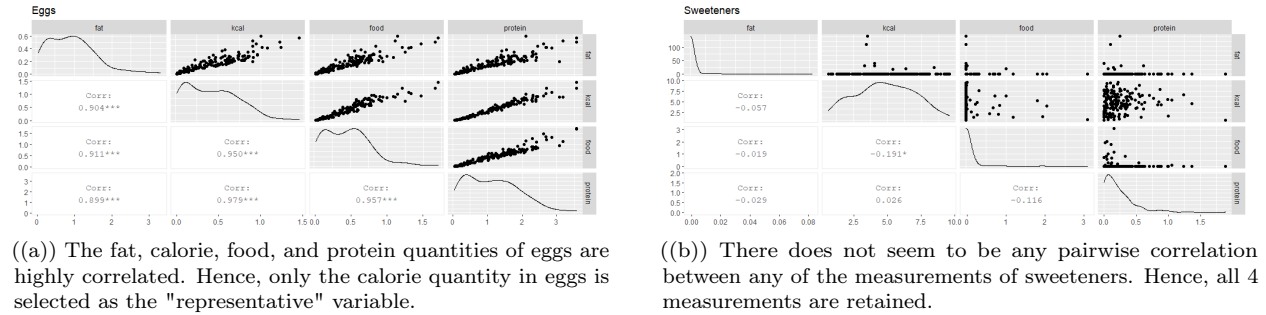


Figure 1: Scatterplot Matrices of Food Categories with Different Measurements

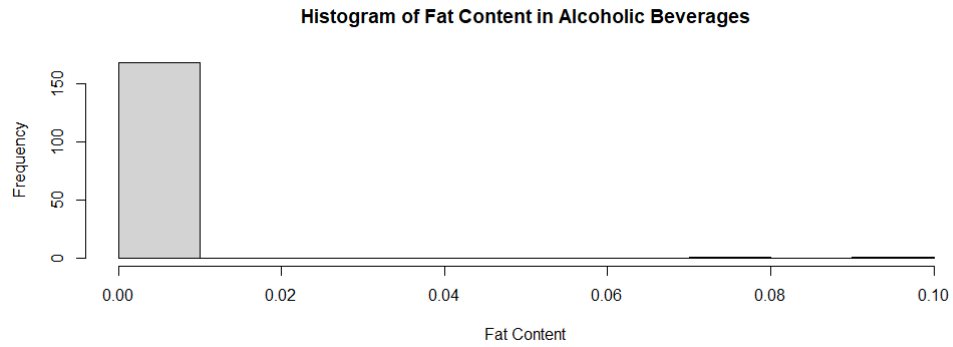
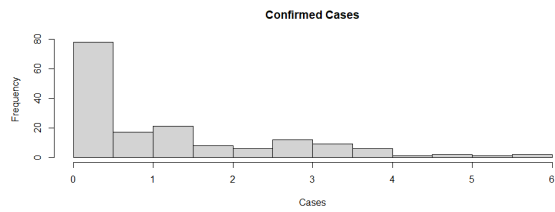


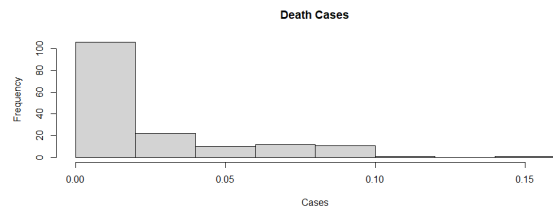
Figure 2: Variables like the fat quantity in alcoholic beverages do not have distributions, and thus can be removed.

Table 1: The variables above have extremely high VIF values, indicating very high intercorrelation with the rest of the predictor variables, and should be subsequently removed.

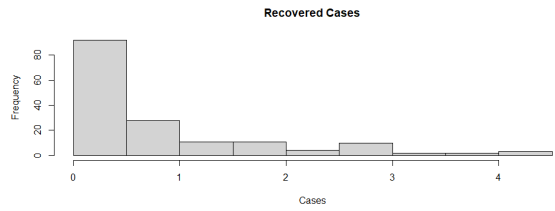
Variable	VIF
Animal Fats (Food Quantity)	2,096,672.147
Vegetables (Food Quantity)	2,094,971.233
Animal Products (Calories)	238,996.825
Vegetable Products (Calories)	237.921.552



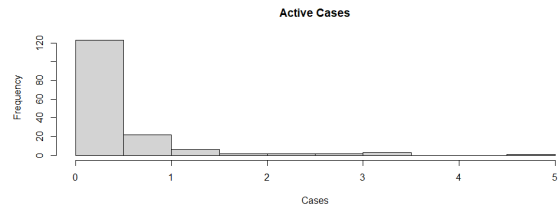
((a)) Histogram of Confirmed Cases



((b)) Histogram of Death Cases

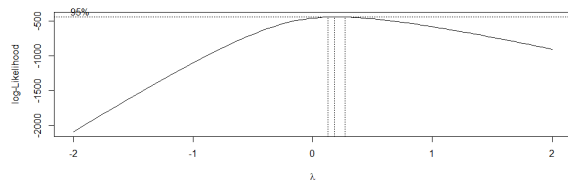


((c)) Histogram of Recovered Cases

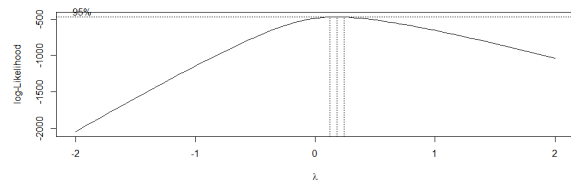


((d)) Histogram of Active Cases

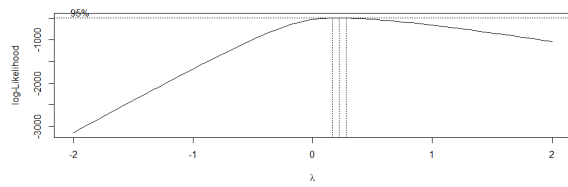
Figure 3: Histograms of Response Variables



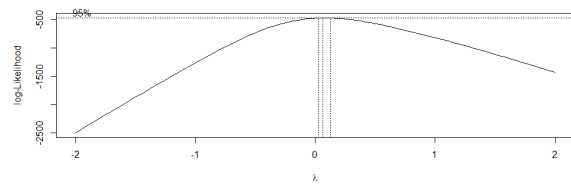
((a)) Box-Cox of Confirmed Cases



((b)) Box-Cox of Death Cases

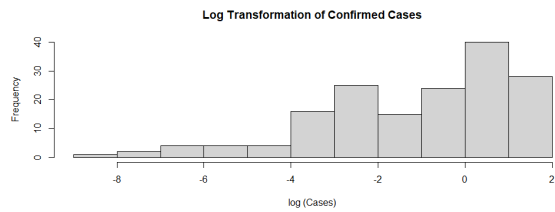


((c)) Box-Cox of Recovered Cases

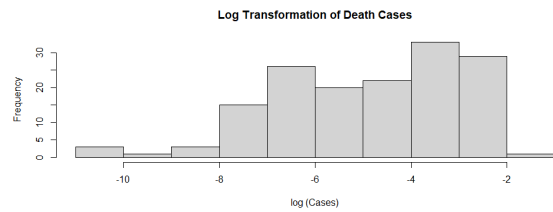


((d)) Box-Cox of Active Cases

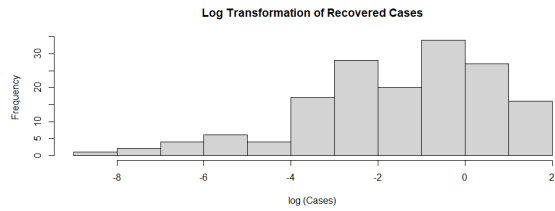
Figure 4: Based on the above plots, it can be seen that the log transformation is the best way to go for all of the response variables.



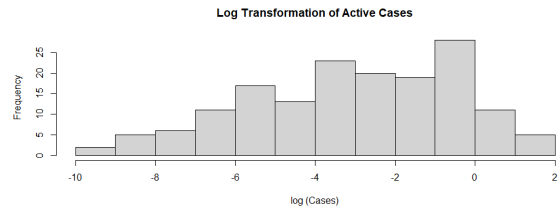
((a)) Histogram of log transformed Confirmed Cases



((b)) Histogram of log transformed Death Cases



((c)) Histogram of log transformed Recovered Cases



((d)) Histogram of log transformed Active Cases

Figure 5: After carrying out the log transformation on each of the response variable, their histograms are significantly closer to being normally distributed. Note: extremely small values like 10^{-11} were added into each data point of death, recovered and active cases due to the existence of zero values.

Table 2: Top 5 highly correlated predictor variables with their respective response variables. Note: The letter codes are P - protein, Q - food quantity, C - calories, F - fat.

Confirmed		Deaths		Recovered		Active	
Misc. (P)	0.59	Misc. (P)	0.45	Stimulants (F)	0.45	Misc. (P)	0.5
Treenuts (P)	-0.52	Eggs (C)	0.43	Misc. (P)	0.42	Animal Prod. (Q)	0.45
Eggs (C)	0.49	Treenuts (P)	-0.42	Treenuts (P)	-0.39	Animal Fat (F)	0.43
Animal Fat (F)	0.47	Animal Fat (F)	0.39	Eggs (C)	0.38	Treenuts (P)	-0.42
Stimulants (F)	0.44	Cereal (F)	-0.35	Offals (P)	-0.32	Veg. Prod. (P)	0.42

Table 3: Analysis of the AIC and BIC models to determine the best model

	AIC Model	BIC Model
p	16	10
SSE	150.12	189.24
MSE	1.53	1.82
R^2	0.73	0.66
R_a^2	0.69	0.63
C_p	16	29.54
$Press_p$	218.35	242.87

(a): Confirmed cases models

	AIC Model	BIC Model
p	34	4
SSE	856.33	3774.95
MSE	10.7	34.32
R^2	0.82	0.2
R_a^2	0.74	0.17
C_p	35.5	253.29
$Press_p$	3978.05	4017.27

(b): Death cases models

	AIC Model	BIC Model
p	5	4
SSE	2203.47	2281.23
MSE	20.22	20.74
R^2	0.3	0.27
R_a^2	0.27	0.25
C_p	5	6.85
$Press_p$	2893.88	2943.18

(c): Recovered cases models

	AIC Model	BIC Model
p	13	5
SSE	302.65	466.12
MSE	3.06	4.35
R^2	0.78	0.66
R_a^2	0.75	0.65
C_p	-41.25	-33.08
$Press_p$	539	711.68

(d): Active cases models

Table 4: Comparison of $MSPE$ values between AIC and BIC models

	AIC Model	BIC Model
$MSPE_{valid}$	7.51	5.65
$MSPE_{train}$	1.32	1.66

(a): Confirmed cases models

	AIC Model	BIC Model
$MSPE_{valid}$	67.15	12.43
$MSPE_{train}$	7.51	33.11

(b): Death cases models

	AIC Model	BIC Model
$MSPE_{valid}$	10.25	8.93
$MSPE_{train}$	19.33	20.01

(c): Recovered cases models

	AIC Model	BIC Model
$MSPE_{valid}$	4.2	5.66
$MSPE_{train}$	2.7	4.16

(d): Active cases models

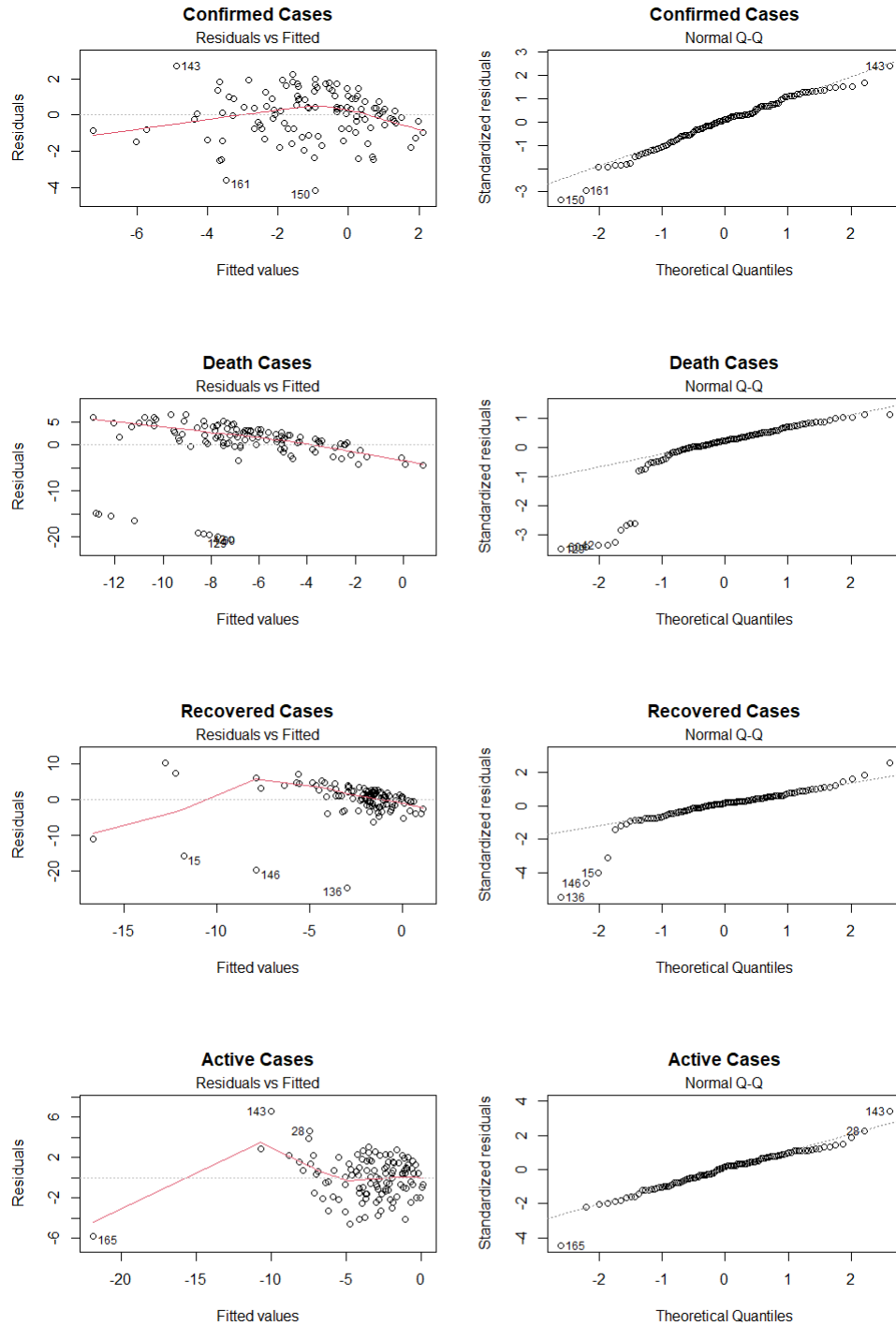


Figure 6: The plots on the left are the residuals vs fitted values plots for each response variable using the BIC models, whereas the plots on the right are their respective QQ plots. Bar a few possible outliers, it is obvious to see that the active cases and the confirmed cases models satisfy the normal error model assumptions.

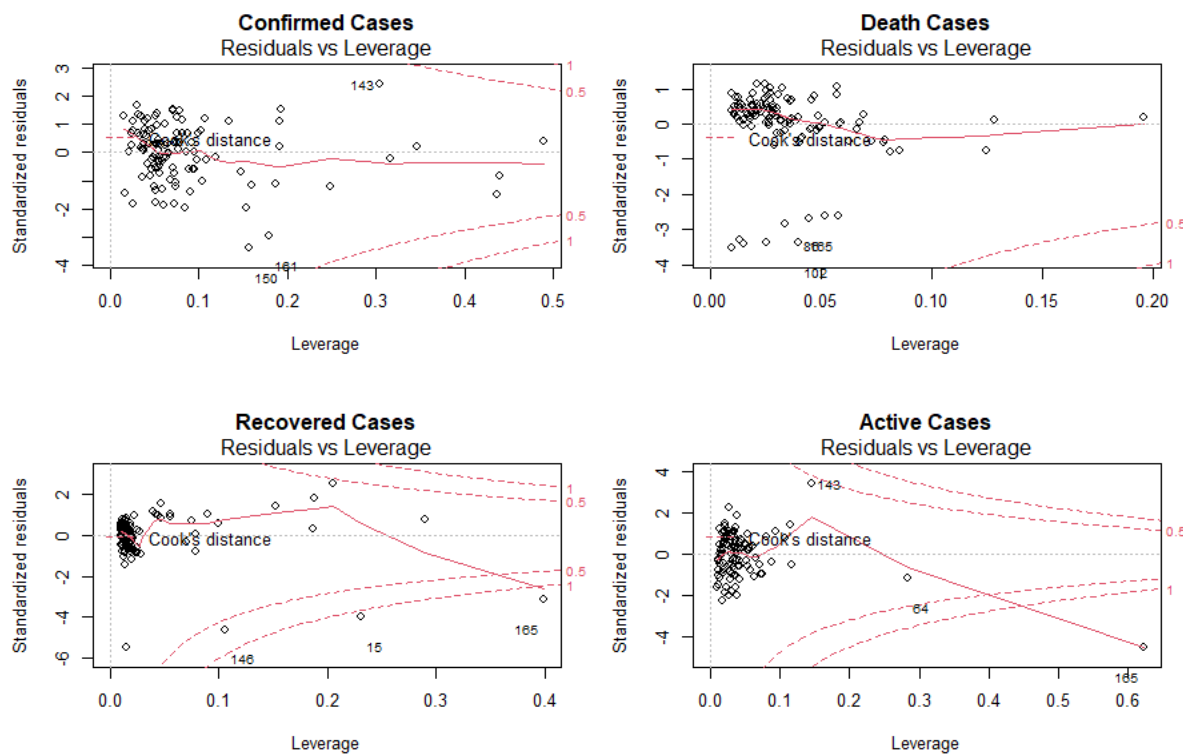
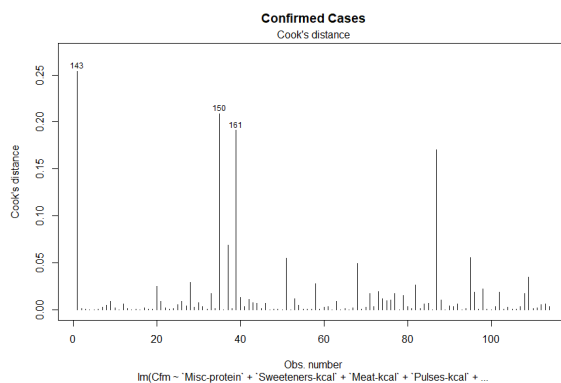
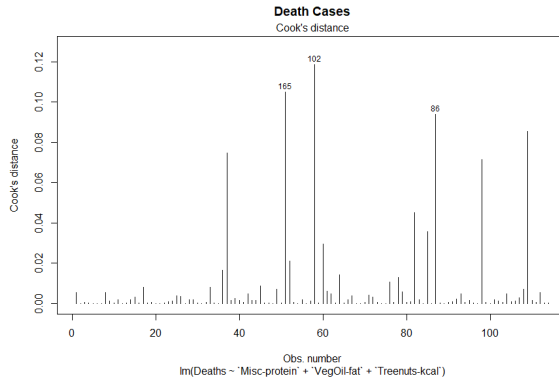


Figure 7: Residual vs Leverage plots for all the response variables



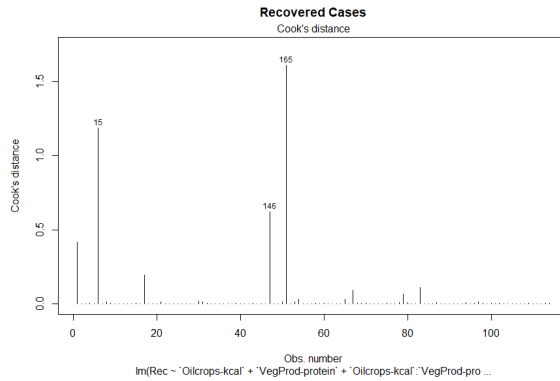
Index	Country	Cook's
143	Sri Lanka	0.25
150	Thailand	0.21
161	Tanzania	0.19
86	Laos	0.17

Figure 8: Significantly large values of Cook's distance in the confirmed cases model



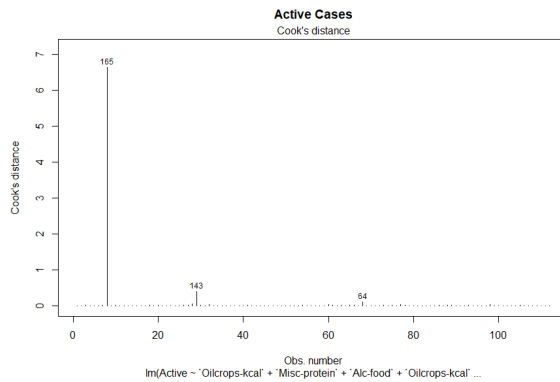
Index	Country	Cook's
102	Mongolia	0.118
165	Vanuatu	0.105
86	Loas	0.094
25	Combodia	0.086
42	Dominica	0.075
151	Timor-Leste	0.071

Figure 9: Significantly large values of Cook's distance in the death cases model



Index	Country	Cook's
165	Vanuatu	1.61
15	Belgium	1.19
146	Sweden	0.62
143	Sri Lanka	0.41

Figure 10: Significantly large values of Cook's distance in the recovered cases model



Index	Country	Cook's
165	Vanuatu	6.62
143	Sri Lanka	0.39

Figure 11: Significantly large values of Cook's distance in the active cases model

4.2 Appendix 2: R Codes

```

1 #####
2 # load all data
3 fat = read.table("Fat_Supply_Quantity_Data.csv", header = TRUE, sep = ",")

```

```

4 kcal = read.table("Food_Supply_kcal_Data.csv", header = TRUE, sep = ",")
5 food = read.table("Food_Supply_Quantity_kg_Data.csv", header = TRUE, sep = ",")
6 protein = read.table("Protein_Supply_Quantity_Data.csv", header = TRUE, sep = ",")
7
8 #####
9 # renaming variables
10 dfAsString = function(x){
11   deparse(substitute(x))
12 }
13 colRename = function(x, name) {
14   names(x) = c("Country", "Alc", "AnimalProd", "AnimalFats",
15               "AqProd", "Cereal", "Eggs", "Seafood", "Fruits",
16               "Meat", "Misc", "Milk", "Offals", "Oilcrops",
17               "Pulses", "Spices", "StRt", "Stim", "SugarCrops",
18               "Sweeteners", "Treenuts", "VegProd", "VegOil", "Veg",
19               "Obesity", "Undernourished", "Cfm", "Deaths", "Rec",
20               "Active", "Pop", "Unit")
21   names(x)[2:24] = paste(names(x)[2:24], name, sep = "-")
22   return(x)
23 }
24 fat = colRename(fat, dfAsString(fat))
25 kcal = colRename(kcal, dfAsString(kcal))
26 food = colRename(food, dfAsString(food))
27 protein = colRename(protein, dfAsString(protein))
28
29 #####
30 ## exploratory data analysis
31 # check correlation between measurements of same food categories
32 library(GGally)
33 for (i in 2:24){
34   tempDf = as.data.frame(cbind(fat = fat[,i], kcal = kcal[,i], food = food[,i], protein =
35                               protein[,i]))
36   title = strsplit(names(fat)[i], "-")[1]
37   print(ggpairs(tempDf, title = title,
38                 upper = list(continuous = "points"), lower = list(continuous = "cor"),
39                 progress = FALSE))
40   print(colSums(cor(tempDf)))
41 }
42 # remove unnecessary columns
43 redDf = cbind(fat[c(1,2,4,6,9,11,12,18,19,20,23)],
44               kcal[c(3,5,6,7,8,10,12,14:17,20:24)],
45               food[c(2:4,13,18,19,20,22,24)],
46               protein[c(4,9,11:31)])
47 dim(redDf)
48
49 # remove X variables with improper distribution
50 # quick view of histograms for all X variables
51 for (i in 2:52){
52   hist(redDf[,i], xlab = names(redDf)[i], main = names(redDf)[i])
53 }
54
55 # remove unnecessary columns
56 drops = c("Alc-fat", "SugarCrops-fat", "Sweeteners-fat", "AqProd-kcal", "Sweeteners-food", "
57           SugarCrops-protein")
58 redDf = redDf[,!(names(redDf) %in% drops)]
59
60 # remove X variables with extremely high VIF
61 VIF = diag(solve(cor(redDf[,2:46])))
62 ord = order(VIF, decreasing = TRUE)
63 VIF[ord]
64
65 # remove unnecessary columns
66 drops = c("AnimalFats-food", "Veg-food", "AnimalProd-kcal", "VegProd-kcal")
67 redDf = redDf[,!(names(redDf) %in% drops)]

```

```

67 dim(redDf)
68
69 # checking correlation between Y variables
70 cor(na.omit(redDf)[,45:48])
71 dim(na.omit(redDf))
72
73 # best correlation between X and Y variables
74 bestCorDf = data.frame(Ord = 1:10)
75 for (i in 45:48){
76   corrXY = cor(na.omit(redDf[,c(2:42,i)]))[,42]
77   ord = order(abs(corrXY), decreasing = TRUE)
78   bestCorDf = cbind(bestCorDf,
79                     names(corrXY[ord[2:11]]),
80                     corrXY[ord[2:11]])
81 }
82 names(bestCorDf) = c("Ord", "varCfm", "corCfm", "varD", "corD", "varRec", "corRec", "varAct",
83                     "corAct")
84 row.names(bestCorDf) = NULL
85 bestCorDf
86
87 # distribution of Y variables
88 hist(redDf$Cfm, xlab = "Cfm", main = "Confirmed Cases")
89 hist(redDf$Death, xlab = "Death", main = "Death")
90 hist(redDf$Rec, xlab = "Rec", main = "Recovered")
91 hist(redDf$Active, xlab = "Active", main = "Active")
92
93 # box cox transformation
94 library(MASS)
95 fit = lm(Cfm ~ ., data = redDf[,c(2:42, 45)]); boxcox(fit)
96 fit = lm(I(Deaths+.00001) ~ ., data = redDf[,c(2:42, 46)]); boxcox(fit)
97 fit = lm(I(Rec+.00001) ~ ., data = redDf[,c(2:42, 47)]); boxcox(fit)
98 fit = lm(I(Active+.00001) ~ ., data = redDf[,c(2:42, 48)]); boxcox(fit)
99
100 # log transformation of Y variables
101 hist(log(fat$Cfm), xlab = "log(Cfm)", main = "log(Confirmed Cases)")
102 hist(log(fat$Death), xlab = "log(Death)", main = "log(Death)")
103 hist(log(fat$Rec), xlab = "log(Rec)", main = "log(Recovered)")
104 hist(log(fat$Active), xlab = "log(Active)", main = "log(Active)")
105
106 #####
107 ## model fitting
108 # omit NA's and form 4 separate data frames
109 cfmDf = na.omit(redDf[,c(1:42, 45)])
110 deathDf = na.omit(redDf[,c(1:42,46)])
111 recDf = na.omit(redDf[,c(1:42,47)])
112 actDf = na.omit(redDf[,c(1:42,48)])
113
114 # check dimensions of each dataframe
115 dim(cfmDf)
116 dim(deathDf)
117 dim(recDf)
118 dim(actDf)
119
120 # transformation of Y variables
121 cfmDf$Cfm = log(cfmDf$Cfm)
122 deathDf$Deaths = log(deathDf$Deaths + 1e-12)
123 recDf$Rec = log(recDf$Rec + 1e-12)
124 actDf$Active = log(actDf$Active + 1e-12)
125
126 # split data into train and validation
127 set.seed(10)
128 ind1 = sample(1:dim(cfmDf)[1], size = 0.7*dim(cfmDf)[1])
129 ind2 = sample(1:dim(actDf)[1], size = 0.7*dim(actDf)[1])
130 # train
131 cfmDf.train = cfmDf[ind1,]

```

```

131 deathDf.train = deathDf[ind1,]
132 recDf.train = recDf[ind1,]
133 actDf.train = actDf[ind2,]
134 #valid
135 cfmDf.valid = cfmDf[-ind1,]
136 deathDf.valid = deathDf[-ind1,]
137 recDf.valid = recDf[-ind1,]
138 actDf.valid = actDf[-ind2,]
139
140 # check distribution of response variables in train and validation sets
141 hist(cfmDf.train$Cfm, main = "Train")
142 hist(cfmDf.valid$Cfm, main = "Valid")
143 hist(deathDf.train$Deaths, main = "Train")
144 hist(deathDf.valid$Deaths, main = "Valid")
145 hist(recDf.train$Rec, main = "Train")
146 hist(recDf.valid$Rec, main = "Valid")
147 hist(actDf.train$Active, main = "Train")
148 hist(actDf.valid$Active, main = "Valid")
149
150 # model fitting for confirmed cases
151 cfm.fit1 = lm(Cfm ~ 1, data = cfmDf.train[,-1])
152
153 # fitting AIC model
154 cfm.fit2 = stepAIC(cfm.fit1,
155                   scope = list(upper = lm(Cfm ~ ., data = cfmDf.train[,-1]),
156                                lower = cfm.fit1),
157                   direction = "both", trace = 0, k = 2)
158 summary(cfm.fit2)
159 cfm.fit4 = stepAIC(cfm.fit1,
160                   scope = list(upper = lm(Cfm ~ ('Misc-protein' + 'Sweeteners-kcal' + 'Meat
161                                           -kcal' +
162                                           'Pulses-kcal' + 'Stim-protein' + '
163                                           Oilcrops-kcal' + 'Fruits-protein' +
164                                           'Eggs-kcal' + 'Alc-food' + 'Treenuts -
165                                           kcal' + 'Spices-kcal' +
166                                           'Sweeteners-protein')^2,
167                                data = cfmDf.train[, -1]),
168                                lower = cfm.fit1),
169                   direction = "both", trace = 0, k = 2)
170 summary(cfm.fit4)
171 # fitting BIC model
172 cfm.fit3 = stepAIC(cfm.fit1,
173                   scope = list(upper = lm(Cfm ~ .^2, data = cfmDf.train[,-1]),
174                                lower = cfm.fit1),
175                   direction = "both", trace = 0, k = log(nrow(cfmDf.train)))
176 summary(cfm.fit3)
177
178 # model fitting for death cases
179 death.fit1 = lm(Deaths ~ 1, data = deathDf.train[,-1])
180 # fitting AIC model
181 death.fit2 = stepAIC(death.fit1,
182                     scope = list(upper = lm(Deaths ~ ., data = deathDf.train[,-1]),
183                                    lower = death.fit1),
184                     direction = "both", trace = 0, k = 2)
185 summary(death.fit2)
186 death.fit4 = stepAIC(death.fit1,
187                     scope = list(upper = lm(formula = Deaths ~ ('Misc-protein' + 'Treenuts -
188                                           kcal' + 'Meat-kcal' +
189                                           'Eggs-kcal' + 'Spices -
190                                           kcal' + 'Pulses-protein' + 'Stim-protein' +
191                                           'Oilcrops-kcal' + 'Alc -
192                                           food' + 'Veg-protein' + 'Veg-kcal' +
193                                           'Misc-fat')^2,
194                                    data = deathDf.train[, -1]),
195                                    lower = death.fit1),

```



```

190         direction = "both", trace = 0, k = 2)
191 summary(death.fit4)
192 # fitting BIC model
193 death.fit3 = stepAIC(death.fit1,
194                     scope = list(upper = lm(Deaths ~ .^2, data = deathDf.train[,-1]),
195                                   lower = death.fit1),
196                     direction = "both", trace = 0, k = log(nrow(deathDf.train)))
197 summary(death.fit3)
198
199 # model fitting for recovery cases
200 rec.fit1 = lm(Rec ~ 1, data = recDf.train[,-1])
201 # fitting AIC model
202 rec.fit2 = stepAIC(rec.fit1,
203                   scope = list(upper = lm(Rec ~ .^2, data = recDf.train[,-1]),
204                                   lower = rec.fit1),
205                   direction = "both", trace = 0, k = 2)
206 summary(rec.fit2)
207 # fitting BIC model
208 rec.fit3 = stepAIC(rec.fit1,
209                   scope = list(upper = lm(Rec ~ .^2, data = recDf.train[,-1]),
210                                   lower = rec.fit1),
211                   direction = "both", trace = 0, k = log(nrow(recDf.train)))
212 summary(rec.fit3)
213
214 # model fitting for active cases
215 act.fit1 = lm(Active ~ 1, data = actDf.train[,-1])
216 # fitting AIC model
217 act.fit2 = stepAIC(act.fit1,
218                   scope = list(upper = lm(Active ~ ., data = actDf.train[,-1]),
219                                   lower = act.fit1),
220                   direction = "both", trace = 0, k = 2)
221 summary(act.fit2)
222 act.fit4 = stepAIC(act.fit1,
223                   scope = list(upper = lm(formula = Active ~ ('Oilcrops-kcal' + 'Misc-
224                                           protein' + 'Meat-kcal' +
225                                           'Treenuts-protein' + '
226                                           Pulses-kcal' + 'Alc-food' + 'Sweeteners-protein')^2,
227                                   data = actDf.train[, -1]),
228                                   lower = act.fit1),
229                   direction = "both", trace = 0, k = 2)
230 summary(act.fit4)
231 # fitting BIC model
232 act.fit3 = stepAIC(act.fit1,
233                   scope = list(upper = lm(Active ~ .^2, data = actDf.train[,-1]),
234                                   lower = act.fit1),
235                   direction = "both", trace = 0, k = log(nrow(actDf.train)))
236 summary(act.fit3)
237
238 #####
239 ## model diagnostics
240 # define functions
241 # compute sse, mse, R2, adjusted R2, Cp, Pressp statistics
242 modelDiag = function(obj, sigma2){
243   n = length(obj$residuals)
244   p = length(obj$coefficients)
245   h = influence(obj)$hat
246   sse = sum(obj$residuals^2)
247   mse = sse/obj$df.residual
248   r2 = summary(obj)$r.squared
249   adjr2 = summary(obj)$adj.r.squared
250   cp = sse/sigma2 - (n - 2*p)
251   pressp = sum((obj$residuals/(1-h))^2)
252   return (c(p = p, SSE = sse, MSE = mse, R2 = r2, Adj.R2 = adjr2, Cp = cp, PRESSp = pressp))
253 }

```

```

253 getMspe = function(fit, data){
254   actual = data[,ncol(data)]
255   pred = predict(fit, data)
256   res = actual - pred
257   m = nrow(data)
258   return (c(MSPE = sum(res^2)/m))
259 }
260
261 modelValid = function(data, fit1, fit2){
262   avg.sse1 = sum(fit1$residuals^2)/nrow(fit1$model)
263   avg.sse2 = sum(fit2$residuals^2)/nrow(fit2$model)
264   return (data.frame(AIC.valid = c(getMspe(fit1, data), 'SSE/n' = avg.sse1),
265                      BIC.valid = c(getMspe(fit2, data), 'SSE/n' = avg.sse2)))
266 }
267
268 # internal and external validation for confirmed cases
269 # confirmed cases
270 cfm.fitfull = lm(formula = Cfm ~ 'Misc-protein' + 'Oilcrops-kcal' + 'Stim-protein' +
271                  'Fruits-protein' + 'Pulses-kcal' + 'Meat-kcal' + 'Eggs-kcal' +
272                  'Sweeteners-kcal' + 'Alc-food' + 'Treenuts-kcal' + 'Oilcrops-kcal': 'Eggs-
273                  kcal' +
274                  'Eggs-kcal': 'Sweeteners-kcal' + 'Meat-kcal': 'Alc-food' +
275                  'Fruits-protein': 'Pulses-kcal' + 'Misc-protein': 'Alc-food', data = cfmDf.
276                  train[, -1])
277 cfm.fit4.diag = modelDiag(cfm.fit4, sum(cfm.fitfull$residuals^2)/cfm.fitfull$df.residual) #
278   train model
279 cfm.fit3.diag = modelDiag(cfm.fit3, sum(cfm.fitfull$residuals^2)/cfm.fitfull$df.residual)
280 cfm.diagComp = data.frame(AIC.train = cfm.fit4.diag, BIC.train = cfm.fit3.diag)
281 cfm.diagComp
282
283 cfm.diagComp.v = modelValid(cfmDf.valid[, -1], cfm.fit4, cfm.fit3) # valid model
284 cfm.diagComp.v
285
286 # internal and external validation for death cases
287 # death cases
288 death.fitfull = lm(Deaths ~ 'Misc-protein' + 'Meat-kcal' + 'Eggs-kcal' +
289                      'Spices-kcal' + 'Pulses-protein' + 'Stim-protein' + 'Oilcrops-kcal' +
290                      'Treenuts-kcal' + 'Veg-protein' + 'Misc-fat' + 'Alc-food' +
291                      'Eggs-kcal': 'Stim-protein' + 'Eggs-kcal': 'Pulses-protein' +
292                      'Meat-kcal': 'Pulses-protein' + 'Meat-kcal': 'Treenuts-kcal' +
293                      'Veg-protein': 'Misc-fat' + 'Meat-kcal': 'Veg-protein' + 'Misc-protein':
294                      'Misc-fat' +
295                      'Eggs-kcal': 'Oilcrops-kcal' + 'Spices-kcal': 'Treenuts-kcal' +
296                      'Meat-kcal': 'Misc-fat' + 'Meat-kcal': 'Alc-food' + 'Treenuts-kcal': 'Alc-
297                      food' +
298                      'Eggs-kcal': 'Veg-protein' + 'Pulses-protein': 'Oilcrops-kcal' +
299                      'Pulses-protein': 'Treenuts-kcal' + 'Misc-fat': 'Alc-food' +
300                      'Veg-protein': 'Alc-food' + 'Misc-protein': 'Oilcrops-kcal' +
301                      'Stim-protein': 'Alc-food' + 'Pulses-protein': 'Misc-fat' +
302                      'Stim-protein': 'Veg-protein' + 'Eggs-kcal': 'Misc-fat' + 'VegOil-fat',
303                      deathDf.train[, -1])
304 death.fit4.diag = modelDiag(death.fit4, sum(death.fitfull$residuals^2)/death.fitfull$df.
305                      residual) # train model
306 death.fit3.diag = modelDiag(death.fit3, sum(death.fitfull$residuals^2)/death.fitfull$df.
307                      residual)
308 death.diagComp = data.frame(AIC.train = death.fit4.diag, BIC.train = death.fit3.diag)
309 death.diagComp
310
311 death.diagComp.v = modelValid(deathDf.valid[, -1], death.fit4, death.fit3) # valid model
312 death.diagComp.v
313
314 # internal and external validation for recovered cases
315 # recovery cases
316 rec.fitfull = rec.fit2
317 rec.fit2.diag = modelDiag(rec.fit2, sum(rec.fitfull$residuals^2)/rec.fitfull$df.residual) #

```

```

train model
310 rec.fit3.diag = modelDiag(rec.fit3, sum(rec.fitfull$residuals^2)/rec.fitfull$df.residual)
311 rec.diagComp = data.frame(AIC.train = rec.fit2.diag, BIC.train = rec.fit3.diag)
312 rec.diagComp
313
314 rec.diagComp.v = modelValid(recDf.valid[, -1], rec.fit2, rec.fit3) # valid model
315 rec.diagComp.v
316
317 # internal and external validation for activ cases
318 # active cases
319 act.fitfull = lm(formula = Active ~ 'Misc-protein' + 'Offals-protein' + 'Meat-kcal' +
320                  'Eggs-kcal' + 'Spices-protein' + 'Stim-fat' + 'Meat-kcal': 'Spices-protein'
321                  'Misc-protein': 'Stim-fat' + 'Misc-protein': 'Meat-kcal' + 'Offals-food',
322                  data = actDf.train[, -1])
323 act.fit4.diag = modelDiag(act.fit4, sum(act.fitfull$residuals^2)/act.fitfull$df.residual) #
train model
324 act.fit3.diag = modelDiag(act.fit3, sum(act.fitfull$residuals^2)/act.fitfull$df.residual)
325 act.diagComp = data.frame(AIC.train = act.fit4.diag, BIC.train = act.fit3.diag)
326 act.diagComp
327
328 act.diagComp.v = modelValid(actDf.valid[, -1], act.fit4, act.fit3) # valid model
329 act.diagComp.v
330
331 ## checking for outlying cases
332 # function to calculate Cook's distance
333 getCooksDist = function(obj){
334   e = obj$residuals
335   p = length(obj$coefficients)
336   mse = sum(e^2)/obj$df.residual
337   h = influence(obj)$hat
338   return (c(e^2*h/(p*mse*(1-h)^2)))
339 }
340
341 # model diagnostics for confirmed cases
342 # residual analysis
343 plot(cfm.fit4, which = 1:2, main = "AIC Model")
344 plot(cfm.fit3, which = 1:2, main = "BIC Model")
345 # outlying cases
346 plot(cfm.fit4, which = 4:6)
347 plot(cfm.fit3, which = 4:6)
348 # calculate cook's distance and determine outliers
349 cooks = getCooksDist(cfm.fit3)
350 cfmOutliers = cbind(cfmDf.train[which(cooks > 4/(cfm.fit3$df.residual)), c(1,43)], CooksDist
= cooks[which(cooks > 4/(cfm.fit3$df.residual))])
351 cfmOutliers[order(cfmOutliers$CooksDist, decreasing = TRUE),]
352 # examine influence of possible outlier
353 cfm.fit3.rmOutlier = lm(formula = cfm.fit3$call$formula, data = cfmDf.train[-c
(143,150,161,86),])
354 per.change = abs((cfm.fit3$fitted.values - predict.lm(cfm.fit3.rmOutlier, cfmDf.train))/cfm.
fit3$fitted.values)*100
355 summary(per.change)
356 # remove outliers
357 cfmDf.train.rmOutlier = cfmDf.train[-c(143,150,161,86),]
358
359 # model diagnostics for death cases
360 # residual analysis
361 plot(death.fit4, which = 1:2, main = "AIC Model")
362 plot(death.fit3, which = 1:2, main = "BIC Model")
363 # outlying cases
364 plot(death.fit4, which = 4:6)
365 plot(death.fit3, which = 4:6)
366 # calculate cook's distance and determine outliers
367 cooks = getCooksDist(death.fit3)
368 deathOutliers = cbind(deathDf.train[which(cooks > 4/(death.fit3$df.residual)), c(1,43)],

```

```

    CooksDist = cooks[which(cooks > 4/(death.fit3$df.residual))]]
369 deathOutliers[order(deathOutliers$CooksDist, decreasing = TRUE),]
370 # examine influence of possible outlier
371 death.fit3.rmOutlier = lm(formula = death.fit3$call$formula, data = deathDf.train[-c
    (102,165,86,25,42,151),])
372 per.change = abs((death.fit3$fitted.values - predict.lm(death.fit3.rmOutlier, deathDf.train)
    )/death.fit3$fitted.values)*100
373 summary(per.change)
374 # remove outlier
375 deathDf.train.rmOutlier = deathDf.train[-c(102,165,86,25,42,151),]
376
377 # model diagnostics for recovered cases
378 # residual analysis
379 plot(rec.fit2, which = 1:2, main = "AIC Model")
380 plot(rec.fit3, which = 1:2, main = "BIC Model")
381 # outlying cases
382 plot(rec.fit2, which = 4:6)
383 plot(rec.fit3, which = 4:6)
384 # calculate cook's distance and determine outliers
385 cooks = getCooksDist(rec.fit3)
386 recOutliers = cbind(recDf.train[which(cooks > 4/(rec.fit3$df.residual)),c(1,43)], CooksDist
    = cooks[which(cooks > 4/(rec.fit3$df.residual))])
387 recOutliers[order(recOutliers$CooksDist, decreasing = TRUE),]
388 # examine influence of possible outlier
389 rec.fit3.rmOutlier = lm(formula = rec.fit3$call$formula, data = recDf.train[-c
    (165,15,146,143),])
390 per.change = abs((rec.fit3$fitted.values - predict.lm(rec.fit3.rmOutlier, recDf.train))/rec.
    fit3$fitted.values)*100
391 summary(per.change)
392 # remove outliers from dataset
393 recDf.train.rmOutlier = recDf.train[-c(165,15,146,143),]
394
395 # model diagnostics for active cases
396 # residual analysis
397 plot(act.fit4, which = 1:2, main = "AIC Model")
398 plot(act.fit3, which = 1:2, main = "BIC Model")
399 # outlying cases
400 plot(act.fit4, which = 4:6)
401 plot(act.fit3, which = 4:6)
402 # calculate cook's distance and determine outliers
403 cooks = getCooksDist(act.fit3)
404 actOutliers = cbind(actDf.train[which(cooks > 4/(act.fit3$df.residual)),c(1,43)], CooksDist
    = cooks[which(cooks > 4/(act.fit3$df.residual))])
405 actOutliers[order(actOutliers$CooksDist, decreasing = TRUE),]
406 # examine influence of possible outlier
407 act.fit3.rmOutlier = lm(formula = act.fit3$call$formula, data = actDf.train[-c(165,143),])
408 per.change = abs((act.fit3$fitted.values - predict.lm(act.fit3.rmOutlier, actDf.train))/act.
    fit3$fitted.values)*100
409 summary(per.change)
410
411 #####
412 ## final models: refit models without outliers with BIC criterion
413 # confirmed cases
414 cfmDf.rmOutlier = rbind(cfmDf.train.rmOutlier, cfmDf.valid)
415 cfm.fit.final = lm(cfm.fit3$call$formula, data = cfmDf.rmOutlier)
416 summary(cfm.fit.final)
417 anova(cfm.fit.final)
418
419 # death cases
420 deathDf.rmOutlier = rbind(deathDf.train.rmOutlier, deathDf.valid)
421 death.fit.final = lm(death.fit3$call$formula, data = deathDf.rmOutlier)
422 summary(death.fit.final)
423 anova(death.fit.final)
424
425 # recovered cases

```

```
426 recDf.rmOutlier = rbind(recDf.train.rmOutlier, recDf.valid)
427 rec.fit.final = lm(rec.fit3$call$formula, data = recDf.rmOutlier)
428 summary(rec.fit.final)
429 anova(rec.fit.final)
430
431 # active cases
432 actDf.rmOutlier = rbind(actDf.train, actDf.valid)
433 act.fit.final = lm(act.fit4$call$formula, data = actDf.rmOutlier)
434 summary(act.fit.final)
435 anova(act.fit.final)
```