

## Quality Assessment of high energy biscuits Served in School Feeding Programme in Poverty-Prone Areas in Bangladesh

Md. Anwar Hossain<sup>1\*</sup> and Nitai Chakraborty<sup>2</sup>

<sup>1</sup>Planning and Development Division, Bangladesh Council of Scientific and Industrial Research

Dr. Qudrat-I-Khuda Road, Dhanmondi, Dhaka-1205, Bangladesh

<sup>2</sup>Department of Statistics, Qazi Motahar Hossain Building

University of Dhaka, Dhaka 1000, Bangladesh

### ABSTRACT

*This work purposes to examine the quality of industrially treated fortified high energy biscuits in Bangladesh as served in schools in poverty-prone areas in Bangladesh. Datasets were collected from the Institute of Food Science and Technology (IFST), Bangladesh Council of Scientific and Industrial Research (BCSIR), Dhaka. These are collected with the method of Single Stage Cluster Sampling. The sample referred to as Fortified High Energy Biscuit was subjected to laboratory test of the level of some qualitative parameters such as moisture, protein, fat, sugar, carbohydrate, vitamin A, Iron, etc. This data complies with the World Food Programme (WFP) standard in the framework of predictive and classification tools. In this research, we considered a 306 dataset from a previously physiochemical analysis of high energy biscuits in IFST, BCSIR, Dhaka. Several regression techniques were considered multivariate linear regression, and binary logistic regression. The logistic regression gives stronger results in terms of the interpretation, outperforming the linear regression as well as other models. Such models are useful for biscuits quality evaluation and improving biscuits production. Furthermore, similar techniques can help in the quality testing of different food products.*

**Keywords:** Institute of Food Science and Technology (IFST), Bangladesh Council of Scientific and Industrial Research (BCSIR), Single Stage Cluster Sampling, World Food Programme (WFP), Multivariate Linear Regression, Binary Logistic Regression.

### 1. INTRODUCTION

Fortified High Energy Biscuits are biscuits (small baked bread or cakes) that are supplemented with a premix of vitamins and minerals. This ready-to-eat food takes an interest in the covering pressing needs within the intense stage of a crisis circumstance amid which the populace isn't able to cook due to a need to get to essential offices (clean water, cooking accessories, etc.). Their use is also extended to a complement food ratio (use as snacks) to provide vitamins and minerals in regions/populations where diet is subjected to nutritional deficiencies. biscuits can be used also to prevent micronutrients deficiency in young children and school-age children [1].

**1.1 List of Biscuits Suppliers Companies:** Lists of biscuits processing factories here under are as follows:

| S.L. No. | Suppliers of Biscuits- Bangladesh |
|----------|-----------------------------------|
| 1.       | New Olympia Biscuit Factory       |

| S.L. No. | Suppliers of Biscuits- Bangladesh  |
|----------|--|
| 2.       | Masafi Bread & Biscuit Industries Ltd.                                       |
| 3.       | Resco Biscuit & Bread factory (PVT) Ltd.                                     |
| 4.       | Central Marketing Company (CMC)<br>(Alauddin Food & Chemical Industries Ltd) |
| 5.       | Mona Food Industries   |
| 6.       | Olympic Industries Limited   |
| 7.       | Romania Food & Beverages Ltd.  |

Source: *Bangladesh Milling assessment additional info*[2].

Main ingredients: The biscuits product shall be manufactured from fresh and good quality raw materials free from foreign materials, substances hazardous to health including any contamination from toxic or noxious seeds, extreme moisture, insect damage, and fungal contamination and shall comply with all relevant food safety and quality laws, regulations and standards for each material (if used) as followings:

- Wheat flour shall be harmless and appropriate for human consumption.
- Sugar shall be free from heavy metals. Small amounts of metals may represent a hazard to human health.
- Shortening must be prepared from oil that must be free from trans fatty acids and must contain only antioxidants that comply with Codex and relevant regulations: in case of using palm oil, must conform to codex regulations, and in case using butter, must conform to codex for Butter.
- Skimmed milk powder must conform to standard levels. The maximum level of aflatoxin M1 should be less than 0.5 mcg/kg and the melamine maximum level is 2.5 mg/kg [3].

## 1.2 Analytical Requirements

As per contractual agreement, WFP can appoint an inspection company that will check, based on testing plan below, that the food matches requirements of this specification. Additional tests may be defined in case further quality assessment is required. The following sampling and analytical plans are currently utilized by WFP and shared only for suppliers' information. Suppliers should follow their own food safety and quality management plan. WFP reserves the rights to change these plans at any time.

**Table 1: List of compulsory tests and reference method**

| No | Tests                               | Requirements  |
|----|-------------------------------------|---|
| 1  | Moisture content                    | Max 4.5 %   |
| 2  | Organoleptic (smell, taste, colour) | Typical colour, pleasant smell and palatable taste. |
| 3  | Broken biscuits                     | Max. 5.0 % broken (by weight)                       |
| 4  | Total Protein                       | Min. 10g/100g                                       |
| 5  | Total Fat                           | Min. 15.0 g/100g                                    |
| 6  | Crude fibre                         | Max. 2.3 g/100g                                     |

| No | Tests  | Requirements                       |
|----|--|------------------------------------|
| 7  | Peroxide value   | Max. 10 meq/kg <sup>11</sup> fat   |
| 8  | Vitamin A-Retinol  | 500 – 850 mcg/100g                 |
| 9  | Iron   | 10-17 mg/100g                      |
| 10 | Aerobic mesophilic bacteria                                  | Max. 10,000 cfu/g                  |
| 11 | Coliforms  | Max. 10 cfu/g                      |
| 12 | Escherichia coli   | Absent in 10 g                     |
| 13 | Salmonella   | Absent in 25 g                     |
| 14 | Staphylococcus aureus  | <10 cfu/g                          |
| 15 | Bacillus cereus  | Max. 10 cfu/g                      |
| 16 | Yeasts and moulds  | Max. 100 cfu/g                     |
| 17 | GMO (only if required by contract or country of destination) | Negative (< 0.9 % of GMO material) |

This requirement is for foods at the time of purchase [3].

This study is intended to establish the acceptability for human consumption of high-energy biscuits. To find out which parameters of analysis biscuits are significant in the context of quality. To conclude meaning interpretation of this biscuit's products. To find out meaningful patterns in the biscuits data and enables us to identify the possible characteristics in the data.

## 2. MATERIALS AND METHODS

### 2.1 Classification

Cataloging may be an information mining highlight that entrusts objects to target classifications or classes inside a set. The arrangement objective is to anticipate the objective class precisely in the information for every function. A grouping model might be utilized, for instance, to order advanced candidates with little, medium, or high credit chances. Arrangement errands start with an informational collection that knows the class tasks. Characterization is discrete and doesn't infer requests. Nonstop, skimming point respects will prescribe an objective number rather than a clear-cut one. A prescient model that has a mathematical objective uses a relapse calculation, not a calculation for order. The clearest kind of issue with order is a double grouping. The objective quality in coordinated characterization has reasonable two potential qualities: tall praise score or moo praise assessment, for the event. Multiclass targets have numerous qualities: moo, medium, tall, and darken FICO appraisals, for instance. In the model build strategy (preparing), an arrangement calculation discovers connections between the indicator esteems and the objective qualities. Various calculations for the arrangement utilize explicit strategies to recognize connections. These connections are plotted in a model that would then be able to be applied to another arrangement of information in which the class tasks are obscure [4].

Characterization models are assessed by contrasting the normal qualities in a bunch of test information against realized objective qualities. Scoring a measure of classification results in in-class assignments and probabilities for each particular event. For example, the likelihood of each classification for each customer can also be predicted by a model which classifies customers as low, medium, or high. Thus, the objective of the proposed chapter is to foresee the quality of the biscuits based on physicochemical tests through machine learning models. The upcoming sections precisely narrate the classification steps adopted by them in prediction [4].

## 2.2 The Binary logistic model

Binary logistic regression analysis is a statistical method for exploring and modeling the relationship between a random dichotomous variable and one or more random independent variables (continuous or categorical) in retrospective data [5].

Logistic regression estimates the probability of an outcome. Events are coded as binary variables with a value of 1 representing the occurrence of a target outcome, and a value of zero representing its absence. OLS can also model binary variables using linear probability models [6]. OLS may give predicted values beyond the range (0,1), but the analysis may still be useful for classification and hypothesis testing. The normal distribution and homogeneous error variance assumptions of OLS will likely be violated with a binary dependent variable, especially when the probability of the dependent event varies widely. Both models permit persistent, ordinal, and/or categorical independent factors [7].

In logistic regression, the dependent variable is frequently dichotomous, that is, it can take the value 1 with a probability of success  $\theta$  or the value 0 with a probability of failure  $1 - \theta$ . Bernoulli (or binary) variables are the name for this sort of variable. Applications of logistic regression have also been extended to circumstances where the dependent variable is of more than two cases, known as multinomial or polytomous regression, however this is not as common and is not treated in this treatment [8].

The independent or predictor variables in logistic regression can have any shape, as previously stated. That is, logistic regression makes no assumptions about the independent variables' distribution. Within each group, they do not have to be normally and linearly distributed, or of equal variance. In logistic regression, the relationship between the predictor and response variables is not a linear function; instead, the logistic regression function, which is the logit transformation of the predictor and response variables, is utilized:

$$\theta = \frac{e^{(\alpha + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_k x_k)}}{1 + e^{(\alpha + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_k x_k)}}$$

Where  $\alpha$  = the constant of the equation and,  $\beta$  = the coefficient of the predictor variables.

An alternative form of the logistic regression equation is:

$$\text{logit}[\theta(x)] = \log \left[ \frac{\theta(x)}{1 - \theta(x)} \right] = \alpha + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_k x_k$$

The purpose of logistic regression is to use the most parsimonious model to correctly predict the category of result for individual cases. To do this, a model is built that incorporates all predictor factors that can help predict the response variable. Several options are available during the model construction process. Variables can be included into the model in the researcher's preferred order, or stepwise regression can be used to verify the model's fit after each coefficient is added or removed [9].

In the exploratory phase of research, stepwise regression is used, but it is not recommended for theory testing [6]. Theory testing is the process of putting a priori hypotheses or ideas about the relationships between variables to the test. The purpose of exploratory testing is to find relationships rather than making a priori assumptions about the relationships between the variables.

The preferred approach of exploratory analyses appears to be backward stepwise regression, in which the study starts with a full or saturated model and variables are removed from the model in an iterative process. After each variable is eliminated, the model's fit is checked to confirm that it still sufficiently matches the data. The analysis is complete when no more variables can be removed from the model.

Logistic regression has two key applications. The first is group membership prediction. The findings of the analysis are in the form of an odds ratio since logistic regression evaluates the likelihood of success over the probability of failure. For example, logistic regression is frequently used in epidemiological research to determine the likelihood of acquiring cancer after controlling for other risk factors. Logistic regression also reveals the correlations and strengths between the variables (for example, smoking 10 packs a day puts you at a higher risk of cancer than working in an asbestos mine).

Several alternative strategies are used to test the importance of coefficients before including or excluding them from the model. Each of these points will be discussed further down.

**2.3 Wald Test:** The statistical significance of each coefficient ( $\beta$ ) in the model is determined using the Wald test. The z statistic is calculated via a Wald test, and it is as follows:

$$z = \frac{\hat{B}}{SE}$$

After that, the z value is squared, providing a Wald statistic with a chi-square distribution. However, various authors have discovered issues with the Wald statistic's use. [6] advises that standard error is exaggerated for big coefficients, decreasing the Wald statistic (chi-square) value. According to [10], the likelihood-ratio test is more reliable than the Wald test for small sample sizes.

**2.4 Likelihood-Ratio Test:** The likelihood-ratio test compares the maximum value of the likelihood function for the full model ( $L_1$ ) to the maximum value of the likelihood function for the simpler model ( $L_0$ ). The likelihood-ratio test statistic is equal to the following:

$$-2 \log \left( \frac{L_0}{L_1} \right) = -2 [\log(L_0) - \log(L_1)] = -2(L_0 - L_1)$$

The chi-squared statistic is obtained by transforming the likelihood functions logarithmically. When creating a model using backward stepwise elimination, this is the suggested test statistic to utilize. [9, 10].

## 2.5 Data

High Energy Biscuits analysis data has been collected from the Institute of Food Science and Technology (IFST), BCSIR with permission of the authority. The variable name of the study is Protein, Fat, Carbohydrate, Sugar, etc. The relevant data have been collected from the analytical service report of biscuits which are previously analyzed over the year 2007-2012 by respective scientists of IFST. Data collection methods were the non-participant observation of organizations included in the study. The documented investigation included hard-copy issues of reports of expository records. We applied the data collection method as Single Stage Cluster Sampling.

## 2.6 Software

Statistical Analysis System R version 4.1.0 were used to analyze the data.

### 3. RESULTS AND DISCUSSION

#### 3.1 High Energy Biscuits (HEB) Quality Predictor

We worked on predicting the Biscuits quality for this project. Limit of unacceptable biscuits analysis parameters is represented with 0 and the limit of acceptable biscuits analysis parameters is represented as 1.

#### 3.2 Explanatory Data Analysis (EDA)

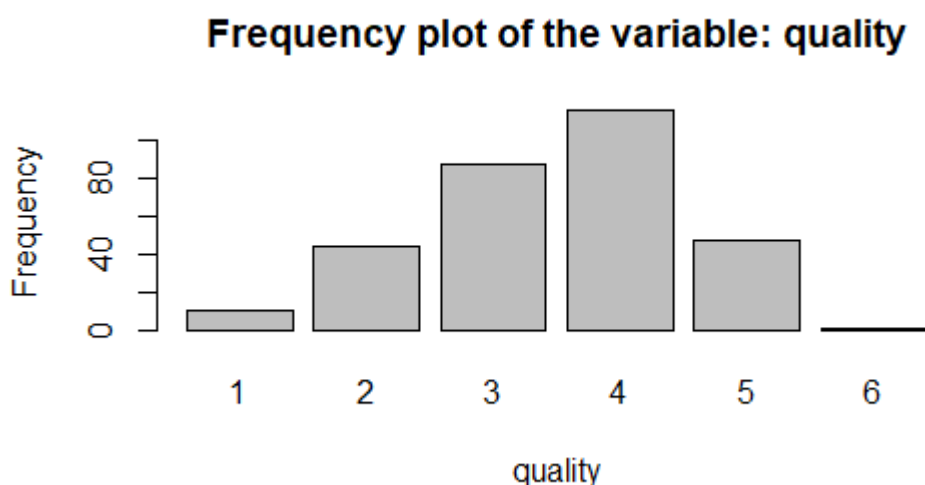
Now, we proceed to develop an EDA on the data to find essential insights and to determine specific relationships between the variables.

First, we developed a descriptive analysis where collected the five-number summary statistics of the data.

**Table 2. The physicochemical data analysis of high energy biscuits.**

| Attribute (units)      | Min.  | 1st Qu. | Median | Mean   | 3rd Qu. | Max.   |
|------------------------|-------|---------|--------|--------|---------|--------|
| Moisture (%)           | 0.960 | 2.430   | 2.955  | 2.957  | 3.450   | 7.810  |
| Protein (%)            | 1.31  | 9.34    | 10.17  | 10.26  | 11.14   | 14.44  |
| Fat (%)                | 7.33  | 13.73   | 14.74  | 14.72  | 15.97   | 21.29  |
| Sugar (%)              | 9.10  | 11.95   | 13.48  | 13.58  | 14.88   | 24.07  |
| Total Carbohydrate (%) | 63.31 | 68.65   | 70.44  | 70.40  | 72.06   | 78.63  |
| Vitamin A (µg/100g)    | 102.7 | 232.2   | 271.0  | 524.7  | 524.7   | 5785.0 |
| Iron (mg/100g)         | 1.910 | 9.283   | 10.725 | 11.708 | 12.660  | 89.000 |

The table above shows the 5-digit sum of each variable in the data. That is, I used the function to get the minimum and maximum values of the numeric variables, the first and third quartiles, and the mean and median. In addition, the summary shows the frequency of the level of the dependent variable. Next, we developed a univariate analysis that consists of examining each variable individually. First, we analyzed the dependent variables. To analyze the result variables, I created a bar chart that visualizes the frequency counts at the category level. We also created a frequency count table to know the exact amount and percentage of values in different layers of each category.



**Fig.1 Frequency plot to analyze the dependent variable**

Our main goal is to identify which of these variables have a significant effect on biscuits quality.

### 3.3 Preliminary Graphical Analysis

Apart from the linear regression model in this study we consider also the Logistic Regression Model. The framework divides each dataset into "bad" and "good" biscuits. The bad is till the quality score less than 3 else is a good one.

In fig.2 there are the plot accordingly to this division.

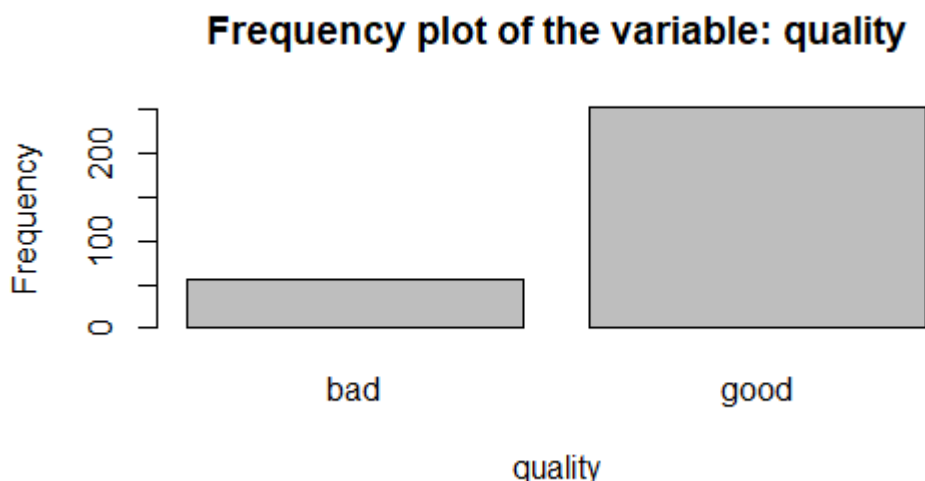


Fig.2: Frequency plot to division of the dependent variable of quality

As we can see in these figures the “good” quality is 82% and 18% of bad quality biscuits. In this section, our primary objectives are to find out which machine learning algorithm will enable the most accurate prediction of biscuit quality from its physicochemical properties? What physicochemical properties of biscuits have the highest impact on their quality.

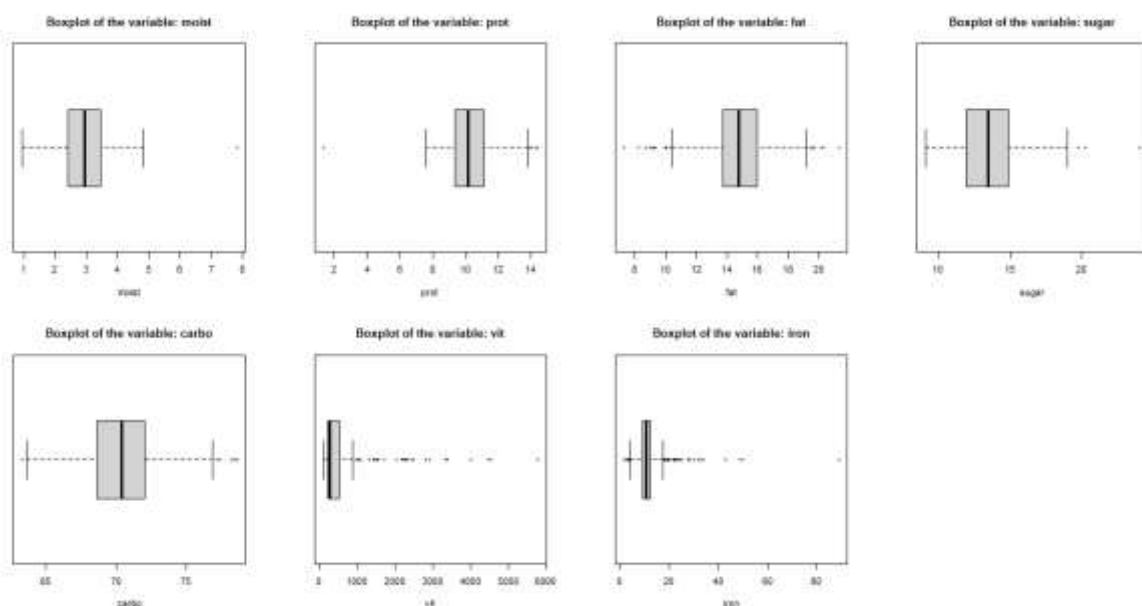


Fig.3: Boxplots to analyzed the numeric independent variables

As we can see, the boxplots show where are the mean, median, and quartile measurements located for each variable, as well as the range of values each variable has.

By analyzing the boxplots, we concluded that all the variables have outliers. Furthermore, the variables “Vitamin A” and “Iron” are the variables that have the greatest number of outliers. As we can see, there is a concentration of values near the mean and median, which is reflected by a very slim interquartile range (IQR).

This information will come in handy at the data preparation step when we proceed to assess the outlier values.

### 3.4 Principal Component Analysis

We perform a Principal Components (multivariate) Analysis to detect collinearity or correlation among the variables before starting regression analysis. Identifying variables that are highly collinear—which can make one of the variables almost redundant in some cases—can help us select the best possible binary logistic regression model and not only. In Fig.4 we can see these directions of the variables.

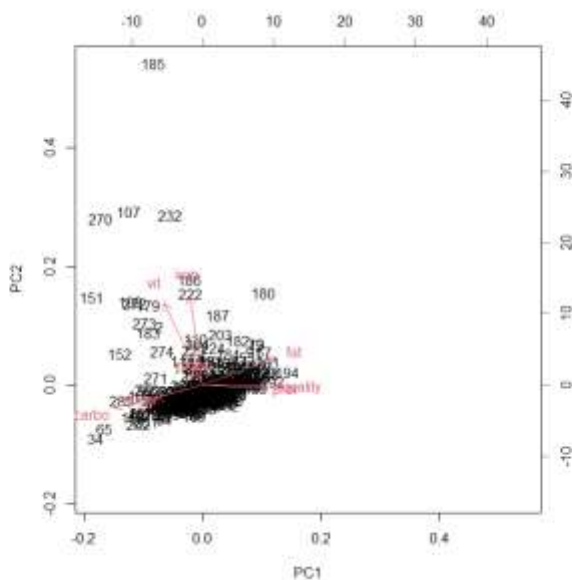


Fig.4 Loading Plot

It is obviously shown that quality and protein go very strongly together in the same direction, this means “collinearity”, on the contrary the sugar and total carbohydrate pointing on the other side mean that decreasing these parameter’s values against quality.

### Principal Components Analysis

Call: principal(r = data, nfactors = 1)

Standardized loadings (pattern matrix) based upon correlation matrix

| Attribute (units) | PC1   | h2    | u2   | com |
|-------------------|-------|-------|------|-----|
| Moisture (%)      | -0.07 | 0.005 | 0.99 | 1   |
| Protein (%)       | 0.67  | 0.443 | 0.56 | 1   |



|                               |       |       |      |   |
|-------------------------------|-------|-------|------|---|
| <b>Fat (%)</b>                | 0.74  | 0.548 | 0.45 | 1 |
| <b>Sugar (%)</b>              | -0.49 | 0.242 | 0.76 | 1 |
| <b>Total Carbohydrate (%)</b> | -0.89 | 0.786 | 0.21 | 1 |
| <b>Vitamin A (µg/100g)</b>    | -0.38 | 0.148 | 0.85 | 1 |
| <b>Iron (mg/100g)</b>         | -0.13 | 0.016 | 0.98 | 1 |
| <b>quality</b>                | 0.80  | 0.643 | 0.36 | 1 |

PC1

SS loadings 2.83

Proportion Var 0.35

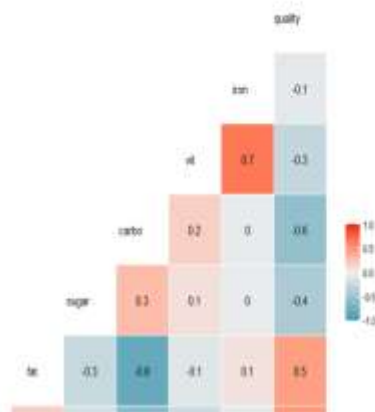
Mean item complexity = 1

Test of the hypothesis that 1 component is sufficient.

The root mean square of the residuals (RMSR) is 0.17  
with the empirical chi square 494.05 with prob  $< 5.1e^{-92}$

Fit based upon off diagonal values = 0.71

These results of principal component analysis confirm the above loading plot visual observations. Fig. 5 shows the correlations visually.



**Fig.5 correlation matrix**

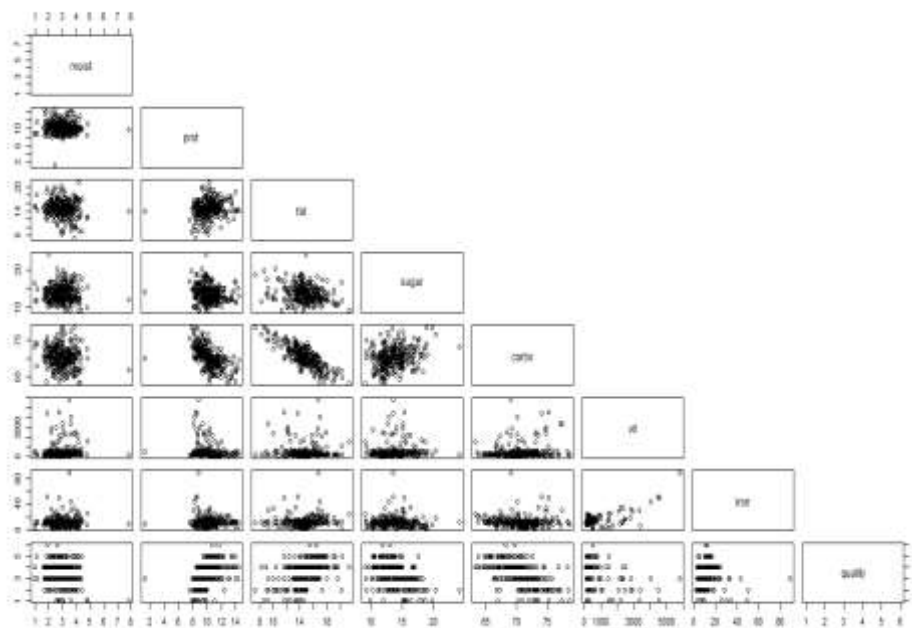


Fig.6 scatter-plot

### 3.5 Linear Regression Model

As mentioned earlier, linear regression is the first attempt to fit a model and predict the quality values of a dataset. To show that our goal is quality, we perform a linear regression on the variable "quality" for all other variables. If you follow the standard method of regression analysis, you will find that the R<sup>2</sup> is very low. Therefore, consider a separate model for datasets that use high-energy cookies. The result is from the cookie dataset.

#### 3.5.1 Linear regression in biscuits datasets

If the response variable quality is assumed continuous the R code regression model is:

```
lm(formula = quality ~ ., data = data)
```

and gives the following residuals and coefficients:

**Residuals:**

| Min      | 1Q       | Median   | 3Q      | Max     |
|----------|----------|----------|---------|---------|
| -1.93586 | -0.52408 | -0.02174 | 0.43450 | 2.05095 |

**Coefficients:**

|             | Estimate               | Std. Error            | t value | Pr(> t )   |
|-------------|------------------------|-----------------------|---------|------------|
| (Intercept) | 1.091e <sup>+01</sup>  | 5.522e <sup>+00</sup> | 1.975   | 0.04918 *  |
| moisture    | -2.236e <sup>-01</sup> | 8.052e <sup>-02</sup> | -2.777  | 0.00583 ** |
| protein     | 1.624e <sup>-01</sup>  | 5.826e <sup>-02</sup> | 2.788   | 0.00564 ** |
| fat         | 4.107e <sup>-02</sup>  | 6.205e <sup>-02</sup> | 0.662   | 0.50854    |

|                           | <i>Estimate</i>        | <i>Std. Error</i>     | <i>t value</i> | <i>Pr(&gt; t )</i>       |
|---------------------------|------------------------|-----------------------|----------------|--------------------------|
| <i>sugar</i>              | -1.239e <sup>-01</sup> | 2.111e <sup>-02</sup> | -5.867         | 1.18e <sup>-08</sup> *** |
| <i>total carbohydrate</i> | -1.036e <sup>-01</sup> | 5.657e <sup>-02</sup> | -1.831         | 0.06806                  |
| <i>vitamin a</i>          | -2.448e <sup>-04</sup> | 8.732e <sup>-05</sup> | -2.803         | 0.00540 **               |
| <i>iron</i>               | 6.390e <sup>-03</sup>  | 7.830e <sup>-03</sup> | 0.816          | 0.41513                  |

---  
 Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 0.7738 on 298 degrees of freedom  
**Multiple R-squared: 0.4697,** Adjusted R-squared: 0.4572  
 F-statistic: 37.7 on 7 and 298 DF, p-value: < 2.2e<sup>-16</sup>

As we seen from the results the most critical coefficient the R-squared is 0.4697. The residuals comport very well converging to 0 with mean: -3.45073e<sup>-17</sup>.

**Prediction table:**

| <i>Pred.var.</i> | <i>1</i> | <i>2</i> | <i>3</i> | <i>4</i> | <i>5</i> | <i>6</i> |
|------------------|----------|----------|----------|----------|----------|----------|
| <i>1</i>         | 2        | 3        | 6        | 0        | 0        | 0        |
| <i>2</i>         | 1        | 15       | 26       | 2        | 0        | 0        |
| <i>3</i>         | 0        | 2        | 57       | 28       | 0        | 0        |
| <i>4</i>         | 0        | 1        | 26       | 72       | 16       | 0        |
| <i>5</i>         | 0        | 0        | 8        | 31       | 7        | 1        |
| <i>6</i>         | 0        | 0        | 0        | 2        | 0        | 0        |

**Accuracy:**

[1] 0.5  
 So, the final percentage of right prediction is about **50.00%**

**3.6 Logistic Regression Model**

To know the results of the linear regression model and get better predictive performance, let's look at another powerful model, logistic regression. Of course, the main difference from linear regression is that it predicts the quality category, not the score. Is it bad or good? If the response variable quality is assumed continuous the R code regression model is:  
`glm(quality ~ ., family = binomial(link = logit), data = data)`  
 and gives the following residuals and coefficients:

**Deviance Residuals:**

| Min      | 1Q      | Median  | 3Q      | Max     |
|----------|---------|---------|---------|---------|
| -3.03796 | 0.05639 | 0.17921 | 0.40496 | 1.80122 |

**Coefficients:**

|                    | Estimate   | Std. Error | z value | Pr(> z )     |
|--------------------|------------|------------|---------|--------------|
| (Intercept)        | 61.1683953 | 20.5519418 | 2.976   | 0.002918 **  |
| moisture           | -1.2761109 | 0.3405191  | -3.748  | 0.000179 *** |
| protein            | 0.3712902  | 0.1959897  | 1.894   | 0.058167     |
| fat                | -0.2842379 | 0.2397227  | -1.186  | 0.235743     |
| sugar              | -0.4295342 | 0.0947178  | -4.535  | 5.76e-06 *** |
| total carbohydrate | -0.6873646 | 0.2183213  | -3.148  | 0.001642 **  |
| vitamin a          | -0.0007059 | 0.0003322  | -2.125  | 0.033578 *   |
| iron               | 0.0363539  | 0.0297636  | 1.221   | 0.221926     |

---

Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1

(Dispersion parameter for binomial family taken to be 1)

Null deviance: 288.25 on 305 degrees of freedom

Residual deviance: 165.06 on 298 degrees of freedom

AIC: 181.06

Number of Fisher Scoring iterations: 6

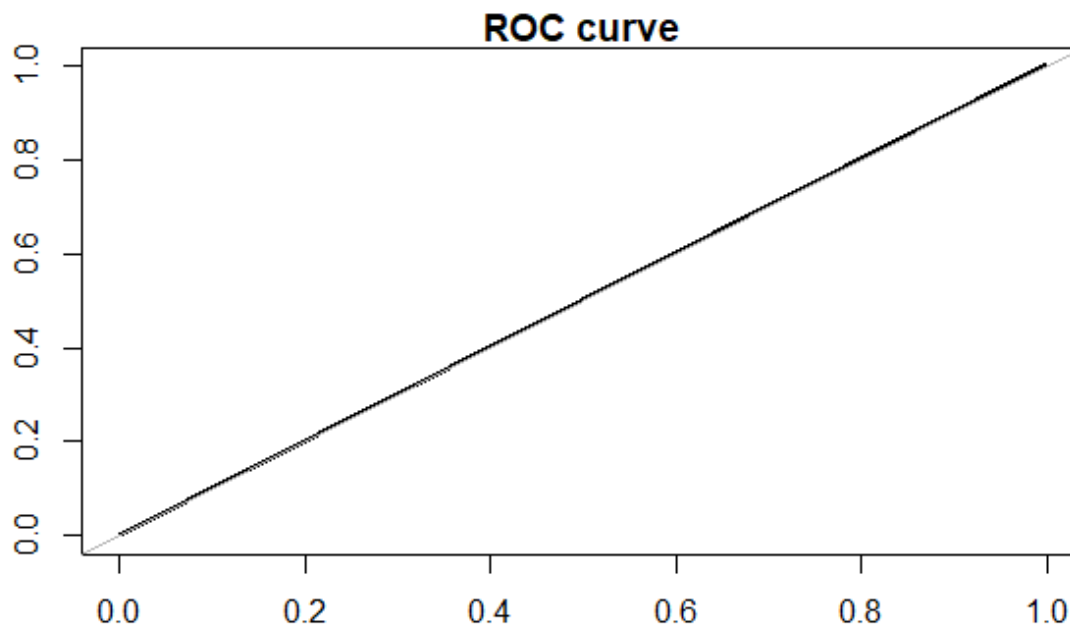
**Prediction table:**

| Pred.var. | bad | good |
|-----------|-----|------|
| Bad       | 1   | 54   |
| Good      | 0   | 251  |

**Accuracy:**

[1] 0.8235294

So, the final percentage of right prediction is about **82.35%**



**Fig.7 ROC curves for logistic regression model**

By analyzing the results, we declared that the most significant variable for this model is “sugar”, followed by the variable “vitamin A”.

Further, we effectuated an investigation to know the performance and impact of these components on the biscuit’s quality.

As we analyzed, the logistic regression model explains the actual theory facts. The investigation of the essential components for a good biscuit quality involved the variables obtained as necessary in the model. For this reason, the variables “sugar” and “vitamin A” are very significant to the model because these elements will be an essential component in indicating if a biscuit has a good or bad quality.

#### 4. CONCLUSION AND RECOMMENDATION

Comparison of linear models

Unfortunately, our model does not meet the criteria for a linear regression model based on a statistical approach with good criteria.  $R^2$  represents only a 54.33% variance of the data, even after removing the outliers. The predictive power of the model measured by RMSE is 0.8242499 and the MAE is 0.6409561. We learned this data is not useful for linear regression. There is also a question as to why this issue occurs. The data has no linear relationship between the predictor and the target. The dataset author also said, "For privacy and logistics reasons, only physicochemical (input) and quality assessment (output) variables are available (for example, data on which species, biscuits brand, biscuits production price). No, etc.) "in the dataset. The target size is actually a normal classification task, and the higher the number of grades, the higher the quality of the biscuit. However, regression can actually solve this problem because it assumes that the target is a continuous number. The quality of the cookie is not determined by the variables present in the data. R-squared achieves only 54.33% variance of the data, and perhaps the other 73% have a greater impact.

After retrieving the results of various machine learning algorithms, we found that the logistic regression model showed higher accuracy in predicting the quality of cookies. With 96.69% accuracy, this model was able to correctly predict the value of 234. This means that the model misclassification error is 3.31%. On the other hand, an analysis of the ROC curve showed that the performance of the model was not as good as expected. The ROC curve was labeled as a failure curve by evaluating the area under the curve (AUC = 50%). In other words, the model is unable to identify the various classes, indicating that the model is performing poorly. For this reason, we conclude that while the model

has excellent accuracy in predicting test set values, it is a disastrous rate to quickly identify a true positive value. In addition, analysis of other ROC curves revealed that the decision tree performed best and the area under the curve was 86.9%. That is, the decision tree model did not have the highest accuracy of any logistic model, but it performs better than logistic regression in identifying different classes of dependent variables. We also use a logistic regression model to identify which physicochemical properties have the greatest impact on biscuit quality. We have distinguished that "sugar" and "vitamin A" are the variables that have the most significant impact on the model. Changes to any of these variables will have a significant impact on the model's results.

## Acknowledgements

I would like to thank Authority of BCSIR for gives an opportunity to use of analytical service data and Ministry of Science and Technology to allocate research fund for this purpose.

## REFERENCES

- [1] Value N, Fsc OF, Food S. Annex 15 NUTRITIONAL VALUE OF FSC STANDARDISED FOOD PACKAGES & Nutritional Value of Proposed FSC Standardised Food Packages ( Immediate and Long-term ) Food Basket for First 7 Days ( Immediate ) Specifications : Below are specifications for the following commodities : Bangladesh Biscuit Palm Olein Oil.
- [2] *Bangladesh Milling assesment additional info.*
- [3] Technical Specifications for HIGH ENERGY BISCUITS (HEB). 2021; 1–9.
- [4] Gupta M, Vanmathi C. A Study and Analysis of Machine Learning Techniques in Predicting Wine Quality.
- [5] Logistic Regression Analysis - an overview | ScienceDirect Topics, <https://www.sciencedirect.com/topics/nursing-and-health-professions/logistic-regression-analysis> (accessed 22 November 2021).
- [6] Menard S. *Applied logistic regression analysis: Sage university series on quantitative applications in the social sciences*. Thousand Oaks, CA: Sage, 1995.
- [7] Pohlman JT, Leitner DW. A comparison of ordinary least squares and logistic regression.
- [8] Tabachnick BG, Fidell LS. *Using Multivariate Statistics*, New York: HarperCollins College Publishers.
- [9] Interval C, Ratio O. *Logistic Regression*. 2016; 6–8.
- [10] Agresti A. *An Introduction to Categorical Data Analysis*. NewYork: John Wiley & Sons, Inc, 1996.

\*Corresponding author. E-mail: [anwarbcsir@yahoo.com](mailto:anwarbcsir@yahoo.com)