



CHRIST

(DEEMED TO BE UNIVERSITY)

BANGALORE • INDIA

MST471 -Categorical Data Analysis

CAC1 Project

Food Preferences and Consumption Patterns

By

Raghotham R C - 2348141

Sowmiya S - 2348145

Swetha R - 2348149

Thulasi V - 2348152

Vidya Sri L - 2348154

4MSTAT

Department of Statistics and Data Science

CHRIST (Deemed to be University), Bengaluru – 560 029

Team Work Bifurcation

- **Sowmiya S(2348145):**
 - Analyze Dietary Preferences: Examine the association between gender and dietary preferences, using odds ratio and relative risk to quantify differences.
 - Assess Health Consciousness Impact: Determine how health consciousness influences food preferences through Chi-Square tests and visualizations.
- **Raghotham R C(2348141):**
 - Evaluate Meal Preferences: Explore the link between meal preferences and frequency of eating out, applying likelihood ratio tests and heatmaps.
 - Apply Multinomial Logistic Regression: Investigate shopping behaviors related to cooking frequency using multinomial logistic regression.
- **Thulasi V(2348152):**
 - Predict Dietary Preferences: Fit and evaluate logistic regression models to predict dietary preferences based on various predictors.
 - Formatting: Responsible for formatting the final report and analysis.
- **Swetha R(2348149):**
 - Fit Log-Linear Model: Test the adequacy of the log-linear model in representing relationships among categorical variables related to food preferences.
 - Interpretation: Contributed to interpreting the results and insights from the analysis.
- **Vidya Sri L(2348154):**
 - Data Collection: Contributed to gathering data for the project.
 - Formed the Questionnaire.
 - Formatted the final report, EDA and Visualization.

Exploring Dietary Preferences and Food Consumption Patterns through Categorical Data Analysis

MST471 - CDA Group Project

2024-08-28

Introduction:

Dietary preferences and food consumption patterns are key aspects of human behavior that have significant implications for public health, nutrition, and the food industry. These preferences are shaped by a complex interplay of demographic factors, health consciousness, cultural influences, and individual lifestyle choices. Understanding these patterns is crucial for designing effective public health interventions, developing targeted marketing strategies, and addressing the diverse needs of consumers.

This study aims to explore the intricate relationships between dietary preferences and various influencing factors such as gender, health consciousness, meal preferences, and eating habits. By employing categorical data analysis techniques, the research seeks to uncover the underlying associations that drive food choices. Specifically, the study will analyze how gender affects dietary preferences, assess the impact of health consciousness on food choices, examine the relationship between meal preferences and eating out frequency, and develop predictive models to forecast dietary preferences. Additionally, advanced statistical methods, including multinomial logistic regression and log-linear modeling, will be used to delve into more complex relationships among multiple categorical variables related to food consumption.

Objectives

1. **Analyze Dietary Preferences:** Examine the association between gender and dietary preferences, using odds ratio and relative risk to quantify differences.
2. **Assess Health Consciousness Impact:** Determine how health consciousness influences food preferences through Chi-Square tests and visualizations.
3. **Evaluate Meal Preferences:** Explore the link between meal preferences and frequency of eating out, applying likelihood ratio tests and heatmaps.
4. **Predict Dietary Preferences:** Fit and evaluate logistic regression models to predict dietary preferences based on various predictors.
5. **Apply Multinomial Logistic Regression:** Investigate shopping behaviors related to cooking frequency using multinomial logistic regression.
6. **Fit Log-Linear Model:** Test the adequacy of the log-linear model in representing relationships among categorical variables related to food preferences.

Formulation of Problem

The study aims to explore and analyze dietary preferences and food consumption patterns using categorical data analysis. The focus is on understanding how various factors such as gender, health consciousness, meal preferences, and eating habits influence dietary choices among respondents. The key issues include:

1. **Dietary Preferences by Gender:** Assessing how gender affects preferences for vegetarian versus non-vegetarian diets.
2. **Health Consciousness and Food Choices:** Investigating the relationship between health consciousness levels and preferred food types.
3. **Meal Preferences and Eating Out Frequency:** Evaluating whether meal preferences are associated with how often individuals eat out.
4. **Dietary Preferences Prediction:** Developing models to predict dietary preferences based on demographic and behavioral factors.
5. **Multinomial and Log-Linear Modeling:** Employing advanced statistical techniques to understand the complex relationships among multiple categorical variables related to food consumption.

Questionnaire

The following questionnaire was designed to collect comprehensive data on respondents' food preferences, eating habits, and spending patterns. It includes various types of questions to categorize responses and assess key aspects of food consumption.

1. **What is your age group?**
 - Under 18
 - 18-28
 - 29-38
 - 39-48
 - 49 and above
2. **What is your gender?**
 - Male
 - Female
3. **What is your dietary preference?**
 - Vegetarian
 - Non-Vegetarian
4. **What is your favorite type of cuisine?**
 - Italian
 - Chinese
 - North Indian
 - South Indian
5. **How often do you eat out?**
 - Daily
 - Weekly

- Monthly
 - Occasionally
6. **Which meal of the day do you prefer the most?**
- Breakfast
 - Lunch
 - Dinner
 - Snacks
7. **How important is health to you when choosing food?**
- Important
 - Neutral
 - Does not care
8. **What is your preferred type of food?**
- Fast Food
 - Home-Cooked
 - Organic
 - Street Food
9. **How much do you typically spend on food per week?**
- Under 200
 - 200-500
 - 600-1000
 - 1000 and above
10. **Where do you primarily shop for food?**
- Supermarket
 - Local Market
 - Organic Store
 - Online
11. **How often do you cook meals at home?**
- Daily
 - Several times a week
 - Weekly

- Occasionally

12. How often do you try new foods or cuisines?

- Very Often
- Often
- Occasionally
- Never

Data Collection:

Description:

Data was collected using an online survey tool, distributed among participants from various age groups. The survey automatically recorded respondents' email addresses to ensure authenticity.

Sample Size:

A total of **131** respondents participated in the survey, providing diverse insights into their food preferences and consumption habits.

Methodology:

The study utilizes a structured questionnaire to collect data on respondents' food preferences, eating habits, and spending patterns. The survey was distributed online, targeting a diverse sample of participants from various age groups, resulting in a total of 131 responses. The dataset was pre-processed to remove unnecessary columns, and various statistical analyses were performed to meet the study's objectives.

- **Contingency Tables and Tests:** We created contingency tables to explore the relationships between categorical variables, such as gender and dietary preferences. The odds ratio and relative risk were calculated to quantify differences, while Chi-Square tests were employed to assess the significance of associations, such as between health consciousness and preferred food types.
- **Visualization:** Heatmaps were generated to visually represent the relationships between variables, providing a clearer understanding of how different factors, such as health consciousness, influence food preferences.
- **Logistic Regression:** A logistic regression model was fitted to predict dietary preferences based on selected predictors, such as gender and favorite cuisine. Model selection techniques, including forward selection and backward elimination, were used to identify the most important predictors.
- **Multinomial Logistic Regression:** This analysis was used to investigate shopping behaviors related to cooking frequency, providing insights into how often people shop at different places based on how frequently they cook meals at home.
- **Log-Linear Modeling:** A log-linear model was fitted to test the adequacy of representing relationships among categorical variables related to food preferences, allowing us to understand the interactions between multiple factors.

Data Import and Initial Processing:

```
library(readxl)
#Load data
library(readxl)
data<- read_excel("D:/Tri_4/CDA/CDA_Project (Responses)final.xlsx")
data <- data[, -c(1, 2)]
# Print first few rows of data
head(data)
```

```
## # A tibble: 6 x 12
##   age      gender diet_pref fav_cuisine eat_out_freq meal_pref health_conscious
##   <chr>    <chr>  <chr>      <chr>      <chr>      <chr>      <chr>
## 1 18-28    Male    Non - Veg~ South Indi~ Daily        Lunch      Neutral
## 2 Under 18 Female Non - Veg~ South Indi~ Weekly        Lunch      Important
## 3 18-28    Female Non - Veg~ Italian    Weekly        Lunch      Important
## 4 18-28    Female Non - Veg~ South Indi~ Weekly        Dinner     Important
## 5 29-38    Female Vegetarian South Indi~ Occasionally Dinner     Neutral
## 6 18-28    Female Non - Veg~ South Indi~ Monthly      Lunch      Important
## # i 5 more variables: pref_typefood <chr>, spend_week <chr>, shop_place <chr>,
## #   cook_freq <chr>, new_food <chr>
```

The dataset is loaded from an Excel file. The first two columns are dropped,because they are unnecessary for the analysis.

Analysis:

```
# Summary statistics for numerical variables
summary(data)
```

```
##      age      gender      diet_pref      fav_cuisine
## Length:131      Length:131      Length:131      Length:131
## Class :character Class :character Class :character Class :character
## Mode :character Mode :character Mode :character Mode :character
## eat_out_freq      meal_pref      health_conscious      pref_typefood
## Length:131      Length:131      Length:131      Length:131
## Class :character Class :character Class :character Class :character
## Mode :character Mode :character Mode :character Mode :character
## spend_week      shop_place      cook_freq      new_food
## Length:131      Length:131      Length:131      Length:131
## Class :character Class :character Class :character Class :character
## Mode :character Mode :character Mode :character Mode :character
```

```
# Frequencies of categorical variables
table(data$age)
```

```
##
##      18-28      29-38 39 and above      Under 18
##      108          7          11          5
```

```
table(data$gender)
```

```
##  
## Female    Male  
##      69     62
```

```
table(data$diet_pref)
```

```
##  
## Non - Vegetarian    Vegetarian  
##           94           37
```

```
table(data$fav_cuisine)
```

```
##  
##      Chinese      Italian North Indian South Indian  
##           7           6           17           101
```

```
table(data$eat_out_freq)
```

```
##  
##      Daily      Monthly Occasionally      Weekly  
##           27           15           26           63
```

```
table(data$meal_pref)
```

```
##  
## Breakfast    Dinner    Lunch    Snacks  
##           22           51           49           9
```

```
table(data$health_conscious)
```

```
##  
## Does not care    Important    Neutral  
##           10           73           48
```

```
table(data$pref_typefood)
```

```
##  
## Fast Food Home-Cooked    Organic Street Food  
##           16           92           12           11
```

```
table(data$spend_week)
```

```
##  
## 1000 above    200-500    600-1000    Under 200  
##           16           60           27           28
```



```
table(data$shop_place)
```

```
##
##   Local Market      Online Organic Store      Supermarket
##           55           24           7           45
```

```
table(data$cook_freq)
```

```
##
##           Daily      Occasionally Several times a week
##           43           45           29
##      Weekly
##           14
```

```
table(data$new_food)
```

```
##
##      Never Occasionally      Often      Very Often
##      12           69           32           18
```

Demographics:

69 females and 62 males. Majority are non-vegetarian (94).

Cuisine Preferences:

South Indian is the most popular cuisine (101). Other favorites include North Indian (17), Italian (6), and Chinese (7).

Meal Preferences:

Most preferred meals are dinner (51) and lunch (49). Fewer respondents prefer breakfast (22) and snacks (9).

Eating Out Frequency:

Weekly (63) is the most common frequency. Other frequencies: Occasionally (26), Daily (27), and Monthly (15).

Food Preferences:

Home-cooked food is highly preferred (92). Less preference for fast food (16) and street food (11). Organic food has a moderate preference (12).

Shopping Habits:

Local Market (55) is the most popular shopping place. Followed by Supermarket (45), Online (24), and Organic Store (7).

Food Importance:

Most consider food important (73). A smaller number are indifferent (10).

Analysis of Gender and Dietary preferences:

```
# Create contingency table with gender as the response variable
contingency_table_gen <- table(data$diet_pref, data$gender)
```

```
# Display the contingency table
print(contingency_table_gen)
```

```
##
##               Female Male
## Non - Vegetarian    42   52
## Vegetarian         27   10
```

The contingency table reveals the following dietary preferences:

- Females: 42 prefer non-vegetarian, and 27 prefer vegetarian.
- Males: 52 prefer non-vegetarian, and 10 prefer vegetarian.

```
# Define the values from the contingency table
a <- 42 # Females preferring non-vegetarian
b <- 27 # Females preferring vegetarian
c <- 52 # Males preferring non-vegetarian
d <- 10 # Males preferring vegetarian
```

```
# Calculate Odds Ratio
odds_ratio <- (a / b) / (c / d)
print(paste("Odds Ratio:", odds_ratio))
```

```
## [1] "Odds Ratio: 0.299145299145299"
```

The odds ratio of approximately 0.299 indicates that the odds of females preferring non-vegetarian food compared to vegetarian food are about 30% of the odds for males. This suggests that females are significantly less likely to prefer non-vegetarian food compared to males.

```
# Calculate Relative Risk
risk_female_non_veg <- a / (a + b)
risk_male_non_veg <- c / (c + d)
relative_risk <- risk_female_non_veg / risk_male_non_veg
print(paste("Relative Risk:", relative_risk))
```

```
## [1] "Relative Risk: 0.725752508361204"
```

The relative risk of approximately 0.727 further supports this, showing that females are about 73% as likely as males to prefer non-vegetarian food.

Analysis of Health consciousness and Preferred foodtype:

```
# Create a contingency table with health_conscious and pref_typefood
contingency_table_hf <- table(data$health_conscious, data$pref_typefood)

# Display the contingency table
print(contingency_table_hf)
```

```
##
##               Fast Food Home-Cooked Organic Street Food
## Does not care      3         5         1         1
## Important          6        56         9         2
## Neutral            7        31         2         8
```

```

# Perform the Chi-Square test
chi_square_test <- chisq.test(contingency_table_hf)

## Warning in chisq.test(contingency_table_hf): Chi-squared approximation may be
## incorrect

# Display the results
print(chi_square_test)

##
## Pearson's Chi-squared test
##
## data:  contingency_table_hf
## X-squared = 13.839, df = 6, p-value = 0.03149

# Critical value for alpha = 0.05 and df = 6
critical_value <- qchisq(0.95, df = 6)
print(critical_value)

## [1] 12.59159

```

The Chi-Square test results for the contingency table with `health_conscious` and `pref_typefood` show a test statistic of 13.839 with 6 degrees of freedom, and a p-value of 0.03149. This p-value is below the typical alpha level of 0.05, indicating that there is a statistically significant association between health consciousness and preferred food type. The critical value for the Chi-Square distribution with 6 degrees of freedom at a 0.05 significance level is approximately 12.592. Since the test statistic 13.839 exceeds this critical value, we reject the null hypothesis.

In other words, the differences observed in the frequency of preferred food types across different levels of health consciousness are unlikely to be due to chance, suggesting a meaningful relationship between health consciousness and food preference.

Visualization:

```

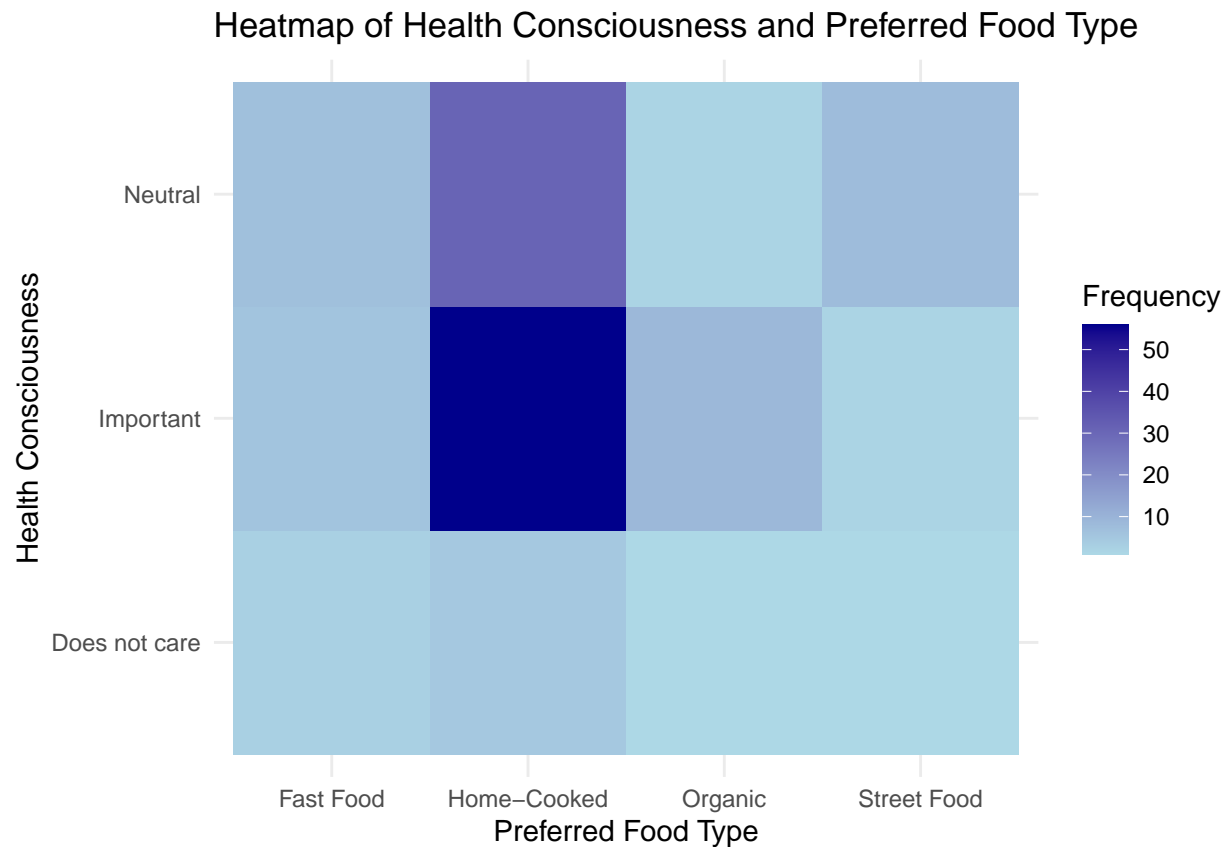
# Convert the contingency table to a data frame
contingency_df <- as.data.frame(contingency_table_hf)

# Load ggplot2 package
library(ggplot2)

## Warning: package 'ggplot2' was built under R version 4.3.2

# Create a heatmap with a color gradient from light to dark
ggplot(contingency_df, aes(x = Var2, y = Var1, fill = Freq)) +
  geom_tile() +
  scale_fill_gradient(low = "lightblue", high = "darkblue") +
  labs(title = "Heatmap of Health Consciousness and Preferred Food Type",
       x = "Preferred Food Type",
       y = "Health Consciousness",
       fill = "Frequency") +
  theme_minimal()

```



The heatmap visually supports this finding by showing how the frequency of each food preference varies with health consciousness, with some categories (like “Important” in “Home-Cooked”) being more pronounced than others.

Analysis of Meal preference and eating outside frequently:

```
#Create contingency table
table = table(data$meal_pref, data$eat_out_freq)
table
```

```
##
##           Daily Monthly Occasionally Weekly
## Breakfast      4      1             2    15
## Dinner        14      5            11    21
## Lunch          6      6            11    26
## Snacks         3      3             2     1
```

```
#Likelihood Ratio Test
library(DescTools)
```

```
## Warning: package 'DescTools' was built under R version 4.3.3
```

```
likelihood_ratio_test <- GTest(table)
print(likelihood_ratio_test)
```

```
##
## Log likelihood ratio (G-test) test of independence without correction
##
## data:  table
## G = 15.477, X-squared df = 9, p-value = 0.07863
```

```
# Calculate the degrees of freedom
df <- (nrow(table) - 1) * (ncol(table) - 1)

# Calculate the chi square tabulated value
p_value <- qchisq(0.95, df)
print(paste("Chi square tabulated value ", p_value))
```

```
## [1] "Chi square tabulated value 16.9189776046204"
```

Here, the likelihood ratio test (G-test), is used to determine whether there is a significant association between meal preference and frequency of eating out.

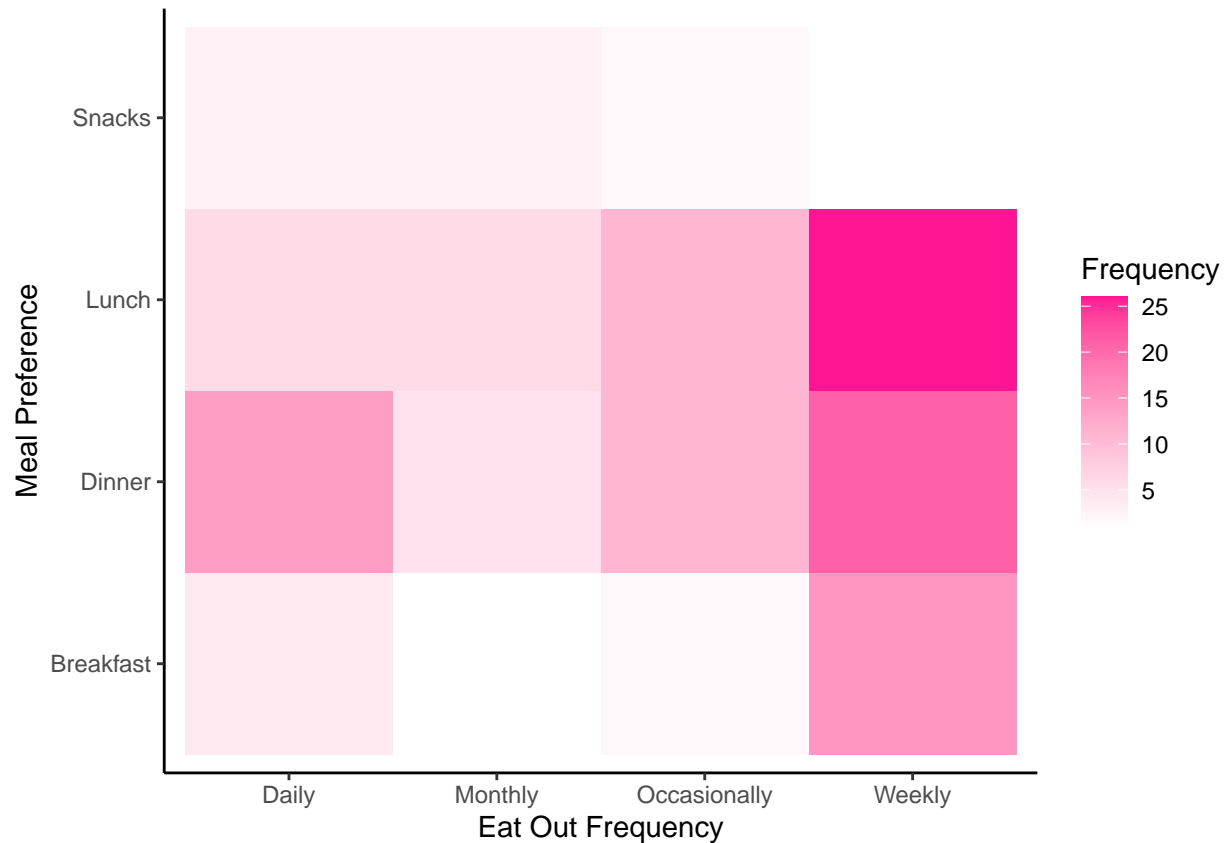
The p-value is greater than the typical significance level of 0.05, and Calculate G is less than table value, therefore we fail to reject the null hypothesis which suggests that there is no significant association between meal preference and frequency of eating out at the 5% level.

###Visualization:

```
# Convert the contingency table to a data frame
table_df <- as.data.frame(table)

# Rename the columns
colnames(table_df) <- c("meal_pref", "eat_out_freq", "freq")

# Create a heatmap using ggplot2
ggplot(table_df, aes(x = eat_out_freq, y = meal_pref, fill = freq)) +
  geom_tile() +
  scale_fill_gradient(low = "white", high = "deeppink") +
  labs(x = "Eat Out Frequency", y = "Meal Preference", fill = "Frequency") +
  theme_classic()
```



The heatmap shows that “Lunch” is the most common meal for eating out, especially on a weekly basis, as indicated by the darkest pink shade. “Dinner” is also popular, particularly on an occasional and weekly basis. Eating out daily is less common across all meal types. “Snacks” and “Breakfast” have the lowest frequencies, with lighter shades across all eating-out frequencies. Overall, weekly dining out, especially for lunch, is the most frequent behavior.

Logistic Regression and Model Evaluation on Dietary Preferences:

##Data Transformation and Initial Model Fitting:

```
data$diet_pref <- ifelse(data$diet_pref == "Vegetarian", 1, 0)
model=lm(data$diet_pref~.,data=data)
```

The diet_pref column is converted into a binary variable where “Vegetarian” is coded as 1 and others as 0. A linear model is fitted with diet_pref as the response variable and all other variables as predictors.

###Test for Normality:

##Shapiro-Wilk Test:

- H_0 :The residuals are normally distributed
- H_1 :The residuals are not normally distributed.

```
r=rstudent(model)
shapiro.test(r)
```

##

```
## Shapiro-Wilk normality test
##
## data:  r
## W = 0.94288, p-value = 3.125e-05
```

Here $p\text{-value} = 3.125e-05 < 0.05$, we reject H_0 . Hence the residuals are not normally distributed.

Model Selection

Forward Selection

Forward selection helps in identifying the most important predictors for `diet_pref` by incrementally adding variables that improve model performance.

```
fitstart=lm(data$diet_pref~1,data = data)
model <- lm(data$diet_pref ~ ., data =data)
fwd=step(fitstart,direction = 'forward',scope = formula(model))
```

```
## Start:  AIC=-207.1
## data$diet_pref ~ 1
##
##           Df Sum of Sq  RSS    AIC
## + gender      1   1.72774 24.822 -213.91
## + fav_cuisine   3   2.17643 24.373 -212.30
## + spend_week    3   1.69757 24.852 -209.76
## + eat_out_freq  3   1.41083 25.139 -208.25
## <none>                26.550 -207.10
## + new_food      3   0.96976 25.580 -205.97
## + health_conscious 2   0.40881 26.141 -205.13
## + meal_pref     3   0.78911 25.761 -205.05
## + pref_typefood  3   0.44348 26.106 -203.31
## + shop_place    3   0.30528 26.244 -202.62
## + age           3   0.16263 26.387 -201.91
## + cook_freq     3   0.06792 26.482 -201.44
##
## Step:  AIC=-213.91
## data$diet_pref ~ gender
##
##           Df Sum of Sq  RSS    AIC
## + fav_cuisine   3   1.83157 22.990 -217.96
## + spend_week    3   1.16304 23.659 -214.20
## <none>                24.822 -213.91
## + health_conscious 2   0.49400 24.328 -212.55
## + eat_out_freq  3   0.81109 24.011 -212.27
## + new_food      3   0.72441 24.098 -211.79
## + meal_pref     3   0.51163 24.310 -210.64
## + pref_typefood  3   0.32241 24.500 -209.63
## + shop_place    3   0.20435 24.617 -209.00
## + age           3   0.19776 24.624 -208.96
## + cook_freq     3   0.04080 24.781 -208.13
##
## Step:  AIC=-217.96
```

```
## data$diet_pref ~ gender + fav_cuisine
##
##           Df Sum of Sq    RSS    AIC
## <none>                22.990 -217.96
## + spend_week         3   0.92108 22.069 -217.31
## + eat_out_freq        3   0.76403 22.226 -216.38
## + health_conscious    2   0.38126 22.609 -216.15
## + new_food            3   0.54980 22.441 -215.13
## + meal_pref           3   0.52252 22.468 -214.97
## + pref_typefood       3   0.44846 22.542 -214.54
## + shop_place          3   0.22933 22.761 -213.27
## + age                 3   0.19431 22.796 -213.07
## + cook_freq           3   0.04479 22.945 -212.21
```

Gender and favourite cuisine are the most important predictors for diet_pref.

Backward Elimination

Backward elimination provides another method to refine the model by removing predictors that do not contribute meaningfully to the prediction of diet_pref.

```
bwd=step(model,direction = "backward")
```

```
## Start:  AIC=-191.86
## data$diet_pref ~ age + gender + fav_cuisine + eat_out_freq +
##      meal_pref + health_conscious + pref_typefood + spend_week +
##      shop_place + cook_freq + new_food
##
##           Df Sum of Sq    RSS    AIC
## - age         3   0.02901 18.895 -197.66
## - cook_freq    3   0.11621 18.982 -197.05
## - eat_out_freq 3   0.39817 19.264 -195.12
## - new_food     3   0.43443 19.300 -194.88
## - shop_place   3   0.64646 19.512 -193.44
## - pref_typefood 3   0.66963 19.535 -193.29
## - meal_pref    3   0.67898 19.545 -193.23
## - health_conscious 2   0.50646 19.372 -192.39
## <none>                18.866 -191.86
## - spend_week    3   0.89221 19.758 -191.81
## - gender        1   0.39994 19.266 -191.11
## - fav_cuisine    3   1.38246 20.248 -188.59
##
## Step:  AIC=-197.66
## data$diet_pref ~ gender + fav_cuisine + eat_out_freq + meal_pref +
##      health_conscious + pref_typefood + spend_week + shop_place +
##      cook_freq + new_food
##
##           Df Sum of Sq    RSS    AIC
## - cook_freq      3   0.12716 19.022 -202.78
## - eat_out_freq    3   0.41164 19.306 -200.83
## - new_food        3   0.45686 19.352 -200.53
## - shop_place      3   0.67447 19.569 -199.06
```



```

## - pref_typefood      3    0.68783 19.583 -198.97
## - meal_pref          3    0.75698 19.652 -198.51
## - health_conscious   2    0.51315 19.408 -198.15
## <none>                18.895 -197.66
## - spend_week         3    0.91639 19.811 -197.45
## - gender              1    0.39486 19.290 -196.95
## - fav_cuisine         3    1.45284 20.348 -193.95
##
## Step:  AIC=-202.78
## data$diet_pref ~ gender + fav_cuisine + eat_out_freq + meal_pref +
##   health_conscious + pref_typefood + spend_week + shop_place +
##   new_food
##
##           Df Sum of Sq    RSS    AIC
## - new_food      3    0.38897 19.411 -206.13
## - eat_out_freq   3    0.44946 19.471 -205.72
## - pref_typefood   3    0.58527 19.607 -204.81
## - shop_place     3    0.63639 19.658 -204.47
## - health_conscious 2    0.49001 19.512 -203.45
## - meal_pref      3    0.81691 19.839 -203.27
## - spend_week     3    0.85281 19.875 -203.03
## <none>           19.022 -202.78
## - gender         1    0.43234 19.454 -201.83
## - fav_cuisine     3    1.43919 20.461 -199.22
##
## Step:  AIC=-206.13
## data$diet_pref ~ gender + fav_cuisine + eat_out_freq + meal_pref +
##   health_conscious + pref_typefood + spend_week + shop_place
##
##           Df Sum of Sq    RSS    AIC
## - eat_out_freq   3    0.44866 19.860 -209.13
## - pref_typefood   3    0.54663 19.957 -208.49
## - shop_place     3    0.56392 19.975 -208.38
## - health_conscious 2    0.36281 19.774 -207.70
## <none>           19.411 -206.13
## - meal_pref      3    0.92596 20.337 -206.02
## - spend_week     3    1.07143 20.482 -205.09
## - gender         1    0.48327 19.894 -204.91
## - fav_cuisine     3    1.60607 21.017 -201.71
##
## Step:  AIC=-209.13
## data$diet_pref ~ gender + fav_cuisine + meal_pref + health_conscious +
##   pref_typefood + spend_week + shop_place
##
##           Df Sum of Sq    RSS    AIC
## - shop_place     3    0.46119 20.321 -212.13
## - pref_typefood   3    0.59149 20.451 -211.29
## - health_conscious 2    0.51774 20.377 -209.76
## - meal_pref      3    0.91912 20.779 -209.21
## <none>           19.860 -209.13
## - gender         1    0.52231 20.382 -207.73
## - spend_week     3    1.23779 21.097 -207.21
## - fav_cuisine     3    1.49768 21.357 -205.61
##

```

```

## Step: AIC=-212.13
## data$diet_pref ~ gender + fav_cuisine + meal_pref + health_conscious +
##   pref_typefood + spend_week
##
##           Df Sum of Sq   RSS   AIC
## - pref_typefood    3   0.72676 21.047 -213.52
## - meal_pref        3   0.80377 21.125 -213.04
## - health_conscious  2   0.49348 20.814 -212.98
## <none>                20.321 -212.13
## - spend_week       3   1.03440 21.355 -211.62
## - gender           1   0.63012 20.951 -210.12
## - fav_cuisine      3   1.51776 21.838 -208.69
##
## Step: AIC=-213.52
## data$diet_pref ~ gender + fav_cuisine + meal_pref + health_conscious +
##   spend_week
##
##           Df Sum of Sq   RSS   AIC
## - meal_pref        3   0.68239 21.730 -215.34
## - health_conscious  2   0.42215 21.470 -214.92
## - spend_week       3   0.94104 21.988 -213.79
## <none>                21.047 -213.52
## - fav_cuisine      3   1.39866 22.446 -211.09
## - gender           1   0.83455 21.882 -210.43
##
## Step: AIC=-215.34
## data$diet_pref ~ gender + fav_cuisine + health_conscious + spend_week
##
##           Df Sum of Sq   RSS   AIC
## - health_conscious  2   0.33935 22.069 -217.31
## - spend_week       3   0.87917 22.609 -216.15
## <none>                21.730 -215.34
## - fav_cuisine      3   1.41241 23.142 -213.09
## - gender           1   1.01179 22.742 -211.38
##
## Step: AIC=-217.31
## data$diet_pref ~ gender + fav_cuisine + spend_week
##
##           Df Sum of Sq   RSS   AIC
## - spend_week      3   0.92108 22.990 -217.96
## <none>                22.069 -217.31
## - fav_cuisine     3   1.58961 23.659 -214.20
## - gender          1   0.97216 23.041 -213.67
##
## Step: AIC=-217.96
## data$diet_pref ~ gender + fav_cuisine
##
##           Df Sum of Sq   RSS   AIC
## <none>                22.990 -217.96
## - fav_cuisine     3   1.8316 24.822 -213.91
## - gender          1   1.3829 24.373 -212.30

```

Gender and favourite cuisine are the most important predictors for diet_pref selected through backward selection.

Logistic Regression Model

```
glm=glm(data$diet_pref~ gender + fav_cuisine,data=data,family = binomial)
summary(glm)
```

```
##
## Call:
## glm(formula = data$diet_pref ~ gender + fav_cuisine, family = binomial,
##      data = data)
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)      0.02355    0.79459   0.030  0.97636
## genderMale       -1.20226    0.45014  -2.671  0.00757 **
## fav_cuisineItalian -0.71669    1.17532  -0.610  0.54200
## fav_cuisineNorth Indian  0.85404    0.94676   0.902  0.36702
## fav_cuisineSouth Indian -0.77279    0.82848  -0.933  0.35093
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##      Null deviance: 155.95  on 130  degrees of freedom
## Residual deviance: 138.44  on 126  degrees of freedom
## AIC: 148.44
##
## Number of Fisher Scoring iterations: 4
```

A generalized linear model (GLM) is fitted using gender and fav_cuisine as predictors, with diet_pref as the binary outcome. The logistic regression analysis shows that gender is a significant predictor of dietary preference. Specifically, males are significantly less likely to be vegetarian compared to females (Estimate = -1.20226, $p = 0.00757$). The decrease from 155.95 (null deviance) to 138.44 (residual deviance) suggests that the model with gender and fav_cuisine improves the fit compared to a model with no predictors.

Model Evaluation

```
# Load necessary libraries
library(caret)      # For confusionMatrix and metrics
```

```
## Warning: package 'caret' was built under R version 4.3.3
```

```
## Loading required package: lattice
```

```
##
```

```
## Attaching package: 'caret'
```

```
## The following objects are masked from 'package:DescTools':
```

```
##
```

```
##      MAE, RMSE
```

```

library(pROC)      # For ROC curve and AUC

## Warning: package 'pROC' was built under R version 4.3.2

## Type 'citation("pROC")' for a citation.

##
## Attaching package: 'pROC'

## The following objects are masked from 'package:stats':
##
##     cov, smooth, var

# Predict probabilities
pred_probs <- predict(glm, data, type = "response")

# Convert probabilities to binary predictions
pred <- ifelse(pred_probs > 0.5, 1, 0)

# Calculate confusion matrix
conf_matrix <- confusionMatrix(as.factor(pred), as.factor(data$diet_pref))

# Extract accuracy, sensitivity, specificity
accuracy <- conf_matrix$overall['Accuracy']
sensitivity <- conf_matrix$byClass['Sensitivity']
specificity <- conf_matrix$byClass['Specificity']

# Print accuracy, sensitivity, and specificity
print(paste("Accuracy:", round(accuracy, 4)))

## [1] "Accuracy: 0.771"

print(paste("Sensitivity:", round(sensitivity, 4)))

## [1] "Sensitivity: 0.9574"

print(paste("Specificity:", round(specificity, 4)))

## [1] "Specificity: 0.2973"

# Calculate F1 Score
precision <- conf_matrix$byClass['Precision']
recall <- sensitivity
f1_score <- 2 * (precision * recall) / (precision + recall)

print(paste("F1 Score:", round(f1_score, 4)))

## [1] "F1 Score: 0.8571"

```

```
# Print AUC value
roc_curve <- roc(data$diet_pref, pred_probs)
```

```
## Setting levels: control = 0, case = 1
```

```
## Setting direction: controls < cases
```

```
auc_value <- auc(roc_curve)
print(paste("AUC:", round(auc_value, 4)))
```

```
## [1] "AUC: 0.6998"
```

The model shows good accuracy (77.1%) and is excellent at identifying vegetarians (sensitivity: 95.74%). However, it struggles to correctly identify non-vegetarians (specificity: 29.73%). The F1 Score (0.8571) reflects a good balance between precision and recall, while the AUC (0.6998) indicates moderate overall discrimination ability.

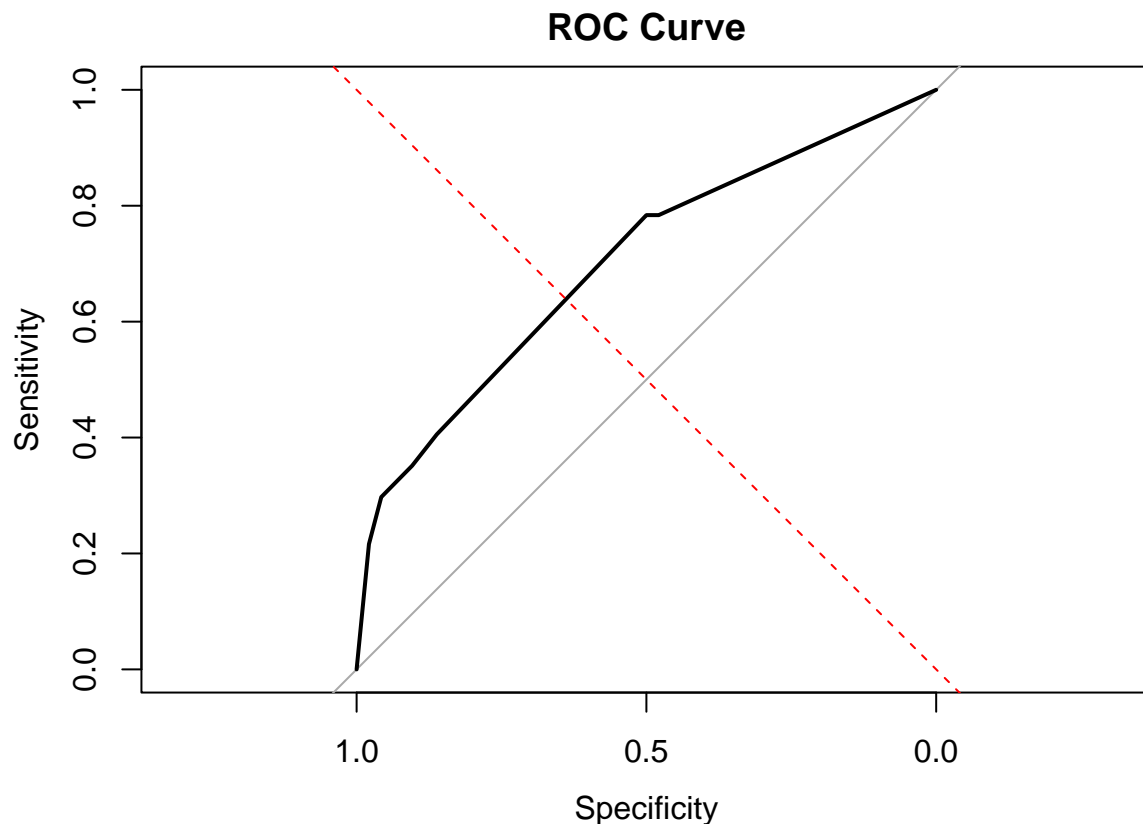
ROC curve

```
# Plot ROC curve
roc_curve <- roc(data$diet_pref, pred_probs)
```

```
## Setting levels: control = 0, case = 1
```

```
## Setting direction: controls < cases
```

```
auc_value <- auc(roc_curve)
plot(roc_curve, main = "ROC Curve")
abline(a = 0, b = 1, lty = 2, col = "red") # Add a diagonal line
```



The ROC curve shows that the model has moderate performance. The curve is close to the diagonal, indicating that the model's ability to distinguish between classes is not strong. This suggests the model is only somewhat better than random guessing, with room for improvement.

The model is very effective at detecting vegetarians but struggles with accurately identifying non-vegetarians. The high sensitivity and F1 score suggest it's well-tuned for situations where identifying vegetarians is more critical, but the low specificity indicates it may produce many false positives for vegetarians.

Fitting of Multinomial Logistic Regression:

```
table1 = table(data$shop_place, data$cook_freq)
table1
```

```
##
##           Daily Occasionally Several times a week Weekly
## Local Market      22         23              8         2
## Online             8          7              7         2
## Organic Store      2          2              1         2
## Supermarket       11         13             13         8
```

```
# Load the necessary libraries
```

```
library(nnet)
```

```
# Convert the shop_place and cook_freq variables to factors
```

```
data$shop_place <- factor(data$shop_place)
```

```
data$cook_freq <- factor(data$cook_freq)
```

```

# Perform multinomial logistic regression analysis
multinom_shop_place <- multinom(shop_place ~ cook_freq, data = data)

## # weights: 20 (12 variable)
## initial value 181.604561
## iter 10 value 150.405605
## final value 150.364156
## converged

# Summarize the results
summary(multinom_shop_place)

## Call:
## multinom(formula = shop_place ~ cook_freq, data = data)
##
## Coefficients:
## (Intercept) cook_freqOccasionally cook_freqSeveral times a week
## Online -1.0115807 -0.17801567 0.8779944
## Organic Store -2.3981398 -0.04437648 0.3185701
## Supermarket -0.6931075 0.12257508 1.1785653
## cook_freqWeekly
## Online 1.011652
## Organic Store 2.398215
## Supermarket 2.079430
##
## Std. Errors:
## (Intercept) cook_freqOccasionally cook_freqSeveral times a week
## Online 0.4128589 0.5973178 0.6620474
## Organic Store 0.7386327 1.0436177 1.2925406
## Supermarket 0.3692701 0.5067150 0.5816177
## cook_freqWeekly
## Online 1.081875
## Organic Store 1.243213
## Supermarket 0.872571
##
## Residual Deviance: 300.7283
## AIC: 324.7283

# Extract the coefficients for each category
coef_shop_place <- coef(multinom_shop_place)

# Calculate the odds ratios for each category
or_shop_place <- exp(coef_shop_place)

print("Odds Ratios for shop_place:")

## [1] "Odds Ratios for shop_place:"

print(or_shop_place)

## (Intercept) cook_freqOccasionally cook_freqSeveral times a week

```

## Online	0.36364371	0.8369293	2.406069
## Organic Store	0.09088686	0.9565937	1.375160
## Supermarket	0.50001985	1.1304040	3.249709
##	cook_freqWeekly		
## Online	2.750141		
## Organic Store	11.003518		
## Supermarket	7.999911		

The multinomial logistic regression shows that people who cook more frequently (especially weekly) are much more likely to shop at supermarkets and organic stores. Those who cook less frequently, like occasionally or several times a week, are more inclined to shop online. The odds ratios indicate a strong preference for supermarkets and organic stores among regular cooks, with the highest likelihood seen in weekly cooks.

Fitting of Log-Linear Model:

Null Hypothesis (H0): The log-linear model accurately represents the relationships among the categorical variables in the data. In other words, the model fits the data well and there are no significant deviations from what is expected.

Alternative Hypothesis (H1): The log-linear model does not accurately represent the relationships among the categorical variables.

```
library(MASS)

# Create a contingency table and set dimension names
contingency_table <- table(data$age, data$pref_typefood, data$spend_week, data$new_food)
dimnames(contingency_table) <- list(
  age = levels(data$age),
  pref_typefood = levels(data$pref_typefood),
  spend_week = levels(data$spend_week),
  new_food = levels(data$new_food)
)

# Fit the log-linear model using loglm
loglin_model <- loglm(~ age + pref_typefood + spend_week + new_food, data = contingency_table)
loglin_model
```

```
## Call:
## loglm(formula = ~age + pref_typefood + spend_week + new_food,
##       data = contingency_table)
##
## Statistics:
##              X^2   df    P(> X^2)
## Likelihood Ratio 138.3906 243 0.99999999
## Pearson          274.7579 243 0.07895611
```

The log-linear model fit well according to the statistics. The Likelihood Ratio Chi-Square value is 138.39 with a p-value close to 1(>0.05), indicating a good model fit. The Pearson Chi-Square value is 274.76 with a p-value of 0.079 (>0.05), suggesting no significant deviation between observed and expected frequencies. Overall, these results imply that the log-linear model adequately represents the relationships among the categorical variables of the data.

Conclusion:

The study successfully explored dietary preferences and food consumption patterns using categorical data analysis. Key findings include:

Gender and Dietary Preferences: Females were found to be significantly less likely to prefer non-vegetarian food compared to males, with both odds ratio and relative risk analyses supporting this conclusion.

Health Consciousness and Food Choices: A significant association was found between health consciousness and preferred food type, indicating that individuals who value health are more likely to prefer home-cooked or organic foods.

Meal Preferences and Eating Out Frequency: The likelihood ratio test suggested no significant association between meal preference and the frequency of eating out, indicating that these factors may be independent of each other.

Predictive Modeling: Logistic regression analysis identified gender and favorite cuisine as significant predictors of dietary preferences. The model demonstrated good accuracy in predicting vegetarians but struggled with specificity in identifying non-vegetarians.

Multinomial and Log-Linear Modeling: The multinomial logistic regression revealed strong preferences for supermarkets and organic stores among frequent cooks, while the log-linear model adequately represented the relationships among various categorical variables related to food consumption.

Overall, the findings provide valuable insights into the factors influencing dietary preferences and food consumption patterns, offering potential applications in public health, marketing, and consumer behavior research. The study's methodology and statistical analyses contribute to a deeper understanding of the complex interactions between demographic, behavioral, and psychological factors in shaping food choices.

Further Recommendation:

Based on the findings of this study, several recommendations can be made to better address the diverse dietary preferences and food consumption patterns observed:

Targeted Health Promotion Campaigns: Public health initiatives should focus on increasing awareness about healthy food choices, especially among groups that show lower health consciousness. Tailored campaigns could encourage healthier eating habits, particularly for individuals who prefer fast food or street food.

Supporting Sustainable Eating Practices: Considering the growing interest in health-conscious and organic food, policies that support sustainable farming practices and make organic products more accessible and affordable could align with consumer preferences and promote overall well-being.

Developing Mobile Apps for Dietary Guidance: Technology companies could develop or enhance mobile apps that provide personalized dietary advice based on user input regarding their food preferences, health consciousness, and eating habits. This could help users make informed food choices in real-time.