

# Nutrition Project Report

*Dong Luo*

*22 October 2018*

## Contents

|          |   |          |
|----------|---|----------|
| <b>1</b> | <b>Executive summary</b>                      | <b>1</b> |
| <b>2</b> | <b>Methodologies</b>                          | <b>1</b> |
| 2.1      | Proportion of the diets . . . . .             | 1        |
| 2.2      | Comparing Demographics by Diet Type . . . . . | 3        |
| <b>3</b> | <b>limitations</b>                            | <b>6</b> |

## 1 Executive summary

Obesity affects approximately one third of the Australian population aged 18 years and over (Australian Bureau of Statistics, 2015), posing major health risks to individuals. Obesity has been demonstrated to be related equally to genetics and environmental factors including diet. A person who is dieting has been defined as one who consumes macronutrient groups in an amount lying outside of the acceptable macronutrient distribution ranges (National Health and Medical Research Council, 2017). Here we are interested in whether consuming a diet outside of normal macro-nutrient ranges may be related to Body Mass Index BMI, and thus obesity, and whether a person's demographics, for example sex, age, waist size, and socioeconomic status is predictive of a person's choice of diet. To investigate this, we developed three research questions.

1. What are the most common diet types among adults in the sample, categorised by the macro- nutrients fat, proteins and carbohydrates.
2. Investigated whether people on different diets have different characteristics and demographics, for example socio- economic status, age, sex, Basal Metabolic Rate (BMR), and BMI.
3. To see if we could predict obesity measured using BMI based on eight variables, including BMR, energy intake, sex, and time spend sedentary.

Throughout the analysis, we focussed on the macro- nutrients protein, fat and carbohydrates. These were because three common types of diets trends have been identified (Kossoff, Turner, Doerrer, Cervenka, Henry, 2016) which are all combinations of different levels of each protein, fat, and carbohydrates. These diets are the Keto diet (high fat, low carbohydrates, medium protein, which made up of 18% of our sample), the Atkins diet (high protein, low carbohydrates, medium fat, 7% of our sample), and the Dash diet (high carbohydrates, low fat, 1.7% of our sample).

## 2 Methodologies

### 2.1 Proportion of the diets

We excluded the observations where age was less than 18. This was because our research questions focussed on obesity which we determined using the BMI. However BMI does not apply to children, so including under 18s may lead to misleading results.

We cut each of the continuous variables `CHOPER1`, `FATPER1`, and `PROPER1`, which are the percentages of total energy coming from carbohydrates, fat, and protein into three distinct levels, low, medium, and high. Table 1 shows the percentages of total energy used to divide the groups.

Table 1: Division of Macro- nutrients by Total Energy Percentage

|                   | Carbohydrates | Fat      | Protein  |
|-------------------|---------------|----------|----------|
| <b>Low (%)</b>    | [0,45]        | [0,20]   | [0,15]   |
| <b>Medium (%)</b> | (45,65]       | (20,35]  | (15,25]  |
| <b>High (%)</b>   | (65,100]      | (35,100] | (25,100] |

We divided the data this way because we were interested in the mean proportion of each diet, so that we could identify the most popular diets within our sample.

Table 2: Proportion of the top 6 diet types

| proportion | fat    | carb   | protein |
|------------|--------|--------|---------|
| 0.1902491  | medium | medium | medium  |
| 0.1785904  | high   | low    | medium  |
| 0.1383148  | medium | low    | medium  |
| 0.1267621  | medium | medium | low     |
| 0.0714361  | medium | low    | high    |
| 0.0557499  | high   | low    | low     |

We fitted log-linear models to see if there was any independence underlying the three-way contingency table.

### 2.1.1 Structural Zeroes

We noticed that there were some cells with value zero in our table. We treated the zeroes as structural zeroes (impossible combinations) since the variables **CHOPER1** (carbohydrate), **FATPER1** (fat), and **PROPER1** (protein) in the original dataset stand for the proportions. This means some of the diets listed in the table, such as high carbohydrate, high fat, high protein were not possible as the sum of the three proportions would exceed 100. We removed the cells of structural zeroes from the and model the incomplete table.

### 2.1.2 Log-Linear Models

There were 20 out of 27 cells with positive entries. The null model will have  $20 - 1 = 19$  degrees of freedom and the additive model will have 13 degrees of freedom.

We started with the additvie model, which stands for the complete independence.

We found the deviance for the additive model (including all three factors) was 3472.4957366. We compared it with the models with one two way interaction. There were three such models, and their residual deviance is shown as in Table 3.

The model with **carb:fat** interaction had the lowest deviance. The difference of deviance between this model and the additive model was 1827.5168456. The two models are nested and when we compared them, the  $H_0$  was the additive model (the smaller model).

$$M1(H_0) : \hat{\mu} = \beta_0 + \alpha_i + \beta_j + \gamma_k,$$

Table 3: Deviances of models with one interaction term

|                   | CF       | CP      | FP       |
|-------------------|----------|---------|----------|
| Residual Deviance | 1644.979 | 2828.39 | 3356.749 |

where  $\alpha_i$ ,  $\beta_j$ , and  $\gamma_k$  denote the  $i^{th}$ ,  $j^{th}$  and  $k^{th}$  group of carbohydrate, fat, and protein levels. The alternative hypothesis is

$$M2(H_A) : \hat{\mu} = \beta_0 + \alpha_i + \beta_j + \gamma_k + (\alpha\beta)_{ij}$$

Under  $H_0$ , the difference in deviance followed a  $\chi^2$  distribution whose degrees of freedom equals the difference in residual degrees of freedom of the two models.

If our table was complete, we would expect the difference in degrees of freedom to be  $(3 - 1) \times (3 - 1) = 4$ , since each factor has 3 levels.

Here the difference in degrees of freedom is 3. This is because there was a structural zero in the marginal table of carbohydrate and fat (we cannot have high carbohydrate and high fat at the same time), so we were only adding three parameters when fitting the model. A coefficient of an interaction level in the `glm` function output will be NA.

The  $p$ -value for the test is close to 0, so we would reject the null hypothesis and prefer the model with `carb:fat` interaction. This model with one interaction term stands for the block independence.

We use the same procedure to test the models with more interaction terms, and the result of the test is summarised in Table 4.

Table 4: Deviances for Poisson log-linear models

| Type                                  | Model    | Deviance | d.f. |
|---------------------------------------|----------|----------|------|
| <b>Completely Independence Model</b>  | C+F+P    | 3472     | 13   |
|                                       | P+CF     | 1645     | 10   |
| <b>Block Independence Model</b>       | F+CP     | 2828     | 10   |
|                                       | C+FP     | 3357     | 9    |
|                                       | CF+CP    | 668      | 7    |
| <b>Conditional Independence Model</b> | CF+FP    | 1566     | 6    |
|                                       | CP+FP    | 2752     | 6    |
| <b>Uniform Association Model</b>      | CF+CP+FP | 449      | 3    |
| <b>Saturated Model</b>                | CFP      | 0        | 0    |

The letters C, F, and P in the table represent the variables carbohydrate, fat, and protein respectively. Symbols such as FP indicate we are including the interaction term `fat:protein`, and when we include an interaction term we must also include the variables used to compute the interaction. Similarly, when we include the three-way interaction term CFP (`carb:fat:protein`), we include all two-way interactions as well as all of the three variables.

The deviance test suggests that we shall choose the saturated model, which implies

$$\hat{\mu}_{ijk} = y_{ijk}.$$

Therefore, the proportion for each diet can be estimated by the sample proportions.

## 2.2 Comparing Demographics by Diet Type

This section worked out the best diets to model different dependent variables; BMI, BMR, waist size (cm), minutes spent sedentary, sex, SES, and age. The focus of this area was to fit eight key demographics using protein, fat, and carbohydrate diets. Each of the variables was chosen through discussion with NUTM3001 students, and were modelled using linear regression for numerical variables and logistic regression for binary variables.

Prior to beginning model selection, assumptions for normality of errors and homogeneity of variance were assessed using Q-Q plots and box plots of residuals, and a residual versus fitted plot was used visually assess

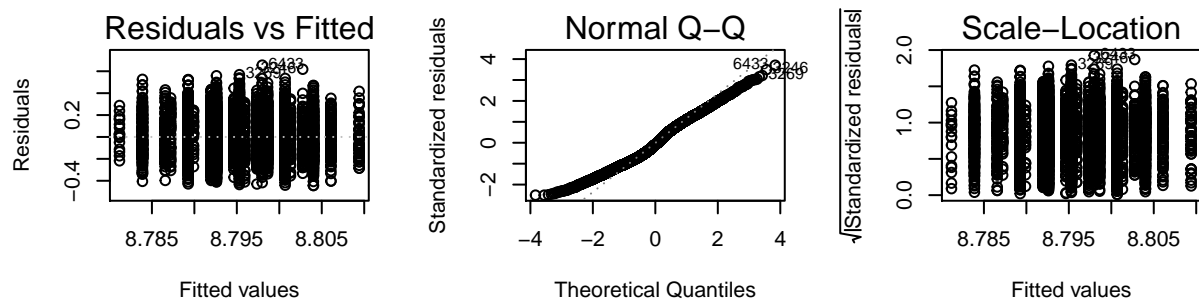


Figure 1: Model Diagnostics

homogeneity of variance and independence. An example of the graphs used has been provided in Figure 1 showing the diagnostic plots for the dependent variable BMR, excluding the boxplot of residuals. For variables BMR and sedentary the normality assumption was not met. We decided to apply a log transformation to BMR and a square root transformation to minutes sedentary to satisfy normality.

Once the the assumptions were met, we began model selection. We used residual deviance tests and AIC Stepwise model selection to find the “best” model for continuous dependent variables and the F test and AIC Stepwise model selection for binary dependent variables. For five of the dependent variables the “best” model was found to be the same using AIC and the residual deviance or F test. For BMR, SES, and minutes spent sedentary the models did not match. When the models did not match, the model selected by the F test was used as the “best” model because the AIC has been claimed to typically prefer overly complex models.

For each dependent variable using the F or deviance tests for nested models, we tested the null hypothesis that the variable was best modelled by a smaller model beginning with the null model, compared to the alternative hypothesis that the variable was best explained by a model with a single extra term. Linear and logistic regression tested the null hypothesis that adding an extra term will not give a better fit of the variable, compared to the alternative hypothesis that adding an extra term will increase the goodness of fit on the model.

A 5% significance level was chosen as a threshold for the inclusion of the model variables. The resulting “best” models explaining each of the dependent variables is found in Table 1.

Table 5: Best models from model selection

|                              | Deviance | AIC   | F- test |
|------------------------------|----------|-------|---------|
| BMI                          | P        | P     |         |
| BMR                          |          | C+PF  | P       |
| Waist (cm)                   |          | P+C   | P+C     |
| Mins sedentary               |          | P+F+C | P+F+C   |
| Sex (proportion of male)     | C+PF     | C+PF  |         |
| SES (proportion of high SES) | C        | C+PF  |         |
| Age                          |          | P+F+C | P+F+C   |

### 2.2.1 Demographics for diets

Table 6 presents the key characteristics of each of the three diets investigated. For BMR, waist measurement, minutes spent sedentary, and age, the mean of each group was reported. For BMI, sex, and SES respectively, the probability of a person being overweight, male, or in the normal decile Index of Relative Socio-Economic Disadvantage has been reported.

Table 6: Demographics for 8 key variables

|                | Atkins   | Keto     | Low Fat, High Carb |
|----------------|----------|----------|--------------------|
| BMI            | 0.358    | 0.282    | 0.282              |
| BMR            | 6724.316 | 6713.304 | 6897.837           |
| Waist (cm)     | 94.510   | 93.264   | 94.489             |
| Mins sedentary | 2259.420 | 2388.410 | 2054.230           |
| Sex            | 0.416    | 0.384    | 0.375              |
| SES            | 0.237    | 0.286    | 0.197              |
| Age            | 50.497   | 48.333   | 44.031             |

### 2.2.2 A couple of key things to note:

- BMI: We found the proportion of people who were obese was highest in the high protein diet. A recent CSIRO Protein Balance Report suggests that a high protein diet can improve body composition (Noakes, 2018). Whilst this is the opposite of our findings, it is possible that individuals with a larger waist circumference may have been consuming a high protein diet with the prospective goal of losing weight. It is unknown whether the diet has been successful in this regard as the AHS only records one point in time. Another consideration may be that people on high protein diets could be focussed on muscle development which could result in a higher body weight, and consequently a higher BMI. More data would need to be collected around subjects exercise activities and purpose for diet choices to test these theories.
- Waist: People on high protein diets had the largest mean waist measurement (95.204cm) compared to any of the other diets tested, including Atkins, Keto, and the low carbohydrate/ high fat diet. A person's waist size can show whether a person is carrying excess fat around their middle and can be an indicator of the level of internal fat deposits covering internal organs. In conjunction with BMI, they can help determine a person's risk of health issues, such as stroke or heart disease.
- Sedentary: The high fat low carbohydrate diet (keto diet) had the highest mean sedentary minutes of 2388.410, which equates to 39.8 hours of being sedentary over two days. This diet is often advertised as an effective diet for athletes, however this sample was not targeted towards athletes. It may have been interesting to look further into this comparing people who said that they were on a diet compared to not on a diet who fell into the Keto diet, and to look into their respective minutes spent sedentary and minutes spent exercising.
- Sex: Within our dataset 46% of subjects were male yet we found that the probability of being male on a high carbohydrate, medium fat, and medium protein diet was 60%.
- SES: Despite having a lower proportion of normal SES participants in the study, there was a high proportion who had a high carbohydrate diet but normal fat and protein. This is a surprising result as much literature supports low SES rather than normal and high SES having high carb diets since carbohydrates are cheaper than protein, and might be consumed excessively in order to meet dietary protein needs (Brooks, Simpson & Raubenheimer, 2010). However some literature comments on the inconsistencies of results related to protein intake in the context of lower SES (Darmon, N., Drewnowski, A., 2008).

### 3 limitations

Poisson distribution assumption, sampling zeroes.

Brooks, Simpson, R.C., and D. Raubenheimer. 2010. “The Price of Protein: Combining Evolutionary and Economic Analysis to Understand Excessive Energy Consumption.” *Obesity Reviews* 11 (12): 887–94.

Cassell, Flom, D.L. 2007. “Stopping Stepwise: Why Stepwise and Similar Selection Methods Are Bad, and What You Should Use.”

Darmon, A., N. & Drewnowski. 2008. “Does Social Class Predict Diet Quality?” *The American Journal of Clinical Nutrition* 87 (5).

Kossoff, Turner, E.H. 2016. “The Ketogenic and Modified Atkins Diet: Treatments for Epilepsy and Other Disorders.” *Obesity Reviews* 11 (12). Springer Publishing Company: 887–94.

Noakes, M. n.d. “Protein Balance: New Concepts for Protein in Weight Management.” CSIRO, Australia.

Thayer, J.D. 1990. “Implementing Variable Selection Techniques in Regression.” *Annual Meeting of the American Educational Research Association*, 16–20.