## 1 IDENTIFYING SMALL GROUPS OF FOODS THAT CAN PREDICT ACHIEVEMENT OF KEY DIETARY

# 2 RECOMMENDATIONS: DATA MINING OF THE UK NATIONAL DIET AND NUTRITION SURVEY, 2008-

- 3 **12**
- 4 Running head: Data mining National Diet & Nutrition Survey
- 5 Philippe J. <u>GIABBANELLI<sup>1,2</sup></u>, Jean <u>ADAMS<sup>1</sup></u>
- 6 <sup>1</sup>UKCRC Centre for Diet and Activity Research (CEDAR), MRC Epidemiology Unit, University of
- 7 Cambridge School of Clinical Medicine, Institute of Metabolic Science, Cambridge CB2 0QQ, United
- 8 Kingdom
- 9 <sup>2</sup>Department of Computer Science, Northern Illinois University, DeKalb IL, USA.
- 10 Note. PJG was based at affiliation (a) when this work was performed; he is now based at affiliation
- 11 (b).
- 12 \*Correspondence concerning this article should be addressed to: Dr Jean Adams, UKCRC Centre for
- 13 Diet and Activity Research (CEDAR), MRC Epidemiology Unit, University of Cambridge School of
- 14 Clinical Medicine, Institute of Metabolic Science, Cambridge CB2 0QQ, United Kingdom;
- 15 jma79@medschl.cam.ac.uk
- 16

# 17 **RESEARCH ETHICS**

- 18 Ethical approval for the National Diet and Nutrition Survey (NDNS) has obtained from the
- 19 Oxfordshire A Research Ethics Committee and all participants provided informed consent to take
- 20 part in the survey. Further ethical approval was not required for this secondary analysis of
- 21 anonymised data.
- 22

# 23 ACKNOWLEDGEMENTS

Thanks to Simon Wheeler at the MRC Epidemiology Unit, Cambridge University for help withinterpreting NDNS data.

#### 27 FINANCIAL SUPPORT

- 28 The work was undertaken by the Centre for Diet and Activity Research (CEDAR), a UKCRC Public
- 29 Health Research Centre of Excellence. Both authors gratefully acknowledge funding from the British
- 30 Heart Foundation, Cancer Research UK, Economic and Social Research Council, Medical Research
- 31 Council (MRC), the National Institute for Health Research (NIHR), and the Wellcome Trust, under the
- 32 auspices of the UK Clinical Research Collaboration (reference MR/K023187/1).
- 33

## 34 CONTRIBUTIONS

- 35 JA conceived the original idea for this work. PJG designed and conducted the data analysis, and
- 36 produced all figures and tables. Both authors interpreted the results and drafted the manuscript.
- 37

## 38 **COMPETING INTERESTS**

- 39 None.
- 40

## 41 ABSTRACT

## 42 Background

- 43 Many dietary assessment methods attempt to estimate total food and nutrient intake. If the
- 44 intention is simply to determine whether participants achieve dietary recommendations, this leads
- 45 to much redundant data. We used data mining techniques to explore the number of foods that
- 46 intake information was required on to accurately predict achievement, or not, of key dietary
- 47 recommendations.

## 48 Methods

- 49 We built decision trees for achievement of recommendations for fruit & vegetables, sodium, fat,
- saturated fat, and free sugar using data from the UK National Diet and Nutrition Survey (NDNS,
- 51 2008-12). Decision trees describe complex relationships between potential predictor variables (age,
- 52 sex, and all foods listed in the NDNS database) and outcome variables (achievement of each of the
- 53 recommendations).

## 54 Results

- 4156 individuals were included in the analysis. Information on consumption of 113 out of 3911 (3%)
- 56 foods, plus age and sex was required to accurately categorise individuals according to all five

- 57 recommendations. The best trade-off between decision tree accuracy and number of foods included
- 58 occurred at between 11 (for fruit and vegetables) and 32 (for fat, plus age) foods, achieving an
- 59 accuracy of 73% (for fat) to 83% (for fruit and vegetables), with similar values for sensitivity and
- 60 specificity.

## 61 Conclusions

- 62 Using information on intake of 113 foods, it is possible to predict with 73-83% accuracy whether
- 63 individuals achieve key dietary recommendations. Substantial further research is required to make
- 64 use of these findings for dietary assessment.

## 65 Keywords

66 Data mining; diet; dietary assessment; dietary pattern analysis; nutrition

#### 67 **INTRODUCTION**

- 68 The intention of many dietary assessment methods is to capture information on all foods consumed,
- 69 or at least those believed to make the largest contribution to total intake,<sup>(1)</sup> in order to estimate
- total nutrient intake. For some purposes, this detailed estimation of total nutrient intake may lead to
- 71 collection of much redundant data. This is particularly the case when assessing adherence with
- 72 policy targets and messages such as 'five-a-day' portions of fruit and vegetables.
- 73 The collection of substantial redundant information places unnecessary burden on research
- 74 participants, and unnecessarily uses scarce research resources. To take a first step to overcoming
- this problem, we applied data mining techniques to explore how many, and which, foods
- information was required on to accurately predict achievement, or not, of key dietary
- 77 recommendations.

#### 78 Data mining, an overview

- 79 Unlike traditional statistical approaches such as multiple regression, data mining allows multiple,
- 80 non-linear, relationships and interaction effects to be efficiently captured.<sup>(2; 3)</sup> Several data mining
- tools exist. In this study, we use 'classifiers'. A classifier is a function that labels individuals on an
- 82 outcome (e.g. achieving a dietary recommendation or not) based on a group of predictor variables
- 83 (e.g. how much of each individual food was consumed). The analysis package is first provided with a
- 84 'training set' of individual-level data in which both the outcome and the predictor variables are
- 85 known, and uses this to learn how the predictor variables are related to the outcome. This produces
- 86 the classifier function, which can then be used to infer the outcome in a new case based on just the
- 87 predictor variables. Finally, the accuracy of the classifier is evaluated on a new 'testing set' of data.
- 88 There are numerous ways to build classifiers. We used 'decision trees'.<sup>(2; 4; 5)</sup> Decision trees provide a
- 89 graphical illustration of a classifier composed of a number of predictor variables. A decision tree
- 90 involves repeated 'cuts' of the data according to the level of included predictor variables to identify
- 91 groups of individuals who are similar in terms of the outcome variable of interest. This produces a
- 92 decision tree where the path from the root to the outcome corresponds to successive 'cuts', or
- 93 divisions, of the population.
- 94 Figure 1 provides a simplified, hypothetical example of a decision tree where the intention is to
- 95 identify whether or not individuals achieve the recommended intake of fruit and vegetables (the
- 96 outcome) using information on consumption of carrots and white bread (the two predictor
- 97 variables). Figure 1a shows the decision tree based on the 'cuts' represented in Figure 1b. Figure 1b
- 98 is a simple graphical plot of consumption of both carrots and white bread with all individuals labelled
- 99 according to whether or not they achieve the recommended intake of fruit and vegetables. There

appear to be five 'clusters' of participants in Figure 1b in terms of meeting fruit and vegetable

- 101 recommendations. A series of 'cuts' can isolate these clusters. The first cut (labelled 'A' in both
- 102 Figure 1a and Figure 1b) divides the population according to consumption of carrots. The next two
- 103 cuts (labelled 'B' and 'C') then divide the resulting two groups according to consumption of white
- bread. Finally, a fourth cut (labelled 'D') divides those with a medium carrot and medium white

105 bread intake according to a more fine-grained assessment of carrot intake.

- 106 To build decision trees with different numbers of predictor variables, the minimum number of
- 107 individual cases that can be further divided by a subsequent 'cut' is varied. If a small group of
- 108 individuals can be further sub-divided, a sizable tree including many predictor variables can result.
- 109 However, if limits are placed on the minimum size of group that can be further sub-divided, a smaller
- 110 decision tree, including fewer predictor variables, results. In the current study, we make use of this
- 111 feature to explore the effect of including more or fewer predictor variables on the accuracy of
- 112 decision trees.
- 113 A small number of studies have applied data mining techniques to nutritional data. These have
- 114 primarily focused on dietary pattern analysis, exploring which dietary components are predictive of a
- range of health outcomes.<sup>(6) (7) (8) (9)</sup> However, we are not aware of any other uses of data mining to
- identify which foods are predictive of achievement, or not, of key dietary recommendations.

#### 117 Aims

- 118 Our aim was: to use data mining techniques to determine the number of foods that intake
- 119 information was required on to accurately predict achievement, or not, of dietary recommendations
- 120 for intake of fruits & vegetables, free sugars, sodium, fat, and saturated fat.

## 121 METHODS

- 122 We built decision trees for achievement of key dietary recommendations using data from the first
- 123 four years of the rolling programme of the UK's national dietary surveillance dataset: the National
- 124 Diet and Nutrition Survey (NDNS).

#### 125 Data source

- 126 The NDNS is an annual cross-sectional survey assessing the diet, nutrient intake and nutritional
- 127 status of the general population aged 18 months and upwards living in private households in the
- 128 UK.<sup>(10)</sup> Since 2008, an annual 'rolling programme' has been in place, allowing data to be combined
- 129 over years. We used data from years 1-4 of this programme, collected in 2008-12.
- 130 The NDNS aims to collect data from a sample of 1,000 respondents per year: at least 500 adults
- 131 (aged 19 years and older) and at least 500 children (aged 1.5 to 18 years). Households across the UK

- are selected to take part in the NDNS using a multi-stage probability design. In each wave, a random
- 133 sample of primary sampling units is selected for inclusion. These are small geographical areas that
- allow more efficient data collection by enabling it to be geographically focused. Within these
- primary sampling units, private addresses are randomly selected for inclusion. If, on visiting, it is
- 136 found that more than one household lives at a particular address, one is randomly selected for
- 137 inclusion. Within participating households, up to one adult and one child are randomly selected to
- 138 take part as 'respondents'. Data collection includes completion of four-day estimated food diary –
- 139 where participants estimate the weight of foods consumed using food labels and household
- 140 measures.<sup>(11)</sup>
- NDNS data were obtained from the UK Data Archive an online resource that makes research data
  available to the UK research community.
- 143 Inclusion and exclusion criteria
- 144 NDNS participants were included in the analysis if they completed three or four days of the
- 145 estimated food diary. As recommendations for fruit and vegetable intake only apply to those aged
- 146 11 years or older, children aged less than 11 years were excluded from this component of the
- 147 analysis.

## 148 Outcomes of interest – achievement of dietary recommendations

- 149 Information on which foods were consumed, and how much participants estimated was consumed,
- 150 was combined with nutritional information to determine mean daily intake of fruit and vegetables
- 151 (80g portions), and sodium (mg); and mean daily percentage of energy derived from fat, saturated
- 152 fat, and free sugars for each individual. This information was then used to determine whether or not
- 153 each individual met international, or UK, recommendations for these variables.
- 154 We used UK recommendations or fruit and vegetable and sodium intake, as these have been graded
- according to age. It is recommended that individuals aged 11 years and older consume at least five
- 156 80g portions of fruit and vegetables per day. This includes a maximum of one portion of juice, with
- additional juice portions not counted. For sodium, current UK recommendations are that those aged
- 158 11 years and older consume no more than 2400mg per day; children aged 7-10 years, no more than
- 159 2000mg; children aged 4-6 year, no more than 1200mg; and children aged 1-3 years, no more than
- 160 800mg.<sup>(12)</sup>
- 161 The World Health Organization recommends population food and nutrient intake goals for the
- avoidance of diet related diseases. These state that no more than 30% of energy should be derived
- 163 from fat, no more than 10% from saturated fatty acids, and no more than 10% from free sugars.<sup>(13)</sup>

#### 164 **Predictor variables of interest – foods consumed**

In total, 3911 different foods (including drinks) have been recorded in NDNS food diaries. We used
total estimated weight (in grams) of each individual food eaten by each individual as potential
predictor variables. Age and sex were also included as potential predictor variables. The use of
including markers of socio-economic position (education, income, and social class) as potential
predictor variables was explored but these were found to add no additional increase in accuracy
over and above age, sex and individual foods. Decision trees reported here do not include any socioeconomic predictor variables.

#### 172 Data analysis

- 173 Our analysis scripts and detailed decision trees are available at https://osf.io/znv82. In all cases
- 174 except sodium, the proportion of individuals achieving the recommendations was substantially less
- than 50%; for sodium substantially more than 50% of individuals achieved the recommendations
- 176 (Table 1). As detailed in Supplementary File 1, this imbalance in outcome variables can lead to low-
- 177 quality classifiers. To correct this, we pre-processed the data using the Synthetic Minority Over-
- 178 sampling TEchnique (SMOTE),<sup>(14)</sup> which creates new cases for the group which accounted for less
- 179 than 50% of participants by interpolating between existing cases that lie together. WEKA software<sup>(15)</sup>
- 180 was then used to build decision trees using the J48 algorithm and error pruning.
- For each outcome of interest we built a series of decision trees with different numbers of predictor 181 182 variables by varying the minimum number of individual cases that could be further divided. For each 183 of the decision trees built, we calculated the number of predictor variables used and overall 184 accuracy in correctly classifying individuals. We used the standard 10-fold cross-validation procedure<sup>(16)</sup> in which the entire eligible NDNS dataset was split into 10 approximately equally sized 185 186 parts. Nine parts were used in turn as training sets, and the remaining 10th part was used as testing 187 set. The ability of decision trees to correctly identify those who achieved the recommendations 188 (sensitivity) and those who did not (specificity) was also calculated. Adaptive sampling was used to 189 identify the maximum overall accuracy that could be achieved, as well as the optimum trade-off 190 between minimising number of predictor variables and maximising overall accuracy.

## 191 **RESULTS**

- 192 Overall, 91% of households eligible for inclusion agreed to take part in the first four waves of NDNS.
- 193 Within these, 56% (2083 adults and 2073 children; 4156 participants in total) of individuals selected
- to take part completed three or four days of the estimated food diary and were included in the
- analysis for sodium, free sugars, fat and saturated fat. Of these 4156 participants, 2967 (71.4%) were

aged 11 years or older and included in the analysis for fruit and vegetables. There were no missingdata on sex or age.

The distributions of age and sex in the analytical sample compared to the UK population as a whole are shown in Table 1. As the NDNS sample contains relatively equal numbers of children aged 18 years or younger, and adults, distributions are provided separately for adults and children in this table. The main differences between the age and sex distributions in the analytical sample and UK population were that the analytical sample had a higher proportion of adult women and a lower proportion of young adults (aged 19-29 years) than the UK population.

204 Figure 2 shows the overall accuracy of decision trees for each of the five outcomes plotted against 205 the number of predictor variables in decision trees. Overall accuracy ranged from 69% (fat; 10 206 predictor variables) to 84% (fruit and vegetables; 50 predictor variables) depending on the outcome 207 of interest and number of predictor variables included. For all guidelines but sodium, the 208 relationship between the number of predictor variables and the accuracy was best described using a 209 logarithmic trend model (p<0.01 in all cases). Thus, increasing the number of predictor variables 210 from around 10 to 30 improved the accuracy by a maximum of around five percentage points, but 211 beyond this adding even a large number of additional predictor variables yielded only a very small 212 additional improvement. We were unable to fit any function to the relationship between accuracy 213 and number of predictor variables for sodium.

214 Table 2 provides information on the decision tree for each outcome that represented the best trade-215 off between accuracy and number of predictor variables. Information on the most accurate possible 216 tree for each outcome is also shown in Table 2. Between 11 (for fruit and vegetables) and 33 (for fat) 217 predictor variables provided the best trade-off to identify whether individuals achieved each of the 218 recommendations, achieving overall accuracy of 73% (for fat) to 83% (for fruit and vegetables). 219 Adding further predictor variables beyond this improved accuracy by a maximum of 2% (for 220 saturated fat) and less than 1% (for all other outcomes). Sensitivity and specificity were similar to 221 overall accuracy for fruit and vegetables and free sugars (and saturated fat when the maximum 222 number of predictor variables were included). However, specificity was higher than sensitivity for fat 223 (and saturated fat), but the reverse was seen for sodium. Predictor variables in decision trees with 224 the best trade-off between accuracy and number of predictor variables accounted for between 13% 225 (for fat) and 31% (for free sugars) of total intake of relevant outcome variables.

Predictor variables used in decision trees with the best trade-off between accuracy and number of
predictor variables are shown in Table 3. In total, 113 foods (out of a total 3911 [3%] recorded as

228 consumed), age and sex were included in the decision trees for all five outcomes. Overall, there was

- 229 little overlap in predictor variables across outcomes. Age and two foods were included as predictor
- variables in the decision trees for three outcomes. A further six foods were included as predictor
- variables in the decision trees for two outcomes. The remaining 104 foods were included as
- 232 predictor variables in only one decision tree.

#### 233 DISCUSSION

#### 234 Summary of results

- This is the first work we are aware of using data mining techniques to explore the number of foods
- that information is required on to predict achievement of dietary recommendations. In total,
- 237 information on consumption of 113 of 3911 foods (3%), plus age and sex was required to accurately
- 238 categorise individuals according to all five dietary recommendations (fruit & vegetables, free sugars,
- sodium, fat, and saturated fat). The best trade-off between decision tree accuracy and number of
- foods included was achieved at between 11 (for fruit and vegetables) and 32 (for fat, plus age) foods.
- 241 These decision trees had an overall accuracy of 73% (for fat) to 83% (for fruit and vegetables), with
- similar values for sensitivity and specificity. Few individual foods were present in the decision tree
- 243 for more than one dietary recommendation, although age was present in three.
- 244 Strengths and limitations of methods
- 245 We used data from a population-based sample meaning our findings are likely to be generalizable
- across the UK and to other countries with similar dietary profiles. However, diets vary
- 247 internationally<sup>(17)</sup> and our results may not be more widely generalizable. The analytical sample had a

slightly higher proportion of adult women and lower proportion of younger adults (aged 19-29

- 249 years) than the UK population as a whole.
- 250 The data used were collected using 'estimated' food diaries where portion sizes were estimated
- but not weighed. These are considered to be one of the more accurate methods of measuring
- dietary intake,<sup>(18)</sup> meaning that both the predictor and outcome variables are likely to be valid.
- 253 However, even estimated food diaries have their limitations, particularly in terms of participant
- burden and under-reporting of energy intake.<sup>(19; 20)</sup> Doubly labelled water has been used to estimate
- total energy expenditure in a subsample of NDNS participants and compare this to reported energy
- 256 intake from food diaries. This reveals that reported energy intake is 12-34% lower than estimated
- total energy expenditure, depending on the age of participants.<sup>(11)</sup> This mismatch may be due to
- 258 intentional or unintentional misreporting; participants changing their food intake in response to
- 259 recording it; or a variety of other reasons. However, misreporting is unlikely to affect all foods and
- 260 nutrients equally. For example, participants may be more likely to misreport confectionary than

vegetable intake. For this reason, misreporting is not adjusted for in NDNS and we have not adjustedfor misreporting here.

Data mining using decision trees is computationally and statistically efficient. For example, inclusion of all 3911 foods consumed by NDNS participants in regression models with achievement of dietary recommendations as outcomes would be computationally, and statistically, demanding and unlikely to produce satisfactory results. Decision trees also produce transparent, and intuitively understandable, outputs (ours are provided at https://osf.io/znv82).<sup>(21)</sup>

- 268 Many of food included in the analysis had very skewed distributions. Indeed, the vast majority of 269 foods in the database (3618) were eaten by less than 150 people. Decision trees seek to maximize 270 information gain at each step, rather than working with the distribution as a whole as in traditional 271 regression analysis. If an item is very discriminatory and helps differentiate between those who do 272 and do not meet a particular guideline then it will be included, even if it is only consumed by a small 273 number of people. Conversely, if an item is eaten by almost everyone but is not discriminatory, then 274 it would be unlikely to be included. There was no overall trend between the proportion of participants who ate a food and the chance that that food was included in a decision tree (data not 275 276 shown).
- 277 We used adaptive sampling to identify decision trees that achieved the best trade-off between 278 accuracy and number of predictor variables included. Thus, instead of systematically calculating the 279 accuracy of all decision trees including all possible number of predictor variables, we focused on 280 identifying the relationship between accuracy and number of predictor variables (logarithmic in 281 most cases), where the optimum trade-off between accuracy and number of predictor variables 282 occurred (i.e. where the logarithmic curve flattened out). This means we cannot be absolutely sure 283 that we have identified the decision trees with the best trade-off between accuracy and number of 284 predictor variables in all cases. However, given the very small additional improvements in accuracy 285 achieved by the most accurate, versus best trade-off, decision trees, we are certainly likely to have 286 identified the near-best trade-off decision trees.
- We used estimated dietary records as our 'gold standard' tool for determining whether or not
  individuals achieved recommendations. Further work will be required to compare the accuracy of
  our decision trees to other methods of estimating who achieves dietary recommendations, such as
  food frequency questionnaires.

## 291 Interpretation and implications of findings and areas for future work

Our findings indicate that information on only a small number of foods is required to determinewhether individuals achieve five important dietary recommendations. If such binary outcomes are

- the key outcome of interest, then more detailed dietary assessment methods, may inappropriatelyuse scarce research resources and be unnecessarily burdensome to participants.
- 296 Whilst our results suggest that information on only a limited number of foods needs to be captured 297 when assessing whether guidelines are met, substantial further research will be needed before these 298 findings could be applied in the form of a new dietary assessment instrument. Firstly, it would be 299 helpful to replicate our analyses in a different, but comparable, sample. We have not done is as we 300 are not aware of a comparable UK population-representative sample in whom diet diaries have been 301 collected. Our decision trees used information on exact intake of 113 foods over 3-4 days. Assessing 302 exact intake of a small number of foods may be no less burdensome for participants than assessing 303 estimated intake of all foods using a food diary. Future work could compare the accuracy of decision 304 trees based on exact intake of 113 foods, approximate intake of these foods (e.g. using the ordinal categories often used in food frequency questionnaires), and exact and approximate intake of foods 305 306 at the food group, rather than individual food, level. Acceptability to research participants and 307 resource implications of collecting the data required in all cases should also be compared.
- 308 Our analysis focused on which foods can be used to predict whether or not individuals achieve 309 dietary recommendations. But it is not necessarily the case that it is the foods included in the 310 decision tress which cause people to achieve the recommendations or not. Only a maximum of 32% 311 of total intake of relevant nutrients or foods were accounted for by predictor variables in decision 312 trees with the best trade-off between accuracy and number of predictor variables. Thus, decision trees did not particularly include foods that account for the majority of intake of nutrients and foods 313 314 of interest – as might be expected in a food frequency questionnaire. The complex relationships 315 between individual foods included in our decision trees and the dietary recommendations they are 316 associated with may offer further useful insights and could be studied further.

#### 317 CONCLUSION

- 318 We used data mining techniques to explore the number of foods that consumption information was
- required on to accurately predict achievement, or not, of five key dietary recommendations.
- 320 Information on consumption of 11-32 foods (plus age and sex) was sufficient to identify with 73-83%
- 321 accuracy whether individuals achieved individual dietary recommendations. In total, information on
- 322 113 foods was required to predict achievement of all five recommendations studied. This method
- 323 could be used to develop a new dietary assessment questionnaire.

## 324 **REFERENCES**

- 1. Willett WC, Sampson L, Stampfer MJ et al. (1985) Reproducibility and validity of a
- 326 semiquantitative food frequency questionnaire. *Am J Epidemiol* **122**, 51-65.
- 2. Crutzen R, Giabbanelli P (2013) Using Classifiers to Identify Binge Drinkers Based on Drinking
   Motives. *Subst Use Misuse*.
- 3. Dierker L, Rose J, Tan X *et al.* (2010) Uncovering multiple pathways to substance use: a
  comparison of methods for identifying population subgroups. *J Prim Prev* **31**, 333-348.
- 4. McKenzie DP, McFarlane AC, Creamer M *et al.* (2006) Hazardous or harmful alcohol use in Royal
  Australian Navy veterans of the 1991 Gulf War: identification of high risk subgroups. *Addict Behav*31, 1683-1694.
- 5. Hillemacher T, Frieling H, Wilhelm J *et al.* (2012) Indicators for elevated risk factors for alcoholwithdrawal seizures: an analysis using a random forest algorithm. *J Neural Transm* **119**, 1449-1453.
- 6. Lazarou C, Karaolis M, Matalas A-L *et al.* (2012) Dietary patterns analysis using data mining
  method. An application to data from the CYKIDS study. *Comput Methods Programs Biomed* **108**, 706-
- 338 714.
- 339 7. Kastorini C-M, Papadakis G, Milionis HJ et al. (2013) Comparative analysis of a-priori and a-
- posteriori dietary patterns using state-of-the-art classification algorithms: A case/case-control study.
   Artif Intell Med 59, 175-183.
- 8. Thangamani D, Sudha P (2014) Identification Of Malnutrition With Use Of Supervised Datamining
   Techniques -Decision Trees And Artificial Neural Networks. *International Journal Of Engineering And Computer Science* 3, 8236-8241.
- 9. Einsele F, Sadeghi L, Ingold R *et al.* (2015) A Study about Discovery of Critical Food Consumption
  Patterns Linked with Lifestyle Diseases using Data Mining Methods. *Proceedings of the International Conference on Health Informatics*, 239-245.
- 10. Bates B, Lennox A, Swan G (editors) (2010) National Diet and Nutrition Survey: Headline results
- *from Year 1 of the Rolling Programme (2008/2009).* London: Foods Standards Agency and
  Department of Health.
- 11. Bates B, Lennox A, Prentice A *et al.* (editors) (2014) *National Diet and Nutrition Survey Results from Years 1, 2, 3 and 4 (combined) of the Rolling Programme (2008/2009 2011/2012).* London:
- 353 Public Health England.
- 12. Scientific Advisory Committee on Nutrition (2003) *Salt and Health*. London: The Stationary Office.
- 355 13. World Health Organisation (2003) Diet, nutrition and the prevention of chronic diseases: report
   356 of a joint WHO/FAO expert consultation. WHO Technical Report Series **916**.
- 14. Chawla N, Bowyer K, Hall L *et al.* (2002) SMOTE: Synthetic Minority Over-sampling Technique.
   *Journal of Artificial Intelligence Research* 16, 321-357.
- 359 15. Bouckaert RR, Frank E, Hall MA *et al.* (2010) WEKA Experiences with a Java Open-Source
  360 Project. *Journal of Machine Learning Research* 11, 2533-2541.

- 361 16. Kuncheva L (2004) *Fundamentals of pattern recognition Combining pattern classifiers: Methods* 362 *and algorithms.* Hoboken, New Jersey: John Wiley & Sons.
- 17. Imamura F, Micha R, Khatibzadeh S *et al.* Dietary quality among men and women in 187
  countries in 1990 and 2010: a systematic assessment. *The Lancet Global Health* **3**, e132-e142.
- 18. Bingham S, Gill C, Welch A *et al.* (1994) Comparison of dietary assessment methods in nutritional
  epidemiology: weighed records v. 24 h recalls, food-frequency questionnaires and estimated-diet
  records. *Br J Nutr* **72**, 619-643.
- 19. Poslusna K, Ruprich J, de Vries JH *et al.* (2009) Misreporting of energy and micronutrient intake
  estimated by food records and 24 hour recalls, control and adjustment methods in practice. *Br J Nutr* **101 Suppl 2**, S73-85.
- 20. Burrows TL, Martin RJ, Collins CE (2010) A systematic review of the validity of dietary assessment
  methods in children when compared with the method of doubly labeled water. *J Am Diet Assoc* 110,
  1501-1510.
- 21. Crutzen R, Giabbanelli PJ, Jander A *et al.* (2015) Identifying binge drinkers based on parenting
- dimensions and alcohol-specific parenting practices: building classifiers on adolescent-parent paired
   data. *BMC Public Health* 15, 747.
- 377

## FIGURE TITLES AND LEGENDS

Figure 1. Schematic illustration of a decision tree (left, Figure 1a.) and how this is formed through repeated 'cuts' of the data (right, Figure 1b)

Figure 1a. Schematic illustration of a decision tree

Figure 1b. Schematic illustration of how a decision tree is formed through repeated 'cuts' of the data

Figure 2. Overall accuracy (with 95% confidence margins) of decision trees against number of predictor variables included

	Adults aged 19y or older		Children aged <19y	
Variable	Analytical sample (n=2083)	UK population	Analytical sample (n=2073)	UK population
Female, n(%)	1182 (56.8)	25,198,773 (51.5)	1007 (48.6)	6,955,262 (48.8)
Age (adults)				
19-29y, n(%)	296 (14.2)	9,447,071 (19.3)		
30-39y, n(%)	390 (18.7)	8,319,926 (17.0)		
40-49y, n(%)	425 (20.4)	9,268,735 (18.9)		
50-59y, n(%)	363 (17.4)	7,708,532 (15.8)		
60-64y, n(%)	181 (8.7)	3,807,975 (7.8)		
65y+, n(%)	428 (20.6)	10,377,127 (21.2)		
Age (children)				
0-4y, n(%)			499 (24.1)	3,913,953 (27.5)
5-9y, n(%)			583 (26.4)	3,516,615 (24.7)
10-14y, n(%)			547 (26.4)	3,669,326 (25.7)
15-18y, n(%)			444 (21.4)	3,152,919 (22.1)

# Table 1. Comparison of analytical sample to UK population

## Table 2. Prevalence of achieving and not achieving dietary recommendations and accuracy of decision trees to predict this

	Fruit& vegetables	Free sugars	Sodium	Fat	Saturated fat
N (%) achieving recommendation without over-sampling	656 (22.1%)	1472 (35.4%)	2524 (60.7%)	1045 (25.1%)	795 (19.1%)
SMOTE over-sampling %*	252% (YES)	85% (YES)	54% (NO)	197% (YES)	322% (YES)
N achieving recommendation after over-sampling	<u>2309</u> *	<u>2679</u>	2524	<u>3103</u>	<u>3354</u>
N not achieving recommendation after over-sampling	2311*	2684	<u>2513</u>	3111	3361
Decision tree with the best trade-off between accuracy and number of predictor variables					
Overall accuracy	83.1%	76.5%	75.9%	72.4%	79.7%
Sensitivity	82.5%	76.1%	81.9%	66.3%	75.8%
Specificity	83.8%	76.9%	69.8%	78.4%	83.6%
Npredictor variables	11	28	28	33	28
% of all relevant food/nutrient (g) accounted for by predictor variables	21.0%**	31.2%	13.4%	13.0%	27.4%
Most accurate decision tree					
Overall accuracy	83.6%	77.0%	76.1%	72.9%	81.7%
Sensitivity	83.9%	75.7%	80.7%	69.3%	81.4%
Specificity	83.3%	78.3%	71.5%	76.4%	81.9%
N predictor variables	50	64	49	123	156
% of all relevant food/nutrient accounted for by predictor variables	30.8%**	38.6%	25.4%	29.5%	42.7%

\*After over-sampling using the SMOTE method (see Appendix); the prevalence affected by over-sampling is underlined

\*\*Percent of all fruit and vegetables (g) recorded, not just those contributing to 5-a-day portions (specifically, fruit juice can only contribute a maximum of one 5-a-day portion)

	Dietary recommendation outcome			2	Food name		
Fat	Free sugars	Fruit & veg	Sodium	Saturated fat			
Yes			Yes	Yes	Age		
Yes					Alcoholic soft drinks spirit based		
		Yes			Almonds kernel only: ground almonds		
	Yes				Apple juice unsweetened cartons pasteurised		
	Yes				Apple juice unsweetened UHT		
	105	Vec			Apples eating raw flesh & skin only		
Voc		105			Avocado poar floch only		
163			Voc		Reconsider the short only		
			Vec		Decon respers back grilled real and lat		
			res		Bacon rashers back not smoked grined extra trim		
			Yes		Baked beans in tomato sauce with pork sausages		
Yes		Yes			Bananas raw flesh only		
				Yes	Beefburger and onion grilled		
Yes					Black pudding fried		
	Yes				Blackcurrant juice drink ready to drink not low calorie		
	Yes				Boiled sweets barley sugar butterscotch glacier mints hard candy		
			Yes		Bread white crusty		
			Yes	Yes	Bread white toasted		
Yes					Bread, 50% white and 50% wholemeal flours		
			Yes		Bread, white sliced, not fortified		
			Yes		Brown sauce bottled		
			Yes		Brussels sprouts-fresh hoiled		
Voc			105		Butter beans dried boiled		
Voc				Voc	Butter salted		
res				Vec	Dutter unselted		
				res	Butter unsalted		
	Yes				Carbonated beverages no juice not low calorie canned		
Yes	Yes			Yes	Carbonated beverages no juice not low calorie not canned		
		Yes			Celery, fresh raw		
Yes					Chapati brown no fat		
Yes				Yes	Cheese cheddar any other or for recipes		
				Yes	Cheese cheddar English		
			Yes		Cheese soft full fat. Philadelphia type		
Yes					Chicken fried in olive oil		
				Yes	Children's fromagefrais fruit with added vitamin D		
				Yes	Chocolate brownie no nuts purchased		
				Yes	Chocolate covered caramels Cadburys caramel		
	Yes				Chocolate Swiss roll with buttercream purchased		
	Yes				Cola cherry cola canned not low calorie		
	Yes				Cola not canned not low calorie not caffeine free		
Vec	105				Coleslaw nurchased not low calorie		
Voc					Cookies and hiscuits with chocolate		
163				Voc	Connette tune ice cream checelate er put based		
	Vec			163	Conherry fruit juice drink e.g. Ocean Spray		
	res			Vaa	Cranberry fruit juice unink e.g. Ocean Spray		
				res	Cream double		
	Yes				Cream egg		
				Yes	Croissants plain not filled		
	Yes				Drinking chocolate instant dry weight		
			Yes		Fat spread (62-72% fat) not polyunsaturated		
	Yes				Fruit gums winegums		
	Yes				Fruit juice drink carbonated not low calorie not canned		
	Yes				Fruit juice drink with 5% fruit juice ready to drink		
				Yes	Fully coated chocolate biscuits with biscuit filling		
Yes					Garlic bread. Lower fat		
			Yes		Ham unspecified not smoked not canned		
			Yes		Hamburger Big Mac McDonalds		
	Yes				High juice ready to drink not blackcurrant or low calorie		
	Yes				Ice Iollies		
	Yes				laffa Cakes		
				Yes	Kit Kat		

# Table 3. Predictor variables (individual foods, age and sex) included in decision trees for predicting achievement of five dietary recommendations

Yes					Lager not canned e.g. Heineken
Yes					Lager not canned e.g. Skol
Yes					Lamb scrag and neck stewed lean only
	Yes				Lemonade not low calorie not canned
				Yes	Light spreadable butter (60% fat)
	Yes				Lucozade sport isotonic drink not carbonated
Yes				Yes	Mayonnaise (retail)
			Yes	Yes	Milk chocolate bar
	Yes				Milk shake thick style takeaway
Yes					Milk skimmed after boiling
				Yes	Milk whole pasteurised winter
				Yes	Milk whole summer pasteurised
Yes					Mushrooms fried in olive oil
			Yes		Naan bread plain
		Yes			Oatcakes
	Yes				Olive oil
		Yes			Onions boiled
	Yes				Orange juice unsweetened UHT
Yes					Oven ready chips
			Yes		Papadums/poppadoms fried in vegetable ghee
Yes					Pasta noodles boiled
			Yes		Pasta noodles egg boiled
Yes					Pasta spaghetti boiled white
			Yes		Peanut butter crunchy not wholenut
		Yes			Pears eating raw flesh & skin only no core
Yes					Pepperami
				Yes	Petit Filousfromagefrais
			Yes		Potato cakes (scones) purchased
Yes					Potatoes new boiled skins eaten
			Yes		Potatoes old baked flesh & skin
				Yes	Potatoes old mashed & butter
			Yes		Prawns boiled flesh only
			Yes		Reduced fat spread (41-62%) not polyunsaturated
			Yes		Ribena original blackcurrant drink concentrate
	Yes				Robinsons fruit shoot
			Yes		Rolls white crusty
Yes		Yes		Yes	Sausage roll flaky pastry purchased
Yes					Sausages, pork, grilled
			Yes		Sausages, premium pork, grilled
Yes					Scrambled eggs with skimmed milk and no fat
Yes					Semi-sweet biscuit
			Yes		Sex
	Yes				Soya alternative to milk sweetened plain
		Yes			Spinach fresh raw
				Yes	Spreadable butter (75-80% fat)
	Yes				Sugar white
				Yes	SupernoodlesBatchelorsas served
				Yes	Swiss roll individual chocolate coated purchased
		Yes			Tomatoes raw
			Yes		Turkey slices unsmoked prepack or deli
	Yes				Water for concentrated soft drinks not diet
Yes					White chocolate buttons mice
Yes					Whole milk after boiling
Yes					Wine white dry not canned
	Yes	X		Yes	Yogurt twinpot with cereal/crumble
		Yes			Yogurt, Greek style, cows, natural, whole milk
			Yes		Yorkshire pudding frozen