

Identification of predictive factors of the degree of adherence to the Mediterranean diet through machine-learning techniques

Alba Arceo ^{Equal first author, 1}, **Carlos Fernandez-Lozano** ^{Corresp., Equal first author, 2, 3}, **Salvador Pita-Fernández** ¹, **Sonia Pértega-Díaz** ¹, **Alejandro Pazos** ^{2, 3}

¹ Clinical Epidemiology and Biostatistics Research Group,, Instituto de Investigación Biomédica de A Coruña (INIBIC), Complejo Hospitalario Universitario de A Coruña (CHUAC), SERGAS, Universidade da Coruña, A Coruña, Spain

² Department of Computer Science and Information Technologies, Faculty of Computer Science, CITIC-Research Center of Information and Communication Technologies, Universidade da Coruña, A Coruña, Spain

³ Grupo de Redes de Neuronas Artificiales y Sistemas Adaptativos. Imagen Médica y Diagnóstico Radiológico (RNAS-IMEDIR). Instituto de Investigación Biomédica de A Coruña (INIBIC). Complejo Hospitalario Universitario de A Coruña (CHUAC), SERGAS, Universidade da Coruña, A Coruña, Spain

Corresponding Author: Carlos Fernandez-Lozano
Email address: carlos.fernandez@udc.es

Food consumption patterns have undergone changes that in recent years have resulted in serious health problems. Studies based on the evaluation of the nutritional status have determined that the adoption of a food pattern-based primarily on a Mediterranean diet has a preventive role, as well as the ability to mitigate the negative effects of certain pathologies. A group of more than 500 adults, aged over 40 years from our cohort in Northwestern Spain was surveyed. Under our experimental design, ten experiments were run with four different machine-learning algorithms and the predictive factors most relevant to the adherence of a Mediterranean diet were identified. A feature selection approach was explored and under a null hypothesis test, it was concluded that only 16 measures were of relevance, suggesting the strength of this observational study. Our findings indicate that the following factors have the highest predictive value in terms of the degree of adherence to the Mediterranean diet: basal metabolic rate, mini nutritional assessment questionnaire total score, weight, height, bone density, waist-hip ratio, smoking habits, age, EDI-OD, circumference of the arm, activity metabolism, subscapular skinfold, subscapular circumference in cm, circumference of the waist, circumference of the calf and brachial area.

1 Identification of predictive factors of the 2 degree of adherence to the Mediterranean 3 diet through machine-learning techniques

4 Alba Arceo-Vilas¹, Carlos Fernandez-Lozano^{2,3}, Salvador Pita¹, Sonia
5 Pérttega-Díaz¹, and Alejandro Pazos^{2,3}

6 ¹Clinical Epidemiology and Biostatistics Research Group, Instituto de Investigación
7 Biomédica de A Coruña (INIBIC), Complejo Hospitalario Universitario de A Coruña
8 (CHUAC), SERGAS, Universidade da Coruña, A Coruña, Spain.

9 ²Department of Computer Science and Information Technologies, Faculty of Computer
10 Science, CITIC-Research Center of Information and Communication Technologies,
11 Universidade da Coruña, A Coruña, Spain

12 ³Grupo de Redes de Neuronas Artificiales y Sistemas Adaptativos. Imagen Médica y
13 Diagnóstico Radiológico (RNASA-IMEDIR). Instituto de Investigación Biomédica de A
14 Coruña (INIBIC). Complejo Hospitalario Universitario de A Coruña (CHUAC), SERGAS,
15 Universidade da Coruña, A Coruña, Spain

16 Corresponding author:

17 Carlos Fernandez-Lozano

18 Email address: carlos.fernandez@udc.es

19 ABSTRACT

20 Food consumption patterns have undergone changes that in recent years have resulted in serious health
21 problems. Studies based on the evaluation of the nutritional status have determined that the adoption
22 of a food pattern-based primarily on a Mediterranean diet has a preventive role, as well as the ability to
23 mitigate the negative effects of certain pathologies. A group of more than 500 adults, aged over 40 years
24 from our cohort in Northwestern Spain was surveyed. Under our experimental design, ten experiments
25 were run with four different machine-learning algorithms and the predictive factors most relevant to the
26 adherence of a Mediterranean diet were identified. A feature selection approach was explored and under
27 a null hypothesis test, it was concluded that only 16 measures were of relevance, suggesting the strength
28 of this observational study. Our findings indicate that the following factors have the highest predictive
29 value in terms of the degree of adherence to the Mediterranean diet: basal metabolic rate, mini nutritional
30 assessment questionnaire total score, weight, height, bone density, waist-hip ratio, smoking habits, age,
31 EDI-OD, circumference of the arm, activity metabolism, subscapular skinfold, subscapular circumference
32 in cm, circumference of the waist, circumference of the calf and brachial area.

33 INTRODUCTION

34 The economic development, urbanisation and industrialisation worldwide have changed individuals'
35 eating habits and lifestyles, such as smoking, excessive consumption of alcohol, sedentary lifestyle and
36 stress, leading to a nutritional transition which its principle cost in the health sector, is the appearance of
37 non-transmissible chronic diseases.

38 A consequence of the alteration of dietary patterns is what has been called "epidemic obesity", defined
39 by the World Health Organization (WHO) as the first non-viral epidemic of the 21st century, with 500
40 million obese people worldwide (Finucane et al., 2011; Krzysztoszek et al., 2015) affecting more than
41 50% of the adult population in Spain (López-Sobaler et al., 2016; Anta et al., 2013; Rodríguez Rodríguez,
42 E; Lopez Plaza, B; Lopez Sobaler, M; Ortega, 2011).

43 The assessment of nutritional status of a population is one of the best indicators of the health status of
44 the said population, being a methodology that must include three important aspects: a global assessment,
45 a study of the dimension and a study of body composition (Ravasco et al., 2010). With adequate

interpretation of the findings, appropriate therapeutic measures should be taken to correct deviations from normality.

In the context of nutrition and public health, the Mediterranean diet (MD) has been forged over the centuries, being characterised by cereal, olive oil, low saturated fats and meat, moderate consumption of dairy and a regular and moderate intake of wine, being a lifestyle in accordance with geographic, climatological, orographic, cultural and environmental conditions within the countries and regions that surround the Mediterranean Sea (Pérez C, 2011).

There is an increasing interest in the study of the preventive role of MD and also as a treatment for various pathologies associated with chronic inflammation, such as metabolic syndrome, diabetes mellitus, cardiovascular disease (CVD), neurodegenerative diseases, breast cancer and psycho-organic deterioration, leading to greater longevity and better quality of life (Dussaillant et al., 2016; Chrysoshoou et al., 2004; Trichopoulou, 2004; Serra-Majem et al., 2006; Estruch et al., 2013; Sofi et al., 2014; Della Camera et al., 2017). Moreover, the importance of MD has also been identified as a potential element contributing to the prevention of breast cancer (Shapira, 2017) or in patients carrying the BRCA mutation (Bruno et al., 2017). In 2010, UNESCO declared this diet an Intangible Cultural Heritage of Humanity (UNESCO, 2010).

Numerous studies have been published over the past decades, showing the relationship between MD intake and CVD (Martínez-González et al., 2015; Widmer et al., 2015), and meta-analyses that relate it to general health status (Sofi et al., 2014). In the Greek cohort EPIC (European Prospective Investigation into Cancer and Nutrition Study) a 2-point increase in adherence to this diet was associated with a 33% reduction in CVD mortality (Sofi et al., 2014). Additionally, the analysis of a sub-cohort of 2,700 individuals over 60 years old, with a history of myocardial infarction showed that a greater adherence to MD had an 18% drop in overall mortality (Lack et al., 2003). Other studies have confirmed these associations, including the follow-up of a Spanish cohort of 13,600 adults with coronary heart disease. After 5 years, it was observed that 2 points of increase in adherence to MD were associated with a 26% decrease in coronary risk (et al Trichopoulou A, Bamia C, Norat T, Overvad K, Schmidt EB, Tjonneland A, 2007).

Eating disorders are linked to a distorted perception of one's own body image, as well as to body dissatisfaction. The importance of a study on body dissatisfaction is due to the fact that recent investigations have confirmed that alterations in body image have a causal participation in an eating disorder, rather than being secondary to it (Míguez Bernárdez et al., 2011). Body image is considered a qualitative approximation to the nutritional status of the individual (Sámano et al., 2015) and can be determining for their nutritional management (Martínez-González et al., 2011).

One of the main fields of application of Machine-Learning (ML) techniques since its origins is in the field of Biomedicine, finding previously published studies in related areas such as biomedical image (Fernandez-Lozano et al., 2016b), characterisation of different types of carcinomas (Kim et al., 2017), measurement of activity in genetic networks (Hu et al., 2016), deformable models for image comparison (Rodriguez et al., 2014), gene selection, and classification of microarray data (Díaz-Uriarte and Alvarez de Andrés, 2006), to name a few.

Moreover, due to the great versatility of ML techniques, they have been used in a wide variety of application areas, to discover hidden patterns in the datasets: identification and authentication of tequilas (Pérez-Caballero et al., 2017), wearable sensor data fusion (Kanjó et al., 2018), predicting the outcomes of organic reactions (Skoraczynski et al., 2017), animal behaviour detection (Pons et al., 2017) or to measure the visual complexity of images (Machado et al., 2015). In particular, ML techniques have proven to be able to uncover unimaginable relationships in very diverse fields of application, such as image or voice recognition, sentiment analysis or language translation (Li et al., 2015; Perez-de Viñaspre and Oronoz, 2015).

The main objective of this work is the development of ML models for the prediction of the degree of adherence to the Mediterranean diet. To this end, information on different anthropometric and socio-demographic variables, nutritional status and self-perception of body image is used in order to identify which of the variables have a greater influence and are key in the adherence to a healthy diet such as MD, allowing our patients to improve their quality of life and to reduce the negative effects of well-known and related diseases.

Taking into account all of the above, the experimental methodology proposed in the development of this study is based on the collection and generation of data to be analysed with our cohort in Galicia

(Spain), as well as on the use of ML techniques. The purpose is to extract and explain the underlying information in the data and determine which of these variables are the most important to classify people as having either a good or poor adherence to the MD. As mentioned before, there are several health benefits related to this particular food diet, especially for: chronic inflammation, metabolic syndrome, diabetes mellitus, CVD, neurodegenerative diseases, cancer and psycho-organic deterioration, moreover leading to greater longevity and better quality of life. Thus, this study is relevant for understanding how to measure the degree of adherence, in order to ensure the aforementioned benefits.

The structure of the article is as follows: in the Materials and methods section, the subjects are presented, the variables are measured for each of them. Next, the machine learning and feature selection techniques are described, along with the experimental design followed to ensure that the results are reproducible and representative of the studied problem. In the next section, the results are presented and discussed, and the final section of the article includes the conclusions of the work.

MATERIALS AND METHODS

The present study was structured as follows. Initially, a population from our cohort was selected to carry out the study; the population was grouped into two categories: with high and low degree of adherence to the MD. Once the set of the population on which the study will be carried out has been identified, the information is collected from each of the users of the health system. The type of study carried out will be described below, as well as the sample size will be justified and all measurements collected will be explained in detail. Once the dataset is generated, it will be analyzed with four different ML techniques and a feature selection phase will be applied for dimensionality reduction.

Population and data description

This is an observational prevalence study, conducted in Northwestern Spain (municipality of Cambre, A Coruña, Spain), which included a randomly selected population aged 40 years and over. The sampling population consisted of individuals residing in Cambre, identified through the National Health System card census. In Spain, the National Health System has universal coverage, and almost all Spanish citizens are beneficiaries of public healthcare services.

The sample size was calculated taking into account the total population of the municipality ($n = 12,446$). After stratification by age and gender, ($n = 503$) persons were selected to participate in the study. Sample size was estimated using the single proportion formula, with 95% confidence Interval. A sample size of ($n = 503$) subjects was estimated based on an adherence to mediterranean diet rate of 50%. Precision was set at 4.3% and percentage of losses at 10%. Population data is shown in Table 1.

Table 1. Population data of the municipality of Cambre (A Coruña) for the year 2012 and sample data according to age and sex.

Age groups	Population			Sample		
	Total	Men	Women	Total	Men	Women
40 – 44	2465	1202(26.9%)	1263(27.8%)	33	19(12.9%)	14(13.2%)
45 – 49	2231	1110(24.8%)	1121(24.7%)	85	52(35.4%)	33(31.1%)
50 – 54	1763	857(19.2%)	906(19.9%)	54	32(21.8%)	22(20.8%)
55 – 59	1383	702(15.7%)	681(14.9%)	33	18(12.3%)	15(14.1%)
60 – 64	1170	598(13.4%)	572(12.6%)	48	26(17.7%)	22(20.8%)
Total (40 – 64)	9012	4469	4543	253	147	106
65 – 69	1027	497(33%)	530(27.5%)	94	57(38%)	37(37%)
70 – 74	688	337(22.4%)	351(18.2%)	77	46(30.7%)	31(31%)
75 – 79	777	326(21.6%)	451(23.4%)	46	28(18.7%)	18(18%)
80 – 84	511	198(13.1%)	313(16.2%)	24	12(8%)	12(12%)
85 – more	431	148(9.8%)	283(14.7%)	9	7(4.7%)	2(2%)
Total (65 and more)	3434	1506	1928	250	150	100

A personal interview was arranged with each individual. After obtaining their permission and written consent, a trained nurse proceeded to the measurement of anthropometric variables and to the collection of the necessary data to cover the questionnaires. The patients who could not go to the health center

due to personal or displacement reasons and those who suffer from a cognitive impairment, making it impossible for them to perform the study, were excluded. The study received written approval from the Regional Ethics Committee for Clinical Research (code 2012/390 CEIC Galicia).

The information described below was collected from each selected subject: socio-demographic variables: age, gender, level of education, marital status and relationships of coexistence; prevalence of arterial hypertension and smoking: the systolic and diastolic blood pressure of each patient was recorded at the beginning and at the end of the visit, obtaining the prevalence of arterial hypertension; the smoking habit was recorded according to self-reported information. Anthropometric variables: the anthropometric parameters allow us to know the state of the protein and caloric reserves, besides providing guidance to the health professional about the consequences of the imbalances in these reserves.

All measurements were made during the same session, to avoid variations in the environmental or biological conditions. For the measurement of weight and size, the person was barefoot and with light clothing; an MB-201T plus Asimed scale-rod was used with an accuracy of 100 grams (weight) and 1 mm (size). BMI was obtained by means of the $BMI_{ratio} = \frac{weight(kg)}{height(m)^2}$, and grouped according to the WHO classification of $BMI < \frac{18.5kg}{m^2}$: low weight; 18.5 to $24.99 \frac{kg}{m^2}$: normal weight; 25 to $29.99 \frac{kg}{m^2}$: overweight, and $\geq 30 \frac{kg}{m^2}$: obesity.

The waist and hip circumference was measured with an inelastic tape measure with the patient standing upright, the abdomen relaxed, the upper limbs hanging at the sides, and with the feet and knees joined together. The waist circumference was measured by taking the mid-point between the lower costal margins and the iliac crests, as it is considered a risk factor for cardiovascular disease when it is wider than 80 cm in women and wider than 94 cm in men, and a very high risk if it exceeds 88 cms and 102 cms, respectively (Alberti et al., 2009).

The hip circumference was measured as the maximum circumference around the buttocks. Based on these two values, the waist-hip ratio was calculated using the cut-off points proposed by the WHO, where normal levels of 0.8 are found in women and 1 in men, higher values indicating abdominal visceral obesity, which is associated with increased cardiovascular risk (Jover E, 1997).

The calf circumference was measured in the widest section of the ankle-knee distance (cuff area) showing a good correlation with fat-free mass and muscle strength (cols Rolland Y, Lauwers-Cances V, 2003; Barbosa Murillo et al., 2007; Bonnefoy M, Jauffret M, Kostka T, 2002). The measurement was carried out with an inextensible tape measure in cm.

Subscapular skin fold, this fold measures truncal obesity. The measurement is made one centimeter below the lower angle of the scapula, following the natural furrow of the skin. The scapula protrudes when the arm is carefully placed behind the back and the lower angle can be located this way. The measurement of the fold will be diagonally over an angle of 45° to ensure the correct thickness measurement. The plicometer forceps should be applied 1 cm in the inferolateral position to the thumb and finger that lifts the fold.

To assess the amount of subcutaneous adipose tissue, the skin folds were measured in millimetres in the tricipital, bicipital, subscapular and suprailiac areas. A digital caliper Trim meter was used, including a double layer of skin and underlying adipose tissue, always avoiding the muscle. The tricipital skinfold was measured longitudinally, at the back of the non-dominant upper limb, at the midpoint between acromion and olecranon, with the limb relaxed, parallel to the axis of the arm; the bicipital fold was measured at the same point as the tricipital, but on the under arm.

The circumference of the arm was measured with an inextensible anthropometric tape measure in cm. The measurement was taken at the midpoint of the non-dominant arm, in the same place where the tricipital skinfold was measured and without compression with the anthropometric tape.

Once the data of the different measurements were obtained, the mid-arm muscle circumference was found, with which the skeletal muscle mass of the patients (protein compartment) was known and expressed in cm. The arm muscle area indicated that the muscle compartment was based on the brachial circumference and tricipital skinfold measurements. The fat area of the arm indicated that the patient's fatty compartment used the total brachial area and the muscular area of the arm. The Adipose Muscular Index, which evaluates the nutritional status from the adipose and muscular areas of the arm, was also calculated, being essentially applied in the assessment of obesity.

For the determination of body fat percentage by electrical bio-impedance, a Beurer and BG55 model bio-impedance meter was used, with a maximum capacity of 150 kg and a precision of 0.1% for body fat, body water and muscle percentage, and 100 g for body weight, according to the information provided by

the manufacturer.

These methods are based on physical principles, such as the different ability of conduction or resistance that the tissues show to the passage of an electric current, with greater conductivity of the lean tissues than the fatty ones (Norman et al., 2007). Thus, by means of bio-impedance, the following values were obtained: weight, fat mass, liquid mass, muscle mass, bone density, basal metabolic rate (BMR) and activity metabolism. Data on socio-demographic variables, such as age, gender (male/female), cohabitation (with whom the live), prevalence of current smokers, ex smokers (patient stopped smoking more than 12 months before entering the study) and non-smokers were estimated. Additionally, blood pressure was recorded.

Adherence to the Mediterranean diet

Consumption of a characteristic food pattern of MD is associated with numerous health benefits. These benefits are attributed to bioactive compounds that exert synergistic effects and decrease the risk for development of chronic diseases.

In order to assess the quality of dietary habits (adherence to a Mediterranean dietary pattern), the Mediterranean diet adherence test was used (Insituto de Salud Carlos III, 2009). It is a questionnaire consisting of fourteen quick questions that allow us to understand whether participants' usual diet can be considered as following the parameters of the MD. Each question answered affirmatively adds a point. It is considered that a person correctly follows the Mediterranean diet when their score is equal to or greater than nine points.

The assessment of nutritional status was determined using the Mini Nutritional Assessment (MNA) questionnaire (Insituto de Salud Carlos III, 2009). It is a validated method which, through eighteen short questions, evaluates anthropometric measures, dietary habits, lifestyle, pharmacological treatments and mobility, and performs a subjective evaluation of health and nutritional status. The total value of the MNA scale is thirty points, a score < 17 being considered malnutrition, there is a risk of malnutrition between 17-23,5, and well nourished subjects obtain scores of twenty-four points and higher.

A measure of subjective weight is included by asking: "I consider that my weight is: A) higher than normal, B) normal, C) lower than normal", following the model proposed in (Espina et al., 2001). Based on the answer, the population is classified into three groups: "fairly subjective weight" those who believe to be at an ideal weight, "more subjective kilograms" for those who believe that they are overweight and "less subjective kilograms" for those who think they weigh less than they should.

Two of the eleven sub-scales of Garner's Eating Disorder Inventory (EDI-2) (Garner, 1998) were used to study body image: Body Dissatisfaction (EDI-IC) and Obsession for Thinness (EDI-OD), as they evaluate aspects directly related to perceptual alterations. The body dissatisfaction sub-scale (EDI-CI) measures the dissatisfaction of the subject with the general shape of their body or with those parts of the body that most concern those with eating disorders (stomach, hips, thighs, buttocks). The thinness obsession sub-scale (EDI-O) measures concern about weight, diets and fear of getting fat.

This questionnaire was validated in Spain by (Corral, S., González, M., Pereña, J. y Seisdedos, 1998). The fourteen items of these two sub-scales were mixed in the questionnaire to avoid the subjects guessing the construct being evaluated. All items were answered and corrected according to the form proposed in the questionnaire manual. The Mediterranean Diet Adherence test was used to determine the degree of adherence to the Mediterranean Diet, being a short specific questionnaire of fourteen items validated for the Spanish population and used by the Mediterranean Diet Prevention Group (PREDIMED) (Martínez-González et al., 2015).

Machine Learning and statistical analysis

The authors tested different ML techniques for solving this problem, using cross-validation techniques to avoid over-training, while ensuring that the generalised capability of the model is the best possible, as well as different runs of the experiments to check the behaviour of the techniques. Thus, all experiments were repeated ten times to check the stability of the results and the observed deviation between the experiments was small, as shown in the results section. In particular, a tenfold cross validation was used to divide the dataset in such a way that nine random partitions were used to train and one to validate the results, each time taking a different subset for validation.

In order to compare the performance of the ML techniques, the Area Under the Receiver Operating Characteristic Curve (AUROC) was used. This is a combined measurement which, besides being

independent of the threshold used, includes both Type I and type II errors, ensuring that it is not conditioned by differences in the total number of cases of each class (Fawcett, 2006).

An experimental design was employed (Fernandez-Lozano et al., 2016a), allowing us to divide the data using a cross-validation technique which ensured that the performance results obtained, as mentioned above, were not skewed. That is, they were adjusted to the data, and researchers are able to identify which of the hyper-parameters are most suitable to find the best model with each ML techniques, according to its particular hyper-parameter configuration. To this end, the programming environment R (R Core Team, 2016) and the package mlr (Bischl et al., 2016) were used, which also allowed us to perform the considered experimental design. In addition, another of the objectives pursued by this study was to find as few variables as possible that would yield a performance value as high as possible, preferably at least equal to that obtained using all available variables. This is basically a feature selection approach where the main aims are the following: avoid overfitting and improve model performance, to provide faster and more cost-effective models, and moreover to gain a deeper insight into the underlying processes that generated the data as mentioned in (Saey et al., 2007). There are three approaches in ML to perform this process and the use of a filter approximation was chosen, for its velocity and independence of the classifier (Saey et al., 2007). In general, performing this feature selection process helps to reduce inherently the present noise in such datasets.

The final step of our experimental analysis was a null hypothesis test for choosing the best model in order to ensure whether the performance of a particular ML technique is statistically better than the others or not. In our case, as there were more than two repeated measures, an ANOVA or a Friedman test should be considered. In particular, three different conditions should be checked: normality, independence and homoscedasticity. If our results fulfil the three conditions, a parametric test is applicable, and the ANOVA one should be considered, otherwise the non-parametric version, the Friedman test. Finally, a post hoc procedure had to be used in order to correct the p-values for multiple testing.

Machine Learning techniques for classification problems

A large number of experiments were carried out in an attempt to identify the best ML model able to solve the problem and to ensure that the results are reproducible, real and obtained under equal conditions. In addition, the search space was explored for the best possible parameters for each technique in the same way, so that all techniques could have the same possibilities of exploration across the same subsets of data and avoid the over fit that could occur. In particular, the following well-known state-of-the-art techniques were implemented: Random Forest (Breiman, 2001), Support Vector Machines (Cortes and Vapnik, 1995; Vapnik, 1995), Elastic Net (Tibshirani, 1994; Zou and Hastie, 2005) and weighted k-Nearest Neighbours (Hechenbichler and Schliep, 2004).

Random Forest (RF) (Breiman, 2001) is a state-of-the-art ML technique that was used in multiple domains with good results. One of its main strengths is that the results obtained are very easy to understand, it is based on very simple concepts and in general, although it is applied with little experience in the parametrisation of hyper-parameters, good results are obtained. This technique combines multiple decision trees, each of them tuned over a subset of bootstrapped data. In this way, RF combines each of the individual predictions of the decision trees into a global prediction that, in general, is more successful than any of the simple ones. Of all the possible variables in the dataset, a number were randomly selected (with replacement) and a number of trees were constructed based on the set of examples used for the training phase and obtained from the previously selected subset. When there are classification problems, it is recommended to use a square root number of the total number of variables existing in the dataset. To explore the solution space in the best way possible, in our experiments we used a parameter domain that was adjusted by a grid search and that, for a number of trees (1000), we explored randomly selected values of variables (2-6). In addition, values that varied (1-4) according to the size of the terminal nodes of the tree were explored.

Support Vector Machines (SVM) (Cortes and Vapnik, 1995; Vapnik, 1995) is also one of the ML techniques that have been commonly used in different domains in recent years and have obtained good results. In fact, along with RF, it is one of the algorithms considered state-of-the-art, easy to understand, and the results obtained are verifiable. In problems that occur during a study, the main objective of SVM is to find the hyperplane that best separates the examples between high and low degree of adherence to the Mediterranean diet and at the same time to maximise the distance of separation between both examples and the hyperplane. That is, it attempts to find the separation hyperplane that generalises in the best

possible way (Burges, 1998). To achieve this goal, SVM introduces a particular mathematical concept known as kernel: it is a mathematical function that allows the conversion of the input space into a higher dimension, which is used to transform a non-separable linear problem into one that is separable. There are different kernel functions, which in general could be interpreted as a measure of similarity between two objects (60), and one of the most used is Gaussian Radial Basis (RBF), because basically any surface can be obtained with this function (61). In this case, the domain of the parameters used to search for the best model consists of a grid search of two different parameters. The first one (parameter C) is directly related to the model and is used as a balance between the classification errors and the simplicity of the decision surface, while the second (gamma parameter) is the free parameter of the Gaussian function and in particular, SVM is very sensitive to changes in this parameter. For both parameters, and according to the usual practice, values were evaluated in potencies of two between -12 and 12. To better understand this technique, the following reading materials are recommended (Burges, 1998; Vert et al., 2004; Cristianini and Shawe-Taylor, 2000).

Elastic Net (ENET) (Tibshirani, 1994; Zou and Hastie, 2005) is based on lasso (penalised least squares method) and was specifically developed to solve some of the limitations encountered for this technique (56). On the one hand, a grid search was performed on two different parameters, the alpha penalty parameter was searched (it has values in the range of 0 to 1) and in particular the following 0, 0.15, 0.25, 0.35, 0.5, 0.65, 0.75, 0.85 and 1. On the other hand, the best value of the lambda parameter was used, as recommended by the authors of the technique, from values less than or equal to one to negative powers of ten, in particular the following values were used: 0.0001, 0.001, 0.01, 0.1, 1.

Finally, a simple k-Nearest Neighbour (KNN) (Hechenbichler and Schliep, 2004) assigned, through a decision rule, an unclassified example belonging to a class by frequency of occurrence to its k-most similar neighbouring examples. Then, in accordance to the distance of Minkowski for each of the examples and following the maximum accumulated kernel densities the weighted k-Nearest neighbour are identified (Hechenbichler and Schliep, 2006; Samworth, 2012). In particular, neighbouring values of less than or equal to nine were used. Therefore, this particular and improved implementation of a k-Nearest Neighbour used kernel functions to measure the degree of similarity of the examples, as previously mentioned in the case of the SVMs.

RESULTS

The dataset has a total of 38 variables employed to characterise the differences underlying in the data between high and low adherence to the MD. The data has been standardised using the z-score formula to have a mean equal to zero and a standard deviation equal to 1. Four different ML techniques were used to verify the results obtained, in an attempt to identify the technique that provides the best-performing results. Initially, the analysis of the complete set of study variables is carried out. It can be seen in Figure 1.a and b how the techniques present a fairly stable behavior in the prediction. Even a simple a priori technique such as KNN obtains the best results of the entire experimental phase, indicating that almost all variables contain relevant information. In any case, in order to understand whether there is noise or contradictory or correlated information that may be hindering the learning process of the algorithms, a phase of dimensionality reduction will then be carried out.

Additionally, a process of feature selection was carried out to reduce the number of variables as much as possible, so that the results could remain similar without statistical differences, if not better, for those obtained using all variables. Our approach is a filter feature selection using a T-test to quantify the correlation between each feature and the class (high or low adherence to the MD) before the training process. Three subsets of 4, 16 and 32 features were evaluated of the original ordering according to the highest p-value from the T-test. The average AUROC results of the execution of the ten 10-fold cross-validation experiments are shown in Figure 1.

As the number of features increases, there is a clear growing tendency in performance and obtained results in AUROC with 16 and 32 features are very close to those obtained with the full dataset. In any case, a study should be conducted on whether the differences are statistically significant between the subsets of 16, 32 variables and the full dataset to ensure that the subset with fewer features is statistically the best option. Finally, as shown in Figure 1.a, SVM is the best model in three out of the four datasets, and manages to reach values closest to 0.94 in AUROC.

However, as previously mentioned, a single mean measure is not enough and it is necessary to analyse the behaviour of the models during the whole experimental phase and to verify how stable they are, as

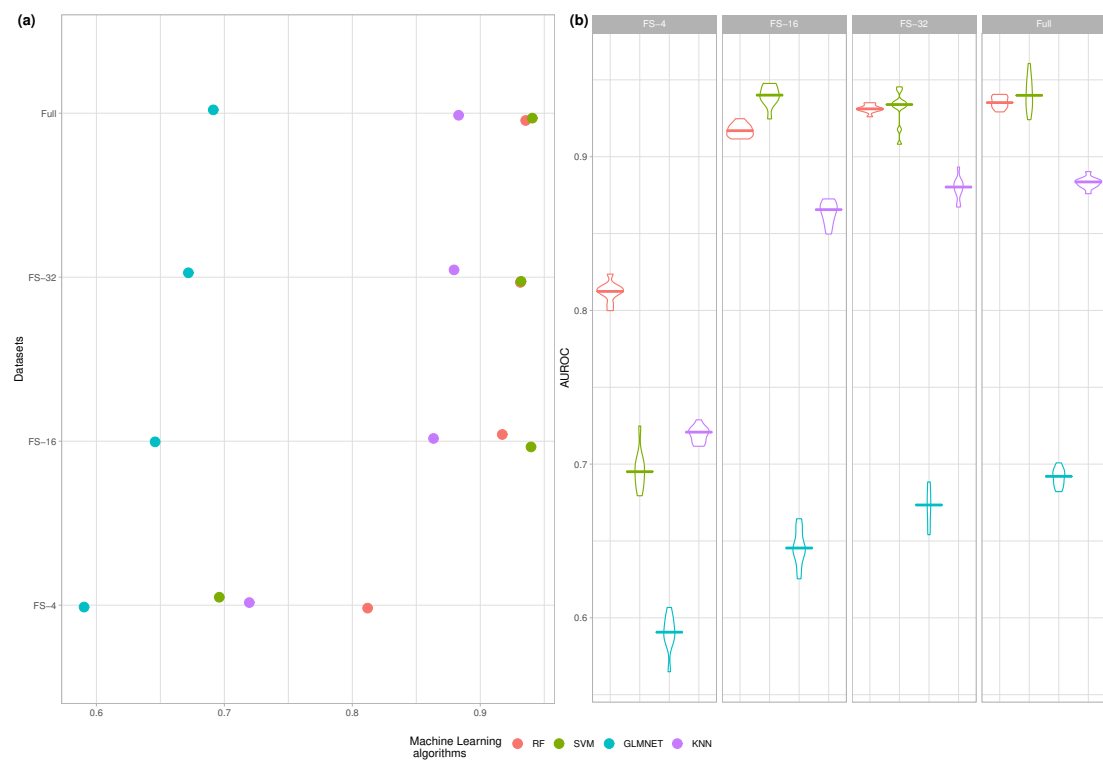


Figure 1. Summary of the performance (AUROC) of the four ML techniques (RF, SVM, GLMNET and KNN) for each one of the subsets of features. a) Average of the experiments for each size analyzed and b) boxplot of the results in order to check the behaviour of the techniques through the learning process.

shown in Figure 1.b.

Figure 1.b shows that if the number of variables is very small (4), the models are skewed and there is a higher variability in the performance because there is not enough information in the data to find a good classification model. It is also important to note that the results obtained with 16 and 32 features showed that this variability was significantly reduced until reaching average and standard deviation values very similar to those obtained using all the variables.

As observed in the two previous figures, the best results in AUROC were obtained using SVM. The same results in accuracy are shown in Figure 2.

To check whether the difference between the three winning models (SVM with 16, 32 and all variables) is significant or not, a null hypothesis test was applied. Following the experimental methodology proposed in (Fernandez-Lozano et al., 2016a) for the normality condition we used the Shapiro-Wilk test (Shapiro and Wilk, 1965), with a confidence level $\alpha = 0.05$. with the null hypothesis that our results follow a normal distribution. The null hypothesis was not rejected with values $W = 0.96179$ and $p - value = 0.3438$, therefore our results did follow a normal distribution.

Next, a Bartlett test (Bartlett, 1937) was performed, with a confidence level $\alpha = 0.05$ and with the null hypothesis that our data were heterocedastic. The test result indicates that the null hypothesis should not be rejected with a value of Bartlett's K-squared 2.3128 with 2 degrees of freedom and $p - value = 0.3146$. The result of both tests indicates that a parametric ANOVA test should be conducted, with a confidence level $\alpha = 0.05$ assuming the null hypothesis that our results are statistically equal. The results of the ANOVA test indicates that we fail to reject the null hypothesis and the three ML models are statistically equal with an adjusted $p - value = 0.1124$. Consequently, a 16-feature model should be considered (BMR, MNA total score, weight, height, bone density, waist-hip ratio, smoker, age, EDI-OD, circumference of the arm, activity metabolism, subscapular skin fold, subscapular circumference in cm, circumference of the waist, circumference of the calf, brachial area) as the best-performing one, and half of the initial features that are not relevant for the SVM were removed.

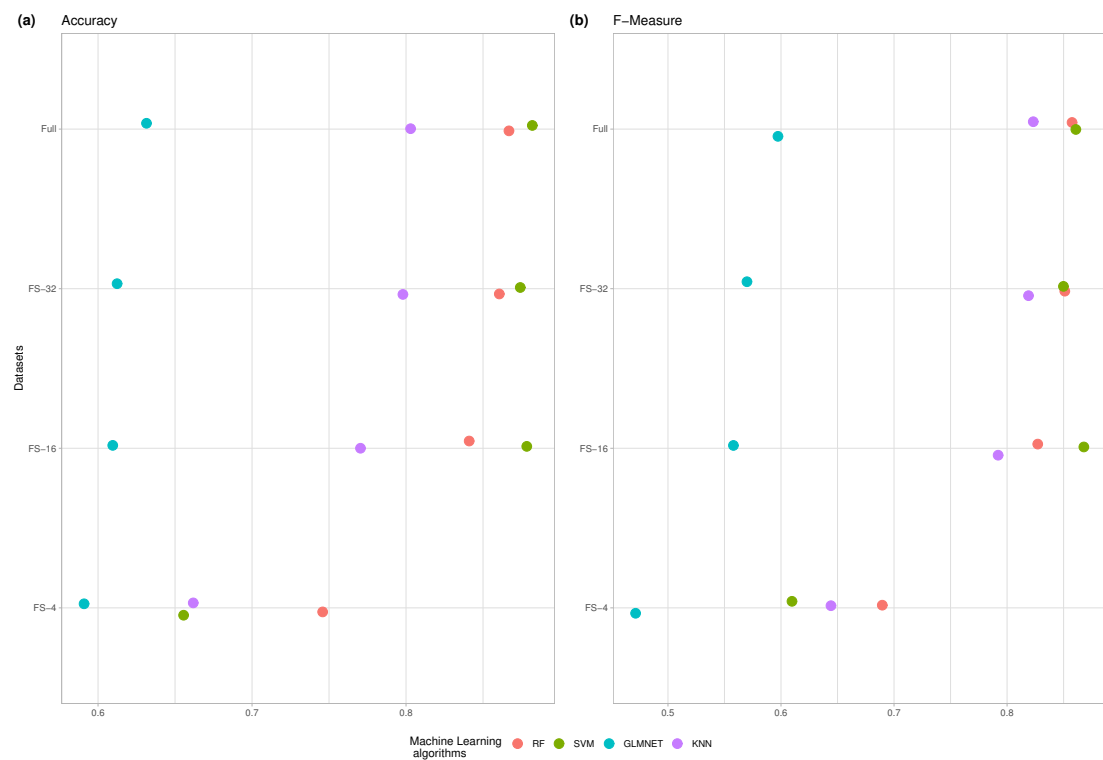


Figure 2. Summary of the average performance of the experiments: a) (Accuracy) and b) (F-measure) of the four ML techniques (RF, SVM, GLMNET and KNN) for each one of the subsets of features.

DISCUSSION

To check whether our results are relevant and are in accordance with what has been previously published, the state-of-the-art articles published on the topic were reviewed, in an attempt to identify the degree of adherence to the variables most related to a MD. The search results led to previous studies that also found the variables identified by SVM as the most important, in particular BMR (Cutillas et al., 2013; Careau, 2017; Srivastava et al., 2017; Bonfanti et al., 2014), MNA total score (Farias et al., 2016; Zaragoza Martí et al., 2015; Abreu-Reyes et al., 2017), weight and height (de la Montaña Miguélez et al., 2012; Bach-faig et al., 2008; Ortega Anta and López Sobaler, 2014; Travé and Castroviejo, 2011), bone density (Romero Pérez and Rivas Velasco, 2014; Savanelli et al., 2017; Melaku et al., 2017; Štefan et al., 2017) or waist-hip ratio (Downer et al., 2016; Estruch et al., 2016; Bertoli et al., 2015) and if the patient is a smoker (Zaragoza Martí et al., 2015; Marventano et al., 2017; Grao-Cruces et al., 2015). Therefore, the results were contrasted, ensuring the ability of ML techniques to identify underlying patterns in the data. According to the feature selection process, the remaining predictors are not relevant for all the ML techniques.

CONCLUSIONS

The first model based on ML that was proposed for the prediction of the degree of adherence to the Mediterranean diet depended on information related to different anthropometric variables, socio-demographic variables, nutritional status and self-perception of body image.

Initially, experiments with four different ML methods were performed and feature selection techniques were applied to reduce the dimensionality of the problem. SVM is the best-performing model according to the experimental design after a null hypothesis test, and our study found that using a feature selection approach, the number of features could be drastically reduced to 16 (less than half of the initial number) achieving an equivalent performance value in AUROC. The best model obtained was an SVM with an RBF kernel as a decision function. The importance of each one of the predictors cannot be studied because a nonlinear SVM is like a black box and the internal mapping function is unknown. Furthermore, the

weight vector cannot be explicitly computed.

Finally, our results are in accordance with the findings of previous publications and have primarily served to establish new factors related to the degree of adherence to the Mediterranean diet.

REFERENCES

- Abreu-Reyes, J. A., Álvarez-Luis, D., Arteaga-Hernández, V., Sánchez-Mendez, M., and Abreu-González, R. (2017). Mediterranean diet adherence by patients with primary open angle glaucoma. *Archivos de la Sociedad Española de Oftalmología (English Edition)*, 92(8):353–358.
- Alberti, K. G. M. M., Eckel, R. H., Grundy, S. M., Zimmet, P. Z., Cleeman, J. I., Donato, K. A., Fruchart, J. C., James, W. P. T., Loria, C. M., and Smith, S. C. (2009). Harmonizing the metabolic syndrome: A joint interim statement of the international diabetes federation task force on epidemiology and prevention; National heart, lung, and blood institute; American heart association; World heart federation; International . *Circulation*, 120(16):1640–1645.
- Anta, R. M. O., Lopez-Solaber, A. M., Perez-Farinos, N., Ortega Anta, R. M., López-Solaber, A. M., and Pérez-Farinos, N. (2013). Associated factors of obesity in Spanish representative samples. *Nutr Hosp*, 28(5):56–62.
- Bach-faig, A., Buckland, G., Faig, A. B., and Majem, L. S. (2008). prevención de la obesidad . Una revisión de la Revisión Eficacia de la dieta mediterránea en la prevención de la obesidad . Una revisión de la bibliografía. (January).
- Barbosa Murillo, J. A. P., Rodríguez M, N. G., Hernández H De Valera, Y. M., Hernández H, R. A., and Herrera M, H. A. (2007). Masa muscular, fuerza muscular y otros componentes de funcionalidad en adultos mayores institucionalizados de la Gran Caracas-Venezuela. *Nutricion Hospitalaria*, 22(5):578–583.
- Bartlett, M. S. (1937). Properties of Sufficiency and Statistical Tests. *Proceedings of the Royal Society of London. Series A, Mathematical and Physical Sciences*, 160(901):268–282.
- Bertoli, S., Leone, A., Vignati, L., Bedogni, G., Martínez-González, M. Á., Bes-Rastrollo, M., Spadafranca, A., Vanzulli, A., and Battezzati, A. (2015). Adherence to the Mediterranean diet is inversely associated with visceral abdominal tissue in Caucasian subjects. *Clinical Nutrition*, 34(6):1266–1272.
- Bischl, B., Lang, M., Kotthoff, L., Schiffner, J., Richter, J., Jones, Z., and Casalicchio, G. (2016). mlr: Machine Learning in R.
- Bonfanti, N., Fernandez, J. M., Gomez-Delgado, F., and Perez-Jimenez, F. (2014). Effect of two hypocaloric diets and their combination with physical exercise on basal metabolic rate and body composition. *Nutricion hospitalaria*, 29(3):635–643.
- Bonnefoy M, Jauffret M, Kostka T, J. J. (2002). Usefulness of calf circumference measurement in assessing the nutritional state of hospitalized elderly people. *Gerontology*, 48(3):162–9.
- Breiman, L. (2001). Random Forests. *Machine Learning*, 45(1):5–32.
- Bruno, E., Manoukian, S., Venturelli, E., Oliverio, A., Rovera, F., Iula, G., Morelli, D., Peissel, B., Azzolini, J., Roveda, E., and Pasanisi, P. (2017). Adherence to Mediterranean Diet and Metabolic Syndrome in BRCA Mutation Carriers. *Integrative Cancer Therapies*, page 153473541772101.
- Burges, C. J. C. (1998). A Tutorial on Support Vector Machines for Pattern Recognition. *Data Mining and Knowledge Discovery*, 2(2):121–167.
- Careau, V. (2017). Energy Intake, Basal Metabolic Rate, and Within-Individual Trade-Offs in Men and Women Training for a Half Marathon: A Reanalysis. *Physiological and Biochemical Zoology*, 90(3):392–398.
- Chrysohoou, C., Panagiotakos, D. B., Pitsavos, C., Das, U. N., and Stefanadis, C. (2004). Adherence to the Mediterranean diet attenuates inflammation and coagulation process in healthy adults: The ATTICA study. *Journal of the American College of Cardiology*, 44(1):152–158.
- cols Rolland Y, Lauwers-Cances V, C. M. y. (2003). Sarcopenia, calf circumference, and physical function of elderly women: a cross-sectional study. *J Am Geriatr Soc*, 51(8):1120–4.
- Corral, S., González, M., Pereña, J. y Seisdedos, N. (1998). *Adaptación española del Inventario de trastornos de la conducta alimentaria*. MADRID.
- Cortes, C. and Vapnik, V. (1995). Support-Vector Networks. In *Machine Learning*, volume 20, pages 273–297.

- 454 Cristianini, N. and Shawe-Taylor, J. (2000). *An Introduction to Support Vector Machines: And Other*
455 *Kernel-based Learning Methods*. Cambridge University Press, New York, NY, USA.
- 456 Cutillas, A. B., Herrero, E., de San Eustaquio, A., Zamora, S., and Pérez-Llamas, F. (2013). Prevalencia
457 de peso insuficiente, sobrepeso y obesidad, ingesta de energía y perfil calórico de la dieta de estudiantes
458 universitarios de la comunidad autónoma de la región de Murcia (España). *Nutrición Hospitalaria*,
459 28(3):683–689.
- 460 de la Montaña Miguélez, J., Cobas, N., Rodríguez, M., Míguez Bernárdez, M., and Castro Sobrino, L.
461 (2012). Adherencia a la dieta mediterránea y su relación con el índice de masa corporal en universitarios
462 de Galicia. *Nutrición clínica y dietética hospitalaria*, ISSN 0211-6057, Vol. 32, Nº. 3, 2012, págs.
463 72–80, 32(3):72–80.
- 464 Della Camera, P. A., Morselli, S., Cito, G., Tasso, G., Cocci, A., Laruccia, N., Travaglini, F., Del Fabbro,
465 D., Mottola, A. R., Gacci, M., Serni, S., Carini, M., and Natali, A. (2017). Sexual health, adherence to
466 Mediterranean diet, body weight, physical activity and mental state: factors correlated to each other.
467 *Rivista Urologia*, page 0.
- 468 Díaz-Uriarte, R. and Alvarez de Andrés, S. (2006). Gene selection and classification of microarray data
469 using random forest. *BMC bioinformatics*, 7:3.
- 470 Downer, M. K., Gea, A., Stampfer, M., Sánchez-Tainta, A., Corella, D., Salas-Salvadó, J., Ros, E., Estruch,
471 R., Fitó, M., Gómez-Gracia, E., Arós, F., Fiol, M., De-la Corte, F. J. G., Serra-Majem, L., Pinto, X.,
472 Basora, J., Sorlí, J. V., Vinyoles, E., Zazpe, I., and Martínez-González, M.-Á. (2016). Predictors of
473 short- and long-term adherence with a Mediterranean-type diet intervention: the PREDIMED
474 randomized trial. *The international journal of behavioral nutrition and physical activity*, 13(1):67.
- 475 Dussailant, C., Echeverría, G., inés UrquíaGa, niColás VelasCo, and atilio RiGotti (2016). Evidencia
476 actual sobre los beneficios de la dieta mediterránea en salud. *artículo de revisión rev Med chileRev*
477 *Med Chile*, 144(144):1044–1052.
- 478 Espina, A., Ortego, M. A., Ochoa de Alda, I. n., Yenes, F., and Aleman, A. (2001). La imagen corporal en
479 los trastornos alimentarios. *Body shape in eating disorders*, 13(4):533–538.
- 480 Estruch, R., Martínez-González, M. A., Corella, D., Salas-Salvadó, J., Fitó, M., Chiva-Blanch, G.,
481 Fiol, M., Gómez-Gracia, E., Arós, F., Lapetra, J., Serra-Majem, L., Pintó, X., Buil-Cosiales, P.,
482 Sorlí, J. V., Muñoz, M. A., Basora-Gallisá, J., Lamuela-Raventós, R. M., Serra-Mir, M., and Ros, E.
483 (2016). Effect of a high-fat Mediterranean diet on bodyweight and waist circumference: a prespecified
484 secondary outcomes analysis of the PREDIMED randomised controlled trial. *The Lancet Diabetes &*
485 *Endocrinology*, 4(8):666–676.
- 486 Estruch, R., Ros, E., Salas-Salvad, J., Covas, M.-I., Dpharm, Corella, D., Ars, F., Gmez-Gracia, E.,
487 Ruiz-Gutiérrez, V., Fiol, M., Lapetra, J., Lamuela-Raventós, R. M., Serra-Majem, L., Pint, X.,
488 Basora, J., Muñoz, M. A., Sorl, J. V., Martnez, J. A., and Martinez-Gonzalez, M. A. (2013). Primary
489 Prevention of Cardiovascular Disease with a Mediterranean Diet. *New England Journal of Medicine*,
490 page 130225030008006.
- 491 et al Trichopoulou A, Bamia C, Norat T, Overvad K, Schmidt EB, Tjonneland A (2007). Modified
492 Mediterranean diet and survival after myocardial infarction: the EPIC-Elderly study. *Eur J Epidemiol*,
493 22(12):871–81.
- 494 Farias, G., Thieme, R. D., Teixeira, L. M., Heyde, M. E., Bettini, S., and Radominski, R. (2016). Nutrición
495 Hospitalaria Trabajo Original. *Nutr. Hosp.*, 33(5):1108–1115.
- 496 Fawcett, T. (2006). An introduction to ROC analysis. *Pattern Recognition Letters*, 27(8):861–874.
- 497 Fernandez-Lozano, C., Gestal, M., Munteanu, C. R., Dorado, J., and Pazos, A. (2016a). A methodology
498 for the design of experiments in computational intelligence with multiple regression models. *PeerJ*,
499 4:e2721.
- 500 Fernandez-Lozano, C., Seoane, J., Gestal, M., Gaunt, T., Dorado, J., Pazos, A., and Campbell, C. (2016b).
501 Texture analysis in gel electrophoresis images using an integrative kernel-based approach. *Scientific*
502 *Reports*, 6.
- 503 Finucane, M., Stevens, G., Cowan, M., Danaei, G., Lin, J. K., Paciorek, C. J., Singh, G. M., Gutierrez,
504 H. R., Lu, Y., Bahalim, A. N., Farzadfar, F., Riley, L. M., and Ezzati, M. (2011). National, regional,
505 and global trends in body-mass index since 1980: systematic analysis of health examination surveys
506 and epidemiological studies with 960 country-. *The Lancet*, 377(9765):557–567.
- 507 Garner, D. (1998). *EDI-2: Inventario de Trastornos de la Conducta Alimentaria. Manual*. MADRID.
- 508 Grao-Cruces, A., Nuviala, A., Fernandez-Martinez, A., and Martinez-Lopez, E.-J. (2015). Relationship

- 509 of physical activity and sedentarism with tobacco and alcohol consumption, and Mediterranean diet in
- 510 Spanish teenagers. *Nutricion hospitalaria*, 31(4):1693–1700.
- 511 Hechenbichler, K. and Schliep, K. (2004). Weighted k-Nearest-Neighbor Techniques and Ordinal
- 512 Classification.
- 513 Hechenbichler, K. and Schliep, K. (2006). Weighted k-nearest-neighbor techniques and ordinal classifica-
- 514 tion. In *Discussion Paper 399, SFB 386*.
- 515 Hu, Z., Mao, J.-H., Curtis, C., Huang, G., Gu, S., Heiser, L., Lenburg, M. E., Korkola, J. E., Bayani, N.,
- 516 Samarajiwa, S., Seoane, J. A., Dane, M. A., Esch, A., Feiler, H. S., Wang, N. J., Hardwicke, M. A.,
- 517 Laquerre, S., Jackson, J., W. Wood, K., Weber, B., Spellman, P. T., Aparicio, S., Wooster, R., Caldas,
- 518 C., and Gray, J. W. (2016). Genome co-amplification upregulates a mitotic gene network activity that
- 519 predicts outcome and response to mitotic protein inhibitors in breast cancer. *Breast Cancer Research*,
- 520 18(1):70.
- 521 Insituto de Salud Carlos III (2009). Estudio predimed. Prevención primaria de la enfermedad Cardiovas-
- 522 cular con la Dieta Mediterránea. pages 1–41.
- 523 Jover E (1997). Índice cintura/cadera. Obesidad y riesgo cardiovascular. *An Med Intern*, 14(1-2).
- 524 Kanjo, E., Younis, E., and Sherkat, N. (2018). Towards unravelling the relationship between on-body,
- 525 environmental and emotion data using sensor information fusion approach. *Information Fusion*, 40.
- 526 Kim, J., Bowlby, R., Mungall, A. J., Robertson, A. G., Odze, R. D., Cherniack, A. D., Shih, J., Pedamallu,
- 527 C. S., Cibulskis, C., Dunford, A., Meier, S. R., Kim, J., Raphael, B. J., Wu, H.-T., Wong, A. M., Willis,
- 528 J. E., Bass, A. J., Derks, S., Garman, K., McCall, S. J., Wiznerowicz, M., Pantazi, A., Parfenov, M.,
- 529 Thorsson, V., Shmulevich, I., Dhankani, V., Miller, M., Sakai, R., Wang, K., Schultz, N., Shen, R.,
- 530 Arora, A., Weinhold, N., Sánchez-Vega, F., Kelsen, D. P., Zhang, J., Felau, I., Demchok, J., Rabkin,
- 531 C. S., Camargo, M. C., Zenklusen, J. C., Bowen, J., Leraas, K., Lichtenberg, T. M., Curtis, C., Seoane,
- 532 J. A., Ojesina, A. I., Beer, D. G., Gulley, M. L., Pennathur, A., Luketich, J. D., Zhou, Z., Weisenberger,
- 533 D. J., Akbani, R., Lee, J.-S., Liu, W., Mills, G. B., Zhang, W., Reid, B. J., Hinoue, T., Laird, P. W.,
- 534 Shen, H., Piazuelo, M. B., Schneider, B. G., McLellan, M., Taylor-Weiner, A., Cibulskis, C., Lawrence,
- 535 M., Cibulskis, K., Stewart, C., Getz, G., Lander, E., Gabriel, S. B., Ding, L., McLellan, M. D., Miller,
- 536 C. A., Appelbaum, E. L., Cordes, M. G., Fronick, C. C., Fulton, L. A., Mardis, E. R., Wilson, R. K.,
- 537 Schmidt, H. K., Fulton, R. S., Ally, A., Balasundaram, M., Bowlby, R., Carlsen, R., Chuah, E., Dhalla,
- 538 N., Holt, R. A., Jones, S. J. M., Kasaian, K., Brooks, D., Li, H. I., Ma, Y., Marra, M. A., Mayo, M.,
- 539 Moore, R. A., Mungall, A. J., Mungall, K. L., Robertson, A. G., Schein, J. E., Sipahimalani, P., Tam,
- 540 A., Thiessen, N., Wong, T., Cherniack, A. D., Shih, J., Pedamallu, C. S., Beroukhi, R., Bullman, S.,
- 541 Cibulskis, C., Murray, B. A., Saksena, G., Schumacher, S. E., Gabriel, S., Meyerson, M., Hadjipanayis,
- 542 A., Kucherlapati, R., Pantazi, A., Parfenov, M., Ren, X., Park, P. J., Lee, S., Kucherlapati, M., Yang,
- 543 L., Baylin, S. B., Hoadley, K. A., Weisenberger, D. J., Bootwalla, M. S., Lai, P. H., Van Den Berg,
- 544 D. J., Berrios, M., Holbrook, A., Akbani, R., Hwang, J.-E., Jang, H.-J., Liu, W., Weinstein, J. N., Lee,
- 545 J.-S., Lu, Y., Sohn, B. H., Mills, G., Seth, S., Protopopov, A., Bristow, C. A., Mahadeshwar, H. S.,
- 546 Tang, J., Song, X., Zhang, J., Laird, P. W., Hinoue, T., Shen, H., Cho, J., Defrictas, T., Frazer, S.,
- 547 Gehlenborg, N., Heiman, D. I., Lawrence, M. S., Lin, P., Meier, S. R., Noble, M. S., Voet, D., Zhang,
- 548 H., Kim, J., Polak, P., Saksena, G., Chin, L., Getz, G., Wong, A. M., Raphael, B. J., Wu, H.-T., Lee,
- 549 S., Park, P. J., Yang, L., Thorsson, V., Bernard, B., Iype, L., Miller, M., Reynolds, S. M., Shmulevich,
- 550 I., Dhankani, V., Abeshouse, A., Arora, A., Armenia, J., Kundra, R., Ladanyi, M., Lehmann, K.-V.,
- 551 Gao, J., Sander, C., Schultz, N., Sánchez-Vega, F., Shen, R., Weinhold, N., Chakravarty, D., Zhang,
- 552 H., Radenbaugh, A., Hegde, A., Akbani, R., Liu, W., Weinstein, J. N., Chin, L., Bristow, C. A., Lu,
- 553 Y., Penny, R., Crain, D., Gardner, J., Curley, E., Mallery, D., Morris, S., Paulauskis, J., Shelton, T.,
- 554 Shelton, C., Bowen, J., Frick, J., Gastier-Foster, J. M., Gerken, M., Leraas, K. M., Lichtenberg, T. M.,
- 555 Ramirez, N. C., Wise, L., Zmuda, E., Tarvin, K., Saller, C., Park, Y. S., Button, M., Carvalho, A. L.,
- 556 Reis, R. M., Matsushita, M. M., Lucchesi, F., de Oliveira, A. T., Le, X., Paklina, O., Setdikova, G., Lee,
- 557 J.-H., Bennett, J., Iacocca, M., Huelsenbeck-Dill, L., Potapova, O., Voronina, O., Liu, O., Fulidou, V.,
- 558 Cates, C., Sharp, A., Behera, M., Force, S., Khuri, F., Owonikoko, T., Pickens, A., Ramalingam, S.,
- 559 Sica, G., Dinjens, W., van Nistelrooij, A., Wijnhoven, B., Sandusky, G., Stepa, S., Crain, D., Paulauskis,
- 560 J., Penny, R., Gardner, J., Mallery, D., Morris, S., Shelton, T., Shelton, C., Curley, E., Juhl, H., Zornig,
- 561 C., Kwon, S. Y., Kelsen, D., Kim, H. K., Bartlett, J., Parfitt, J., Chetty, R., Darling, G., Knox, J.,
- 562 Wong, R., El-Zimaity, H., Liu, G., Boussioutas, A., Park, D. Y., Kemp, R., Carlotti, C. G., da Cunha
- 563 Tirapelli, D. P., Saggiaro, F. P., Sankarankutty, A. K., Noushmehr, H., dos Santos, J. S., Trevisan, F. A.,

- Eschbacher, J., Dubina, M., Mozgovoy, E., Carey, F., Chalmers, S., Forgie, I., Godwin, A., Reilly, C.,
Madan, R., Naima, Z., Ferrer-Torres, D., Vinco, M., Rathmell, W. K., Dhir, R., Luketich, J., Pennathur,
A., Ajani, J. A., McCall, S. J., Janjigian, Y., Kelsen, D., Ladanyi, M., Tang, L., Camargo, M. C.,
Ajani, J. A., Cheong, J.-H., Chudamani, S., Liu, J., Lolla, L., Naresh, R., Pihl, T., Sun, Q., Wan, Y.,
Wu, Y., Demchok, J. A., Felau, I., Ferguson, M. L., Shaw, K. R. M., Sheth, M., Tarnuzzer, R., Wang,
Z., Yang, L., Zenklusen, J. C., Hutter, C. M., Sofia, H. J., and Zhang, J. (2017). Integrated genomic
characterization of oesophageal carcinoma. *Nature*.
- Krzyszczoszek, J., Wierzejska, E., and Zielińska, A. (2015). Obesity. An analysis of epidemiological and
prognostic research. *Archives of Medical Science*, 11(1):24–33.
- Lack, G., Fox, D., Northstone, K., and Golding, J. (2003). New England Journal. *The New England
journal of medicine*, pages 977–985.
- Li, X., Li, J., and Wu, Y. (2015). A global optimization approach to multi-polarity sentiment analysis.
PLoS ONE, 10(4):1–18.
- López-Sobaler, A. M., Aparicio, A., Aranceta-Bartrina, J., Gil, Á., González-Gross, M., Serra-Majem,
L., Varela-Moreiras, G., and Ortega, R. M. (2016). Overweight and General and Abdominal Obesity
in a Representative Sample of Spanish Adults: Findings from the ANIBES Study. *BioMed research
international*, 2016:8341487.
- Machado, P., Romero, J., Nadal, M., Santos, A., Correia, J., and Carballal, A. (2015). Computerized
measures of visual complexity. *Acta Psychologica*, 160:43–57.
- Martínez-González, M. A., García-López, M., Bes-Rastrollo, M., Toledo, E., Martínez-Lapiscina, E. H.,
Delgado-Rodríguez, M., Vazquez, Z., Benito, S., and Beunza, J. J. (2011). Mediterranean diet and the
incidence of cardiovascular disease: A Spanish cohort. *Nutrition, Metabolism and Cardiovascular
Diseases*, 21(4):237–244.
- Martínez-González, M. A., Salas-Salvadó, J., Estruch, R., Corella, D., Fitó, M., and Ros, E. (2015).
Benefits of the Mediterranean Diet: Insights From the PREDIMED Study. *Progress in Cardiovascular
Diseases*, 58(1):50–60.
- Martínez-González, M. A., Salas-Salvadó, J., Estruch, R., Corella, D., Fitó, M., and Ros, E. (2015).
Benefits of the mediterranean diet: Insights from the predimed study. *Progress in Cardiovascular
Diseases*, 58(1):50 – 60. Preventive Cardiology Update: Controversy, Consensus, and Future Promise.
- Marventano, S., Godos, J., Platania, A., Galvano, F., Mistretta, A., and Grosso, G. (2017). Mediter-
ranean diet adherence in the Mediterranean healthy eating, aging and lifestyle (MEAL) study cohort.
International Journal of Food Sciences and Nutrition, pages 1–8.
- Melaku, Y. A., Gill, T. K., Taylor, A. W., Adams, R., and Shi, Z. (2017). Association between nutrient
patterns and bone mineral density among ageing adults. *Clinical Nutrition ESPEN*.
- Míguez Bernárdez, M., De la Montaña Miguélez, J., González Carnero, J., and González Rodríguez,
M. (2011). Concordancia entre la autopercepción de la imagen corporal y el estado nutricional en
universitarios de Orense. *Nutricion Hospitalaria*, 26(3):472–479.
- Norman, K., Smoliner, C., Valentini, L., Lochs, H., and Pirlich, M. (2007). Is bioelectrical impedance
vector analysis of value in the elderly with malnutrition and impaired functionality? *Nutrition*,
23(7-8):564–569.
- Ortega Anta, R. M. and López Sobaler, A. M. (2014). Primeras Jornadas UCM-ASEN Avances y
controversias en nutrición y salud. *Nutrición hospitalaria*, 30(2):21–28.
- Pérez C, A. J. (2011). ¿Es posible la dieta Mediterránea en el siglo XXI? *La dieta Mediterránea en el
marco de la nutrición comunitaria: luces y sombras*.
- Pérez-Caballero, G., Andrade, J., Olmos, P., Molina, Y., Jiménez, I., Durán, J., Fernandez-Lozano, C.,
and Miguel-Cruz, F. (2017). Authentication of tequilas using pattern recognition and supervised
classification. *TrAC - Trends in Analytical Chemistry*, 94.
- Perez-de Viñaspre, O. and Oronoz, M. (2015). SNOMED CT in a language isolate: an algorithm for a
semiautomatic translation. *BMC Med Inform Decis Mak*, 15 Suppl 2(Suppl 2):S5.
- Pons, P., Jaen, J., and Catala, A. (2017). Assessing machine learning classifiers for the detection of
animals’ behavior using depth-based tracking. *Expert Systems with Applications*, 86:235–246.
- R Core Team (2016). R: A Language and Environment for Statistical Computing.
- Ravasco, P., Anderson, H., and Mardones, F. (2010). Métodos de valoración del estado nutricional. *Nutr
Hosp*, 25(Supl. 3):57–66.
- Rodríguez, A., Fernandez-Lozano, C., Dorado, J., and Rabuñal, J. R. (2014). Two-dimensional gel

- 619 electrophoresis image registration using block-matching techniques and deformation models. *Analytical*
620 *biochemistry*, 454:53–9.
- 621 Rodríguez Rodríguez, E; Lopez Plaza, B; Lopez Sobaler, M; Ortega, R. (2011). Prevalencia de sobrepeso
622 y obesidad en adultos españoles. *Nutricion Hospitalaria*, 26(2):355–363.
- 623 Romero Pérez, A. and Rivas Velasco, A. (2014). Adherence to Mediterranean diet and bone health.
624 *Nutricion hospitalaria*, 29(5):989–96.
- 625 Saeys, Y., Inza, I., and Larrañaga, P. (2007). A review of feature selection techniques in bioinformatics.
626 *Bioinformatics*, 23(19):2507–2517.
- 627 Sámano, R., Rodríguez-ventura, A., Sánchez-jiménez, B., Godínez, E., Noriega, A., Zelonka, R., Garza,
628 M., and Nieto, J. (2015). Satisfacción de la imagen corporal en adolescentes y adultos mexicanos y
629 su relación con la autopercepción corporal y el índice de masa corporal real. *Nutricion Hospitalaria*,
630 31(3):1082–1088.
- 631 Samworth, R. J. (2012). Optimal weighted nearest neighbour classifiers. *The Annals of Statistics*,
632 40(5):2733–2763.
- 633 Savanelli, M. C., Barrea, L., Macchia, P. E., Savastano, S., Falco, A., Renzullo, A., Scarano, E., Net-
634 tore, I. C., Colao, A., and Di Somma, C. (2017). Preliminary results demonstrating the impact of
635 Mediterranean diet on bone health. *Journal of Translational Medicine*, 15(1):81.
- 636 Serra-Majem, L., Roman, B., and Estruch, R. (2006). Scientific evidence of interventions using the
637 Mediterranean diet: a systematic review. *Nutrition reviews*, 64(2 Pt 2):S27–47.
- 638 Shapira, N. (2017). The potential contribution of dietary factors to breast cancer prevention. *European*
639 *Journal of Cancer Prevention*, page 1.
- 640 Shapiro, S. S. and Wilk, M. B. (1965). An analysis of variance test for normality (complete samples).
641 *Biometrika*, 52(3-4):591–611.
- 642 Skoraczynski, G., Dittwald, P., Miasojedow, B., Szymkuć, S., Gajewska, E. P., Grzybowski, B. A., and
643 Gambin, A. (2017). Predicting the outcomes of organic reactions via machine learning: are current
644 descriptors sufficient? *Scientific reports*, 7(1):3582.
- 645 Sofi, F., Macchi, C., Abbate, R., Gensini, G. F., and Casini, A. (2014). Mediterranean diet and health
646 status: an updated meta-analysis and a proposal for a literature-based adherence score. *Public Health*
647 *Nutrition*, 17(12):2769–2782.
- 648 Srivastava, R., Batra, A., Dhawan, D., and Bakhshi, S. (2017). Association of energy intake and
649 expenditure with obesity: A cross-sectional study of 150 pediatric patients following treatment for
650 leukemia. *Pediatric Hematology and Oncology*, 34(1):29–35.
- 651 Štefan, L., Čule, M., Milinović, I., Sporiš, G., and Juranko, D. (2017). The relationship between adherence
652 to the Mediterranean diet and body composition in Croatian university students. *European Journal of*
653 *Integrative Medicine*, 13(Supplement C):41–46.
- 654 Tibshirani, R. (1994). Regression Selection and Shrinkage via the Lasso.
- 655 Travé, T. D. and Castroviejo, A. (2011). Adherencia a la dieta mediterránea en la población universitaria.
656 26(3):602–608.
- 657 Trichopoulou, A. (2004). Traditional Mediterranean diet and longevity in the elderly: a review. *Public*
658 *health nutrition*, 7(7):943–947.
- 659 UNESCO (2010). The Mediterranean diet inscription on the Representative List of the Intangible Cultural
660 Heritage of Humanity.
- 661 Vapnik, V. N. (1995). *The nature of statistical learning theory*. Springer-Verlag New York, Inc.
- 662 Vert, J. P., Tsuda, K., and Schölkopf, B. (2004). A Primer on Kernel Methods. In *Kernel Methods in*
663 *Computational Biology*, pages 35–70. MIT Press.
- 664 Widmer, R. J., Flammer, A. J., Lerman, L. O., and Lerman, A. (2015). The Mediterranean Diet, its
665 Components, and Cardiovascular Disease. *The American Journal of Medicine*, 128(3):229–238.
- 666 Zaragoza Martí, A., Ferrer Cascales, R., Cabañero Martínez, M. J., Hurtado Sánchez, J. A., and Laguna
667 Pérez, A. (2015). Adherencia a la dieta mediterránea y su relación con el estado nutricional en personas
668 mayores. *Nutrición hospitalaria*, 31(4):1667–74.
- 669 Zou, H. and Hastie, T. (2005). Regularization and variable selection via the elastic net. *Journal of the*
670 *Royal Statistical Society. Series B: Statistical Methodology*, 67(2):301–320.