# Reproducibility of data-driven dietary patterns in two groups of adult Spanish women from different studies

Adela Castelló[1,2,3]*, Virginia Lope[1,2,3], Jesús Vioque[2,4], Carmen Santamariña[5], Carmen Pedraz-Pingarrón[6], Soledad Abad[7], Maria Ederra[8], Dolores Salas-Trejo[9], Carmen Vidal[10], Carmen Sánchez-Contador[11], Nuria Aragonés[1,2,3], Beatriz Pérez-Gómez[1,2,3] and Marina Pollán[1,2,3]  on behalf of DDM-Spain research group†

[1]*Cancer Epidemiology Unit, National Center for Epidemiology, Instituto de Salud Carlos III, Avenida Monforte de Lemos, 5, 28029, Madrid, Spain*
[2]*Consortium for Biomedical Research in Epidemiology & Public Health (CIBERESP), Instituto de Salud Carlos III, Avenida Monforte de Lemos, 5, 28029, Madrid, Spain*
[3]*Cancer Epidemiology Research Group, Oncology and Hematology Area, IIS Puerta de Hierro (IDIPHIM), Calle Manuel de Falla, 1, 28222 Majadahonda, Madrid, Spain*
[4]*Universidad Miguel Hernandez, Crta. Nacional 332, s/n, 03550, Sant Joan D'Alacant, Alicante, Spain*
[5]*Galician Breast Cancer Screening Program, Galician Regional Health Authority, C/Gregorio Hernández, 2, 4, 15011, A Coruña, Spain*
[6]*Castile-León Breast Cancer Screening Program, General Directorate of Public Health, Avenida Sierra de Atapuerca, S/N, 09071, Burgos, Spain*
[7]*Aragón Breast Cancer Screening Program, Aragon Health Service, C/Ronda de Liberación, 1, 44002, Teruel, Zaragoza, Spain*
[8]*Navarre Breast Cancer Screening Program, Public Health Institute, C/ Leire, 15, 31003, Pamplona, Spain*
[9]*Valencian Breast Cancer Screening Program, General Directorate of Public Health, C/ Micer Mascó, 31, 46010, Valencia, Spain*
[10]*Cancer Prevention and Control Unit, Catalonian Institute of Oncology (ICO), Avenida Gran Vía, S/N, km 2.7, 08907, L´Hospitalet de Llobregat, Barcelona, Spain*
[11]*Balearic Islands Breast Cancer Screening Program, Regional Authority for Health & Consumer Affairs, C/ Cecilio Metelo, 18, 07012, Palma de Mallorca, Islas Baleares, Spain*

## Abstract

The objective of the present study was to assess the reproducibility of data-driven dietary patterns in different samples extracted from similar populations. Dietary patterns were extracted by applying principal component analyses to the dietary information collected from a sample of 3550 women recruited from seven screening centres belonging to the Spanish breast cancer (BC) screening network (Determinants of Mammographic Density in Spain (DDM-Spain) study). The resulting patterns were compared with three dietary patterns obtained from a previous Spanish case–control study on female BC (Epidemiological study of the Spanish group for breast cancer research (GEICAM: grupo Español de investigación en cáncer de mama)) using the dietary intake data of 973 healthy participants. The level of agreement between patterns was determined using both the congruence coefficient (CC) between the pattern loadings (considering patterns with a CC≥0·85 as fairly similar) and the linear correlation between patterns scores (considering as fairly similar those patterns with a statistically significant correlation). The conclusions reached with both methods were compared. This is the first study exploring the reproducibility of data-driven patterns from two studies and the first using the CC to determine pattern similarity. We were able to reproduce the EpiGEICAM Western pattern in the DDM-Spain sample (CC = 0·90). However, the reproducibility of the Prudent (CC = 0·76) and Mediterranean (CC = 0·77) patterns was not as good. The linear correlation between pattern scores was statistically significant in all cases, highlighting its arbitrariness for determining pattern similarity. We conclude that the reproducibility of widely prevalent dietary patterns is better than the reproducibility of more population-specific patterns. More methodological studies are needed to establish an objective measurement and threshold to determine pattern similarity.

**Key words: Dietary patterns: Reproducibility: Congruence coefficients: Principal component analyses: Component loadings: Component scores**

---

Diet is a key modifiable risk factor, but the exploration of its role in disease occurrence is complicated because of methodological issues related to the dietary assessment method used[1–3], food and nutrient interactions[4,5] and differences in food consumption across populations[6–8]. Traditionally, nutritionists and researchers have explored the effect of individual dietary factors in disease occurrence. However, some authors advocate the use of dietary patterns instead of individual foods and nutrients, arguing that they may better capture variability in the population's diet, while allowing the evaluation of interactions between dietary factors[9–11].

These patterns can be identified with data-driven methods such as principal component analysis (PCA), factor analysis (FA) and cluster analysis or can be represented by investigator-driven patterns known as dietary quality indices. Investigator-driven patterns assign a set of scores based on individuals' fulfilment of a set of fixed recommendations. Therefore, they are widely applicable, facilitating the exploration of the reproducibility of their association with different diseases in independent populations[12–16]. However, they present the disadvantage of being very disease dependent, given that they are mainly based on existing evidence of the association between diet and CVD[17]. On the other hand, data-driven dietary patterns are more representative of the diet of the specific population from which they have been extracted and independent of the diseases, but many authors argue that the patterns obtained are very population-dependent, and therefore difficult to reproduce in other settings[11,18,19]. The reproducibility of data-driven dietary patterns has been assessed previously by various authors using dietary information obtained with common assessment tools at different moments of time within the same sample[20–23]. However, no previous studies have explored the reproducibility of data-driven dietary patterns extracted from different samples.

The objective of this study was to assess the reproducibility of data-driven dietary patterns in different samples extracted from similar populations. We compared the results from a previous case–control study Epidemiological study of the Spanish group for breast cancer research (GEICAM: grupo Español de investigación en cáncer de mama) on diet and female breast cancer (BC) in Spain[24] with those obtained from a sample of Spanish women attending BC screening programmes (*Determinantes de la Densidad Mamográfica en España* – Determinants of Mammographic Density in Spain (DDM-Spain)), by evaluating the correlation between pattern scores and the congruence between the composition of patterns in both populations.

## Methods

### Study population and data collection

We used information on three dietary patterns obtained from a previous case–control study on female BC (EpiGEICAM study) using the dietary intake data of 973 healthy participants, aged 22–71 years, and recruited from fourteen Spanish provinces during the period 2006–2011[24]. These patterns will be used as a reference to explore their reproducibility in a different sample using data from the DDM-Spain participants. DDM-Spain is a cross-sectional, multicentre study carried out in seven

screening centres belonging to the Spanish BC screening network and located throughout the Spanish peninsula[25,26]. In Spain, all women aged 50–69 years (45–69 years in some regions), regardless of nationality or legal status, are invited to be screened under these government-sponsored programmes every 2 years. Women were randomly selected among all screening attendants and invited to participate on a daily basis until the minimum sample size of 500 for each centre was reached. A total of 3550 women were recruited between 2007 and 2008, with an average participation rate of 74·5 % (range 64·7–84·0 % across centres). Women were interviewed at the screening centres by trained interviewers who collected demographic, anthropometric, physical activity, gynaecologic, obstetric and occupational data, as well as family and personal history (including weight and height at age 18 years). Information on smoking included current status and months since quitting for ex-smokers. Current smokers were defined as women who smoked at the time of mammography or had quit <6 months before. Dietary intake during the preceding year was collected using a validated 117-item FFQ[27,28]. Postmenopausal status was defined as self-reported absence of menstruation in the previous 12 months. Interviewers measured weight, height and waist and hip circumferences twice using the same protocol and identical balance scales, stadiometers and measuring tapes. A third measure was taken when the first two were not equal.

The DDM-Spain study was conducted according to the guidelines laid down in the Declaration of Helsinki, and all procedures involving human subjects were approved by the bioethics and animal welfare committee at the Carlos III Institute of Health. All participants signed a consent form, including permission to publish results from the current research.

### Dietary patterns

The FFQ used in both studies were designed to assess the whole diet, had similar structures and were based on a validated FFQ[27,28]. However, the FFQ of the DDM-Spain study included some additional food items that were not contained in the FFQ of the EpiGEICAM study[25,26]: the FFQ used in the EpiGEICAM study contained ninety-nine items from which eighty-six were used to create the food groups (after excluding the non-energetic and alcoholic beverages), whereas the FFQ from DDM-Spain included 117 items (the same ninety-nine from DDM-Spain plus eighteen additional foods) from which ninety-nine were used to create the food groups (after excluding non-energetic and alcoholic beverages). In both cases, the dietary information collected was grouped into the exact same twenty-six food groups that are summarised in Table 1, where the items only included in the DDM-Spain study are represented in italics.

The EpiGEICAM study identified three dietary patterns over twenty-six food groups: a Western pattern characterised by elevated intakes of high-fat dairy products, processed meat, refined grains, sweets, energetic drinks and other convenience foods and sauces and by low intakes of low-fat dairy products and whole grains; a Prudent pattern defined by high intakes of low-fat dairy products, vegetables, fruits, whole grains and juices; and a Mediterranean pattern represented by a high intake of fish, vegetables, legumes, boiled potatoes, fruits,

The CC between the pattern loadings of a given pattern from EpiGEICAM ($l_{1j}$) and the pattern loadings of a given pattern from DDM-Spain ($l_{2j}$) for each of the $j = 1, \ldots, 26$ food groups were calculated as follows:

$$CC = \frac{\sum_{j=1}^{26} l_{1j} \times l_{2j}}{\sqrt[2]{\left(\sum_{j=1}^{26} l_{1j}^2\right) \times \left(\sum_{j=1}^{26} l_{2j}^2\right)}}.$$

In addition, to follow the same methodology commonly used in studies exploring the reproducibility of dietary patterns, Spearman's correlation coefficients (Corr) between the EpiGEICAM and the DDM-Spain pattern scores were calculated. For that purpose, patterns scores (which reflect the level of compliance of each woman with each one of the dietary patterns) were calculated as the linear combination of consumption of food groups weighted by the pattern loadings from EpiGEICAM Western, Prudent and Mediterranean patterns and from the set of selected patterns resulting from applying PCA to the DDM-Spain data as follows[34]:

$$P_{ki} = \sum_j (L_{kj} \cdot C_{ji}),$$

where $P$ is the pattern score, $L$ the loading score, $C$ the centred food consumption, $k$ the Western, Prudent and Mediterranean patterns from EpiGEICAM and Western, Prudent and Mediterranean patterns from DDM-Spain, $i = 1, \ldots, 3550$ women and $j = 1, \ldots, 26$ food groups.

CC is the preferred measure for component/factor similarity extracted with PCA/FA because its validity is supported by methodological research[31–33]. In addition, a recent study has questioned the ability of using solely Pearson's correlation (Corr) coefficient to assess pattern similarity[35]. However, the majority of studies exploring the reproducibility of dietary patterns base their conclusions on the latter measure, considering any significant correlation as being indicative of pattern similarity regardless of its value[20–23]. In this study, we provide the correlation coefficient for the sake of comparability with published data, but we will base our final conclusion regarding pattern reproducibility on the CC.

To take into account sampling variability in the estimation of pattern loadings using DDM-Spain data, and subsequently in the estimation of the agreement measurements between the patterns identified within the EpiGEICAM and the DDM-Spain studies, we performed a non-parametric bootstrap estimation with 5000 replications. Using sampling replacement, the bootstrap obtained 5000 replicates of the original DDM-Spain data set. PCA was then applied in each replication, and the three principal components that proved to be more similar to those reported in the EpiGEICAM were selected on the basis of the distance between the pattern loadings (more details are given in the online Supplementary Method 1). The 95 % percentile CI for each parameter were represented by percentiles 2·5 and 97·5 of the 5000 bootstrap point estimates' distribution.

Similar analyses were carried out by applying the PCA to food groups from the DDM-Spain study, which included the same exact eighty-six items considered in the EpiGEICAM analysis (online Supplementary Table S1 and Fig. S1).

Analyses were performed using STATA/MP 14.0.

## Results

The anthropometric, reproductive and socio-demographic characteristics of the EpiGEICAM controls[24] and DDM-Spain women are summarised in Table 2. The DDM-Spain study recruited a higher percentage of older and postmenopausal women (77 v. 47 %), women with higher energy intake (on average 656 kJ/d (157 kcal/d) more in the DDM-Spain group), women with higher BMI and a higher percentage of women who practised physical activity with moderate-to-vigorous intensity (76 v. 63 %). On the other hand, these women reported lower intake of alcohol, lower educational level (34 % with primary school or less in DDM and 16 % in EpiGEICAM), lower percentage of family history of BC (7 v. 20 %), lower age at first delivery (43 % of parous women in the DDM had their first child before 25 years of age, whereas this proportion was 26 % in EpiGEICAM) and there was a lower percentage of nulliparous (9 v. 23 %) women. The distribution of age at menarche and smoking appeared to be fairly similar in both studies.

Fig. 1–3 show the comparison between the original loadings from the EpiGEICAM study with their corresponding values in the DDM-Spain study. Western patterns from both studies were characterised by high intakes of high-fat dairy products, refined grains, energetic drinks and convenience food and sauces and low intakes of low-fat dairy products and whole grains. Correlations with the intake of red and/or processed meat and with sweets were also close to the 0·3 threshold. Moreover, the DDM-Spain Western pattern seemed to be negatively correlated with the consumption of white fish, a result that was not observed in EpiGEICAM. Despite these small differences, the elevated CC between patterns (CC = 0·90) indicates a fair similarity between the Western patterns extracted from the EpiGEICAM and the DDM-Spain data (Fig. 1).

We did not identify a pattern among women of the DDM-Spain study that was highly congruent with the EpiGEICAM Prudent pattern. The most similar pattern presented a high consumption of whole grains and juices but failed to correlate with low-fat dairy products, vegetables and fruits (Fig. 2). Something similar was observed with the Mediterranean pattern: several high correlations were observed with some vegetables, legumes, potatoes and nuts. However, the pattern from the DDM-Spain study did not include other typical factors of the Mediterranean diet, such as fish, olive oil and fruits (even if pattern loadings for these food groups were not low), whereas other foods more common in the Prudent diet, such as low-fat dairy products, or in the Western diet, such as sweets, and sugary and convenience foods, were included with high correlations. According to the CC (0·77), the EpiGEICAM and the DDM-Spain Mediterranean patterns cannot be considered similar (Fig. 3).

Finally, had we considered any significant correlation as being indicative of similarity, we would have concluded that all patterns extracted from the EpiGEICAM data were reproducible in the DDM-Spain study.

## Discussion

To the best of our knowledge, this is the first study exploring the reproducibility of data-driven patterns in two different

**Table 2.** Anthropometric, reproductive and socio-demographic characteristics of EpiGEICAM controls and Determinants of Mammographic Density in Spain (DDM-Spain) women
(Mean values and standard deviations; medians and interquartile ranges (IQR); numbers and percentages)

| | EpiGEICAM controls (*n* 973)* | | | DDM-Spain (*n* 3550) | | |
|---|---|---|---|---|---|---|
| | Mean | SD | % v.e. | Mean | SD | % v.e. |
| EpiGEICAM patterns | | | | | | |
| Western pattern | 0·00 | 3·77 | 16 | 0·00 | 2·31 | 16 |
| Prudent pattern | 0·00 | 3·34 | 13 | 0·34 (−2·21–1·92)† | | 15 |
| Mediterranean pattern | 0·00 | 2·70 | 8 | 0·00 | 1·50 | 7 |
| Participants' characteristics | | | | | | |
| Energy intake (kJ/d) | 7937 | 2627 | | 8593·93 | 2012·88 | |
| Energy intake (kcal/d) | 1897 | 628 | | 2054·15 | 481·09 | |
| Alcohol intake (g/d) | | | | | | |
| Median | 2 | | | 0·85 | | |
| IQR | 0·04–7·10 | | | 0–5·68 | | |
| BMI (kg/m²) | 25·36 | 4·28 | | 28·03 | 4·99 | |
| Age (years) | 50·63 | 9·47 | | 56·20 | 5·46 | |
| Age at menarche (years) | 12·44 | 1·52 | | 13 | 12–14 | |
| | *n* | | % | *n* | | % |
| Physical activity in the last year | | | | | | |
| Low | 287 | | 30 | 842 | | 24 |
| Moderate | 368 | | 38 | 1842 | | 52 |
| Vigorous | 246 | | 25 | 866 | | 24 |
| Unknown | 72 | | 7 | – | | – |
| Smoking | | | | | | |
| Never or former +6 months | 645 | | 67 | 2180 | | 61 |
| Smoker or former smoker <6 months | 325 | | 33 | 1370 | | 39 |
| Unknown | 3 | | 0 | – | | – |
| Education | | | | | | |
| Primary school or less | 158 | | 16 | 1204 | | 34 |
| Secondary school | 489 | | 50 | 1978 | | 56 |
| University | 318 | | 33 | 363 | | 10 |
| Unknown | 8 | | 1 | 5 | | 0 |
| Family history of BC | | | | | | |
| No | 782 | | 80 | 3291 | | 93 |
| Yes | 191 | | 20 | 259 | | 7 |
| Age at first delivery (years) | | | | | | |
| <20 | 45 | | 5 | 302 | | 9 |
| 20–24 | 208 | | 21 | 1194 | | 34 |
| 25–29 | 266 | | 27 | 1271 | | 36 |
| >29 | 148 | | 15 | 465 | | 13 |
| Nulliparous | 220 | | 23 | 316 | | 9 |
| Unknown | 86 | | 9 | 2 | | 0 |
| Menopausal status | | | | | | |
| Premenopausal | 513 | | 53 | 816 | | 23 |
| Postmenopausal | 460 | | 47 | 2734 | | 77 |

v.e., Total variability in food group intakes explained by the pattern; BC, breast cancer.
* Descriptive data extracted from the scientific article of Castello *et al.*[24].
† As distribution of the prudent score was skewed, the median and IQR were used to describe this score.

samples extracted from similar populations. We were able to reproduce the Western pattern identified in women from the EpiGEICAM study among women attending BC screening programmes who participated in the DDM-Spain study. However, the reproducibility of the Prudent and Mediterranean patterns cannot be considered good.

The association between dietary patterns and BC has been explored in many studies in different settings. Most of these studies identified a Western/Unhealthy pattern, which shares the most important characteristics with the Western patterns identified in EpiGEICAM and DDM-Spain, such as high consumption of fatty dairy products, red/processed meat, refined grains, sweets and convenience foods[36–41]. However, the Mediterranean and Prudent patterns have often been mixed

under the names of Vegetable, Prudent, Healthy or Mediterranean diet. These patterns are characterised by a high consumption of vegetables and fruits[36–47] that are an important part of the Mediterranean diet, but fail to include other items such as olive oil[36,38–41,44–47], nuts[36–41,43–47], legumes[37,39–41,44,46,47] or fish[38,41], which are key foods to differentiate the so-called Prudent or Healthy patterns from the Mediterranean.

None of the above-mentioned studies have been able to identify both, a Prudent and a Mediterranean pattern in the same population, probably reflecting the difficulty in differentiating them in contexts where the Mediterranean diet is not very prevalent. On the other hand, the higher agreement in the definition of a Western pattern across studies is consistent with the greater reproducibility of this pattern observed in our study.

**Fig. 1.** Pattern loadings of the Western dietary pattern extracted from the EpiGEICAM study[24] (left) and pattern loadings and 95% percentile CI of the Western pattern extracted from Determinants of Mammographic Density in Spain (DDM-Spain) data (right). * Congruence coefficient (CC) and 95% percentile CI between EpiGEICAM and DDM-Spain pattern loadings. † Correlation coefficient (Corr) and 95% percentile CI between EpiGEICAM and DDM-Spain pattern scores. All correlations were significant at a 95% confidence level.



**Fig. 2.** Pattern loadings of the Prudent dietary pattern extracted from the EpiGEICAM study[24] (left) and pattern loadings and 95% percentile CI of the Prudent pattern extracted from Determinants of Mammographic Density in Spain (DDM-Spain) data (right). * Congruence coefficient (CC) and 95% percentile CI between EpiGEICAM and DDM-Spain pattern loadings. † Correlation coefficient (Corr) and 95% percentile CI between EpiGEICAM and DDM-Spain pattern scores. All correlations were significant at a 95% confidence level.

EpiGEICAM Mediterranean / DDM-Spain Mediterranean

| | EpiGEICAM | DDM-Spain | 95% CI |
|---|---|---|---|
| High-fat dairy products | 0·20 | 0·11 | 0·01, 0·26 |
| Low-fat dairy products | −0·01 | 0·48 | −0·20, 0·55 |
| Eggs | 0·16 | 0·18 | 0·09, 0·22 |
| White meat | 0·18 | 0·16 | 0·08, 0·21 |
| Red meat | 0·22 | 0·21 | 0·07, 0·27 |
| Processed meat | 0·26 | 0·25 | 0·13, 0·30 |
| White fish | 0·34 | 0·05 | −0·02, 0·12 |
| Oily fish | 0·44 | 0·21 | 0·12, 0·27 |
| Seafood/shellfish | 0·35 | 0·25 | 0·17, 0·31 |
| Leafy vegetables | 0·40 | 0·22 | 0·13, 0·30 |
| Fruiting vegetables | 0·45 | 0·27 | 0·17, 0·34 |
| Root vegetables | 0·44 | 0·30 | 0·17, 0·38 |
| Other vegetables | 0·42 | 0·22 | 0·14, 0·29 |
| Legumes | 0·34 | 0·32 | 0·21, 0·36 |
| Potatoes | 0·40 | 0·33 | 0·21, 0·41 |
| Fruits | 0·31 | 0·23 | 0·13, 0·29 |
| Nuts | 0·29 | 0·44 | 0·29, 0·50 |
| Refined grains | 0·23 | 0·27 | 0·11, 0·32 |
| Whole grains | 0·06 | −0·03 | −0·08, 0·06 |
| Olives and vegetable oil | 0·34 | 0·22 | 0·14, 0·27 |
| Other edible fats | 0·11 | 0·16 | 0·10, 0·20 |
| Sweets | 0·05 | 0·45 | 0·15, 0·53 |
| Sugary | 0·00 | 0·32 | 0·19, 0·40 |
| Juices | −0·39 | −0·13 | −0·17, −0·09 |
| Energetic drinks | −0·25 | −0·08 | −0·19, 0·00 |
| Convenience food and sauces | 0·24 | 0·34 | 0·14, 0·39 |

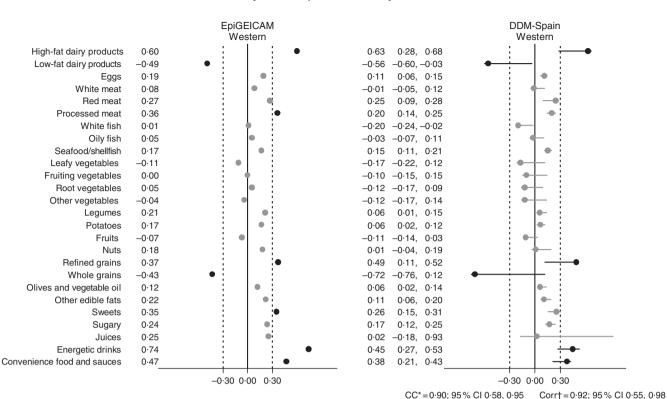CC* = 0·77; 95 % CI 0·65, 0·83    Corr† = 0·74; 95 % CI 0·63, 0·79

**Fig. 3.** Pattern loadings of the Mediterranean dietary pattern extracted from the EpiGEICAM study[24] (left) and pattern loadings and 95 % percentile CI of the Mediterranean pattern extracted from Determinants of Mammographic Density in Spain (DDM-Spain) data (right). * Congruence coefficient and 95 % percentile CI between EpiGEICAM and DDM-Spain pattern loadings. † Correlation coefficient (Corr) and 95 % percentile CI between EpiGEICAM and DDM-Spain pattern scores. All correlations were significant at a 95 % confidence level.
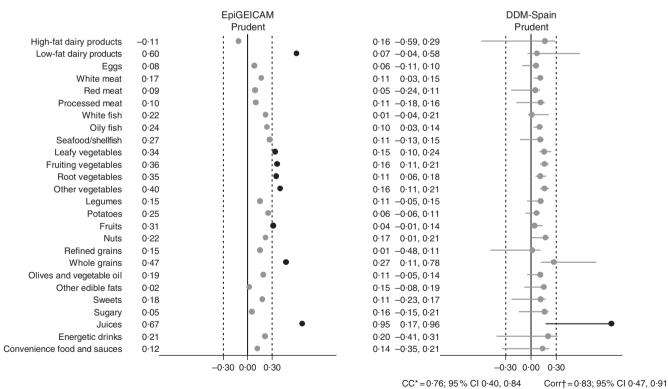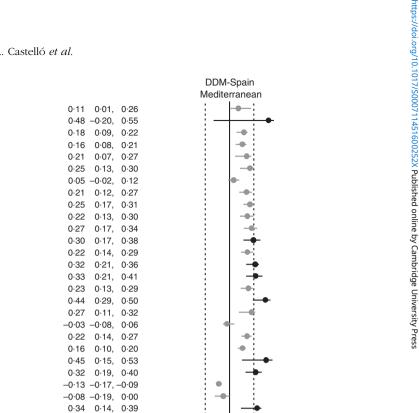
As noted earlier in this study, PCA reduces a set of inter-correlated variables to a group of principal components (dietary patterns in this case) so that the maximum correlation between the variables within components and the minimum correlation among components are obtained[48]. Therefore, the greater the variability in diet, the easier it will be to find clearly differentiated independent patterns. In our study, although EpiGEICAM included women from fourteen Spanish provinces (four of them on the Mediterranean coast), DDM-Spain participants were recruited from screening centres located in seven provinces (three of them located on the Mediterranean coast). Therefore, the greater geographical distribution in the EpiGEICAM study may imply a greater representativeness of all diets across the Spanish territory. In addition, distribution of age among DDM-Spain women was more homogeneous (range = 45–69) than that observed in the EpiGEICAM participants (range = 22–71). As García-Arenzana et al.[49] previously described, older women tend to have healthier dietary habits than younger women, which may have produced a more heterogeneous distribution of dietary habits in the EpiGEICAM study. This heterogeneity might have facilitated the identification of more specific patterns, not only limited to the discrimination of two antagonistic patterns (Western v. Healthy/Prudent/Mediterranean) but also allowing the clear differentiation of patterns with subtle differences, such as the Prudent and Mediterranean patterns.

Regarding the pre-established thresholds for the CC that define the similarity of dietary patterns in both studies, we based our decision on three published pieces of research that evaluated concordance coefficients in light of the subjective opinion of several experienced researchers judging the equivalence between different components[31–33]. Haven and Nesselroade[31,33] argue that values over 0·80 are enough to assume fair similarity between components, whereas Lorenzo-Seva & Berge[32] maintain a more conservative approach setting the cut-off point for fair similarity at 0·85 and preventing a CC below this value from being interpreted as indicative of similarity. All three articles agree on the difficulty in setting up a cut-off point under which patterns should be considered clearly different. Despite the fact that the CC is considered a good measure of agreement between components or factors extracted with PCA or FA[31–33], the existing bibliography evaluating the reproducibility of data-driven dietary patterns does not use this measure and bases its conclusions only on the correlations between pattern scores, considering any significant correlation as being indicative of similarity regardless of its value[20–23], which can be as low as 0·27[23]. In our case, the correlations were significant and high for all three patterns (Fig. 1–3). However, according to the CC, only the Western pattern can be considered fairly similar between studies, which highlights the arbitrariness of the significance of the linear correlation to define pattern similarity and the need to choose an appropriate measure and a concrete threshold for such a measure to determine the level of congruence between patterns. In this regard, we have recently explored the applicability of previously reported dietary patterns in a different setting and we found that, for CC between pattern loadings ≥0·82 or correlations between pattern scores ≥0·57, patterns not only appear to have a very similar composition

but also are similarly associated with BC risk[35]. The same direction of the associations but loss of significance was observed for values of the CC between pattern loadings ≤0·77 and values of the correlation between pattern scores ≤0·52. In the present study, taking into account only the methodological studies published regarding the threshold of the CC for pattern similarity[31–33], we followed the most conservative approach and considered dietary patterns to be fairly similar if CC values were ≥0·85.

A major limitation of the use of dietary patterns is the potential for subjective interpretations by the investigator to be introduced at various stages of the dietary patterns' construction. Subjective decisions that might affect the comparability between studies are as follows: which foods should be included in each of the defined groups, the thresholds chosen to determine the contribution of food groups to the identified dietary patterns and the assignation of a label to each of these patterns[9–11,18,19]. However, we have demonstrated that this limitation can be overcome by a detailed analysis when comprehensive information on food grouping and loadings is provided by Castello et al.[35]. On the other hand, both FFQ from EpiGEICAM and DDM-Spain collected information on ninety-nine identical foods, except for the fact that DDM-Spain included eighteen additional foods that were not included in EpiGEICAM. In addition, the same group of researchers took principal responsibility for the analysis of the data; therefore, food grouping and labelling were very similar in both studies.

Finally, we summarise the main strengths of the present study. As previously mentioned, various studies have assessed the reproducibility of investigator-driven patterns[12–16]. The reproducibility of data-driven dietary patterns extracted from the same sample using the dietary information obtained with different assessment tools or in different time points[20–23] has also been explored. However, to our knowledge, this is the first study assessing the reproducibility of data-driven dietary patterns in different samples from similar populations and the first using the CC to evaluate their similarity. In addition, most of the published studies on reproducibility of data-driven dietary patterns based their conclusions on limited sample sizes that ranged from 124–498[20–22]. Dietary patterns from EpiGEICAM were extracted over 973 healthy women, and for DDM-Spain the sample size was 3550, a size only exceeded by the Newby et al. study[23].

## Conclusions

The reproducibility of widely prevalent dietary patterns such as the Western pattern is better than the reproducibility of patterns more specific to certain populations, such as the Mediterranean. More methodological studies exploring the reproducibility of dietary patterns are needed to establish a more objective threshold for the CC between pattern loadings and their equivalent Corr between pattern scores that define pattern similarity.

## Acknowledgements

## Supplementary material

For supplementary material/s referred to in this article, please visit http://dx.doi.org/10.1017/S000711451600252X

## References

1. Bingham SA, Luben R, Welch A, et al. (2003) Are imprecise methods obscuring a relation between fat and breast cancer? Lancet 362, 212–214.
2. Kelemen LE (2007) GI Epidemiology: nutritional epidemiology. Aliment Pharmacol Ther 25, 401–407.
3. Willett W (2001) Commentary: dietary diaries versus food frequency questionnaires – a case of undigestible data. Int J Epidemiol 30, 317–319.
4. Jacobs DR Jr & Steffen LM (2003) Nutrients, foods, and dietary patterns as exposures in research: a framework for food synergy. Am J Clin Nutr 78, 508S–513S.
5. Messina M, Lampe JW, Birt DF, et al. (2001) Reductionism and the narrowing nutrition perspective: time for reevaluation and emphasis on food synergy. J Am Diet Assoc 101, 1416–1419.
6. Irala-Estevez JD, Groth M, Johansson L, et al. (2000) A systematic review of socio-economic differences in food habits in Europe: consumption of fruit and vegetables. Eur J Clin Nutr 54, 706–714.
7. Sanchez-Villegas A, Martinez JA, Prattala R, et al. (2003) A systematic review of socioeconomic differences in food habits in Europe: consumption of cheese and milk. Eur J Clin Nutr 57, 917–929.
8. Teufel NI (1997) Development of culturally competent food-frequency questionnaires. Am J Clin Nutr 65, 1173S–1178S.
9. Barkoukis H (2007) Importance of understanding food consumption patterns. J Am Diet Assoc 107, 234–236.
10. Hu FB (2002) Dietary pattern analysis: a new direction in nutritional epidemiology. Curr Opin Lipidol 13, 3–9.
11. Jacques PF & Tucker KL (2001) Are dietary patterns useful for understanding the role of diet in chronic disease? Am J Clin Nutr 73, 1–2.
12. George SM, Ballard-Barbash R, Manson JE, et al. (2014) Comparing indices of diet quality with chronic disease

mortality risk in postmenopausal women in the women's health initiative observational study: evidence to inform national dietary guidance. *Am J Epidemiol* **180**, 616–625.

13. Harmon BE, Boushey CJ, Shvetsov YB, *et al.* (2015) Associations of key diet-quality indexes with mortality in the Multiethnic cohort: the dietary patterns methods project. *Am J Clin Nutr* **101**, 587–597.
14. Liese AD, Krebs-Smith SM, Subar AF, *et al.* (2015) The dietary patterns methods project: synthesis of findings across cohorts and relevance to dietary guidance. *J Nutr* **145**, 393–402.
15. McCullough ML (2014) Diet patterns and mortality: common threads and consistent results. *J Nutr* **144**, 795–796.
16. Reedy J, Krebs-Smith SM, Miller PE, *et al.* (2014) Higher diet quality is associated with decreased risk of all-cause, cardio-vascular disease, and cancer mortality among older adults. *J Nutr* **144**, 881–889.
17. Fung TT, McCullough ML, Newby PK, *et al.* (2005) Diet-quality scores and plasma concentrations of markers of inflammation and endothelial dysfunction. *Am J Clin Nutr* **82**, 163–173.
18. Martinez ME, Marshall JR & Sechrest L (1998) Invited commentary: factor analysis and the search for objectivity. *Am J Epidemiol* **148**, 17–19.
19. Slattery ML & Boucher KM (1998) The senior authors' response: factor analysis as a tool for evaluating eating patterns. *Am J Epidemiol* **148**, 20–21.
20. Hu FB, Rimm E, Smith-Warner SA, *et al.* (1999) Reproduci-bility and validity of dietary patterns assessed with a food-frequency questionnaire. *Am J Clin Nutr* **69**, 243–249.
21. Khani BR, Ye W, Terry P, *et al.* (2004) Reproducibility and validity of major dietary patterns among Swedish women assessed with a food-frequency questionnaire. *J Nutr* **134**, 1541–1545.
22. Nanri A, Shimazu T, Ishihara J, *et al.* (2012) Reproducibility and validity of dietary patterns assessed by a food frequency questionnaire used in the 5-year follow-up survey of the Japan Public Health Center-Based Prospective Study. *J Epidemiol* **22**, 205–215.
23. Newby PK, Weismayer C, Akesson A, *et al.* (2006) Long-term stability of food patterns identified by use of factor analysis among Swedish women. *J Nutr* **136**, 626–633.
24. Castelló A, Pollan M, Buijsse B, *et al.* (2014) Spanish Medi-terranean diet and other dietary patterns and breast cancer risk: case-control EpiGEICAM study. *Br J Cancer* **111**, 1454–1462.
25. Lope V, Perez-Gomez B, Sanchez-Contador C, *et al.* (2012) Obstetric history and mammographic density: a population-based cross-sectional study in Spain (DDM-Spain). *Breast Cancer Res Treat* **132**, 1137–1146.
26. Pollan M, Lope V, Miranda-Garcia J, *et al.* (2012) Adult weight gain, fat distribution and mammographic density in Spanish pre- and post-menopausal women (DDM-Spain). *Breast Cancer Res Treat* **134**, 823–838.
27. Vioque J, Navarrete-Munoz EM, Gimenez-Monzo D, *et al.* (2013) Reproducibility and validity of a food frequency questionnaire among pregnant women in a Mediterranean area. *Nutr J* **12**, 26.
28. Willett WC, Sampson L, Stampfer MJ, *et al.* (1985) Reprodu-cibility and validity of a semiquantitative food frequency questionnaire. *Am J Epidemiol* **122**, 51–65.
29. Burt C (1948) Factor analysis and canonical correlations. *Br J Math Stat Psychol* **1**, 95–106.
30. Tucker LR (1951) A method for the synthesis of factor analysis studies, Personnel Research Section Report no. 984. Washington, DC: Department of the Army.
31. Haven S & Berge J (1977) Tucker's coefficient congruence as a measure of factorial invariance: an empirical study. Heymans Bulletin no. 290 EX. Groningen, The Netherlands: University of Groningen.
32. Lorenzo-Seva U & Berge J (2006) Tucker's congruence coef-ficient as a meaningful index of factor similarity. *Methodology* **2**, 54–67.
33. Nesselroade J & Baltes P (1970) On a dilemma of comparative factor analysis: a study of factor matching based on random data. *Educ Psychol Meas* **30**, 935–948.
34. Schulze MB, Hoffmann K, Kroke A, *et al.* (2003) An approach to construct simplified measures of dietary patterns from exploratory factor analysis. *Br J Nutr* **89**, 409–419.
35. Castello A, Buijsse B, Martin M, *et al.* (2016) Evaluating the applicability of data-driven dietary patterns to independent samples with focus on measurement tools for pattern simi-larity. *J Acad Nutr Diet* (In the Press).
36. Agurs-Collins T, Rosenberg L, Makambi K, *et al.* (2009) Dietary patterns and breast cancer risk in women participating in the black women's health study. *Am J Clin Nutr* **90**, 621–628.
37. Cottet V, Touvier M, Fournier A, *et al.* (2009) Postmenopausal breast cancer risk and dietary patterns in the E3N-EPIC prospective cohort study. *Am J Epidemiol* **170**, 1257–1267.
38. Cui X, Dai Q, Tseng M, *et al.* (2007) Dietary patterns and breast cancer risk in the Shanghai breast cancer study. *Cancer Epidemiol Biomarkers Prev* **16**, 1443–1448.
39. Terry P, Suzuki R, Hu FB, *et al.* (2001) A prospective study of major dietary patterns and the risk of breast cancer. *Cancer Epidemiol Biomarkers Prev* **10**, 1281–1285.
40. Velie EM, Schairer C, Flood A, *et al.* (2005) Empirically derived dietary patterns and risk of postmenopausal breast cancer in a large prospective cohort study. *Am J Clin Nutr* **82**, 1308–1319.
41. Wu AH, Yu MC, Tseng CC, *et al.* (2009) Dietary patterns and breast cancer risk in Asian American women. *Am J Clin Nutr* **89**, 1145–1154.
42. Adebamowo CA, Hu FB, Cho E, *et al.* (2005) Dietary patterns and the risk of breast cancer. *Ann Epidemiol* **15**, 789–795.
43. Bessaoud F, Tretarre B, Daures JP, *et al.* (2012) Identification of dietary patterns using two statistical approaches and their association with breast cancer risk: a case-control study in Southern France. *Ann Epidemiol* **22**, 499–510.
44. De Stefani E, Deneo-Pellegrini H, Boffetta P, *et al.* (2009) Dietary patterns and risk of cancer: a factor analysis in Uruguay. *Int J Cancer* **124**, 1391–1397.
45. Demetriou CA, Hadjisavvas A, Loizidou MA, *et al.* (2012) The Mediterranean dietary pattern and breast cancer risk in Greek-Cypriot women: a case-control study. *BMC Cancer* **12**, 113.
46. Hirose K, Matsuo K, Iwata H, *et al.* (2007) Dietary patterns and the risk of breast cancer in Japanese women. *Cancer Sci* **98**, 1431–1438.
47. Zhang CX, Ho SC, Fu JH, *et al.* (2011) Dietary patterns and breast cancer risk among Chinese women. *Cancer Causes Control* **22**, 115–124.
48. Rencher A (2002) Principal component analysis. In *Methods of Multivariate Analysis*, pp. 380–407. New York: John Wiley & Sons, Inc.
49. García-Arenzana N, Navarrete-Munoz EM, Peris M, *et al.* (2012) Diet quality and related factors among Spanish female participants in breast cancer screening programs. *Menopause* **19**, 1121–1129.