

A Statistical Agreement-Based Approach for Difference Testing

F. Calle-Alonso, C. J. Pérez

Faculty of Veterinary Medicine, Biostatistics Unit, University of Extremadura, Avda. de la Universidad s/n, 10003 Cáceres, Spain.

Abstract

A statistical approach based on measures of agreement is proposed for use in a sensory analysis context. This approach considers the idea of using statistical agreement to provide information on the homogeneity of the raters' responses, so that this information can then be used to discriminate between products. It can also be used to measure the expertise level of raters. Although the prime focus is on difference testing by the triangle test (ISO 4120:2008), the proposed methodology can also be applied in other contexts such as the paired comparison test (ISO 5495:2009) or the duo-trio test (ISO 10399:2010), among others. The proposed approach is not a substitute for binomial statistical analysis, but rather it can be used as a complement. It is especially useful when few panelists are available and replications are needed. An experiment that evaluates two types of Iberian dry-cured pork loins through the triangle test is performed to illustrate the applicability of the proposed approach.

Keywords: Binomial model, Beta-binomial model, Difference testing, Measures of agreement, Multiple raters, Sensory analysis, Triangle test.

1. Introduction

2 Sensory analysis can be used to provide subjective information about the
3 acceptance of different products, and is also widely used in determining over-
4 all quality. As is well known, the use of a panel is a very important tool in
5 attempting to describe a product's different and complex features. But it
6 also has some drawbacks, subjectivity and low repeatability, for instance. In
7 order to improve the reliability of the results and avoid these problems, some

8 countries have enacted laws which give legal value to sensorial analysis tests
9 and are aimed at homogenizing the results. The International Standard-
10 ization Organization (ISO) has proposed a standard for sensory analysis to
11 ensure that products and services are safe, reliable, and of good quality (see
12 International Organization for Standardization (TC34/SC12)). This standard
13 has been applied in many different fields: quality control, research and de-
14 velopment, market research, protected designation of origin,... For example,
15 the International Organization for Standardization (2004b) norm is generally
16 used by government agencies to regulate all aspects of the triangle test.

17 As well as the triangle test (International Organization for Standardization
18 (2004b)), others such as the paired comparison test (International Organiza-
19 tion for Standardization (2005)) and the duo-trio test (International Organiza-
20 tion for Standardization (2004a)) use the binomial distribution to discriminate
21 between products. However, the binomial model is not suitable in situations
22 in which there is overdispersion. An extension of the binomial model – the
23 beta-binomial model – is used to fit overdispersed binomial data (see Ennis
24 and Bi (1998)). The information provided by the binomial model can be
25 complemented with that obtained with statistical measures of agreement.

26 Agreement among raters is of great importance for researchers and prac-
27 titioners who describe and evaluate objects and behaviours in a number of
28 fields, including the social and behavioural sciences. Fleiss et al. (1969) and
29 Fleiss Fleiss (1971) were two landmark articles on agreement measures. Since
30 then, statistical agreement has been an active research area whose techniques
31 have been widely used in practice. The most popular measure of agreement
32 is Cohen's kappa. There are, however, many others available, each one with
33 its own particular characteristics that make it interesting to use in differ-
34 ent contexts. Agresti (1996) presented several modeling techniques for the
35 analysis of categorical data, in addition to an invaluable summary of the
36 state-of-the-art. Von Eye and Mun (2005) provided a comprehensive refer-
37 ence book that analyses rater agreement from four different perspectives,
38 including log-linear modeling.

39 Although statistical measures of agreement have been widely used in
40 many fields of knowledge, especially in the biomedical sciences, they have
41 remained almost unexplored in the field of sensory analysis. Nevertheless, a
42 few interesting results can be found in the literature. For example, Wu and
43 Chen (1995) considered the agreement among raters to evaluate the agree-
44 ment of tea sensory data, and Mounchili et al. (2005) considered agreement
45 in a sensory analysis of milk samples. In the present communication, a sta-

46 tistical agreement-based approach is presented for sensory analysis. This
47 approach proposes the use of an efficient measure of agreement for two or
48 more raters in which the response is given on a qualitative scale. It is shown
49 how these measures can provide information about the process of seeking for
50 differences. The idea is based on measuring the homogeneity of the raters'
51 responses, and then using this information to analyse differences between
52 products. The proposed approach is connected to the standard binomial
53 procedure. Also, measures of agreement can be used to qualify novice raters'
54 aptitudes, and mark when they become experts. Although the prime focus
55 is on difference testing using the triangle test, the methodology can also be
56 applied in other difference or similarity tests (see Bi (2011)).

57 There are many examples in the literature showing the importance of
58 sensory analysis in terms of designing, testing, launching, and rethinking
59 food products. For example, the characterization of dry-cured shoulder of
60 pork (Lorenzo et al. (2008)), Iberian dry-cured ham (Martín et al. (2010)),
61 pineapple juice (Silva et al. (2010)), and Gamonedo cheese (Ramos-Guajardo
62 and González-Rodríguez (2011)). In the present study, the differences of two
63 Iberian dry-cured pork loins are evaluated through a triangle test by using
64 measures of agreement.

65 The paper is organized as follows. Section 2 presents a short discussion of
66 the main measures of inter-rater agreement and their application to sensory
67 analysis. Section 3 describes the agreement-based approach and connects
68 it to the standard binomial procedure. Some illustrative examples are also
69 presented. In Section 4, an experiment designed involving a triangle test il-
70 lustrates the applicability of the proposed approach. Finally, the conclusions
71 are presented in Section 5.

72 **2. Measures of agreement**

73 Currently there is no standard measure of agreement used by the scientific
74 community, although Cohen's kappa has a long history of use as an index
75 of inter-rater agreement. However, Cohen's kappa is not always the best
76 choice (see, e.g., Gwet (2002) and Fletcher et al. (2011)). When two raters
77 are involved, there is a wide range of available measures of agreement in
78 the statistical literature. For more than two raters, the number of available
79 measures is dramatically reduced because of the difficulty of interpreting the
80 results. Also, the levels of agreement tend to decrease as the number of raters

81 grows. In the following paragraphs, we shall describe the main measures of
82 agreement from the perspective of the proposal of the present work.

83 Sometimes, measures of agreement can be affected by chance. When the
84 raters are not sure about the correct classification of a product, some guessing
85 may occur. Guessing may be total, if they are not able to distinguish any-
86 thing at all, or partial, if they are guessing only some of the samples. When
87 two raters make their predictions by chance, they sometimes agree. The
88 question is when such agreements should count towards a statistical index of
89 agreement. Theoretically, if the agreement by chance can be estimated, then
90 this effect could be removed from the total agreement to discover the true
91 agreement among raters. This is what chance-corrected measures try to do,
92 but it is not at all clear that the final conclusion is reliable. Indeed, some-
93 times these measures yield paradoxical and counter-intuitive results. The
94 choice between chance-corrected or non-chance-corrected measures has been
95 a topic of some debate (see, e.g., Guggenmoos-Holzmann (2006)). However,
96 the best option is to use the measure of agreement that by definition and
97 meaning best fits the nature of the problem being addressed, regardless of
98 whether or not it includes corrections for chance.

99 Consider m raters and c alternatives on a categorical scale. For the trian-
100 gle, paired comparison, duo-trio, 2-AFC, and 3-AFC tests, the alternatives
101 are right (positive) or wrong (negative) responses, i.e., $c = 2$. For the sake
102 of simplicity, we shall first consider the notation for two panelists ($m = 2$),
103 to subsequently generalize it to $m \geq 2$. For a qualitative variable X ranging
104 over 1, 2 (positive and negative ratings, respectively), n_{ij} will denote the ob-
105 served frequency for rater 1 giving the response $X = i$, and rater 2 giving
106 the response $X = j$. The observed frequencies can be presented in a con-
107 tingency table of dimension 2×2 , or generally $m \times 2$. One has, of course,
108 that $\sum_i \sum_j n_{ij} = n$. The proportion parameters, ρ_{ij} , are estimated from the
109 observed proportions, i.e., $\hat{\rho}_{ij} = n_{ij}/n$.

110 The simplest non-chance-corrected measure is the proportion of overall
111 agreement, defined as $\sum_i \rho_{ii}$ and estimated as $\sum_i n_{ii}/n$. The estimated value
112 will be 0 when there is no agreement at all, and 1 when the agreement is
113 absolute. This measure has been criticized, because it can be high even with
114 hypothetical raters who randomly guess on each case with probabilities equal
115 to the observed base rate. There are other non-chance-corrected measures,
116 such as the Holley and Guildford G coefficient and the Rogot and Goldberg
117 A_1 and A_2 indices (see Gwet (2002)). These measures were defined for only
118 two raters, but they can be extended to three or more. They include the

119 agreement observed for all the possible rater responses, i.e., for the right and
 120 wrong responses jointly. They are unable to distinguish between agreement
 121 on right responses and agreement on wrong ones. Since our interest is in
 122 discriminating products, it is important to know the agreement for the right
 123 and the wrong cases separately, and while G , A_1 , and A_2 do provide some
 124 information, they are really unsuitable for the present discrimination context.

125 In order to address the problem of analysing the agreement based on
 126 only one specific response, there are some non-chance-corrected measures
 127 available for two raters – concordance proportion, Dice index, Goodman and
 128 Kruskal λ_r , and the Jaccard measure, among others (see Shoukri (2004)).
 129 The most interesting measure in this context is Dice index since it is easier
 130 to interpret than the others and leads to more realistic agreement values.
 131 Dice’s index has been widely used in several fields of knowledge (see, e.g.,
 132 Ajmone-Marsan et al. (2001) and LaPara et al. (2002)), but has been left
 133 practically unexplored in that of sensory analysis. One exception is the work
 134 of Mouchili et al. (2005) who applied it to the organoleptic analysis of milk
 135 samples.

136 The proposal that we shall describe in the following section is based on
 137 the use of positive and negative Dice indices to discriminate products and
 138 assess agreement. In the following paragraphs of this section, we shall present
 139 the main results for these indices.

For two raters, the Dice indices are defined as:

$$D_i^{(2)} = \frac{2\rho_{ii}}{\rho_i^{(1)} + \rho_i^{(2)}}, i = 1, 2,$$

140 where $\rho_i^{(1)}$ and $\rho_i^{(2)}$ are the marginal probabilities for each rater, i.e., $\rho_i^{(1)} =$
 141 $\rho_{i1} + \rho_{i2}$ and $\rho_i^{(2)} = \rho_{1i} + \rho_{2i}$. $D_1^{(2)}$ refers to the positive response and $D_2^{(2)}$
 142 to the negative one. Both values are defined in the interval $[0, 1]$, taking
 143 the value 1 when there is total agreement for the i -th alternative, and 0
 144 when there is no agreement at all for that alternative. Graham and Bull
 145 (1998) and Mackinnon (2000) used the delta method to derive formulas for
 146 the asymptotic standard errors of these specific response measures. Alterna-
 147 tively, the standard errors can be estimated by Jackknife or nonparametric
 148 bootstrap (see Severiano et al. (2011)) techniques. We have not found any
 149 closed expressions for sampling distributions of Dice’s indices for hypothesis
 150 testing in the literature. However, simulation-based approaches can be used
 151 to estimate confidence intervals and to test hypotheses.

152 Dice already noted that the similarity measure proposed in an context of
 153 ecology in that work could be extended to three or more species. Warrens
 154 (2008) discussed this extension further in analysing similarity coefficients for
 155 binary data. The Dice index for more than two raters can be defined as:

$$D_i^{(m)} = \frac{m\rho_{ii\dots i}}{\sum_{j=1}^m \rho_i^{(j)}}, \quad (1)$$

156 where $\rho_{ii\dots i}$ is the proportion parameter for the case where all the raters
 157 have chosen the alternative i , and $\rho_i^{(j)}$ is the marginal probability that is
 158 obtained for rater j with alternative i . This generalized index maintains the
 159 same properties as the two-rater one. However, there are no longer any closed
 160 expressions for asymptotic standard errors such as there were in the two-rater
 161 case. Confidence intervals or hypothesis testing must be performed using
 162 simulation-based approaches such as Monte Carlo or resampling methods
 163 (see, e.g., Manly (1997)).

With respect to the chance-corrected measures, it is worth mentioning
 that they have traditionally been far more commonly used than the non-
 chance-corrected ones. They can be presented with the common expression:

$$M(I) = \frac{I_o - I_e}{1 - I_e},$$

164 where I_o and I_e are the observed and the expected values of the non-chance-
 165 corrected index of agreement, respectively.

166 The most extensively used measures of agreement are Cohen's kappa for
 167 two raters (Cohen (1960)) and Fleiss's generalized kappa for several raters
 168 (Fleiss (1971)). They have been applied to the estimation of conjoint agree-
 169 ment (agreement for all possible alternatives) in many fields of knowledge.
 170 In the specific field of sensory analysis, there are only a few published appli-
 171 cations of agreement measures, and most of these use Cohen's kappa. Pons-
 172 Sanchez-Cascado et al. (2006) and Baixas-Nogueras et al. (2003) used this
 173 index to evaluate the agreement between two rejection methods for anchovies
 174 and hake, respectively, and Jenschke et al. (2007) used it to assess the agree-
 175 ment between panelists in a beef tasting experiment. Wu and Chen (1995)
 176 used the multi-rater agreement Kappa to evaluate the agreement of tea sen-
 177 sory data. Cohen's Kappa has also been proposed for use in a complementary
 178 way (see, e.g., Cicchetti and Feinstein (1990a)). A chance-corrected measure
 179 that could be used for discrimination testing is the conditional Kappa, al-
 180 though it does suffer from some drawbacks in that context. In the difference

181 problem, it is unusual for the raters to try to guess – they might do so some-
182 times, but only in very few cases. Thus, the basic logic behind studying
183 a chance-corrected measure such as the conditional kappa is inappropriate
184 here. In addition, since both Dice indices are used, the overall agreement
185 for the two possible responses is determined from the effect of the marginal
186 proportions that are considered in the conditional kappa, so that there is no
187 need to correct for possible effects of chance as has to be done in the condi-
188 tional kappa. Note also that the agreement given by the conditional Kappa
189 is in most case underestimated and that, to distinguish fair agreement, the
190 number of agreed responses needs to be very high. Finally, chance-corrected
191 measures may yield misleading values for binary ratings, such as in the prob-
192 lem to be addressed in the present work (see Guggenmoos-Holzmann (2006)).
193 The generalized Dice index does not have these drawbacks, and thus provides
194 a clear and realistic measure of the agreement that may be useful in discrim-
195 inating products.

196 In the following section, we shall present the proposed approach and
197 connect it to the standard procedures.

198 **3. An agreement-based approach**

199 *3.1. Introduction*

200 The problem of discriminating products is a special case studied in sensory
201 analysis. The three standardized tests defined for this kind of problem are
202 the triangle test, the paired test, and the duo-trio test. Other tests put for-
203 ward in the literature are demonstrating high potential, for example, Tetrad
204 (see Garcia et al. (2012)), A-Not A, 2-AFC, and 2-AFCR (see Van Hout
205 et al. (2011)). Besides, alternative strategies have been recently considered,
206 like Bayesian methodology used by Bi (2011) and Dubnicka (2013). As a
207 consequence, differentiation tests is a very active research topic.

208 For discrimination problems, it is usual for there not to be many panelists
209 available, and for a considerable number of observations to be required for
210 the results to be to significant, this latter usually due to the smallness of the
211 differences between the products. Meyners and Brockhoff (2003) showed that
212 one might be able to add to the power of the test while reducing the number
213 of raters by increasing the total number of assessments with replications. In
214 particular, the test must be repeated several times to provide the necessary
215 amount of information, and then results combined. The question of whether
216 it is permissible to combine results from replicated triangular tests has been

217 extensively discussed by various authors. According to Kunert and Meyn-
218 ers (1999), if the experiment is properly randomized and controlled then the
219 assessments are will be independent and will have a binomially distributed
220 success probability. But they also noted that it is difficult for these assump-
221 tions to be satisfied when replications are performed, and in such a case the
222 choice probability and a measure of heterogeneity should be estimated. The
223 binomial distribution assumes the existence of only one source of variability
224 – that based on the samples. Therefore, when panelists are rating identically
225 from one sample to another, the variance is completely explained by the
226 binomial distribution. But rating identically for all replications is unusual,
227 although it may sometimes be the case.

228 A general problem with discrimination testing is the assumption that all
229 panelists have the same probability of discrimination, that there are only two
230 kinds of raters – non-discriminators and perfect discriminators. The former
231 type always guesses, and the latter always discriminates correctly through
232 all the replications. This assumption is unrealistic, and it is evident that
233 panelist variability needs to be taken into account when collecting replicated
234 observations from the same panelists (see Ennis and Jesionka (2011)). To deal
235 with this difficulty, a beta distribution can be used instead of a binomial to
236 model variation in inter-trial choice probabilities.

237 The beta-binomial model considers the variability among samples as well
238 as the variability among raters (also termed overdispersion), making it pos-
239 sible to combine responses across raters and replications. This increases the
240 power of the test for a small panel size (see Anderson (1988)). The beta-
241 binomial distribution is the natural extension of the binomial. It is based on
242 the binomial with parameter p following a beta distribution with parameters
243 a and b . It is useful to apply the re-parameterization $\mu = a/(a + b)$ and
244 $\gamma = 1/(a + b + 1)$, which are the mean of the binomial parameter p and a
245 scale parameter that measures its variation, respectively (see, e.g, Ennis and
246 Bi (1998)). The scale parameter varies from 0 when there is no overdispers-
247 sion to 1 when there is total overdispersion. Ennis and Bi (1998) provide
248 hypothesis tests to evaluate whether the parameters differ significantly from
249 the quantities of interest. With the proportion μ , one tests whether or not
250 the differentiation was the result of guessing. In testing the overdispersion
251 parameter, one may study whether the appropriate distribution is the bino-
252 mial ($H_0 : \gamma = 0$) or the beta-binomial ($H_1 : \gamma \neq 0$). It is more important,
253 however, to correctly estimate and interpret the parameters and their vari-
254 ances that apply to the problem at hand.

255 The results given by the binomial and beta-binomial models for the same
 256 problem are generally sufficiently different for different conclusions to be
 257 drawn regarding the products being tested. When the sensorial judgement
 258 of these products is fairly easy to do, it is generally advantageous to collect
 259 replicated data and analyse it using the beta-binomial model (see Ligget
 260 and Delwiche (2005)). Even so, the question of overdispersion should be
 261 considered, because the binomial model might be appropriate in some cases.
 262 For example, in a test of the sensory quality of cabbage, Radovich et al.
 263 (2004) found from their use of the beta-binomial model that overdispersion
 264 was not significant in their case, and that the binomial model was better
 265 suited to their problem.

266 The binomial procedure can be complemented with information obtained
 267 from an agreement-based approach. Specifically, the positive and negative
 268 Dice indices can be used to provide information on the homogeneity of the
 269 raters' responses. This can then be used to discriminate between products,
 270 and to provide complementary evidence on their differences. As a spin-off, the
 271 approach also provides information on the "quality" of the raters involved.

272 3.2. Using Dice's indices

Firstly in this subsection, we shall consider the relationship between
 Dice's indices and difference tests. For these latter, let p_0 be the guess-
 ing success probability, where $p_0 = 1/2$ for the 2-AFC, paired, and duo-trio
 methods, and $p_0 = 1/3$ for the 3-AFC and triangle methods. If all raters
 are guessing then the marginal proportions of success for the m independent
 raters are the same, i.e., $\rho_1^{(j)} = p_0, j = 1, 2, \dots, m$. Therefore, the positive
 response Dice index is

$$D_1^{(m)} = \frac{mp_0^m}{\sum_{j=1}^m p_0} = p_0^{m-1}.$$

273 Then, the hypothesis tests

$$\begin{aligned} H_0 : D_1^{(m)} &\leq p_0^{m-1} \\ H_1 : D_1^{(m)} &> p_0^{m-1} \end{aligned} \quad (2)$$

274 show whether the raters are discriminating the positive response more than
 275 would be expected by chance. Analogously, for the negative response Dice
 276 index, one will replace p_0 by $1 - p_0$, and test the hypotheses

$$\begin{aligned} H_0 : D_2^{(m)} &\geq (1 - p_0)^{m-1} \\ H_1 : D_2^{(m)} &< (1 - p_0)^{m-1} \end{aligned} \quad (3)$$

277 showing whether the raters fail in the discrimination less than expected by
278 chance. For example, in the triangle test with two panelists, when $D_1^{(2)}$ and
279 $D_2^{(2)}$ are significantly greater than $1/3$ and less than $2/3$, respectively, one
280 can assume that the panelists are not guessing and are actually revealing
281 differences between the products being evaluated.

282 When both hypothesis tests, (2) and (3), are significant, the raters are
283 indeed discriminating products. In this case, $D_1^{(m)}$ must be large and $D_2^{(m)}$
284 must be small. This means that the raters are mostly giving correct responses
285 (i.e., a high degree of agreement is attained), and therefore they are noticing
286 differences between the products. Otherwise, there is no evidence that the
287 raters are properly discriminating products. This procedure is also applicable
288 using bilateral hypothesis tests.

289 Dice indices have several advantages over other measures in studying the
290 agreement for difference tests. Together, positive and negative Dice indices
291 show the consistency of the raters in the two directions, indicating whether
292 the products are different or, on the contrary, are similar. In addition, some
293 other measures can report low overall agreement while the separate agree-
294 ments for both the positive and negative responses are high. For example,
295 the effect of symmetrically unbalanced marginal totals may lead to a low
296 value of Cohen's kappa (see Cicchetti and Feinstein (1990b)). In this case,
297 the wrong conclusion may be drawn that the products are similar when they
298 are actually different, and the Dice indices could have detected any potential
299 differences. Finally, when just one of the positive or negative agreements is
300 low, most indices tend also to be low because they reward symmetry between
301 agreement and disagreement. in contrast, a low negative dice index with a
302 high positive dice index is indicative of major differentiation.

303 In order to perform the hypothesis tests of (2) and (3), one must know
304 the statistical distributions under the null hypotheses. Sometimes it is not
305 possible to derive a closed form expression for a given sampling distribution,
306 and indeed this seems to be the case here. No sampling distribution has as
307 yet been obtained in a closed form relating to the Dice index. However, sam-
308 pling distributions of interest may be estimated by Monte Carlo simulation.
309 Using this technique, it is possible to generate approximations to the true
310 sampling distributions of the test statistics. The precision of the approxi-
311 mation depends strongly on the number of simulations performed. Monte
312 Carlo estimation of sampling distributions is widely used in many practical
313 settings. For example, the IBM SPSS software package offers it as an op-

314 tion when the data does not satisfy the necessary conditions for asymptotic
 315 methods to be used or the samples are so large that the computation time
 316 required is prohibitive. The procedure provides an unbiased estimate of the
 317 exact p-value (see, e.g., Mehta and Patel (2010)).

318 Once the hypothesis tests have been set, the sampling distribution is es-
 319 timated according to the number of raters m , the number of replications k ,
 320 and the number of simulations. Table 1 presents the critical values for several
 321 significance levels in different scenarios. These critical values were obtained
 322 by Monte Carlo estimating the sampling distributions with 1 000 000 simu-
 323 lations. The estimated sampling distributions for $D_1^{(m)}$ are asymmetrically
 324 distributed with right-side tails, whereas those for $D_2^{(m)}$ are approximately
 325 normal. Since one-sided hypothesis tests are considered, attention must be
 326 paid to the right (left) tail for the sampling distribution of $D_1^{(m)}$ ($D_2^{(m)}$). A
 327 reduced table is presented for illustrative purposes.

		$\hat{D}_1^{(m)}$						$\hat{D}_2^{(m)}$			
		$k = 10$	$k = 20$	$k = 30$	$k = 40$			$k = 10$	$k = 20$	$k = 30$	$k = 40$
$m = 2$	0.01	0.8000	0.6667	0.6316	0.5926	$m = 2$	0.99	0.2000	0.3636	0.4324	0.4681
	0.025	0.7500	0.6316	0.5833	0.5517		0.975	0.2857	0.4211	0.4706	0.5000
	0.05	0.6667	0.5882	0.5455	0.5185		0.95	0.3636	0.4615	0.5128	0.5306
	0.1	0.6000	0.5333	0.5000	0.4800		0.9	0.4444	0.5185	0.5455	0.5652
$m = 3$	0.01	0.5455	0.4091	0.3529	0.3243	$m = 3$	0.99	0.0000	0.1579	0.1935	0.2308
	0.025	0.4615	0.3600	0.3103	0.2857		0.975	0.0000	0.1765	0.2308	0.2647
	0.05	0.4000	0.3158	0.2813	0.2553		0.95	0.1579	0.2368	0.2679	0.2917
	0.1	0.3333	0.2727	0.2368	0.2195		0.9	0.1765	0.2647	0.3103	0.3288
$m = 4$	0.01	0.3636	0.2667	0.2105	0.1905	$m = 4$	0.99	0.0000	0.0000	0.0588	0.1132
	0.025	0.3077	0.2105	0.1860	0.1538		0.975	0.0000	0.0784	0.1067	0.1250
	0.05	0.2667	0.1667	0.1538	0.1404		0.95	0.0000	0.0851	0.1463	0.1569
	0.1	0.2222	0.1429	0.1081	0.0889		0.9	0.0000	0.1509	0.1622	0.1887

Table 1: Critical values for the Monte Carlo estimated sampling distributions.

328 The rejection regions for the two tests generally grow with increasing
 329 number of replications and/or number of panelists. This means that less
 330 positive agreement and more negative agreement are necessary to detect sig-
 331 nificant results, and consequently to reveal product differences. For example,
 332 in the triangle test case, rejecting $H_0 : D_1^{(m)} \leq (1/3)^{m-1}$ is difficult when
 333 there are only 10 replications and 2 panelists. With $\alpha = 0.01$, a positive
 334 agreement value greater than 0.8 is needed, and rejecting the null hypoth-
 335 esis becomes easier as the number of replications increases. For $D_2^{(m)}$, the
 336 lowest rejection values are attained with 4 raters and 10 replications. With

337 more replications and fewer panelists it is easier to reject the null hypothesis
338 $H_0 : D_2^{(m)} \geq (2/3)^{m-1}$.

339 In practice, it is not necessary to use a table of critical values and proba-
340 bilities, since the approximated sampling distributions are already available,
341 and can be used to obtain the necessary p-values. Monte Carlo simulations
342 can also be used to calculate the one-sided confidence intervals.

343 It must be remarked that the proposed approach is not valid for overdis-
344 persed binomial data. In such cases, the beta-binomial model should be used
345 instead because the parameter p can have great variability. This possibly
346 extreme variability directly affects the estimates of the agreement indices by
347 yielding large positive and negative Dice indices. With overdispersion, even
348 when there is a clear difference between products, both $D_1^{(m)}$ and $D_2^{(m)}$ are
349 large, leading to misinterpretations if the proposed approach is used. There-
350 fore, before applying this proposed approach, it is advisable to perform an
351 overdispersion test. In Subsection 3.4, the effect of overdispersion is illus-
352 trated with a simulation based example.

353 3.3. *Pairwise comparisons*

354 Besides the information on discrimination provided by the approach, the
355 problem can also be decomposed into $m(m-1)/2$ two-rater problems, i.e.,
356 performing pairwise comparisons for all the raters. These comparisons pro-
357 vide information on the agreement between raters. This can be useful in
358 evaluating the degree of agreement for each panelist relative to the others,
359 and to determine the level of expertise of novice trainee panelists.

360 These Dice indices can be interpreted in a similar way to that of the
361 general problem. Comparison of all the pairwise results together leads to two
362 possibilities – either all the panelists agree in the same way by pairs or they
363 do not. When all the pairwise comparison hypothesis tests are significant,
364 all the panelists are discriminating products in a similar way.

365 On the one hand, if all the panelists agree in the same way, i.e., both $D_1^{(2)}$
366 and $D_2^{(2)}$ are similar for all the pairs then the panelists have the same level
367 of expertise and roughly the same discriminatory reliability. In particular,
368 they all have approximately the same influence on the general agreement and
369 on the differentiation between products. On the other hand, if the panelists
370 agree differently by pairs, one can identify which of them are the sources of
371 the increase or decrease in the general agreement. In Subsection 3.4, we shall
372 present an illustrative example of the interpretation of pairwise comparisons.

373 The proposed framework for pairwise comparisons is particularly helpful
 374 when novice panelists are being trained by an expert. The expert can be
 375 taken as the gold standard, with the pairwise comparisons representing an
 376 objective form of ranking the panelists by efficacy.

377 *3.4. Illustrative examples*

378 In this subsection, we shall present three examples illustrating some typ-
 379 ical scenarios of difference test problems using the triangle test.

380 **Example 1. Binomial data.** We first considered two scenarios in the
 381 binomial model. In one, the proportion of successes used for the simulation
 382 was taken to be $p = 1/3$, corresponding to agreement by chance, so that the
 383 products should not be discriminated. In the other, we took $p = 2/3$, corre-
 384 sponding to the raters being able to properly discriminate between products.
 385 In each case, 20 000 contingency tables were simulated for 3 raters and 20
 386 replications, fitting the beta-binomial model, and estimating its parameters
 387 and the Dice indices, $D_1^{(3)}$ and $D_2^{(3)}$. The p-values were for the hypothesis
 388 tests of the beta-binomial model (see Ennis and Bi (1998)) and of the Dice
 389 indices (see Subsection 3.2) were also calculated. The averages of the param-
 390 eter estimates and the p-values over the 20 000 simulations are presented in
 391 Table 2.

Param.	$\hat{\mu}$	\hat{p} -value	$\hat{\gamma}$	\hat{p} -value	$\hat{D}_1^{(3)}$	\hat{p} -value	$\hat{D}_2^{(3)}$	\hat{p} -value
$\mu = 1/3$	0.3330	0.6238	0.0442	0.7957	0.1042	0.3875	0.4372	0.4997
$\mu = 2/3$	0.6667	0.0006	0.0435	0.7995	0.4363	0.0311	0.1035	0.0314

Table 2: Estimated parameters and p-values for binomial simulated data.

392 The estimates of the parameter μ agree with the pre-set values used to
 393 generate the data with the binomial model, i.e., $p = 1/3$ and $p = 2/3$. The
 394 hypothesis test is non-significant for the case generated with $\mu = 1/3$, and
 395 significant for the case generated with $\mu = 2/3$. The overdispersion tests are
 396 not significant, and the estimates of γ are close to zero in both cases. This
 397 validates the simulation process.

398 When the probability of success is $p = 1/3$, the positive and negative
 399 agreement indices are close to their expected values $1/9$ and $4/9$, respectively.
 400 According to the estimated p -values, these indices are not significant, and
 401 consequently the products can not be considered to have been discriminated.

402 Finally, when $p = 2/3$, both hypothesis tests are significant, indicating that
 403 the agreement is not by chance, and that the raters properly discriminate
 404 the products.

405 Table 3 presents the results for the pairwise comparisons. All the positive
 406 Dice indices for the pairs are similar and close to the expected values $1/3$
 407 and $2/3$, indicating homogeneity among raters for the positive response. The
 408 same is the case for the negative response with respect to the expected values
 409 $2/3$ and $1/3$. Thus, the panelists agree (disagree) and differentiate (do not
 410 differentiate) in the same way.

$\mu = 1/3$	$\hat{D}_1^{(3)}$	\hat{p} -value	$\hat{D}_2^{(3)}$	\hat{p} -value
Rater 1 vs Rater 2	0.3179	0.4966	0.6588	0.4970
Rater 1 vs Rater 3	0.3165	0.4976	0.6592	0.4976
Rater 2 vs Rater 3	0.3138	0.5036	0.6577	0.4961
$\mu = 2/3$	$\hat{D}_1^{(3)}$	\hat{p} -value	$\hat{D}_2^{(3)}$	\hat{p} -value
Rater 1 vs Rater 2	0.6576	0.0450	0.3149	0.0432
Rater 1 vs Rater 3	0.6581	0.0441	0.3163	0.0437
Rater 2 vs Rater 3	0.6577	0.0441	0.3156	0.0434

Table 3: Estimated Dice indices and p-values for pairwise comparisons.

411 **Example 2. Consequences of overdispersion.**

412 As previously observed, the proposed approach is not valid for overdis-
 413 persed binomial data. The following is a simulation-based example to illus-
 414 trate the effects of overdispersion on the Dice indices. The beta-binomial
 415 model is used to generate the data with different levels of overdispersion:
 416 low, medium, and high ($\gamma = 0.2$, $\gamma = 0.5$, and $\gamma = 0.8$). Low and high
 417 success probabilities were also considered ($\mu = 1/3$ and $\mu = 2/3$). Again
 418 20 000 contingency tables were simulated for 3 raters and 20 replications for
 419 the different scenarios. The results are summarized in Table 4.

420 The estimated values for the parameters of the beta-binomial distribution
 421 agree with the ones set beforehand, validating the simulation process. It can
 422 be seen that the Dice indices increase as the overdispersion increases. When
 423 γ is low, the positive and negative indices are closer to their expected values,
 424 whereas, when overdispersion increases, the agreement indices increase too.

425 Table 5 presents the results for the pairwise comparisons. Again, all the
 426 Dice index values increase as the overdispersion increases.

427 The extreme variability distorts any interpretation of the agreement and

Parameters	$\hat{\mu}$	\hat{p} -value	$\hat{\gamma}$	\hat{p} -value	$\hat{D}_1^{(3)}$	\hat{p} -value	$\hat{D}_2^{(3)}$	\hat{p} -value
$\mu = 1/3, \gamma = 0.2$	0.3339	0.7516	0.1847	0.3503	0.2426	0.1993	0.5599	0.7618
$\mu = 2/3, \gamma = 0.2$	0.6668	0.0008	0.1897	0.3388	0.5637	0.0060	0.2474	0.8239
$\mu = 1/3, \gamma = 0.5$	0.3336	0.7571	0.4804	0.0260	0.5015	0.0328	0.7339	0.9668
$\mu = 2/3, \gamma = 0.5$	0.6750	0.0001	0.4784	0.0322	0.7407	0.0002	0.4909	0.3769
$\mu = 1/3, \gamma = 0.8$	0.3329	0.7662	0.7889	0.0005	0.7905	0.0013	0.8961	0.9994
$\mu = 2/3, \gamma = 0.8$	0.6687	0.0000	0.7860	0.0002	0.9018	$1.1 \cdot 10^{-6}$	0.7869	0.9642

Table 4: Estimated parameters and p-values for beta-binomial simulated data.

$\mu = 1/3, \gamma = 0.2$	$\hat{D}_1^{(3)}$	\hat{p} -value	$\hat{D}_2^{(3)}$	\hat{p} -value
Rater 1 vs Rater 2	0.4461	0.2928	0.7237	0.6694
Rater 1 vs Rater 3	0.4438	0.2954	0.7234	0.6701
Rater 2 vs Rater 3	0.4455	0.2916	0.7249	0.6803
$\mu = 2/3, \gamma = 0.2$	$\hat{D}_1^{(3)}$	\hat{p} -value	$\hat{D}_2^{(3)}$	\hat{p} -value
Rater 1 vs Rater 2	0.7261	0.0178	0.4467	0.1478
Rater 1 vs Rater 3	0.7263	0.0180	0.4479	0.1500
Rater 2 vs Rater 3	0.7274	0.0174	0.4500	0.2001
$\mu = 1/3, \gamma = 0.5$	$\hat{D}_1^{(3)}$	\hat{p} -value	$\hat{D}_2^{(3)}$	\hat{p} -value
Rater 1 vs Rater 2	0.6493	0.0782	0.8276	0.8975
Rater 1 vs Rater 3	0.6488	0.0782	0.8273	0.8968
Rater 2 vs Rater 3	0.6495	0.0777	0.8280	0.8972
$\mu = 2/3, \gamma = 0.5$	$\hat{D}_1^{(3)}$	\hat{p} -value	$\hat{D}_2^{(3)}$	\hat{p} -value
Rater 1 vs Rater 2	0.8354	0.0020	0.6574	0.5331
Rater 1 vs Rater 3	0.8333	0.0028	0.6466	0.4962
Rater 2 vs Rater 3	0.8236	0.0025	0.6307	0.5013
$\mu = 1/3, \gamma = 0.8$	$\hat{D}_1^{(3)}$	\hat{p} -value	$\hat{D}_2^{(3)}$	\hat{p} -value
Rater 1 vs Rater 2	0.8576	0.0058	0.9311	0.9923
Rater 1 vs Rater 3	0.8573	0.0060	0.9306	0.9923
Rater 2 vs Rater 3	0.8572	0.0062	0.9312	0.9998
$\mu = 2/3, \gamma = 0.8$	$\hat{D}_1^{(3)}$	\hat{p} -value	$\hat{D}_2^{(3)}$	\hat{p} -value
Rater 1 vs Rater 2	0.9298	0.0001	0.8556	0.8985
Rater 1 vs Rater 3	0.9293	0.0001	0.8522	0.8971
Rater 2 vs Rater 3	0.9314	0.0001	0.8570	0.9067

Table 5: Estimated Dice indices and p-values for pairwise comparisons.

428 the differentiation. Therefore, the proposed approach must only be applied
429 to non-overdispersed binomial data.

430 **Example 3. Detecting non-accurate raters.**

431 The effect of one or more conflictive raters can be analysed in the pair-
 432 wise comparison framework. Table 6 is the contingency table for a scenario
 433 in which one rater disagrees with the other two. Non-overdispersed bino-
 434 mial data are obtained ($\hat{\gamma} = 0.00001$, and $\hat{p} - value = 1$). Using the bino-
 435 mial procedure, the three raters are seen to discriminate between products
 436 ($\hat{p} = 0.6480$, $\hat{p} - value = 2.1 \cdot 10^{-9}$), but it can not be seen which panelists
 437 differentiate and which do not.

				Rater 3		
				A	F	Total
Rater 1	A	Rater 2	A	4	21	25
		F	2	2	4	
		Total	6	23	29	
	F	A	4	1	5	
		F	1	1	2	
		Total	5	2	7	

Table 6: Contingency table.

438 Table 7 presents the general Dice indices and the pairwise comparisons. It
 439 can be observed that the hypothesis tests for these indices are not simultane-
 440 ously significant, indicating that the three raters do not properly discriminate
 441 between products.

	$\hat{D}_1^{(m)}$	\hat{p} -value	$\hat{D}_2^{(m)}$	\hat{p} -value
Raters 1 vs 2 vs 3	0.1714	0.2294	0.0789	0.1118
Rater 1 vs Rater 2	0.8475	0.0003	0.3077	0.0037
Rater 1 vs Rater 3	0.3000	0.5538	0.1250	$8 \cdot 10^{-5}$
Rater 2 vs Rater 3	0.3902	0.3526	0.1936	0.0972

Table 7: Estimated Dice indices and p-values.

442 The hypothesis tests for raters 1 and 2 are significant, denoting that
 443 they are indeed able to discriminate. In contrast, the comparisons between
 444 rater 3 and the other two (1-3 and 2-3) indicate that the agreement is by
 445 chance because the hypothesis tests are not simultaneously significant. It
 446 is apparent that rater 3 is the only one with low agreement, but that that

447 rater's failures decisively affect the general agreement. Raters 1 and 2 are
448 able to differentiate the given products, but rater 3 is only guessing.

449 The following section illustrates the application of the proposed approach
450 in a real context.

451 **4. Application**

452 An experimental study was performed in order to illustrate the potential
453 of the approach in discriminating between two meat products and to evaluate
454 the inter-rater agreement. The triangle test was used with the guidelines
455 defined by the norm International Organization for Standardization (2004b).
456 We shall first describe the experiment.

457 Two Iberian pork loins of the same quality (Iberian pigs fed partly on
458 fodder and partly on mast) were evaluated. The first is a *Carrefour* house-
459 brand pork loin, and the second is produced by a traditional company, *La*
460 *Flor Piornalega*. This variety of pork loin is obtained from free-range Iberian
461 pigs fed on cereals and mast (acorns) and sacrificed at 12 months. The two
462 pieces considered were dry-cured at specialist sites in Spain under very similar
463 conditions of humidity and altitude (Guijuelo for the *Carrefour* product, and
464 Piornal for the *La Flor Piornalega* product). The two loins were tasted by
465 three panelists, and the results analysed by the present proposed approach.

466 The procedure was as follows. A set of three samples was presented
467 simultaneously to each rater, two of them belonging to the same loin. This
468 step was repeated several times with different sets of samples. The raters
469 had to identify which sample was different in each set presented. There were
470 six sessions, and every rater tasted just six sets per session to avoid sensory
471 fatigue. In total, therefore, each rater dealt with 36 sets. All three raters
472 were novices because the objective was to identify whether any differences
473 were noticeable from a regular consumer's point of view.

474 The samples in the sets were displayed uniformly, and all corresponded to
475 the same two pieces of pork loin. The experiment was performed under the
476 same conditions of temperature and lighting in a standardized tasting room.
477 The patterns followed to display the samples were: CPP, PCC, CCP, PPC,
478 CPC y PCP, with C being *Carrefour* and P *La Flor Piornalega*. To tabulate
479 the results, the two possible responses were *A* if they found the different
480 sample, and *F* if they failed. The results are presented in Table 8. Note
481 that the first rater obtained 30 correct responses and only 6 were incorrect,

482 the second rater obtained 28 correct responses and 8 incorrect, and the third
 483 rater obtained 27 correct responses and 9 incorrect.

				Rater 3		
				A	F	Total
Rater 1	A	Rater 2	A	20	4	24
			F	4	2	6
			Total	24	6	30
	F		A	2	2	4
			F	1	1	2
			Total	3	3	6

Table 8: Contingency table for the experimental results.

484 First, we shall approach the discrimination problem by following the stan-
 485 dard methodology. Depending on the properties of the data, there are two
 486 possibilities. If there is no variation among trials then the binomial model
 487 considered in the norm International Organization for Standardization (2004b)
 488 can be applied. Otherwise, one should use the beta-binomial model (see En-
 489 nis and Bi (1998)). In order to choose the model, an overdispersion analysis
 490 is applied. The maximum likelihood estimates for the beta-binomial param-
 491 eters are $\hat{\mu} = 0.8151$ and $\hat{\gamma} = 0.0921$, and the 95% two-sided confidence
 492 intervals are (0.7355, 0.8948) and (0.0000, 0.3222), respectively. The scale
 493 parameter estimate is close to zero, and the corresponding hypothesis test

$$\begin{aligned}
 H_0 : \gamma &= 0 \\
 H_1 : \gamma &\neq 0,
 \end{aligned}
 \tag{4}$$

494 provides a non-significant result with $p - value = 0.3946$. There is no evi-
 495 dence that γ is different from zero, and hence neither of overdispersion being
 496 present. The binomial model can thus be applied.

497 For the binomial model, p must be estimated using the number of correct
 498 differentiations, x_c , and the number of experiments, $k \cdot m = 36 \cdot 3 = 108$. In
 499 this experiment the estimated probability of success is $\hat{p} = 0.8148$, with a
 500 95% two-sided confidence interval equal to (0.7424, 1). The hypothesis test

$$\begin{aligned}
 H_0 : p &= 1/3 \\
 H_1 : p &\neq 1/3,
 \end{aligned}
 \tag{5}$$

501 yields $p - value = 2.2 \cdot 10^{-16}$. Thus, the panelists are not guessing, and they
502 are discriminating between products.

503 The present proposal allows the foregoing information to be comple-
504 mented with other aspects, such as how intense the agreement is with respect
505 to the differentiation, or whether the panelists differentiate the samples in a
506 similar way.

507 With respect to the agreement, it can be observed that the most frequent
508 result is the agreement among the three raters (A, A, A), which occurs 20
509 times out of a total of 36. When there are more than two raters, it becomes
510 more difficult to differentiate all the samples simultaneously. In this case,
511 the three raters found the different sample for the same sets 56% of the time.
512 It is also remarkable that there was only one jointly failed differentiation.
513 The proportion of agreement $\sum_i n_{iii}/n = 21/36 = 0.58$ summarizes this
514 information. This is quite a high proportion result for three raters, but it is
515 interesting to distinguish whether these agreements come from right or from
516 wrong differentiations between the two products.

517 The generalized Dice index of agreement is used to evaluate the overall
518 conditional agreement for each response. The Dice indices are $\hat{D}_1^{(3)} = 0.71$ for
519 the correct responses, and $\hat{D}_2^{(3)} = 0.13$ for the incorrect ones. The 95% one-
520 sided confidence intervals are (0.3158, 0.8182) and (0,0.5400) respectively.
521 These indices lead to the conclusion that the three raters differentiate the two
522 samples quite well, because the positive agreement among them is high and
523 the negative agreement is low. In order to formalize this result, the proposed
524 one-sided hypothesis tests are applied. Monte Carlo simulations were used
525 to generate the Dice index distributions for 3 raters and 36 replications.
526 Figure 1 shows the distribution of the Dice indices under the null hypotheses
527 $H_0 : D_1^{(3)} \leq 1/9$ and $H_0 : D_2^{(3)} \geq 4/9$, respectively.

528 The first hypothesis test provides a $p - value = 0$, i.e., none of the
529 1 000 000 values generated from the statistical distribution is greater than
530 0.71. The second hypothesis test gives a $p - value = 0.000426$. Both tests are
531 significant, indicating that, simultaneously, the positive agreement is greater
532 than expected by chance and the negative agreement is less than expected
533 by chance. This means that differences between the two products are indeed
534 found. This result reinforces that previously obtained with the binomial
535 model. Moreover, this approach yields information about the degree of dis-
536 crimination, which, in this case, is high.

537 A pairwise comparison provides information about whether the raters

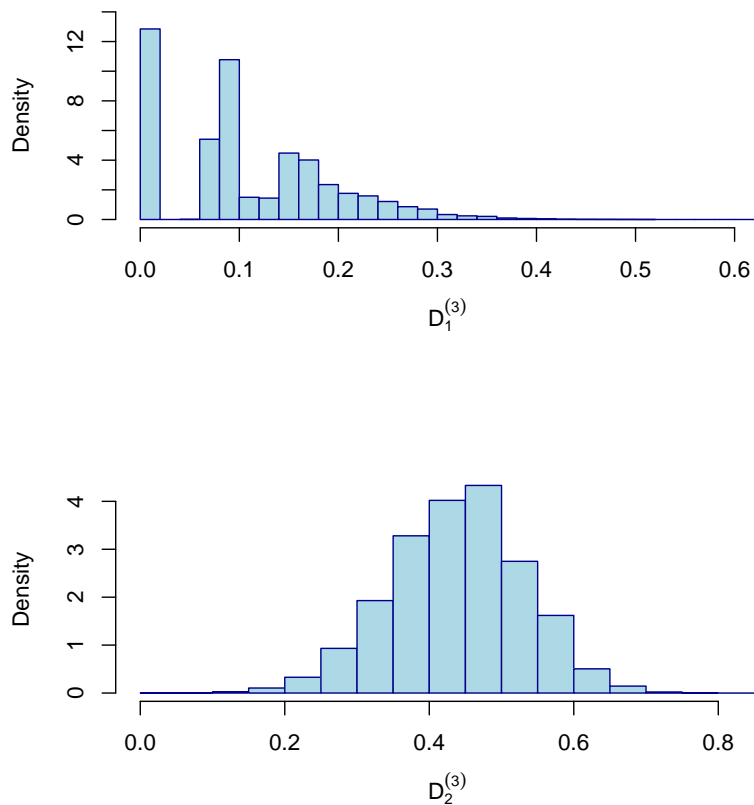


Figure 1: Distribution of Dice indices under the null hypothesis for 3 raters and 36 replications.

538 are discriminating the samples in a similar way or whether one or more of
 539 them are influencing the agreement more. Table 9 presents the pairwise Dice
 540 indices for the three raters. It also gives the p-values according to the tests
 541 defined in Section 3.

542 Note that the positive Dice indices are very similar (from 0.80 to 0.84)
 543 and are high. The indices for the negative response are also similar (from
 544 0.29 to 0.40). This emphasizes that the raters all behave in a similar way,
 545 i.e., there is no rater performing the test in a better or worse way than the
 546 others. The three raters seem to have a similar capability to differentiate,

	$D_1^{(2)}$	p-value	$D_2^{(2)}$	p-value
Rater 1 vs Rater 2	0.83	0.0004	0.40	0.0016
Rater 1 vs Rater 3	0.84	0.0002	0.29	0.0149
Rater 2 vs Rater 3	0.80	0.0005	0.35	0.0049

Table 9: Dice indices and p-values for pairwise comparisons.

547 and they obtained a high degree of differentiation. When performing the
548 hypothesis tests, all the p-values are very small (the largest is 0.0005), re-
549 jecting the possibility that the panelists are guessing when comparing them
550 pairwise. According to these results, the panelists found significant differ-
551 ences between the products studied, both all together and individually. The
552 degree of differentiation was very high.

553 5. Conclusions

554 A novel approach to sensory analysis discrimination tests has been de-
555 scribed. It is based on the generalized positive and negative Dice agree-
556 ment indices, which are used to develop two hypothesis tests. Monte Carlo
557 simulation is used to obtain the distribution under the null hypothesis and
558 the corresponding p-values. The approach provides information on the dis-
559 crimination between products and its strength. Pairwise comparison is used
560 to examine the influence of each rater on the discrimination process. This
561 framework can also be used to train novice panelists by comparing them with
562 experts.

563 The present proposal is not a substitute for the traditional method based
564 on the binomial distribution, but complements it by providing additional in-
565 formation. The applicability of the approach was illustrated by way of some
566 examples, and an experiment was performed using the triangle test scheme to
567 differentiate between two meat products. The results reinforced those given
568 by using the classical binomial model, and showed that the degree of differen-
569 tiation was quite high. Moreover, all the raters were good at discriminating
570 the products, and none was better or worse than the others in this task.

571 The proposed approach is especially interesting when the standardized
572 binomial method is not recommended, i.e., when few raters are involved and
573 replications are needed. It is also recommendable when the interest is on
574 rating novice panelists being trained by an expert, because it allows them to
575 be ranked by skill.

576 **Acknowledgments**

577 We thank Julia Calvarro of the Food Hygiene Area (University of Ex-
578 tremadura) for her assistance in the tasting experiment and Jacinto Martín
579 for his insightful comments. We also thank two anonymous referees for com-
580 ments and suggestions which have improved the content and readability of
581 the paper. This research was partially funded by the *Ministerio de Economía*
582 *y Competitividad, Spain* (Project MTM2011-28983-C03-02), *Gobierno de Ex-*
583 *tremadura, Spain* (Project GRU10110), and *European Union* (European Re-
584 gional Development Funds).

585 **References**

- 586 AGRESTI, A., 1996. An introduction to categorical data analysis. Wiley.
- 587 AJMONE-MARSAN, P., NEGRINI, R., CREPALDI, P., MILANESI, E.,
588 GORNI, C., VALENTINI, A., CICOGNA, M., 2001. Assessing genetic
589 diversity in Italian goat populations using AFLP markers. *Animal Genetics*
590 *32 (5)*, 281–288.
- 591 ANDERSON, D. A., 1988. Some models for overdispersed binomial data.
592 *Australian Journal of Statistics* *30 (2)*, 125–148.
- 593 BAIKAS-NOGUERAS, S., BOVER-CID, S., VECIANA-NOGUS, T.,
594 NUNES, M., VIDAL-CAROU, M., 2003. Development of a quality index
595 method to evaluate freshness in Mediterranean hake (*Merluccius merluc-*
596 *cius*). *Journal of Food Science* *68 (3)*, 1067–1071.
- 597 BI, J., 2011. Similarity tests using forced-choice methods in terms of Thursto-
598 nian discriminial distance, d' . *Journal of Sensory Studies* *26 (2)*, 151–157.
- 599 CICCETTI, D. V., FEINSTEIN, A. R., 1990a. High agreement but low
600 kappa: I. the problem of two paradoxes. *Journal of Clinical Epidemiology*
601 *43 (6)*, 551–558.
- 602 CICCETTI, D. V., FEINSTEIN, A. R., 1990b. High agreement but low
603 kappa: II. resolving the paradoxes. *Journal of Clinical Epidemiology* *43 (6)*,
604 551–558.
- 605 COHEN, J., 1960. A coefficient of agreement for nominal scales. *Educational*
606 *and Psychological Measurement* *20 (1)*, 37–46.

- 607 DUBNICKA, S. R., 2013. A bayesian approach to analyzing replicated pref-
608 erence tests. *Journal of Sensory Studies*.*28 (3)*, 171–187.
- 609 ENNIS, D. M., BI, J., 1998. The beta-binomial model: Accounting for inter-
610 trial variation in replicated difference and preference tests. *Journal of Sen-
611 sory Studies* *13 (4)*, 389–412.
- 612 ENNIS, J., JESIONKA, V., 2011. The power of sensory discrimination meth-
613 ods revisited. *Journal of Sensory Studies* *26 (5)*, 371–382.
- 614 FLEISS, J. L., 1971. Measuring nominal scale agreement among many raters.
615 *Psychological Bulletin* *76 (5)*, 378–382.
- 616 FLEISS, J. L., COHEN, J., EVERITT, B. S., 1969. Large sample standard
617 errors of kappa and weighted kappa. *Psychological Bulletin* *72 (5)*, 323–
618 327.
- 619 FLETCHER, I., MAZZI, M., NUEBLING, M., 2011. When coders are re-
620 liable: The application of three measures to assess inter-rater reliabil-
621 ity/agreement with doctor-patient communication data coded with the
622 VR-codes. *Patient Education and Counseling* *82 (3)*, 341 – 345.
- 623 GARCIA, K., ENNIS, J. M., PRINYAWIWATKUL, W., 2012. A large-scale
624 experimental comparison of the Tetrad and triangle tests in children. *Jour-
625 nal of Sensory Studies* *27 (4)*, 217–222.
- 626 GRAHAM, P., BULL, B., 1998. Approximate standard errors and confidence
627 intervals for indices of positive and negative agreement. *Journal of Clinical
628 Epidemiology* *51 (9)*, 763 – 771.
- 629 GUGGENMOOS-HOLZMANN, I., 2006. How reliable are chance-corrected
630 measures of agreement? *Statistics in Medicine* *12 (23)*, 2191–2205.
- 631 GWET, K., 2002. Kappa statistic is not satisfactory for assessing the extent
632 of agreement between raters. *Statistical Methods For Inter-Rater Reliabil-
633 ity Assessment* *1*, 1–6.
- 634 INTERNATIONAL ORGANIZATION FOR STANDARIZATION.
635 TC34/SC12 - Sensory Analysis.
- 636 INTERNATIONAL ORGANIZATION FOR STANDARIZATION, 2004a.
637 ISO 10399:2004. Sensory analysis. Methodology. Duo-Trio test.

- 638 INTERNATIONAL ORGANIZATION FOR STANDARIZATION, 2004b.
639 ISO 4120:2004. Sensory analysis. Methodology. Triangular test.
- 640 INTERNATIONAL ORGANIZATION FOR STANDARIZATION, 2005.
641 ISO 5495:2005. Sensory analysis. Methodology. Paired comparison test.
- 642 JENSCHKE, B. E., HODGEN, J. M., MEISINGER, J. L., HAMLING, A. E.,
643 MOSS, D. A., AHNSTRM, M. L., ESKRIDGE, K. M., CALKINS, C. R.,
644 2007. Unsaturated fatty acids and sodium affect the liver-like off-flavor in
645 cooked beef. *Journal of Animal Science* 85 (11), 3072–3078.
- 646 KUNERT, J., MEYNER, M., 1999. On the triangle test with replications.
647 *Food Quality and Preference* 10 (6), 477 – 482.
- 648 LAPARA, T. M., NAKATSU, C. H., PANTEA, L. M., ALLEMAN, J. E.,
649 2002. Stability of the bacterial communities supported by a seven-stage
650 biological process treating pharmaceutical wastewater as revealed by PCR-
651 DGGE. *Water Research* 36 (3), 638 – 646.
- 652 LIGGET, R. E., DELWICHE, J. F., 2005. The beta-binomial model: vari-
653 ability in overdispersion across methods and over time. *Journal of Sensory*
654 *Studies* 20 (1), 48–61.
- 655 LORENZO, J. M., GARCÍA FONTÁN, M. C., FRANCO, I., CARBALLO,
656 J., 2008. Biochemical characteristics of dry-cured lacón (a Spanish tradi-
657 tional meat product) throughout the manufacture, and sensorial properties
658 of the final product. Effect of some additives. *Food Control* 19 (12), 1148
659 – 1158.
- 660 MACKINNON, A., 2000. A spreadsheet for the calculation of comprehensive
661 statistics for the assessment of diagnostic tests and inter-rater agreement.
662 *Computers in Biology and Medicine* 30 (3), 127–134.
- 663 MANLY, B. F. J., 1997. Randomization, Bootstrap and Monte Carlo Method
664 in Biology. Chapman and Hall.
- 665 MARTÍN, A., BENITO, M., ARANDA, E., RUIZ-MOYANO, S.,
666 CÓRDOBA, J., CÓRDOBA, M., 2010. Characterization by volatile com-
667 pounds of microbial deep spoilage in Iberian dry-cured ham. *Journal of*
668 *Food Science* 75 (6), 360–365.

- 669 MEHTA, C. R., PATEL, N. R., 2010. IBM SPSS Exact Tests. IBM.
- 670 MEYNER, M., BROCKHOFF, P., 2003. The design of replicated difference
671 tests. *Journal of Sensory Studies* 18 (4), 291–324.
- 672 MOUNCHILI, A., WICHTEL, J., BOSSET, J., DOHOO, I., IMHOF, M.,
673 ALTIERI, D., MALLIA, S., STRYHN, H., 2005. HS-SPME gas chromatographic
674 characterization of volatile compounds in milk tainted with off-
675 flavour. *International Dairy Journal* 15 (12), 1203 – 1215.
- 676 PONS-SANCHEZ-CASCADO, S., VIDAL-CAROU, M. C., NUNES, M. L.,
677 VECIANA-NOGUES, M. T., 2006. Sensory analysis to assess the freshness
678 of Mediterranean anchovies (*Engraulis encrasicolus*) stored in ice. *Food
679 Control* 17 (7), 564 – 569.
- 680 RADOVICH, T. J. K., KLEINHENZ, M. D., DELWICHE, J. F., LIGGETT,
681 R. E., 2004. Triangle tests indicate that irrigation timing affects fresh cab-
682 bage sensory quality. *Food Quality and Preference* 15 (5), 471 – 476.
- 683 RAMOS-GUAJARDO, A., GONZÁLEZ-RODRÍGUEZ, G., 2011. Hypoth-
684 esis testing with fuzzy data: An application to quality control of cheese.
685 In: 2011 11th International Conference on Intelligent Systems Design and
686 Applications (ISDA). pp. 1335 –1340.
- 687 SEVERIANO, A., CARRICO, J. A., ROBINSON, D. A., RAMIREZ, M.,
688 PINTO, F. R., 2011. Evaluation of jackknife and bootstrap for defin-
689 ing confidence intervals for pairwise agreement measures. *PloS one* 6 (5),
690 e19539.
- 691 SHOUKRI, M., 2004. Measures of interobserver agreement. Chapman and
692 Hall.
- 693 SILVA, F., DUARTE, M., CAVALCANTI-MATA, M. E., 2010. Nova
694 metodologia para interpretação de dados de análise sensorial de alimen-
695 tos. *Engenharia Agrícola* 30 (5), 967–973.
- 696 VAN HOUT, D., HAUTUS, M. J., LEE, H.-S., 2011. Investigation of test
697 performance over repeated sessions using signal detection theory: Com-
698 parison of three nonattribute-specified difference tests 2-AFCR, A-NOT A
699 and 2-AFC. *Journal of Sensory Studies* 26 (5), 311–321.

- 700 VON EYE, A., MUN, E., 2005. Analyzing rater agreement: Manifest variable
701 methods. Lawrence Erlbaum Associates.
- 702 WARRENS, M. J., 2008. Similarity coefficients for binary data : properties
703 of coefficients, coefficient matrices, multi-way metrics and multivariate co-
704 efficients. Ph.D. thesis, Leiden University.
- 705 WU, H., CHEN, L., 1995. Sensory analysis in quality control—the agreement
706 among raters. *Botanical bulletin of Academia Sinica* 36, 121–133.