

1 True and error analysis instead of test of correlated
2 proportions: Can we save lexicographic semiorder
3 models with error theory?

4 Michael H. Birnbaum¹

5 ¹California State University, Fullerton

6 ¹mbirnbaum@fullerton.edu

7 May 4, 2022

8 **Abstract**

9 This article criticizes conclusions drawn from the standard test of correlated proportions
10 when the dependent measure contains error. It presents a tutorial on a new method of
11 analysis that uses a fairly general error model called the true and error model of choice. This
12 method allows the investigator to separate measurement of error from substantive conclusions
13 about effects of the independent variable but it requires replicated measures of the dependent
14 variable. The method is illustrated with hypothetical examples and with empirical data from
15 tests of Lexicographic semiorder (LS) models as descriptive models of risky decision making.
16 LS models imply a property known as interactive independence. Data from two previous
17 studies are re-analyzed to test interactive independence. The new analyses yielded clear
18 answers: interactive independence can be rejected; therefore, lexicographic semiorders can
19 be rejected as descriptive models, even if a flexible error model is allowed. The new methods
20 of analysis can be applied to situations in which the test of correlated proportions has been
21 used in the past, where it is possible to obtain replicated measures.

22 **keywords**

23 test of correlated proportions; true and error theory; choice theory; lexicographic semiorder;
24 risky decision making

25 **acknowledgments**

26 Thanks are due to Julien Rouvere for helpful suggestions on the manuscript.

27 **running head**

28 True and error analysis

1 Introduction

This paper criticizes conclusions drawn from a statistical test that has been and continues to be widely used in psychology and economics, and it presents new methods that can address the criticisms. The test of correlated proportions (McNemar, 1947) is used to test the statistical significance of a difference between response proportions obtained in a within-subjects design. The new methods are based on models known as true and error (TE) models, which are analogous to, but not the same as, models used in classical test theory (Novick, 1966; Spearman, 1904). These models extended models of Lichtenstein and Slovic (1971), combined with constraints provided by replications (Birnbaum, 2004, p. 59-60). This paper presents new techniques developed and refined in recent articles (Birnbaum, 2013; 2019; Birnbaum & Bahra, 2012a, 2012b; Birnbaum & Wan, 2020; Birnbaum, Schmidt, & Schneider, 2017; Birnbaum & Quispe-Torreblanca, 2018).

Examples are presented to show how the new methods can lead to different conclusions from those reached by older ones. Hypothetical data show how the test of correlated proportions can be significant and yet the results can be attributed to random error, and how two proportions can be exactly equal so the test will be nonsignificant, and yet there is a significant difference between conditions when analyzed by deeper methods.

The following is a classic method to compare rival theories: One theory implies that two situations are equivalent and the other implies that there is a systematic difference. The experimenter manipulates situations as an independent variable and measures responses as a dependent variable. Suppose Conditions 1 and 2 of Table 1 are two situations that are theoretically equivalent, and the two possible responses of the dependent variable are S and R . The entries in Table 1 represent hypothetical frequencies of the responses in the two conditions.

Suppose the hypothetical data in Table 1 came from a between-subjects experiment in

Table 1: Hypothetical data for a test between two conditions

Independent Variable	Dependent Variable		Row Totals
	<i>S</i>	<i>R</i>	
Condition 1	65	35	100
Condition 2	35	65	100
Column Totals	100	100	

Note: The $\chi^2(1) = 18, p < 0.01$.

54 which there were 200 participants, 100 randomly assigned to each condition. Table 1 shows
 55 that in Condition 1, 65 of 100 participants responded *S*, whereas in Condition 2, only 35
 56 made this response. The Fisher exact test (for small *n*) or the standard Chi-Square test
 57 of independence can be used to assess whether data in a table like this are likely to have
 58 occurred given the null hypothesis that the probability to respond *S* is the same in both
 59 conditions. In this case, $\chi^2(1) = 18, p < 0.01$, so an experimenter would reject theories that
 60 implied no difference in response probabilities between these conditions in favor of theories
 61 that would allow these results. (Upper and lower case, $P(S)$ and $p(S)$ are used here to denote
 62 the obtained proportion and inferred probability of an observed response, respectively.)

63 Now suppose the data in Table 1 arose from a within-subjects experiment in which 100
 64 participants experienced both Conditions 1 and 2 (with suitable counterbalancing). The
 65 analysis of within-subjects data is a bit more complicated, because it involves not only
 66 the marginal response proportions in the two conditions, but also the correlation (non-
 67 independence, or contingency) between the responses by the same people in the two condi-
 68 tions. The test of correlated proportions (McNemar, 1947), developed for this situation, is
 69 described in the next section.

Table 2: Hypothetical data for a within-subjects test between two conditions

Response in Condition 1	Response in Condition 2		Row Totals
	S'	R'	
S	29	36	65
R	6	29	35
Column Totals	35	65	100

Note: The test of correlated proportions compares equality of frequencies of SR' against RS' ; i.e., 36 versus 6.

70 1.1 Test of correlated proportions

71 Table 2 is a cross-tabulation that reveals the contingency between responses by the same
72 people in the two conditions. The row and column sums of Table 2 are the same as the row
73 entries of Table 1. To distinguish responses in the two conditions, let S' and R' designate the
74 responses in Condition 2 corresponding to S and R of Condition 1, respectively. If people
75 responded S if and only if S' , then off-diagonal entries would be zero. In Table 1, responses
76 are not perfectly correlated, nor are responses in the two conditions independent, which
77 would require that $p(SS') = p(S)p(S')$, where $p(SS')$ is the probability of the conjunction.¹
78 Instead, responses are positively correlated. A majority made the same responses in both
79 conditions ($29 + 29 = 58$) but 36 people switched from S to R' and 6 switched in the opposite
80 direction.

81 The test of correlated proportions tests the hypothesis, H, that the probability of re-
82 sponding S in the first condition is the same as the probability of responding S' in the sec-
83 ond condition; i.e., H: $p(S) = p(S')$. The proportions are certainly different, since $P(S) =$
84 $0.65 \neq 0.35 = P(S')$. Asking if marginal probabilities are equal is equivalent to asking if
85 the two types of response reversals are equally probable (McNemar, 1947); this equivalence
86 holds for both observed proportions and probabilities: $p(S) = p(S') \Leftrightarrow p(SR') = p(RS')$.

¹In Table 2, $P(SS') = 0.29 \neq P(S)P(S') = (0.65)(0.35) = 0.2275$.

87 In other words, we can ignore cases where the participant made the same responses in both
88 conditions and examine only cases where a person switched responses. The hypothesis that
89 marginal response probabilities are equal, $H: p(R) = p(R')$, that an equal number switch in
90 either direction, can be tested for these data by a binomial distribution with $n = 36 + 6 =$
91 42 trials, where $p = 0.50$. The binomial probability to observe 36 or more SR' reversals out
92 of 42 preference reversals is about one in a million, so we would reject the hypothesis H that
93 the marginal response probabilities are equal.²

94 As n grows large, the binomial can be approximated by the normal distribution and
95 one can compare a calculated z value with the standard normal distribution. The mean and
96 standard deviation for a binomial are $\mu = np$ and $\sigma = \sqrt{np(1-p)}$. With $p = 0.5$ and $n = 42$
97 for Table 1, $\mu = 21$ and $\sigma = 3.24$, so $z = (36 - 21)/3.24 = 4.63$, an extremely improbable
98 value, leading to the same conclusion as the binomial calculation. This standard formula
99 for z is often called "Conlisk's z -test" in the economics literature after Conlisk (1989); it is
100 equivalent to McNemar's (1947) Chi-Square test. Whether calculated by exact binomial, by
101 the normal approximation (Conlisk's z test), or via the equivalent Chi-Square (McNemar,
102 1947; Lichtenstein & Slovic, 1971), these calculations are all tests of correlated proportions.

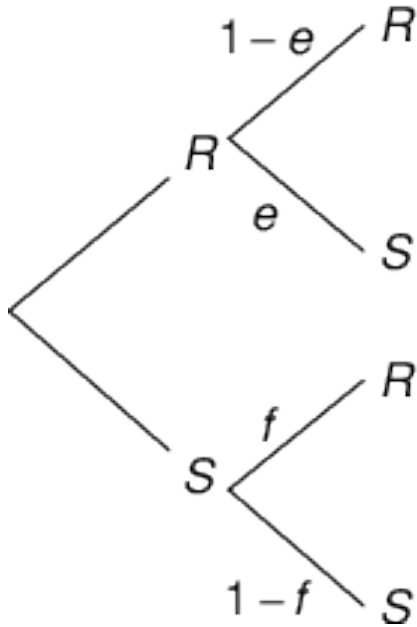
103 Note that in this example, the marginal proportion, $P(S) = 0.65$, is significantly greater
104 than 0.5 by a binomial test, but the marginal proportion, $P(S') = 0.35$, is significantly less
105 than 0.5. This case seems a strong one for concluding that the response probabilities, $p(S)$
106 and $p(S')$, are not equal.

107 A person applying these methods for Table 2 can conclude that we should reject hypoth-
108 esis H that the response probabilities are the same and might *therefore reject a theory that*
109 *true preferences are equivalent and the observed results are due to random response errors.*
110 However, the last part of this argument, in italics, does not follow if we allow a plausible

²The binomial calculation assumes that participants respond independently of each other, which is not controversial when people are tested separately.

Choice Problem 1

True Error Response



Choice Problem 2

True Error Response

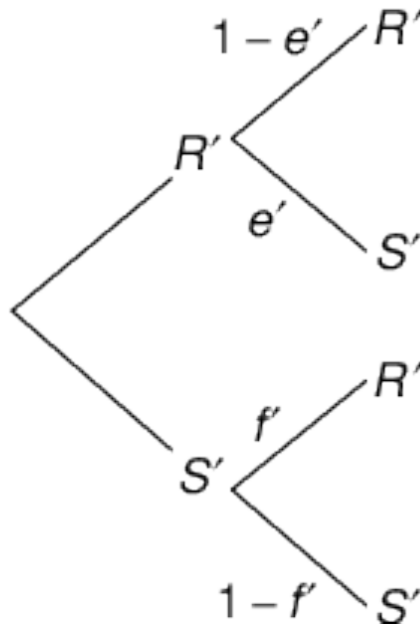


Figure 1: True and Error models for two choice problems: In TE4, all four error terms are free; TE2, assumes $e = f$ and $e' = f'$; TE1 assumes $e = f = e' = f'$. After Birnbaum & Quispe-Torreblanca (2018).

111 theory of error to intervene between true preferences and observed responses in the depen-
 112 dent measure. The next section shows that the results in Table 2 are compatible with the
 113 null hypothesis, H_0 , that the two conditions produced the same true preferences, and that
 114 random errors are responsible for the observed difference in response proportions.³

³The test of correlated proportions was discussed by Lichtenstein and Slovic (1971) and by Conlisk (1989). Although these authors had acknowledged limitations of the test, it became the standard method for analyzing paradoxes of choice in both psychology and economics. For example, a recent review by Blavatsky, Ortmann, & Panchenko (2022) summarizes strength and direction of evidence regarding the Allais paradox in terms of Conlisk z values from 81 experiments. As will be shown here, significant z values do not rule out the theory that the "paradox" is produced by random error.

Table 3: Implications of Null Hypothesis, $H_0: p_{SR'} = p_{RS'} = 0$.

Response		Responses in Problem 2	
Problem 1	S'		R'
S	$p_{RR'}(e)(e') + p_{SS'}(1-f)(1-f')$		$p_{RR'}(e)(1-e') + p_{SS'}(1-f)(f')$
R	$p_{RR'}(1-e)(e') + p_{SS'}(f)(1-f')$		$p_{RR'}(1-e)(1-e') + p_{SS'}(f)(f')$

Note: $p(R) = p_{RR'}(1-e) + p_{SS'}(f)$; $p(R') = p_{RR'}(1-e') + p_{SS'}(f')$.

1.2 True and Error Models of Choice

The true and error model in Figure 1 was developed in the context of choice theory, where the dependent measure is a choice response by a person who is asked to read descriptions of alternatives and to state which of the two alternatives she or he would prefer. For example, "would you rather have $S = \$45$ for sure or $R =$ a fifty-fifty gamble to win either \$100 or \$1?" Such choices are known as decisions under risk, because the consequences and probabilities are known to the decision maker. In this literature, the notations, S and R , are often used to designate "safe" and "risky" options, respectively.

When presented on multiple occasions with the same choice problem, the same person does not always make the same response. Humans might make errors; i.e., a person might truly prefer $S \succ R$ (where \succ denotes "truly preferred to"), and yet the person responds " R " or vice versa. How might people make errors in such an experiment? A person might mis-read the problem, might mis-remember or mis-aggregate the information, might mis-remember her or his evaluations or decisions, or might accidentally push the wrong button to signal the response.

For the rest of this article, the examples will refer to preferences between risky prospects, as in the experiments reanalyzed in this paper, but the reader should keep in mind that the methods described here are also applicable to many other situations in which the dependent measure of an experiment contains error as in Figure 1.

134 In research on risky decision making, the "conditions" are related choice problems de-
 135 signed to be equivalent, according to one theory of human decision making and expected
 136 to differ, according to a rival theory. That is, one can deduce from one theory that $S \succ R$
 137 if and only if $S' \succ R'$, where \succ denotes "is truly preferred to". This theory implies that
 138 except for error, a person should prefer S and S' in Conditions 1 and 2, or prefer R and R' ,
 139 respectively. That is, if the responses contained no error, all of the data in Table 2 would
 140 fall on the diagonal.

141 Figure 1 depicts possible errors in two choice problems. In Choice Problem 1 (left side of
 142 Figure 1), if a person truly prefers R , she or he might erroneously respond S with probability
 143 e . If the person truly prefers S , he or she might respond R with probability f . In Choice
 144 Problem 2 (right), the corresponding errors occur with probabilities e' and f' , respectively.
 145 The errors are assumed to be mutually independent and to have probabilities less than $1/2$.

146 Let p_S denote the probability that a person truly prefers S , which is distinguished from
 147 $p(S)$, the probability that a person responds "S". In general, a person might have any of
 148 four true preference patterns: SS' , SR' , RS' , or RR' , which have probabilities of $p_{SS'}$, $p_{SR'}$,
 149 $p_{RS'}$, and $p_{RR'}$, respectively.

150 According to H0: $p_{SR'} = p_{RS'} = 0$, no person ever has opposite true preferences in the
 151 two choice problems. This definition is not the same as H, which is that $p(SR') = p(RS')$,
 152 that the probabilities of the two types of observed preference reversals are equal. Assuming
 153 H0, it follows that the probabilities that a person would show each response pattern are as
 154 given in Table 3. In other words, any off-diagonal entry is due to error, according to H0.
 155 For the rest of this paper, "H0" will refer to this null hypothesis, which is different from H,
 156 which is the null hypothesis of the test of correlated proportions.

157 Table 3 shows that H0 does not imply H, nor does H imply H0: According to H0, the
 158 probability of the two types of response reversals need not equal each other. For example,
 159 if $p_{RR'} = p_{SS'} = 0.5$; $e' = f = 0.1$, and $e = f' = 0.4$, then the null hypothesis, H0, is

160 compatible with the data of Table 2; one can reproduce the frequencies in Table 2 from the
161 null hypothesis in Table 3 using these parameters. Thus, H0 can be satisfied and H violated.

162 Therefore, no one should reject H0 based on rejection of H in the test of correlated propor-
163 tions. Similarly, this example also shows that simply because one proportion is significantly
164 greater than 0.5 and the other is significantly less than 0.5, one cannot reject H0 that the
165 two experimental conditions induced the same true preferences, because these results can
166 also be reproduced using the same parameters.

167 Although H0 (Table 3) is perfectly compatible with Table 2, other theories are also
168 compatible with those data, including H1, the theory that $e = f = e' = f' = 0.1$, $p_{SS'} =$
169 $p_{RR'} = 0.313$, $p_{SR'} = 0.375$, and $p_{RS'} = 0$. Indeed, the values in Table 2 can be reproduced
170 by many other such theories in which H0 is false. If we knew by some method (but not
171 simply by assumption or faith) that all error rates are equal, then the data in Table 2 would
172 indicate a violation of H0. The data in Table 2 have only three degrees of freedom (since
173 the four entries sum to the number of participants), and the model of Figure 1 allows 7
174 parameters: $e, e', f, f', p_{SS'}, p_{SR'}$, and $p_{RS'}$, so there are many possible solutions. In other
175 words, there are multiple ways to describe the data and one cannot determine which of them
176 is more likely true.

177 It should therefore be clear that with the experimental design as in Table 2 and the test
178 of correlated proportions, we cannot properly test H0 and therefore cannot answer questions
179 we wish to address. Fortunately, we can estimate errors and test theories, if we do a better
180 experiment that includes replications and we analyze the pattern information in the data,
181 as shown in the next two sections.

182 **1.3 Estimating Error from Replications**

183 One can estimate error rates from variation of response by the same person to the same
184 choice problem in the same brief session (Birnbaum, 2004, p. 59-60). To replicate, one

Table 4: TE analysis of replication of a single choice problem.

Response in		Responses in Replicate 2	
Replicate 1	S	R	
S	$p_R(e)(e) + (1 - p_R)(1 - f)(1 - f)$	$p_R(e)(1 - e) + (1 - p_R)(1 - f)(f)$	
R	$p_R(1 - e)(e) + (1 - p_R)(f)(1 - f)$	$p_R(1 - e)(1 - e) + (1 - p_R)(f)(f)$	

Note: $p(R) = p_R(1 - e) + (1 - p_R)(f)$

185 presents each choice problem twice to each participant, suitably separated, counterbalanced,
 186 and embedded among other choice problems.

187 In the simplest design for *individual* true and error theory (*i*TET), one individual serves
 188 in many sessions, and within each session, each choice problem is replicated twice (Birnbaum
 189 & Bahra, 2012a, 2012b). In the simplest design for *group* true and error theory (*g*TET), each
 190 of many participants serve in one session each, and each choice problem is replicated twice
 191 in the session. The key assumption in either form of TE model is that preference reversals
 192 to the same choice problem by the same participant in the same brief experimental session
 193 are due to random error.

194 In studies of an individual, *i*TET allows that the person may have different true pref-
 195 erences over time. This theory is modelled by the assumption that the person may have
 196 different preferences in different sessions but the same preferences hold within a brief session
 197 (Birnbaum & Wan, 2020). In studies of group data, different people may have different true
 198 preferences. The reanalysis of studies in this paper are cases of *g*TET.

199 Suppose we present one choice problem (e.g., S versus R) twice to the same participants,
 200 suitably embedded among many other trials. A participant can have four possible response
 201 patterns (combinations of expressed preferences) for these two replicated trials: The person
 202 can respond S or R on both occasions (SS or RR response patterns), or can make a reversal
 203 of preferences (SR or RS patterns) between the replicates. According to the TE model of
 204 Figure 1, the probabilities of the four patterns (for Choice Problem 1 replicated) are as given

205 in Table 4.

206 Table 4 shows that responses to replicated choice problems are not expected to be inde-
207 pendent, but instead, there will likely be a positive correlation in which the entries on the
208 diagonal will be more probable, and the off-diagonals will be less probable and equal. Note
209 that the four cells of Table 4 constrain 3 parameters, so we have gained constraint relative to
210 Table 3. But there is even more information available, if we replicate both choice problems
211 and use the information from all 16 response patterns.

212 With two choice problems and two replications each, there are 2 by 2 by 2 by 2 = 16
213 possible response patterns. These pattern data provide not only the constraints required to
214 estimate the parameters, but also to test the model. From the relative frequencies of the 16
215 response patterns, which have 15 degrees of freedom (df), one can estimate four error rates,
216 four probabilities of true preference patterns (which sum to 1 and thus consume 3 df), and
217 there remain 8 df to test the TE model. One can then test H_0 as a special case of the TE
218 model, because it has two fewer parameters, since $p_{SR'} = p_{RS'} = 0$.

219 Tables 5, 6, and 7 contain hypothetical examples of such arrays, in which the row and
220 column marginal sums match the frequencies of Table 2, where H is rejected. However,
221 these examples illustrate cases in which the null hypothesis, H_0 ($p_{SR'} = p_{RS'} = 0$) should be
222 rejected (Table 5), where H_0 can be retained (Table 6), and where the TE model itself can
223 be rejected (Table 7). Table 8 contains hypothetical data in which H is satisfied perfectly
224 and yet TE analysis indicates that H_0 should be rejected.

225 The next section is a tutorial on TE methods, showing how these hypothetical cases are
226 analyzed to reach these conclusions.

227 2 True and Error Analysis

228 According to the TE model of Figure 1, the probability to show the SR' response pattern
229 on two replications, denoted SR', SR' , is as follows:

$$\begin{aligned} p(SR', SR') = p_{SS'}(1 - e)^2(e')^2 + p_{SR'}(1 - e)^2(1 - f')^2 + \\ p_{RS'}(f)^2(e')^2 + p_{RR'}(f)^2(1 - f')^2 \end{aligned} \quad (1)$$

230 where $p(SR', SR')$ is the theoretical probability to observe SR' response pattern on both
231 replications; $p_{SS'}$, $p_{SR'}$, $p_{RS'}$, and $p_{RR'}$, are the probabilities of the four possible true prefer-
232 ence patterns; and the error rates, e , f , e' , and f' , are as defined in Figure 1.

233 Note that in each of the four possible true preference states, there is a pattern of errors
234 that could produce each possible observed response pattern. For example, if a person has the
235 true preference pattern SS' , then that person can respond SR', SR' (SR' on two replications)
236 by making no error on the two presentations of the choice between S and R and by making
237 errors on both presentations of the choice between S' and R' .

238 There are 16 equations (including Equation 1) for the 16 possible response patterns. The
239 16 corresponding observed frequencies (counts) of these response patterns have 15 degrees of
240 freedom (df), because the 16 frequencies sum to the total number of response patterns. In
241 g TET with two replicates in one session, this total is the number of participants; in i TET,
242 where one individual served in a number of sessions, it is the number of sessions for the
243 individual.

244 2.1 Fitting TE Models

245 The free, open-source program, TEMAP2.R, can be used to perform statistical analysis to
246 fit and test the six models.⁴ The program analyzes crosstabulation tables like Tables 5, 6, 7,

⁴TEMAP2.R is freely available in the online supplement to Birnbaum & Quispe-Torreblanca (2018); the URL is:

Table 5: Case 1: Hypothetical Frequencies of Response Patterns; H0 Rejected.

Responses in Replicate 1	Responses in Replicate 2				Total
	SS'	SR'	RS'	RR'	
SS'	21	5	2	1	29
SR'	5	25	1	5	36
RS'	2	1	1	2	6
RR'	1	5	2	21	29

Note: TE fit: $G(8) = 0.42$; TE+H0: $G(10) = 25.44$, H0: $G(2) = 25.02$.

Table 6: Case 2: Hypothetical Frequencies Satisfying H0.

Responses in Replicate 1	Responses in Replicate 2				Total
	SS'	SR'	RS'	RR'	
SS'	15	10	2	2	29
SR'	10	14	2	10	36
RS'	2	2	0	2	6
RR'	2	10	2	15	29

Note: TE: $G(8) = 0.47$; TE+H0: $G(10) = 0.58$; H0: $G(2) = 0.11$.

247 and 8. The program estimates parameters to minimize either the standard χ^2 index of fit or
 248 the G index (sometimes called G^2), which is equivalent to a maximum likelihood solution.⁵

$$G = 2 \sum \sum O_{ij} \ln (O_{ij}/E_{ij}) \quad (2)$$

249 where the summation is over the 16 cells, O_{ij} is the observed frequency (count) in Row i
 250 and Column j , E_{ij} is the corresponding "expected" ("predicted" or "fitted") frequency in
 251 the cell according to the particular TE model.

<http://journal.sjdm.org/vol13.5.html>

⁵Programming for Bayesian analysis of true and error models has been presented by Lee (2018) and by Schramm (2020). In cases studied so far, Bayesian and classical statistical analyses have led to similar solutions and conclusions, although some caution is required for the interpretation of Bayesian posterior probabilities for these nested models (Lee, 2018; Birnbaum, 2019).

252 Each of the 16 "expected" (aka, "fitted" or "predicted") frequencies, E_{ij} , is based on the
253 "best-fit" parameter values estimated from the data to minimize G . Each predicted value
254 is equal to the number of participants in a group analysis, n , multiplied by the model's
255 calculated probability (as in Equation 1).

256 The G index is similar to χ^2 and is asymptotically Chi-Square distributed. Because
257 there are 15 df in the data matrix (which sums to the number of participants), there are
258 $15 - 4 - 3 = 8$ degrees of freedom in the test of the TE model.

259 2.2 Testing Special Cases

260 Because the null hypothesis ($H_0: p_{SR'} = p_{RS'} = 0$) is a special case of TE in which 2 fewer df
261 are consumed, the difference in fit between the TE model and its corresponding H_0 special
262 case is asymptotically Chi-Square distributed with 2 df. That is, we calculate the fit with all
263 parameters free (TE fit) and the fit with the constraint that these parameters are fixed to
264 zero (TE+ H_0), and compute the difference, $G(2) = G(8) - G(10)$, which is the test of H_0 .

265 TEMAP2.R can also be applied in cases with relatively small samples where one might be
266 concerned of the applicability of the asymptotic Chi-Square distribution for G . The program
267 employs Monte Carlo simulation to construct sampling distributions of the test statistics,
268 and it uses bootstrapping to estimate confidence intervals for the fitted parameters.

269 The TE model in Figure 1 is denoted TE4 because there are 4 different error rates. A
270 special case of this model, TE2, assumes $e = f$ and $e' = f'$, and a further special case,
271 TE1, assumes that $e = e' = f = f'$.⁶ As shown here, TE4 can produce data as in Table 2,
272 even though H_0 is true, but TE2 cannot reconcile Table 2 with H_0 . The null hypothesis,
273 $H_0: p_{SR'} = p_{RS'} = 0$ is a special case of each of these three TE models. Birnbaum (2019)
274 presented a figure to show the nesting relationships among these six possible models.

⁶TE1 is similar to the model in Conlisk (1989) that might justify the test of correlated proportions; TE1 is sometimes called the "trembling hand" model. However, TE1 does not necessarily yield the same conclusions as the test of correlated proportions, since it can reject H_0 when H can be accepted.

Table 7: Case 3: Hypothetical Frequencies Violating TE model.

Responses in		Responses in Replicate 2			
Replicate 1	SS'	SR'	RS'	RR'	Total
SS'	10	6	3	10	29
SR'	16	9	1	10	36
RS'	1	1	1	3	6
RR'	2	20	1	6	29

Note: TE fit: $G(8) = 35.20$.

275 2.3 Hypothetical Examples

276 The values in Table 5 were constructed from the TE1 model with the assumptions that
 277 $e = f = e' = f' = 0.1$, $p_{SS'} = p_{RR'} = 0.313$, $p_{SR'} = 0.375$, and $p_{RS'} = 0$. That is, these data
 278 were constructed from the assumption that H_0 is false. The hypothetical values in Table 5
 279 are based on $n = 100$, rounded or adjusted to nearby integers, to produce the same row and
 280 column marginal totals as in Table 2.

281 When Table 5 is fit to the TE4 model (with all parameters free) using TEMAP2.R, the
 282 TE model fits well, $G(8) = 0.42$, as expected, since the data were constructed from a special
 283 case of this model, and the best-fit estimates were approximately those used to generate the
 284 data. However, when $p_{SR'}$ and $p_{RS'}$ were fixed to zero, the constrained TE model (TE4+ H_0)
 285 does not fit well, $G(10) = 25.44$, so the test of H_0 yields, $G(10) - G(8) = 25.02$, which far
 286 exceeds 9.3, which is the critical value of $\chi^2(2)$ with $\alpha = 0.01$. Therefore, if we observed
 287 real data as in Table 5, these TE analyses would lead to rejection of H_0 ($p_{SR'} = p_{RS'} = 0$).
 288 Monte Carlo simulations of the sampling distributions also agree with the conclusions that
 289 TE can be retained and that H_0 can be rejected. Based on 10,000 bootstrapping samples
 290 from Table 5, the 95% confidence interval for $p_{SR'}$ is estimated to range from 0.19 to 0.50,
 291 indicating one can be confident of a substantial violation of H_0 with $p_{SR'} > 0$.

Table 8: Case 4: Marginals Satisfy Test of Correlated Proportions and H0 Rejected.

Responses in Replicate 1	Responses in Replicate 2				Total
	SS'	SR'	RS'	RR'	
SS'	17	4	7	2	30
SR'	4	10	2	4	20
RS'	7	2	4	7	20
RR'	2	4	7	17	30

Note: TE fit: $G(8) = 0.34$; TE+H0: $G(10) = 14.07$; H0: $G(2) = 13.73$.

292 The values in Table 6 were similarly constructed to match the marginal proportions,
 293 except Table 6 was built on the assumption that H0 is true: $p_{SR'} = p_{RS'} = 0$, $p_{RR'} = p_{SS'} =$
 294 0.5 ; $e' = f = 0.1$, and $e = f' = 0.4$. Predictions were rounded or slightly adjusted so that all
 295 entries are integers and the marginal proportions match. Again, TE fits the rounded values
 296 well, $G(8) = 0.47$, but this time, so does the special case of H0, $G(10) = 0.58$, so the test of
 297 H0 is $G(2) = 0.11$. The estimated $p_{SR'}$ was only 0.05, with a 95% bootstrapped confidence
 298 interval from 0 to 0.24. These results indicate that we can retain H0 for Table 6.

299 Therefore, the TE analysis of response patterns in the replicated experiment can distin-
 300 guish cases where H0 should be rejected (Table 5) or retained (Table 6), whereas the test of
 301 correlated proportions would lead to rejection of H in both cases.

302 By comparing Tables 5 and 6, one can gain insight in how the data lead to these different
 303 conclusions, even though the marginal sums are the same. In Table 5 the large frequency of
 304 SR', SR' (in the diagonal entry of Table 5), and low frequencies in the off-diagonal entries
 305 in the SR' row and column indicate that the SR' response pattern is "real". In contrast, in
 306 Table 6, one can see that there are high frequencies on the off-diagonals of switching between
 307 SS' or RR' in one replicate and SR' in the other replicate (there are 40 such cases), but
 308 rarely from these patterns to RS' , indicating that the errors, e and f' , account for the large
 309 marginal frequency of SR' , rather than reversals of true preference from S to R' .

310 The example of Table 6 also illustrates a difference between TE4 and TE2. For Table
311 6, the fit of TE2 is $G(10) = 10.99$ and the fit of the TE2 model with H0 (TE2 + H0) is
312 $G(12) = 51.10$, so $G(2) = 40.1$, which is significant. Thus, an investigator who used only
313 TE2 might conclude that H0 should be rejected in Table 6, whereas H0 can be retained if
314 TE4 is allowed. Although TE2 + H0 allows that the response frequency of SR' need not
315 equal that of RS' , it cannot imply that $P(S) > 0.5$ and $P(S') < 0.5$, but in this example,
316 $P(S) = 0.65$ and $P(S') = 0.35$. Thus, these two additional error parameters in TE4 (beyond
317 those of TE2) can potentially reverse the conclusions that two researchers might draw from
318 the same data if they employed TE4 and TE2.

319 Table 7 was constructed by arbitrarily choosing numbers in the table to produce row and
320 column totals to match those in Tables 5 and 6, but without any model to guide the pattern.
321 Can the TE model fit any such arbitrary data? The answer is, "no." The $G(8)$ for the fit
322 of the TE4 model is 35.2, which exceeds the critical value of $\chi^2(8)$ with $\alpha = 0.01$, which is
323 20.1. It should be clear that there are many ways to construct a 4 by 4 array with 15 df that
324 will significantly violate a model with only 7 degrees of freedom in its parameters.

325 Birnbaum (2019) suggested the following method to generate arbitrary data arrays: ran-
326 domly permute actual data (that is, simply take the same empirical frequencies observed in
327 a real experiment and re-arrange them randomly in the table). He then attempted to fit TE4
328 to each of 70,000 random permutations of empirical data with $n = 107$, where the original
329 data fit TE acceptably, $G(8) = 13.2$ It was found that 99.65% of such random permutations
330 had $G > 20$. The example of Table 7 and Birnbaum's (2019) analysis makes clear that the
331 TE model, like Factor Analysis with two dimensions, will not be expected to fit any arbitrary
332 set of numbers. Like other analytic models, TE models are not only statistical devices, but
333 also empirical theories that may or may not fit actual data.

334 Table 8 was constructed to illustrate a case in which H0 is false, but the test of correlated
335 proportions would conclude that H is perfectly acceptable. An experimenter who examined

336 only the proportions of SR' and RS' choices would find that the two types of reversals are
337 exactly equal, and that $P(S) = P(S')$. However, Table 8 actually contains strong evidence
338 against H_0 , since $G(2) = 13.73$, and the best-fit solution to TE4 indicates that $p_{SR'} = 0.34$,
339 $p_{SS'} = p_{RR'} = 0.33$, $p_{RS'} = 0$, $e = f' = 0.28$, and $e' = f = 0.04$. The bootstrapped 95%
340 confidence interval for $p_{SR'}$ is 0.09 to 0.47. Note that there are 28 cases of reversals in which
341 SS' or RR' appears in one replicate and RS' appears in the other, and there are only 4 cases
342 where RS' repeats in both replicates.

343 Another example might have been constructed to illustrate that H does not imply H_0 .
344 Suppose T1 is true and $p_{SR'} = p_{RS'} > 0$, in which case $p(S) = p(S')$, so H is satisfied
345 perfectly even though H_0 is false. By including replications, such cases can be detected by
346 TE and one can estimate e , $p_{SR'}$ and $p_{RS'}$. So even when the error model that justifies the
347 test of correlated proportions is correct, that test may fail to detect true violations of H_0
348 that might have been detected and their magnitudes assessed by TE methods.

349 In sum, TE analyses of hypothetical cases illustrate that the TE model provides a method
350 for deciding whether H_0 should be retained or rejected in cases where the test of correlated
351 proportions is or is not significant. Case 3 in Table 7 also illustrates that the TE model itself
352 is testable and may not fit the data.

353 The next section applies these TE methods to real data to compare two families of models
354 that make different predictions for a property known as interactive independence that can
355 be violated according to expected utility (and other theories in its class) and which must
356 be satisfied according to lexicographic semiorders (and other theories in its class). These
357 studies had been previously analyzed by means of TE2, which as illustrated in the analysis
358 of Table 6 might lead to rejection of H_0 in cases where TE4 would allow H_0 to reproduce
359 the data.

360 **3 Expected Utility versus Lexicographic Semiorders**

361 This section explores a test between two classes of risky decision making models: interactive
362 and non-interactive. Expected utility theory is an example of an interactive model, and a
363 lexicographic semiorder (LS) is an example of a non-interactive model.

364 Let $A = (x_A, q_A; y_A)$ represent a prospect (a "gamble") with a probability of q_A to win
365 $\$x_A$ and otherwise (with probability $1 - q_A$) receive $\$y_A$, where $x_A \geq y_A$. Two models that
366 describe how people might choose among such gambles are presented in the next subsections.

367 **3.1 Expected Utility Theory**

368 According to expected utility (EU) theory, a person prefers $A = (x_A, q_A; y_A)$ over $B =$
369 $(x_B, q_B; y_B)$ (denoted, $A \succ B$, where \succ represents "is truly preferred to") if and only if the
370 expected utility of A exceeds that of B. For two-branch gambles, EU implies:

$$A \succ B \Leftrightarrow q_A(u(x_A)) + (1 - q_A)(u(y_A)) > q_B(u(x_B)) + (1 - q_B)(u(y_B)) \quad (3)$$

371 where $u(x)$ is the utility function for money. Note that in this theory, increasing the proba-
372 bility to win x multiplies $u(x)$ and $u(y)$, so changing the value of q can be said to "interact"
373 with the effects of the consequences, x and y .

374 **3.2 Lexicographic Semiorders**

375 In the LPH lexicographic semiorder (LPH LS), the decision maker first compares the lower
376 consequences of the two alternatives (y_A, y_B) and if the difference exceeds a threshold param-
377 eter, the prospect with the better lowest consequence is chosen (without considering other
378 attributes); but if the difference does not exceed threshold, the decision maker next compares
379 the probabilities. If the difference in probabilities exceeds a threshold, the alternative with

380 the better probability is chosen; but if the difference does not exceed threshold, the highest
381 consequences are then examined and the prospect with the better highest consequence is
382 chosen. LS models can imply violations of transitivity (Tversky, 1969); that is, it is possible
383 to find A, B , and C , such that $A \succ B, B \succ C$, and $C \succ A$, where \succ indicates true preference
384 in the theory.

385 Another individual might use another LS model to compare gambles: she might use a
386 different order of considering the attributes. For example, a person might examine the highest
387 consequences first, then the lowest, then the probabilities (HLP LS). Different individuals
388 might also use different threshold parameters, which could also produce different preferences.
389 And in EU theory, if different people have different u functions, there could also be individual
390 differences among people with the same choice problem. Thus, under either EU or LS
391 theories, there might be individual differences that produce variability in true preferences
392 among individuals, which will be combined with variability due to random error. These two
393 sources of variation in responses might make it difficult to compare the fit of these models
394 to a given set of data.

395 **3.3 A Test of Interactive Independence**

396 Rather than compare models by asking how "well" they fit data obtained with an arbitrary
397 set of choice problems, it can be useful to conduct experiments that test critical properties.
398 A critical property is a property that can be deduced as a theorem from one theory and
399 which can be violated according to the other theory.

400 Birnbaum (2010) and Birnbaum and Gutierrez (2007, p. 107) devised and reported tests
401 of critical properties that must be satisfied by any mixture of LS models. Among these
402 critical properties is interactive independence, which is the assumption that the effect of a
403 difference between alternatives on one attribute is independent of any other attribute that
404 has the same value in both alternatives. This property must be satisfied by a mixture of LS

405 models (Birnbaum, 2010), but it can easily be violated by expected utility theory as well as
406 by other theories, such the TAX model (Birnbaum, 2008).

407 Interactive independence requires that for all $A = (x_A, p; y_A)$, $B = (x_B, p; y_B)$, $A' =$
408 $(x_A, p'; y_A)$, and $B' = (x_B, p'; y_B)$,

$$A \succ B \Leftrightarrow A' \succ B'. \quad (4)$$

409 Note that p is common to both A and B , which have the same consequences as A' and
410 B' , respectively, except that the (common) probability is now p' instead of p . In the specific
411 test below, $x_A > x_B > y_B > y_A$; because A has greater variance in outcomes it is thus more
412 "risky" compared to B ; the notations, R and S , are used to denote these "risky" and "safe"
413 gambles. Interactive independence can be tested in the following two choice problems:

414 Problem 1: Which do you prefer?

415 $R = (\$7.25, 0.05; \$1.25)$

416 or

417 $S = (\$4.25, 0.05; \$3.25)$

418 Problem 2: Which do you prefer?

419 $R' = (\$7.25, 0.95; \$1.25)$

420 or

421 $S' = (\$4.25, 0.95; \$3.25)$

422 Note that R is a "risky" gamble in which one might win either \$7.25 or \$1.25, and S
423 is a "safer" gamble in which the least one can win is \$3.25, but the most one can win is
424 \$4.25. In this case, the expected value of S is greater than that of R . In the second choice
425 problem, the consequences, S' and R' , are the same as those of S and R , respectively, but
426 the probability to win the higher prize (same in both gambles) is higher than it is in Problem
427 1. In the second problem, R' has the higher expected value than S' .

Table 9: Empirical frequencies in test of interactive independence.

Replicate 1	Responses on Replicate 2				Total
	SS'	SR'	RS'	RR'	
SS'	24	21	0	3	48
SR'	10	190	3	7	210
RS'	0	1	14	2	17
RR'	6	7	3	30	46

Note: Data from Birnbaum & Gutierrez (2007, Exp. 2), $n = 321$.

428 According to interactive independence, $S \succ R$ if and only if $S' \succ R'$. In any LS model or
 429 mixture of LS models, a person can have only two preference patterns, RR' or SS' (Birnbaum,
 430 2010, p. 376, p. 383), so interactive independence must be satisfied, apart from error. Thus,
 431 LS models implies interactive independence, $H_0: p_{SR'} = p_{RS'} = 0$.

432 On the other hand, if probabilities and consequences interact, as they do in EU (and
 433 many other theories), then a person might prefer $S \succ R$ in the Problem 1, and prefer $R' \succ$
 434 S' in Problem 2. This pattern of preferences is denoted SR' and would be indicative of an
 435 interaction. Depending on the utility function in EU theory, a person might have preference
 436 patterns of SR' , SS' or RR' .⁷

437 4 Reanalysis of Birnbaum & Gutierrez (2007)

438 Birnbaum and Gutierrez (2007) searched for violations of transitivity predicted by a lexico-
 439 graphic semiorde model using stimuli similar to those of Tversky (1969), who had argued
 440 that some participants might use a lexicographic semiorde that could produce intransitive
 441 preferences. Transitivity is a critical test between EU and LS theories that must be satis-
 442 fied by EU, but which can be violated by LS. Interspersed among trials testing transitivity,

⁷For example, the SR' pattern is implied for these choice problems when $u(x) = x$; but if $u(x) = x^b$, the RR' pattern is implied when $b \geq 3.82$; and if $u(x) = 1 - e^{-ax}$, the SS' pattern follows when $a \geq 1.02$.

Table 10: Indices of fit, G , of TE models to empirical data in Table 9.

Models	TE4	TE2	TE1
TE	30.8	31.1	38.4
TE + LS	320.1	369.3	771.6
LS	289.3	338.2	733.2

443 Birnbaum and Gutierrez (2007, Experiment 2) included tests of interactive independence de-
 444 scribed above. Problems 1 and 2 were presented twice to each of 321 participants, embedded
 445 in randomized and counterbalanced sequences among many other similar choice problems.⁸

446 Table 9 shows the empirical frequencies (counts) of the number of times that each of
 447 the 16 response patterns was observed in this test of interactive independence (Birnbaum &
 448 Gutierrez, 2007). (Table 9 and this method of analysis were not presented in that paper.)
 449 The most frequent response pattern, shown by 190 participants out of 321, was to repeat
 450 the SR' pattern on both replicates.

451 Table 10 shows the indices of fit, G , from TEMAP2.R for the six models, fit to Table 9.
 452 TE4, TE2, and TE1 models have 8, 10, and 11 df, respectively; corresponding LS models
 453 (TE + H0) have an additional 2 df; critical values of $\chi^2(df)$ for df = 2, 8, 10, and 11 for
 454 $\alpha = 0.05$ level of significance are 5.99, 15.51, 18.31, and 19.68, respectively. The differences
 455 in fit between each TE model and its LS special case are presented in the last row of the
 456 table (LS). (Tests of H0 are tests of interactive independence and therefore tests of LS.) All
 457 of the TE + LS models have G more than 10 times the corresponding values for the TE
 458 models of which they are special cases, and all differences (LS) are significant.

459 There are also violations of the TE models. According to any of the TE models, the matrix
 460 in Table 9 should be symmetric. However, the frequency of $SR'SS'$ is 10, and that of $SS'SR'$

⁸The raw data of both Birnbaum and Gutierrez (2007) and of Birnbaum (2010), as well as other data, are available in the archive at this URL:

<http://psych.fullerton.edu/mbirnbaum/archive.htm>

Table 11: "Predicted" (best-fit) frequencies of repeated pattern SR' ; Empirical = 190

Models	TE4	TE2	TE1
TE full	182.6	173.2	173.1
LS	64.6	63.5	20.1

Table 12: Best-fit estimates of parameters in TE models fit to Table 9.

Model	Parameter							
	$p_{SS'}$	$p_{SR'}$	$p_{RS'}$	$p_{RR'}$	e'	e	f'	f
TE4	20	56	12	13	00	39	22	00
TE2	08	75	05	11	04	08	$=e'$	$=e$
TE1	09	75	05	11	06	$=e'$	$=e'$	$=e'$

Note: Values expressed as percentages; i.e., 05 indicates 0.05.

461 is 21, significantly greater. The TEMAP2.R program calculates best-fit values ("predicted")
 462 corresponding to Table 9. These predictions showed that except for this violation, each of
 463 the TE models gave a fairly good approximation to the values in Table 9.⁹

464 The difference between TE4 and TE2 is not significant, but the small difference between
 465 TE2 and TE1 is significant ($G(1) = 38.4 - 31.1 = 7.3, p < 0.05$).

466 The predictions of the LS models were all quite bad, especially in their best-fit values for
 467 the largest observed frequency in Table 9 (190), for the repeated response pattern, $SR'SR'$.
 468 According to any of the LS models, this pattern only occurs due to errors. Table 11 shows the
 469 best-fit predicted values for the six models. The LS4 model predicts 64.6 for this frequency,
 470 and predictions for the other LS models are even farther below the actual value of 190.
 471 Therefore, LS models fail because they are not able to account for the large number of
 472 people who repeated the SR' pattern.

473 Table 12 shows the estimated parameters of the three TE models, which provide better

⁹See Birnbaum and Quan (2020) for simulation studies of the robustness of TE models with respect to systematic violations in tests of transitivity.

Table 13: Test of interactive independence with $p = 0.01$ and $p' = 0.99$.)

Replicate 1	Responses on Replicate 2			
Series A	SS'	SR'	RS'	RR'
SS'	10	8	0	2
SR'	6	77	1	11
RS'	1	0	2	6
RR'	1	10	2	16
Series B	SS'	SR'	RS'	RR'
SS'	4	12	2	3
SR'	16	84	0	5
RS'	0	0	1	2
RR'	0	7	4	10

Note: Data of Birnbaum (2010, Exp. 3, $n = 153$.)

474 approximations to the data than the LS special cases. (Probabilities are expressed as per-
 475 centages to save space in the table; e.g., 04 indicates 0.04.) The best-fit values indicated
 476 that the percentages of participants with SR' pattern as their true preference pattern were
 477 56%, 75%, and 75%, according to TE4, TE2, and TE1, respectively. The corresponding
 478 95% confidence intervals based on 10,000 bootstrapped samples were 50–81, 70–81, and
 479 70–80, respectively, giving confidence that the majority of the sample violated interactive
 480 independence in the manner predicted by interactive models like expected utility.

481 5 Reanalysis of Birnbaum (2010)

482 Birnbaum (2010, Experiment 3) reported tests of interactive independence in choice problems
 483 of the following type:

484 $R = (\$95, p; \$5)$

Table 14: Test of interactive independence with $p = 0.1$ and $p' = 0.9$.)

Replicate 1	Responses on Replicate 2			
Series A	SS'	SR'	RS'	RR'
SS'	12	9	1	1
SR'	10	48	2	12
RS'	0	0	1	2
RR'	2	14	0	37
Series B	SS'	SR'	RS'	RR'
SS'	17	6	1	1
SR'	12	58	1	13
RS'	3	1	0	1
RR'	0	10	2	27

Note: Data of Birnbaum (2010, Exp. 3, $n = 153$.)

485 or

486 $S = (\$55, p; \$20)$

487 where there were five levels of p (and p'): 0.01, 0.10, 0.50, 0.90, and 0.99. There were
 488 153 participants who responded to each choice problem twice, randomly embedded among
 489 many other trials. There were also two variations (Series A and B) with slightly different
 490 values of the consequences (\$50 and \$15 instead of \$55 and \$20), providing another check
 491 on consistency of the results.

492 Results for both series are shown in Table 13 for $p = 0.01$ and $p' = 0.99$, and in Table 14
 493 for $p = 0.10$ and $p' = 0.90$. The modal response pattern in all four cases is to respond SR'
 494 on both replications: 77 and 84 participants in Series A and B of Table 13 and 48 and 58
 495 participants in Series A and B of Table 14, respectively.

496 Tables 15 and 16 show statistical tests for the six TE models and the tests between each
 497 TE model and its LS (H0) special case. In all 12 cases (4 sets of data by 3 TE models

Table 15: Indices of fit, G , of TE models in tests of interactive independence with $p = 0.01$ and $p' = 0.99$.

Series A	TE4	TE2	TE1
TE	8.9	11.5	11.5
TE + LS	83.8	131.4	291.8
LS	74.9	120.0	280.3
Series B	TE4	TE2	TE1
TE	14.2	15.5	24.7
TE + LS	111.7	160.4	345.0
LS	97.5	144.9	320.3

498 in Tables 15 and 16), the large violations of interactive independence, indicate that the LS
 499 models can be confidently rejected under any of the error models.

500 The differences among the TE models are again smaller than differences between TE and
 501 LS models; however, in Table 16, Series A, TE4 fits significantly better than TE2, and in
 502 Table 15, Series B, TE2 and TE4 fit significantly better than TE1.

503 Table 17 shows the estimated parameters under three error models (TE4, TE2, and TE1)
 504 for the four sets of data. The estimated incidence of violations of interactive independence
 505 ($p_{SR'}$ were substantial in all 12 cases. For example, for TE2 Series A and B, the estimated
 506 incidences are 0.73 and 0.85 when $p = 0.01$, and they are 0.56 and 0.60 when $p = 0.10$.
 507 Bootstrapped estimates of 95% confidence intervals on the parameter estimates agree that
 508 one can reject H_0 with confidence, in favor of the hypothesis that $p_{SR'} > 0$ in all cases.

509 In sum, reanalyses of Birnbaum (2010) and of Birnbaum and Gutierrez (2007) are clear:
 510 violations of interactive independence cannot be attributed to random error as in Figure 1.
 511 Although TE4 analysis has the potential to reverse the conclusions of earlier analyses, like
 512 TE2 or the test of correlated proportions, these reanalyses instead reinforce the conclusions
 513 that had been reached using those methods.

Table 16: Indices of fit, G , of TE models fit to Birnbaum (2010) test of interactive independence with $p = 0.10$ and $p' = 0.90$.

Series A	TE4	TE2	TE1
TE	8.7	18.3	19.1
TE + LS	35.0	91.8	207.7
LS	26.3	73.5	188.6
Series B	TE4	TE2	TE1
TE	6.6	9.2	10.0
TE + LS	42.6	114.3	239.3
LS	36.0	105.1	229.3

514 An important finding of these studies was that most those few participants who appeared
515 to show violations of transitivity also showed systematic violations of interactive indepen-
516 dence. That finding suggests that even for those few participants, we cannot retain LS
517 models as a descriptive theory of the violations of transitivity. Birnbaum and Gutierrez
518 (2007) suggested a rival theory for those cases in terms of an assimilation of subjective val-
519 ues of similar probabilities prior to aggregation by a model with multiplicative interaction
520 between probability and value.

521 6 Discussion

522 These analyses lead to four main conclusions: (1) The test of correlated proportions is not
523 appropriate for testing if two situations are psychologically equivalent, if the dependent
524 measures might contain errors. (2) Investigators should instead employ replications within-
525 subjects and analyze response patterns to assess the error structure. (3) The TE models
526 provide workable methods for estimating error rates and the true response patterns, as well as
527 providing statistical tests of both the substantive issues and of the TE models. (4) Reanalysis

Table 17: Best-fit estimates of parameters in TE models fit to Tables 13 and 14.

Model	Parameter							
$p = 0.01$	$p_{SS'}$	$p_{SR'}$	$p_{RS'}$	$p_{RR'}$	e'	e	f'	f
TE4 Series A	26	58	08	08	11	50	00	00
TE4 Series B	02	72	04	22	00	31	33	14
TE2 Series A	08	73	02	18	09	09		
TE2 Series B	03	85	01	11	06	14		
TE1 Series A	08	73	02	18	09			
TE1 Series B	04	84	01	11	10			
$p = 0.10$	$p_{SS'}$	$p_{SR'}$	$p_{RS'}$	$p_{RR'}$	e'	e	f'	f
TE4 Series A	28	29	03	40	06	45	22	01
TE4 Series B	25	46	00	29	08	26	19	04
TE2 Series A	11	56	00	33	12	10		
TE2 Series B	16	60	00	24	11	08		
TE1 Series A	11	56	00	34	11			
TE1 Series B	16	60	00	24	10			

Note: Values expressed as percentages; i.e., 05 indicates 0.05.

528 of two published experiments via the new methods gives a very clear answer to the question
 529 posed in the title to this paper: LS models cannot be saved by the flexible error theory of
 530 Figure 1. These conclusions are discussed in the next sections.

531 6.1 Test of correlated proportions

532 From the derivations and examples analyzed here, it should be clear that if one allows that
 533 the dependent measure may contain errors as in Figure 1, then one should not use the test of
 534 correlated proportions to decide whether two conditions are or are not equivalent. Similarly,
 535 simply because one condition produces a proportion that is significantly greater than 0.5

536 and another condition produces a proportion significantly less than 0.5, one cannot reject
537 the null hypothesis that the two experimental conditions induced the same true responses.

538 This conclusion can be restated more clearly for algebraic choice theory as follows. A
539 theoretician wishes to test a risky decision making model, which implies that $S \succ R \Leftrightarrow$
540 $S' \succ R'$. Because \succ represents true preference, rather than expressed preference, this theory
541 implies H0, that $p_{SR'} = p_{RS'} = 0$, which implies $p_S = p_{S'}$. The test of correlated proportions,
542 however, tests the null hypothesis, H, that the response proportions are equal, $p(S) =$
543 $p(S')$, which is equivalent to equality of the two types of expressed preference reversals; i.e.,
544 $p(SR') = p(RS')$. In the error theory of Figure 1, only TE1 implies that the two types of
545 observed preference reversals will be equal under H0, but equality does not guarantee that
546 they are both zero, so a test of H is not the same as a test of H0, even when T1 is assumed.
547 Furthermore, preference reversals need not be equal for H0 under either TE2 or TE4. Finally,
548 TE4 does not even require that $S \succ R \Leftrightarrow p(S) > 0.5$; indeed, TE4 can allow cases in which
549 modal response probabilities reverse; i.e., $p(S) > 0.5$ and $p(S') < 0.5$, even when H0 holds—
550 i.e., even when there are no true reversals of preference ($p_{SR'} = p_{RS'} = 0$.) In summary,
551 there is a mismatch in principle between the statistical tests of correlated proportions and
552 the theoretical properties of true preferences one wishes to test.

553 6.2 Need for Replications

554 In studies without replications, as in Table 2, Table 3 shows that one cannot answer questions
555 one wishes to address because one cannot tease out measurement of error from the substantive
556 question of the equivalence of conditions. The data in Table 2 are perfectly compatible with
557 the theory that no one reversed true preferences, but they are also consistent with the theory
558 that people systematically switched from R to S' .

559 Unfortunately, many published studies of interesting problems used the statistical test of
560 correlated proportions and many studies did not even include replications. The conclusions

561 drawn from such studies can therefore be questioned, and those questions cannot be answered
562 by reanalysis. For example, a recent review of evidence on the Allais paradox by Blavatsky,
563 et al. (2022) summarizes 81 experiments using Conlisk's z statistic as an index of strength
564 and direction of the paradox. As shown in this paper, this index is not diagnostic of H_0 ; it can
565 be zero when there is a large asymmetric violation of H_0 or when real but opposite violations
566 balance out; and it can be large in absolute value when H_0 is acceptable. Consequently, it
567 is unclear what conclusions, if any, can be drawn from an analysis based on the z index or
568 its components that does not account for errors of measurement.

569 Because neither H nor H_0 implies the other, it would seem reasonable to reanalyze those
570 studies that included replications and perhaps execute those studies again whose conclusions
571 are important and in doubt. Birnbaum and Quispe-Torreblanca (2018) analyzed the data of
572 Birnbaum, et al. (2017) and concluded that violations of the constant consequence indepen-
573 dence of Allais are indeed "real"; the violations in that study cannot be explained by error
574 as in Figure 1.

575 Birnbaum (2008) summarized a number of "new paradoxes" that rule out both expected
576 utility theory and both versions of prospect theory (Kahneman & Tversky, 1979; Tversky &
577 Kahneman, 1992) as descriptive theories of risky decision making. The "new paradoxes" are
578 critical tests of prospect theory that, like the Allais paradoxes, must be implied with any
579 utility function and weighting function. Many of the early studies of this program of research
580 used the test of correlated proportions (e.g., Birnbaum, 1999b). Birnbaum (2008) replicated
581 many of these paradoxes, including violations of first order stochastic dominance, dissection
582 of the Allais paradox, upper and lower cumulative independence, and violations of restricted
583 branch independence and analyzed them via a simplified version of TE2. However, there is
584 a need to re-run or re-analyze those studies in order to check the possibility that some form
585 of prospect theory might be saved by the more complex error theory of TE4 in Figure 1.

586 6.3 True and Error Model Analysis

587 When replications are included in a study, it becomes possible to fit and test TE models and
588 to test H_0 . The model allows one to estimate not only the error rates in Figure 1, but also the
589 four probabilities of the true preference patterns. These four probabilities (informed by the
590 confidence intervals on them) are crucial to evaluation of the theories under consideration.
591 As shown in the hypothetical examples constructed here in Tables 5–7, the TE analysis of
592 replicated data can properly distinguish cases that are equivalent (have the same marginal
593 proportions) to the test of correlated proportions.

594 The TE model is not only a statistical device or analytic tool, but also a simple descriptive
595 model that can be tested. Like any such model, TE uses simplifying assumptions. For
596 example, in the analyses reported here, the model assumes that each person maintains the
597 same true preferences within the session. If people changed true preferences within a session,
598 it would have the effect of inflating the estimated error terms. Further, the analyses presented
599 here assumed that all people have the same error rates, but we know that there are differences
600 in reliability among people. To handle heterogeneity in error rates, Birnbaum and Gutierrez
601 (2007) subdivided data according to the rates of within-person reliability, and analyzed the
602 reliable and unreliable participants separately, which resulted in a better fit of the TE model
603 to the data so analyzed. It may be that the substantive conclusions are robust with respect
604 to such violations (Birnbaum & Quan, 2020), but this question deserves further study.

605 There might be situations where obtaining replications would be difficult to accomplish,
606 but that is certainly not a valid excuse in studies of decision making, where it is common to
607 collect many responses from each participant.

608 The examples analyzed here all involved within-subjects experiments in which the depen-
609 dent measure could be replicated by the same person in each condition. Between-subjects
610 experiments are simpler to analyze statistically, but theoretically, they are more complicated
611 to analyze than within-subjects studies because there can be different relationships between

612 the subjective value and dependent measures in each group of subjects (Birnbaum, 1982;
613 1999a). In a between-subject studies, 9 can be rated as a "bigger" number than 221, but
614 in within-subject studies 221 is judged "bigger" than 9. Similarly, a "married woman" who
615 is a rape victim is rated more "at fault" than a "divorcee" rape victim in between-subjects
616 studies (Jones & Aronson, 1973; Birnbaum, 1982), However, in within-subjects studies, the
617 divorcee is rated more at fault (Birnbaum, 1982). The differences between within- and
618 between-subjects designs can be reconciled by a theory of how contexts in different groups
619 can be different and confounded with the stimuli in between-subjects studies. Because of
620 such complications, satisfactory TE methods and models have not yet been developed for
621 between-subjects situations.

622 **6.4 Rival Methods**

623 Previous approaches to the analysis of variability of responses in choice studies have been
624 reviewed in a number of papers (Birnbaum, 2004, 2008, 2013; Bhatia & Loomes, 2017;
625 Busemeyer & Townsend, 1993; Carbone & Hey, 2000; Kvam & Busemeyer, 2020; Luce, 1997,
626 2000; Regenwetter, Dana, & Davis-Stober, 2011; Wilcox, 2008). A main theme of these
627 reviews is that because there are multiple sources of possible variability, previous approaches
628 have been unable to separate them without arbitrary assumptions, and those assumptions
629 often interacted with the main purpose of the research, which is to test alternative substantive
630 models of decision making.

631 A rival method to true and error models for the analysis of response proportions in within-
632 subjects studies is the Qtest approach, described in Regenwetter, Davis-Stober, Lim, Cha,
633 Guo, Messner, Popova, & Zwilling (2014) and updated in Zwilling, Cavagnaro, Regenwetter,
634 Lim, Fields, & Zhang (2019). This approach has been applied to cases of individual data with
635 the assumption that repeated responses by the same person are independent and identically
636 distributed (iid); however, that iid assumption has been found to be systematically violated

637 in empirical choice data obtained from individuals (Birnbaum, 2012, 2013, 2022; Birnbaum
638 & Bahra, 2012a, 2012b), including reanalysis of the data of Regenwetter, et al. (2011).

639 If data satisfy iid, then there is no more information in crosstabulation matrices such as
640 Tables 5-9 than in the two marginal, binary response proportions in each case. Although iid
641 can occur in TE models in special cases, such as when there is just a single true preference
642 pattern, iid is not generally implied. But the Qtest approach begins and ends with simple
643 analysis of the binary response proportions and ignores all of the information in tables like
644 Tables 5–7, so it would conclude that those cases are all the same, because $P(S)$ and $P(S')$
645 are the same in all three cases.

646 If Qtest were applied to the hypothetical data in Tables 5–8, it would conclude that one
647 should reject H in Tables 5, 6, and 7 and retain it for Table 8. To my knowledge, the Qtest
648 method has not yet been applied in the situations analyzed here, but if it were, it could
649 be criticized by the same arguments as those directed here against the test of correlated
650 proportions, plus the criticism that it assumes away correlations between repeated responses
651 from the same person.

652 The Qtest method has been applied to tests of transitivity of preference by Regenwetter,
653 et al. (2011) and Cavagnaro & Davis-Stober (2014), among others. Transitivity requires
654 that $S \succ R$ and $R \succ T \Rightarrow S \succ T$. Birnbaum and Wan (2020) have shown that any
655 method, including Qtest, that is based strictly on binary response proportions (ignores the
656 pattern data) cannot be relied upon to distinguish data that have been simulated from
657 either transitive or intransitive models. In contrast, TE methods correctly diagnose the data
658 with respect to the model that simulated the data. TE methods for analysis of the issue
659 of transitivity have been presented in Birnbaum and Bahra (2012b), Birnbaum and Wan
660 (2020), Schramm (2020), and Birnbaum (2022).

661 Instead of forcing the assumption of iid in order to justify off-the shelf statistical tests or
662 to simplify an analysis, it seems preferable to make use of the information provided by the

663 pattern information in the data, which typically violates iid, by means of a model that can
664 describe those patterns.

665 **6.5 Lexicographic Semiorde Models Rejected**

666 The reanalyses of Birnbaum and Gutierrez (2007) and Birnbaum (2010) give a clear answer:
667 the property of interactive independence can be confidently rejected. Because any LS model
668 or mixture of LS models imply interactive independence (Birnbaum, 2010), these models can
669 be rejected as descriptive models of risky decision making, and they cannot be saved by the
670 flexible error theory of Figure 1. Other theories that imply interactive independence, such
671 as the Simplified Additive Difference (SAD) model (Ranyard, et al., 2020) and the priority
672 heuristic (Brandstätter, Gigerenzer, & Hertwig, 2006) can also be rejected as descriptive by
673 these results.

674 Tables 12 and 17 show that the property of interactive independence is violated in a
675 particular way in these tests by more than half of the sample under any of the error theories.
676 However, that allows that some people might actually satisfy the property, so it remains
677 possible that perhaps a subset of people might still satisfy LS models. Indeed, Tversky
678 (1969) concluded that only a small proportion of the people tested showed evidence of the
679 intransitive behavior that could be described by his LS model. Other studies also found
680 that only a small fraction of participants show intransitive behavior (Ranyard, et al, 2020;
681 Birnbaum & Gutierrez, 2007; Birnbaum & Bahra, 2012b; Birnbaum, 2010; Birnbaum, 2020;
682 2022; Butler & Pogrebna, 2018).

683 However, Birnbaum and Gutierrez (2007), Birnbaum and Bahra (2012b) and Birnbaum
684 (2010) found that even among those who appeared to show evidence of intransitive prefer-
685 ences in one design, those same individuals did not show consistency with other predictions
686 of LS models with other choice problems included in the same study. For example, Birn-
687 baum and Bahra (2012b) were unable to find a single case in a sample of 134 participants

688 where one could predict from an LS model of choices among gambles of the form $(x, p; 0)$
689 to choices among $(x, 1/2; y)$ and to choices among $(100, p; y)$, where the levels of $x, p,$ and
690 y were chosen to form interlinked designs. Birnbaum and Gutierrez (2007) and Birnbaum
691 (2010) found that most people who showed evidence of transitivity also showed violations
692 of interactive independence or other properties implied by LS models. These findings imply
693 that some other theory besides LS models, such as the assimilation theory of Birnbaum and
694 Gutierrez (2007), is required in order to account for those few cases that appear to show
695 evidence of intransitive preferences.

696 Because conclusions of a few studies regarding interactive independence and the Allais
697 paradox have not changed as a result of TE reanalysis, one might be tempted to infer that it
698 is safe to assume that previous analytic methods are "good enough" for drawing conclusions
699 about theories of behavior. I think that inference would be a mistake. The algebra shows
700 that the conclusions can be changed by proper experiments and analyses. Further, the few
701 cases selected for reanalysis so far have been cases where the evidence has been quite strong;
702 other sets of data may yield different conclusions. Therefore, I would urge experimenters to
703 employ replications and use the newer methods of analysis in order to avoid drawing false
704 conclusions—conclusions that might be reversed by reanalysis or by a proper experiment with
705 replications.

706 References

707 Bhatia, S. & Loomes, G. (2017). Noisy preferences in risky choice: A cautionary note.
708 *Psychological Review*, 124(5), 678–687. <http://dx.doi.org/10.1037/rev0000073>

709 Birnbaum, M. H. (1982). Controversies in psychological measurement. In B. Wegener
710 (Ed.), *Social attitudes and psychophysical measurement* (pp. 401-485). Hillsdale, N.J.:
711 Lawrence Erlbaum Associates. <https://doi.org/10.4324/9780203780947>

712 Birnbaum, M. H. (1999a). How to show that $9 > 221$: Collect judgments in a between-
713 subjects design. *Psychological Methods*, 4(3), 243-249. [https://doi.org/10.1037/1082-989x.](https://doi.org/10.1037/1082-989x.4.3.243)
714 4.3.243

715 Birnbaum, M. H. (1999b). Testing critical properties of decision making on the Internet.
716 *Psychological Science*, 10(5), 399-407. <https://doi.org/10.1111/1467-9280.00176>

717 Birnbaum, M. H. (2004). Tests of rank-dependent utility and cumulative prospect theory
718 in gambles represented by natural frequencies: Effects of format, event framing, and branch
719 splitting. *Organizational Behavior and Human Decision Processes*, 95(1), 40-65. <https://doi.org/10.1016/j.obhdp.2004.05.004>

720

721 Birnbaum, M. H. (2008). New paradoxes of risky decision making. *Psychological Review*,
722 115, 463-501. <https://doi.org/10.1037/0033-295x.115.2.463>

723 Birnbaum, M. H. (2010). Testing lexicographic semi-orders as models of decision mak-
724 ing: Priority dominance, integration, interaction, and transitivity. *Journal of Mathematical*
725 *Psychology*, 54, 363-386. <https://doi.org/10.1016/j.jmp.2010.03.002>

726 Birnbaum, M. H. (2012). A statistical test of independence in choice data with small
727 samples. *Judgment and Decision Making*, 7(1), 97-109. [http://journal.sjdm.org/11/11605/
728 jdm11605.pdf](http://journal.sjdm.org/11/11605/jdm11605.pdf)

729 Birnbaum, M. H. (2013). True-and-error models violate independence and yet they are
730 testable. *Judgment and Decision Making*, 8, 717-737. [http://journal.sjdm.org/13/13422c/
731 jdm13422c.pdf](http://journal.sjdm.org/13/13422c/jdm13422c.pdf)

732 Birnbaum, M. H. (2019). Bayesian and frequentist analysis of True and Error models.
733 *Judgment and Decision Making*, 14(5), 608-616. [http://www.sjdm.org/journal/19/190822/
734 jdm190822.pdf](http://www.sjdm.org/journal/19/190822/jdm190822.pdf)

735 Birnbaum, M. H. (2020). Reanalysis of Butler and Pogrebna (2018) using true and
736 error mode. *Judgment and Decision Making*, 15(6), 1044-1051. [http://journal.sjdm.org/20/
737 200216/jdm200216.pdf](http://journal.sjdm.org/20/200216/jdm200216.pdf)

738 Birnbaum, M. H. (2022). Testing transitivity of preference in individuals. *Submitted for*
739 *publication, xx*(xx), xxx-xxx.

740 Birnbaum, M. H., & Bahra, J. P. (2012a). Separating response variability from structural
741 inconsistency to test models of risky decision making, *Judgment and Decision Making, 7*,
742 402-426. <http://journal.sjdm.org/12/12315/jdm12315.pdf>

743 Birnbaum, M. H., & Bahra, J. P. (2012b). Testing transitivity of preferences in individ-
744 uals using linked designs. *Judgment and Decision Making, 7*, 524-567. <http://journal.sjdm.org/11/111122/jdm111122.pdf>

746 Birnbaum, M. H., & Gutierrez, R. J. (2007). Testing for intransitivity of preference
747 predicted by a lexicographic semiorder. *Organizational Behavior and Human Decision Pro-*
748 *cesses, 104*, 97-112. <https://doi.org/10.1016/j.obhdp.2007.02.001>

749 Birnbaum, M. H., & Quan, B. (2020). Note on Birnbaum and Wan (2020): True and error
750 model analysis is robust with respect to certain violations of the MARTER model. *Judgment*
751 *and Decision Making, 15*(5), 861-862. <https://sjdm.org/journal/20/200413b/supp.pdf>

752 Birnbaum, M. H., & Quispe-Torreblanca, E. G. (2018). TEMAP2.R: True and error
753 model analysis program in R. *Judgment and Decision Making, 13*(5), 428-440. <http://www.sjdm.org/journal/18/18507/jdm18507.pdf>

755 Birnbaum, M. H., Schmidt, U., & Schneider, M. D. (2017). Testing independence condi-
756 tions in the presence of errors and splitting effects. *Journal of Risk and Uncertainty, 54*(1),
757 61-85. <https://doi.org/10.1007/s11166-017-9251-5>

758 Birnbaum, M. H., & Wan, L. (2020). MARTER: Markov true and error model of drifting
759 parameters. *Judgment and Decision Making, 15*, 47-73. [http://journal.sjdm.org/19/190727/](http://journal.sjdm.org/19/190727/jdm190727.pdf)
760 [jdm190727.pdf](http://journal.sjdm.org/19/190727/jdm190727.pdf)

761 Blavatsky, P., Ortmann, A., & Panchenko, V. (2022). On the experimental robustness
762 of the Allais paradox. *American Economic Journal: Microeconomics, 14*(1): 143-63. <https://www.aeaweb.org/articles?id=10.1257/mic.20190153&&from=f>
763 [/](https://www.aeaweb.org/articles?id=10.1257/mic.20190153&&from=f)

764 Brandstätter, E., Gigerenzer, G., & Hertwig, R. (2006). The priority heuristic: Choices
765 without tradeoffs. *Psychological Review*, *113*(2), 409–432. [https://doi.org/10.1037/0033-295x.](https://doi.org/10.1037/0033-295x.113.2.409)
766 113.2.409

767 Busemeyer, J. R., & Townsend, J. T. (1993). Decision field theory: A dynamic-cognitive
768 approach to decision making in an uncertain environment. *Psychological Review*, *100*(3),
769 432–459. <https://doi.org/10.1037/0033-295X.100.3.432>

770 Butler, D. J., & Pogrebna, G. (2018). Predictably intransitive preferences. *Judgment*
771 *and Decision Making*, *13*(3), 217-236. <https://sjdm.org/journal/17/17912b/jdm17912b.pdf>

772 Carbone, E., & Hey, J. D. (2000). Which error story is best? *Journal of Risk and*
773 *Uncertainty*, *20*(2), 161-176. <https://www.jstor.org/stable/41760978>

774 Cavagnaro, D.R., & Davis-Stober, C. P. (2014). Transitive in our preferences, but tran-
775 sitive in different ways: An analysis of choice variability. *Decision*, *1*, 102-122. <https://www.apa.org/pubs/journals/features/dec-0000011.pdf>

776 //www.apa.org/pubs/journals/features/dec-0000011.pdf

777 Conlisk, J. (1989). Three Variants on the Allais Example. *The American Economic Re-*
778 *view*, *79*, 392-407. <https://EconPapers.repec.org/RePEc:aea:aecrev:v:79:y:1989:i:3:p:392-407>

779 Jones, C., & Aronson, E. (1973). Attribution of fault to a rape victim as a function of
780 respectability of the victim. *Journal of Personality and Social Psychology*, *26*(3), 415–419.
781 <https://doi.org/10.1037/h0034463>

782 Kvam, P. D., & Busemeyer, J. R. (2020). A distributional and dynamic theory of pric-
783 ing and preference. *Psychological Review*, *127*(6), 1053–1078. [https://doi.org/10.1037/](https://doi.org/10.1037/rev0000215)
784 [rev0000215](https://doi.org/10.1037/rev0000215)

785 Lee, M. D. (2018). Bayesian methods for analyzing true-and-error models. *Judgment*
786 *and Decision Making*, *13*(6), 622-635. <https://sjdm.org/journal/18/18507c/jdm18507c.pdf>

787 Lichtenstein, S., & Slovic, P. (1971). Reversals of preference between bids and choices
788 in gambling decisions. *Journal of Experimental Psychology*, *89*, 46–55. [https://doi.org/10.](https://doi.org/10.1037/h0031207)
789 [1037/h0031207](https://doi.org/10.1037/h0031207)

- 790 Luce, R. D. (1997). Some unresolved conceptual problems in mathematical psychology.
791 *Journal of Mathematical Psychology*, 41, 79-87. <https://doi.org/10.1006/jmps.1997.1150>
- 792 Luce, R. D. (2000). *Utility of gains and losses: measurement-theoretical and experi-*
793 *mental approaches*. Mahwah,NJ:Lawrence Erlbaum Associates. [https://doi.org/10.4324/](https://doi.org/10.4324/9781410602831)
794 9781410602831
- 795 McNemar, Q. (1947). Note on the sampling error of the difference between correlated pro-
796 portions or percentages. *Psychometrika*, 12, 153–157. <https://doi.org/10.1007/BF02295996>
- 797 Novick, M. R. (1966). The axioms and principal results of classical test theory. *Journal*
798 *of Mathematical Psychology*, 3(1), 1-18. [https://doi.org/10.1016/0022-2496\(66\)90002-2](https://doi.org/10.1016/0022-2496(66)90002-2)
- 799 Ranyard, R., Montgomery, H., Konstantinidis, E., & Taylor, A. L. (2020). Intransitivity
800 and transitivity of preferences: Dimensional processing in decision making. *Decision*, 7(4),
801 287–313. <https://doi.org/10.1037/dec0000139>
- 802 Regenwetter, M., Dana, J., & Davis-Stober, C. P. (2011). Transitivity of Preferences.
803 *Psychological Review*, 118, 42–56. <https://doi.org/10.1037/a0021150>
- 804 Regenwetter, M., Davis-Stober, C.P., Lim, S.H., Cha, Y.-C., Guo, Y., Messner, W.,
805 Popova, A., & Zwilling, C.(2014). QTEST: Quantitative Testing of Theories of Binary
806 Choice. *Decision*, 1, 2-34. <https://doi.org/10.1037/dec0000007>
- 807 Schramm, P. (2020). The individual true and error model: Getting the most out of
808 limited data. *Judgment and Decision Making*, 15(5), 851-860. [https://sjdm.org/journal/](https://sjdm.org/journal/19/190516/jdm190516.pdf)
809 19/190516/jdm190516.pdf
- 810 Spearman, C. (1904). The proof and measurement of association between two things.
811 *The American Journal of Psychology*, 15 (1), 72-101. <https://doi.org/10.2307/1412159>
- 812 Tversky, A. (1969). Intransitivity of preferences. *Psychological Review*, 76, 31-48. <https://doi.org/10.1037/h0026750>
- 813
- 814 Wilcox, N. T. (2008). Stochastic models for binary discrete choice under risk: A critical
815 primer and econometric comparison. In J. C. Cox, & G. W. Harrison (Eds.), *Risk Aversion*

816 *in Experiments (Research in Experimental Economics; Vol. 12*, pp. 197–292). Bingley, UK:
817 Emerald Group Publishing Limited. [https://doi.org/10.1016/S0193-2306\(08\)00004-5](https://doi.org/10.1016/S0193-2306(08)00004-5)
818 Zwilling, C.E., Cavagnaro, D.R., Regenwetter, M., Lim, S.H., Fields, B., & Zhang, Y.
819 (2019). QTEST 2.1: Quantitative Testing of theories of binary choice using Bayesian in-
820 ference. *Journal of Mathematical Psychology*, *91*, 176-194. [https://doi.org/10.1016/j.jmp.](https://doi.org/10.1016/j.jmp.2019.05.002)
821 2019.05.002