1 *Article*

# 2 3 Guidelines for Assessing Oenological and Statistical Significance of Wine Tasters' Binary Judgments

4 **Dom Cicchetti, PhD.**

5  [1]  Yale University Home Office

6          94 Linsley Lake Road, Box 317

7          North Branford, CT 06471

8  **\*** Correspondence: dom.cicchetti@yale.edu

9 **Abstract:** The purpose of this article is to assess the reliability and accuracy (validity) of
10 hypothetical binary tasting judgments in an enological framework. The model that is utilized
11 allows for the control of a wide array of variables that would be exceedingly difficult to fully
12 control in the typical enological investigation. It is shown that results that are judged to be
13 enologically significant are uniformly judged to be statistically significant as well, whether the level
14 of wine Taster agreement is set at 70% (Fair); 80% (Good), or 90% (Excellent), However, in a
15 number of instances, results that were statistically significant were not enologically significant by
16 standards that are widely accepted and utilized. This finding is consistent with the bio-statistical
17 fact that given a sufficiently large sample size, even the most trivial of results will prove to be
18 statistically significant. Consistent with expectations, multiple patterns of 80% (Good) and 90%
19 (Excellent) agreement tended to be both statistically and enologically significant.

20 **Keywords:** hypothetics; enothetics; reliability; validity; accuracy
21

## 22 1. Introduction

23       The objective of this research report is to present a detailed analysis of the relationship between
24 oenological and statistical significance of research results as they both relate to the reliability and
25 accuracy of wine tasters' hypothetical judgments. Reliability is defined here as the extent to which
26 any given binary wine judgment is interchangeable with that of another wine judge. (e.g., agreement
27 that a wine is of excellent quality). The greater the extent to which this occurs, the higher the level of
28 reliability.

29       The accuracy or validity of a hypothetical binary decision refers to the extent to which any pair
30 of wine tasters renders the same correct judgment for example, they both agree, correctly, that the
31 wine is oaked or unoaked, or that the grape varietal is Syrah rather than Grenache. With respect to
32 oenological   research investigations and scientific investigations more broadly, it is a well-known
33 fact that uncontrolled variables can serve to compromise or call into question the accuracy of the
34 reported findings.

35       In a previous study, a method was introduced, in an oenological context to address this vexing,
36 albeit critical issue. Referred to as hypothetics, or oenothetics in the current research context, the
37 method allows investigators to begin to answer, what findings would occur if it were indeed
38 possible to control for variables that are often very difficult or, in some situations, impossible to
39 control in the typical research study. A distinct advantage of the method is that it can also serve to
40 highlight findings that would have become apparent if it were possible to control relevant variables.
41 For example, in a recent oenological investigation it was shown that overall accuracy is a very poor
42 measure of binary wine judgments, such as whether a wine is oaked or not [1]. Specific measures of
43 judgmental accuracy, such as Sensitivity (Se), Specificity (Sp), Predicted Positive Accuracy (PPA)
44 and Predicted Negative Accuracy (PNA) were found to be much more useful measures of wine
45 judgments than overall accuracy. The bio-statistical importance of such findings has relevance in

46 designing future oenological research investigations and in the design of scientific studies more
47 generally.

## 2. The Role of Chance in Scientific Research

49      With respect to both reliability and accuracy of judgment, it should be noted that in any given
50 inter-taster experiment, whether blind or open, a certain amount of measureable agreement will
51 occur on the basis of chance alone. Therefore, appropriate reliability statistics all present as
52 chance-corrected coefficients. This holds true quite irrespective of whether the statistics were
53 designed for nominal variables, such as binary wine judgments [2]; or the Sensitivity-Specificity
54 model-e.g., in an oenological context [1]; ordinal variables [3] ; or variables that are measured on
55 interval or ratio scales [4-6].
56
57      For binary variables, the level of agreement expected on the basis of chance alone is calculated
58 in the exact same manner as for the venerable and most familiar chi-square(d) statistic; and as
59 applied correctly by Cohen [2] in the development of his kappa statistic, which was recently
60 empirically verified [7].

## 3. Criteria for Assessing Levels of Practical Significance of the Reliability of Wine Judgments

62      There are currently three sets of published guidelines that were developed specifically for
63 assessing the degree of the clinical or practical significance of a binary diagnostic judgment, as
64 opposed to its level of statistical significance. In wine research it would seem useful to refer to the
65 term as oenological significance. Three sets of criteria have been published [8-12]. As one might
66 expect, the term clinical significance has its roots in bio-behavioral research, notably in nosology or
67 diagnostic specialty areas. Practical significance is also synonymous with the phrase *strength of*
68 *agreement* [8] and also with the concept of Effect Size (ES), as introduced by Cohen [13].

## 4. The Landis & Koch (1977); Fleiss (1981) and Cicchetti (1994); Oenological Criteria

70      The Landis & Koch guidelines [8] contain six ordinal categories of increasing gradations of
71 Strength of Agreement. These guidelines would seem particularly useful in an oenological context in
72 which wine experts were teaching less well experienced wine tasters to appreciate some of the
73 nuances of wine judgments and then testing their reliability levels with the wine experts, at specific
74 time points in the training exercise.
75      The Fleiss ~~et al.~~ guidelines [10] consist of three ordinal categories of clinical significance; they
76 would be applicable if the primary emphasis was to tri-chotomize wine judgments into unacceptable
77 (Poor); acceptable (Fair or Good) and highly acceptable (Excellent).
78      And, finally, the Cicchetti guidelines [12] consist of four ordinal categories of clinical
79 significance. It would be most applicable if one were to relate them to clinical diagnoses, as in a
80 nosological investigation of Autism [14] or in the present oenological research context. In comparing
81 the Fleiss, et al. guidelines to those of Cicchetti & Sparrow, the latter make a distinction between Fair
82 and Good, thereby forming four categories rather than three.
83      It should also be noted that because of the demonstrated equivalence between k, kw and the
84 ICC, the criteria apply regardless of the type of variable under investigation. First, Fleiss [16]
85 demonstrated the mathematical equivalence between Cohen's kappa statistic (k) for nominal binary
86 variables and the intra-class correlation coefficient (ICC) for variables deriving from interval scales;
87 and, secondly, Fleiss & Cohen [16] demonstrated the mathematical equivalence between Cohen's
88 weighted kappa coefficient [3] and the ICC [6]. This prompted Fleiss and colleagues to correctly
89 describe these three statistics as belonging to a family of mathematically inter-related coefficients.
90 An analogy in the broader bio-statistical world is the often cited mathematical equivalence between
91 the standard correlation coefficient (r) for interval variables and the Phi coefficient for
92 Nominal-dichotomous variables [17].
93      The three aforementioned sets of clinical/oenological criteria are given in Tables 1A, 1B and 1C.

94      **Table 1A.** The Landis & Koch (1977) Criteria for Assessing Oenological    Significance

| k, kw or ICC: | Strength of Agreement: |
|---|---|
| <0.00 | Poor |
| 0.00-0.20 | Slight |
| 0.21-0.40 | Fair |
| 0.41-0.60 | Moderate |
| 0.61-0.80 | Substantial |
| $\geq$ 0.80 | Almost Perfect |

95      **Table 1B.** The Fleiss (1981) Criteria for Assessing Oenological Significance

| k, $k_w$ or ICC: | Clinical Significance: |
|---|---|
| <0.40 | Poor |
| 0.40-0.74 | Fair to Good |
| > 0.75 | Excellent_ |

96      **Table 1C.** The Cicchetti & Sparrow (1981) Criteria for Assessing Oenological Significance

| k, $k_w$ or ICC: | Clinical Significance: |
|---|---|
| <0.40 | Poor |
| 0.40-0.59 | Fair |
| 0.60-0.74 | Good |
| >0.75 | Excellent |

97

98      In the next section of this report there will be a discussion of the relevance of a very early,
99   seminal, albeit seldom cited, publication, that nonetheless appears to have made a substantial
100   contribution to our knowledge of how best to understand levels of inter-taster agreement, or
101   agreement more broadly. It recalls in me the musically derived phrase referring to an familiar classic
102   as "an oldie but goodie." This contribution was made by a research Sociologist   William Robinson,
103   more than 60 years ago and was published in a prominent research Journal in his field, namely, the
104   Sociological Review [18]. Pertinent to this report, Robinson discovered a simple mathematical
105   relationship between what he referred to as the coefficient of agreement (A) and the
106   intra-class correlation coefficient (ICC).  What makes this additional Agreement statistic most
107   desirable, as we will see, is that it is very easy to compute and because of its mathematical
108   relationship to the ICC, which is a chance-corrected coefficient, A itself becomes a chance-corrected
109   coefficient.

110   **5. The Agreement or A index and its Mathematical Relationship to the ICC**

111      Suppose two wine tasters are asked to rate the quality of each of 200 wines, over a period of one
112   year; and their chance-corrected level of agreement produced an ICC value of 0.62 (Good)- [12] or
113   (Substantial)-[8]. If one desires to interpret the ICC as another agreement coefficient (A), how should
114   one proceed?
115      The mathematical relationship is given by the very simple formula introduced by Robinson [18]
116   as:
117      Agreement (A) = (ICC + 1)/2:                                                    (1)
118      Given our hypothetical ICC value of 0.62, Agreement (A) becomes 1.62/2= 0.81 or 81%.
119      In Table 2, the author shows the conversion of a given k, kw or ICC value into its Agreement
120   (A) equivalent. The relevance this type of thinking has for oenological research is explained in the
121   next section of this report.

122

123

124     **Table 3A.** Revised Landis & Koch Criteria [8] for Assessing Oenological Significance

125

126

| K, Kw or ICC Value: | Percent Agreement | Strength of Agreement: |
|---|---|---|
| <0.00 | < 50 | Poor |
| 0.00-0.19 | 50-59.5 | Slight |
| 0.20-0.39 | 60-69.5 | Fair |
| 0.40-0.59 | 70-79.5 | Moderate |
| 0.60-0.79 | 80-89.5 | Substantial |
| >0.80 | >90 | Almost Perfect |

134 **Table 3B.** Revised Fleiss, Levin & Cho Paik (2003) Criteria for Assessing Oenological
135 Significance

| K, Kw or ICC Value: | Percent Agreement: | Clinical Significance: |
|---|---|---|
| <0.40 | <70 | Poor |
| 0.40-0.79 | 70-89.5 | Fair to Good |
| $\geq$ 0.80 | $\geq$90 | Excellent |

142 **Table 3C.** Revised Cicchetti (1994) Criteria for Assessing Oenological Significance

| K, Kw or ICC Value: | Percent Agreement: | Clinical Significance: |
|---|---|---|
| <0.40 | <70 | Poor |
| 0.40-0.59 | 70-79.5 | Fair |
| 0.60-0.79 | 80-89.5 | Good |
| $\geq$ 0.80 | $\geq$ 90 | Excellent |

150     There is an additional bio-statistical fact that derives from Robinson's scientific contribution:
151 first, when any of the kappa coefficients is at its highest possible level (Case 1 in each of Tables 4, 5
152 and 6), then the level of specific agreement on both Positive and Negative judgments will both be
153 exactly equal to the overall Percentage of Observed agreement (PO).
154     In the context of clinical research, one earlier investigation had as its focus the accuracy of a
155 number of multiple regression techniques and neural networks (NN) for the binary diagnosis of
156 Autism. Each multiple regression technique (Logistic, Linear and Quadratic) produced more
157 accurate diagnostic results than did Neural Networks. Accuracy was assessed using the standard
158 Sensitivity-Specificity model1 whereby: <70%=Poor; 70%-79%=Fair; 80%-89%=Good; and
159 90%-100%=Excellent [14]. The reader will note that the same set of criteria are used by Robert Parker
160 and other putative experts to evaluate the quality of wine.
161     Two pertinent questions arise at this point in the narrative: First, what is the correspondence
162 between ICC values and Agreement across a broad and comprehensive spectrum of values? and
163 second, how does this information relate to the aforementioned sets of criteria defining levels of
164 oenological significance?
165     The answer to the first query appears in Table 2.
166
167
168

169 **Table 2.** The Correspondence between ICC and Percent Agreement [18][1]

| ICC Value: | Percent Agreement: |
|---|---|
| 0.00 (P) | 50 (P) |
| 0.05 (P) | 52.5 (P) |
| 0.10 (P) | 55(P) |
| 0.15 (P) | 57.5 (P) |
| 0.20 (P) | 60 (P) |
| 0.25 (P) | 62.5 (P) |
| 0.30 (P) | 65 (P) |
| 0.35 (P) | 67.5 (P) |
| 0.40 (F) | 70 (F) |
| 0.45 (F) | 72.5 (F) |
| 0.50 (F) | 75 (F) |
| 0.55 (F) | 77.5 (F) |
| 0.60 (G) | 80 (G) |
| 0.65 (G) | 82.5 (G) |
| 0.70 (G) | 85 (G) |
| 0.75 (G) | 87.5 (G) |
| 0.80 (E) | 90 (E) |
| 0.85 (E) | 92.5 (E) |
| 0.90 (E) | 95 (E) |
| 0.95 (E) | 97.5 (E) |
| 1.00 (E) | 100 (E) |

192 **1** Note: Because of the mathematical equivalencies between ICC, Kappa and Weighted Kappa, this
193 relationship holds for each of these three statistics for assessing levels of wine tasters' binary
194 judgments, as well as inter-rater agreement levels more generally. See text for more details. The
195 letters P, F, G, and E can refer, in this context to Poor, Fair, Good and Excellent wine quality as
196 defined by the Robert Parker and similar wine rating scales.

198 The answer to the second question appears next.

199 ## 6. Revising the Criteria for the Oenological Significance of Research Findings

200 In order to produce a correspondence between the aforementioned trifecta
201 of clinical significance criteria with the rating of the quality of wine by the Robert Parker or similar
202 scales, a few minor but oenologically significant changes need to be made in each of the three sets of
203 guidelines. It should be recalled that the Parker scale for rating the quality of wine is already
204 equivalent to the clinical criteria given by the aforementioned investigation by Cicchetti, et al., [14].

205 This minor revision process will be illustrated first with the Landis & Koch guidelines [8].
206 Because of the conceptual similarity between this triad of recommended guidelines, the same logic
207 will apply to the Fleiss guidelines [10] and also those published by Cicchetti [12].   If we now present
208 again the original Landis & Koch guidelines, we have the following**:**

| K, Kw, ICC | Agreement | Strength of Agreement |
|---|---|---|
| <0.00 | <50% | Poor |
| 0.00-0.20 | 50%-60% | Slight |
| 0.21-0.40 | 60.5%-70% | Fair |
| 0.41-0.60 | 70.5%-80% | Moderate |
| 0.61-0.80 | 80.5%-90% | Substantial |
| 0.81-1.00 | 90.5%-100% | Almost Perfect |

Note first that the Parker wine quality rating scale, as Percentages, defines below 70 as Poor; 70-79 as Fair; 80-89 as Good and 90 and above as Excellent. In contrast, each of the acceptable-wine-quality scores in the Landis & Koch guidelines appears at the end of each category rather than at its entry level [8]. By simply subtracting the number one from the Slight, Fair, Moderate and Substantial guidelines; and then combining the first two categories Poor and Slight, the revised Agreement categories become:   50-69=Poor; 70-79=Fair; 80-89=Good; and 90-100=Excellent, which, in this revised format, coincides exactly with the clinical criteria for bio-behavioral diagnoses [14], as well as, with the Parker quality of wine criteria.

Applying the same logic to the Fleiss, et al. guidelines, the Fair to Good category of k, kw, or ICC as 0.40 to 0.74 was changed to 0.40 to 0.79; and the last category was revised to define Excellent at > 0.80 instead of at > 0.75.

Finally, the Cicchetti criteria [12] required that the original category of 60 to 74, representing Good Agreement, be revised to 60-79; and that the final category defining Excellent as >  75 be replaced by > 80.

These revised criteria, with very minor changes, are now in line with both the aforementioned clinical guidelines [14] and the identical set of Parker criteria for judging the quality of wine. These revised criteria appear in Tables 2A, 2B and 2C.

Thus far, the focus has been on clinical, practical, or, in this context oenological significance. This is critical because a research result, oenological or otherwise, must have value beyond its level of statistical significance. It must also have clinical, practical or oenological significance to be worth pursuing further. Thus, the desideratum must be that a given scientific finding should not only occur beyond chance expectation, it must also not be a trivial finding. For a comprehensive discussion of this fundamental issue, the interested reader is referred to the scholarly work of Borenstein [19].

Thus far, the focus has been on the overall levels of inter-taster agreement or the overall level of chance-corrected agreement, again on an overall level. In the next part of this report, the issue of specific category agreement will be pursued.

## 7. Specific Category Agreement Levels

In the binary taster agreement context, one is referring to the agreement on positive and negative taster judgments. For example, let us suppose that the oenological researcher is investigating the reliability level of inter-taster agreement as to whether wines are oaked (+) or unoaked ( - ) and the overall agreement, based upon 100 wines, is 80 %; she wishes to proceed further and asks the question "What is the agreement on the oaked wines and the unoaked wines, treated separately?" Conceptually, overall agreement, as one might expect, is a weighted average of the agreement on positive and negative cases. In order to explain the phenomenon in greater detail, consider the hypothetical results of an oenological wine investigation in which, say, 2 experienced

255    wine Tasters are asked to decide whether 100 wines, evaluated over a period of six months, are
256    oaked or not. Suppose the results, in binary contingency table format, are as follows:

257                                **Taster B:**

258    **Taster A:**        **Oaked( + )**        **Unoaked( - )**        **Totals:**

259    **Oaked ( + )**          60                  20                  80

260    **Unoaked( - )**          0                  20                  20

261    **Total:**               60                  40                  100

262        Summing along the main diagonal, the overall level of Taster agreement is 80%. The agreement
263    on Positive cases is 60/(80+60)/2 or 60/70=85.7%; this is based on an average of (80+60)/2= 70 cases; the
264    agreement on Negative cases, correspondingly, is 20/(20+40)/2 or 20/30=66.7%; this derives from an
265    average of (20+40)/2=the remaining 30 cases. Finally: [(85.7 x .70) + (66.7x.30)] = (60+20)=80%.

### 8. The Sensitivity-Specificity Model in an Oenological Context

267        The relevance of the Sensitivity-Specificity model for studying the accuracy of Tasters' binary
268    judgments about wine was recently investigated [2]. Given its relevance for this report, it seems
269    pertinent to briefly allude to it once again. The five components of the model have their roots in
270    bio-behavioral diagnostic issues.
271        The five components of the Sensitivity-Specificity model are: Overall Accuracy (OA). This refers
272    to the percentage of correct binary judgments summed over both positive and negative cases. Thus if
273    there were Taster agreement on 42 of the Positive cases (the wines are oaked) and a corresponding
274    level of agreement on 38 of the Negative cases (the wine is unoaked), the overall agreement level
275    would be 80%. Sensitivity(Se) measures the percentage of filtered wines that are correctly judged as
276    such. If, of 48 wines known to be filtered, 42 were judged correctly by the Tasters, Se would be
277    calculated as 42/48=87.5%,Specificity (Sp) would indicate the percentage of unfiltered wines that are
278    correctly judged as such. Therefore, if 38 out of 52 wines were judged accurately to be unfiltered, Sp
279    would become 38/52=73%. Predicted Positive Accuracy (PPA) refers to the percentage of wines that
280    the Tasters judge to be filtered that are actually filtered. Thus, if 42 of 56 wines that are judged to be
281    filtered turn out to indeed be filtered, then PPA would become 42/56= 75%.
282        Predicted Negative Accuracy indicates the percentage of wines that the Tasters judge to be
283    unoaked that are actually unoaked. If this were true of 38 of 44 wines, then PPN would become
284    38/44=86%. We now turn to the issue of statistical significance.

### 9. Criteria for Assessing Levels of Statistical Significance

286        There are many statistical tests for establishing the level of statistical significance of a given
287    research finding; common among them are the **t** test, the F test and the Z test, which are all
288    mathematically related to each other.
289        As pertains to the current investigation, the statistical significance of a given kappa value is
290    found by dividing kappa by its standard error (SE)--[20], which produces a Z score, the size of
291    which, is directly translated into a probability (p) value which is interpreted in the usual way, as:
292    <± 1.96= Not Statistically Significant (NS);  ± 1.96 = 0.05; ± 2.58 =0.01; ± 3 =0.003; ±4 = <0.005; and
293    ±5 = <0.0001 [20,21].
294
295
296        Irrespective of which statistic is most appropriate to utilize, the objective is always to determine
297    whether a given research finding (oenological or otherwise) has occurred beyond chance
298    expectation. The standard definition of a chance finding is that it must have occurred at or less than 5

299  times in 100. Although criticized by some, this "Holy Grail" criterion for statistical significance has
300  withstood the test of time as it continues to be defined at the level of 0.05 probability (p).

301      Given the topic investigated here, the focus will be on binary wine tasting judgments, but for
302  reasons already given, the findings will also apply, conceptually, to other types of variables, ordinal,
303  interval or ratio. In two recent investigations, one clinical, the other, oenological, exceedingly high to
304  perfect correlations were found between the reliability and accuracy of binary judgments [1, 7].

305      These results are recast in an oenological format at the following levels of overall wine taster
306  agreement on a hypothetical binary variable, such as, whether a wine was oaked or not, with overall
307  hypothetical Taster agreement levels set at 70%, (Table 4); at 80% (Table 5) or at 90%   (Table 6). In
308  each of these three hypothetical oenological data sets, it was possible to control for a number of
309  variables that would be difficult if not impossible to control in the typical oenological investigation.
310  These variables were controlled at each hypothetical level of overall Taster agreement, whether 70%
311  (Fair), 80 % (Good); or 90% (Excellent), as the following: For OA=70%:

312      The patterns of agreement on Positive and Negative cases were set at: 35-35; 40-30; 45-25; 50-20;
313  55-15; 60-10; 65-5; and 70-0.

314      The numbers of disagreement cases (+ -) and (- + ) were each set at 15. This strategy served two
315  important research purposes: first to control or eliminate hypothetical wine taster bias; more
316  specifically, whenever there was a taster disagreement the first taster was just as likely as the second
317  taster to judge a disagreed upon wine as oaked or unoaked. This same design strategy was utilized
318  for the 80% and 90% condition. It should be noted here that very high levels of inter-taster bias have
319  been demonstrated in the judgments of wine experts such as Jancis Robinson and Robert Parker
320  [22, 23].

321      The outcome variables for each of the 70%, 80% and 90% conditions were the following: The
322  Percentage of agreement expected on the basis of Chance alone (PC); The levels of kappa (k) or
323  chance-corrected agreement; The levels of agreement on both Positive -e.g., or Negative cases, for
324  example, the wine is oaked (+) or the wine is unoaked ( - ).

325      The absolute difference between agreement on Positive and Negative wine Tasting judgments,
326  whereby 0 difference = 100% agreement, and maximum possible disagreement would then be 0%
327  agreement; and The final column in each of the three Tables contains the p values for each kappa
328  value.

329  For OA=80%

330      The patterns of agreement on Positive and Negative cases were set at: 40-40; 45-35; 50-30; 55-25;
331  60-20; 65-15; 70-10; 75-5; and 80-0. The numbers of disagreement cases (+ -) and (- +) were each set at
332  10.

333  For OA=90%:

334      The patterns of agreement on Positive and Negative cases were set at: 45-45; 50-40; 55-35; 60-30;
335  65-25; 70-20; 75-15; 80-10; 85-5; and 90-0. The numbers of disagreement cases (+ -) and (- + ) were each
336  set at 5.

337

**Table 4.** Relationship between the Reliability and Accuracy of Pairs of Hypothetical Tasters Judging Whether a Wine is Oaked (+) or Unoaked ( - ) When the Tasters are in 70% Agreement

| Case: | (++) | (--) | (+-) | (-+) | PC | Kappa[1] | PO+ | PO- | Agreement | PO+/PO- p value[3] |
|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 35 | 35 | 15 | 15 | 50 | 0.40 (F) | 70 (F) | 70 (F) | 100 | <0.0005 |
| 2 | 40 | 30 | 15 | 15 | 50.5 | 0.39 (P) | 67 (P) | 73 (F) | 94 | 0.002 |
| 3 | 45 | 25 | 15 | 15 | 52 | 0.375 (P) | 62.5 (P) | 75 (F) | 87.5 | 0.004 |
| 4 | 50 | 20 | 15 | 15 | 54.5 | 0.34 (P) | 57 (P) | 77 (F) | 80 | 0.01 |
| 5 | 55 | 15 | 15 | 15 | 58 | 0.29 (P) | 50 (P) | 79 (F) | 71 | NS[1] |
| 6 | 60 | 10 | 15 | 15 | 62.5 | 0.20 (P) | 40 (P) | 80 (G) | 60 | NS |
| 7 | 65 | 5 | 15 | 15 | 68 | 0.06 (P) | 25 (P) | 81 (G) | 44 | NS |
| 8 | 70 | 0 | 15 | 15 | 74.5 | -0.18 (P) | 0 (P) | 82 (G) | 18 | NS |

The correlation between the size of kappa and the difference in agreement on Positive and Negative cases is +0.98;

[1] Kappa values are classified as Poor (P), Fair (F), Good (G) or Excellent (E) by the revised Cicchetti criteria in Table 3C. [2] **NS**=not statistically significant at $p \leq 0.05$. [3] Statistical significance is found by dividing kappa by its standard error as derived by Fleiss, Cohen & Everitt, **[20]**. Values of Z are interpreted in the standard manner whereby: $<\pm 1.96 = p$ at the 0.05 level; $\pm 2.58$ is at t 0.01; $\pm 3$ at 0.003; $\pm 4$ at 0.0005; and $\pm 5$ at .0001 **[20, 21]**.

**Table 5.** Relationship between the Reliability and Accuracy of Pairs of Hypothetical Tasters Judging Whether a Wine is Filtered (+) or Not Filtered (-) When the Tasters are in 80% Agreement

| Case: | (++) | (--) | (+-) | (-+) | PC | Kappa | PO+ | PO- | Agreement | PO+/PO- p value[3] |
|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 40 | 40 | 10 | 10 | 50 | 0.60 (G) | 80 (G) | 80 (G) | 100 | <0.0005 |
| 2 | 45 | 35 | 10 | 10 | 50.5 | 0.60 (G) | 82 (G) | 78 (F) | 96 | <0.0005 |
| 3 | 50 | 30 | 10 | 10 | 52 | 0.58 (F) | 83 (G) | 75 (F) | 92 | 0.001 |
| 4 | 55 | 25 | 10 | 10 | 55 | 0.56 (F) | 85 (G) | 71 (F) | 86 | <0.005 |
| 5 | 60 | 20 | 10 | 10 | 58 | 0.52 (F) | 86 (G) | 67 (P) | 81 | <0.005 |
| 6 | 65 | 15 | 10 | 10 | 63 | 0.47 (F) | 87 (G) | 60 (P) | 73 | <0.005 |
| 7 | 70 | 10 | 10 | 10 | 68 | 0.38 (P) | 88 (G) | 50 (P) | 62 | 0.01 |
| 8 | 75 | 5 | 10 | 10 | 74.5 | 0.22 (P) | 88 (G) | 33 (P) | 45 | NS |
| 9 | 80 | 0 | 10 | 10 | 82 | -0.11 (P) | 89 (G) | 0 (P) | 11 | NS |

The correlation between the size of kappa and the difference in agreement on Positive and Negative cases is +0.99; [1] Kappa values are classified as Poor (P), Fair (F), Good (G) or Excellent (E) by the revised Cicchetti criteria in Table 3C. [2] **NS**=not statistically significant at $p \leq 0.05$. [3] Statistical significance is found by dividing kappa by its standard error as derived by Fleiss, Cohen & Everitt, (1969). Values of Z are interpreted in the standard manner whereby: $<\pm 1.96 = p$ at the 0.05 level; $\pm 2.58$ is at t 0.01; $\pm 3$ at 0.003; $\pm 4$ at 0.0005; and $\pm 5$ at .0001. **[20, 21]**.

**Table 6.** Relationship between the Reliability and Accuracy of Pairs of Hypothetical Tasters Judging Whether a Wine is Filtered (+) *or* Not Filtered ( - ) When the Tasters are in 90% Agreement

| Case: | (++) | (--) | (+-) | (-+) | PC | Kappa | PO⁺ | PO⁻ | (PO⁺-PO⁻) | *p*-value |
|---|---|---|---|---|---|---|---|---|---|---|
| 1 → | 45 | 45 | 5 → | 5 | 50 | 0.80 (E) | 90 → | 90 → | 100 | <0.0005 |
| 2 → | 50 | 40 | 5 | 5 | 51 | 0.80 (E) → | 89 → | 91 → | 98 → | <0.0005 |
| 3 → | 55 | 35 | 5 | 5 | 52 | 0.79 (G) | 88 | 92 | 96 | <0.0005 |
| 4 → | 60 | 30 | 5 | 5 | 55 | 0.78 (G) | 86 → | 92 | 94 → | <0.0005 |
| 5 → | 65 | 25 | 5 | 5 | 58 | 0.76 (G) → | 83 → | 93 → | 90 → | <0.0005 |
| 6 → | 70 | 20 | 5 | 5 | 63 | 0.73 (G) | 80 → | 93 → | 87 → | <0.0005 |
| 7 → | 75 | 15 | 5 | 5 | 68 | 0.69 (G) → | 75 → | 94 → | 81 → | <0.0005 |
| 8 → | 80 | 10 | 5 | 5 | 75 | 0.61 (G) | 67 | 94 | 73 → | <0.0005 |
| 9 → | 85 | 5 | 5 | 5 | 82 | 0.44 (F) | 50 | 94 → | 56 → | 0.001 |
| 10 | 90 | 0 | 5 | 5 | 90.5 | -0.05 (P) | 0 → | 95 → | 5 → | NS |

[1]The correlation between the size of kappa and the difference in agreement on Positive and Negative cases is +1.00; [1] Kappa values are classified as Poor (P), Fair (F), Good (G) or Excellent (E) by the revised Cicchetti criteria in Table 3C.   [2] **NS**=not statistically significant at $p \leq 0.05$. .[3] Statistical significance is found by dividing kappa by its standard error as derived by Fleiss, Cohen & Everitt **[20]**. Values of Z are interpreted in the standard manner whereby: $<\pm$ 1.96=*p* at the 0.05 level; $\pm$ 2.58 is at t 0.01; $\pm$ 3 at 0.003; $\pm$ 4 at 0.0005; and $\pm$5 at 0.0001 **[20, 21].**

The advantage of the hypothetical information revealed in these three tables is that they allow for a degree of experimental control that is seldom or almost never possible in the typical oenological study or in clinical research more generally. The method of Hypothetics, or Oenothetics in this context, allows the research scientist to produce the results that would have occurred if the actual experiments they represent were feasible. The general findings will precede those occurring on a case by case basis, separately for the 70%, 80% and 90% condition.

## 10. Overall Results: Correlations between the Reliability and Accuracy of Wine Tasters' Hypothetical Binary Judgments

As we examine the results deriving from Tables 4, 5 and 6, it should be noted that the correlations between reliability and overall accuracy or validity of Tasters' hypothetical binary wine judgments is exceptionally high, that is, almost perfect to completely perfect. This holds true whether the overall Taster agreement levels were expressed at 70% (Fair); 80% (Good) or 90% (Excellent). The three correlation are, respectively, + 0.98, +0.99 and +1.00.

An advantage of using the standard correlation coefficient to measure the relationship between the reliability and accuracy/validity of hypothetical Tasters' judgments is that it provides a familiar and easy-to-interpret result. A major disadvantage is that the correlation coefficient is an omnibus statistic that provides no information about reliability and accuracy of judgment on a case by case basis, as would be true of individual kappa coefficients or the components of the Sensitivity-Specificity model in whatever clinical or other research context.

## 11. Hypothetical Results on a Case by Case Basis

With respect to the hypothetical data in Table 4 (patterns of 70% agreement), the Case 1 result indicates both oenological and statistical significance; and Cases 5 through 8 indicate results that are neither oenologically nor statistically significant; however, Cases 2, 3 and 4 produce findings that are statistically significant but not oenologically significant.

The results for the 80% condition, as spread in Table 5, show that the first nine Cases yield results that are both oenologically and statistically significant, while the tenth Case indicates a result that is neither oenologically nor statistically significant.

393   The data for the 90% condition indicates that the first nine Cases produced results that are both
394   oenologically and statistically significant while the result for the tenth Case was neither
395   oenologically nor statistically significant.
396   Taken as a whole these results are consonant with two research results that occur in both the
397   oenological and clinical world of science: first, given an appropriate sample size even the most trivial
398   of results will be statistically significant; and secondly, the greater the level of agreement, the more
399   likely the result is apt to be both statistically significant and of material importance, whether
400   oenological, clinical, or otherwise.
401   One way to provide more specific information on a case by case basis is to summarize the data
402   from Tables 4, 5 and 6 in a single table as follows: The hypothetical information for the 70%
403   condition was based upon 8 cases; the 80% condition on 9 cases; and the 90% condition was based
404   upon an additional 10 cases. These sum to 27 cases in all.
405   If one now recasts the data into a 2 x 2 or binary Table, it will then be possible to perform the
406   kappa statistic, to measure the level of hypothetical Taster reliability as well as to obtain the 5
407   accuracy components of the Sensitivity-Specificity model. The recast data appear in Table 7.

408   **Table 7.** Illustrating the Relationship between the Reliability and Accuracy of Wine Tasters'
409   Hypothetical Binary Judgments of Whether a Wine is Oaked (+) or Not Oaked ( - ), Expressed in
410   Percentages

411   _____

|            | **Taster 2:** |        |            |
|------------|---------------|--------|------------|
| **Taster 1:** | **( + )**  | **( - )** | **Totals:** |
| **( + )**  | 12            | 4      | 16         |
| **( - )**  | 0             | 11     | 11         |
| **Totals:** | 12           | 15     | 27         |

### 12. Summary and Conclusions

419   Utilizing a new methodology called Hypothetics, or Oenothetics, in a wine tasting
420   investigation, a model was introduced that makes it possible to control for a large number of
421   variables that are often most difficult to control in the typical oenological study, beverage study or
422   more generally. Because of this level of control, the method allows for findings and insights that are
423   often not possible using available standard methodologies and standard data analytic strategies. In
424   this fundamental sense, the hypothetical results that were obtained appear to have heuristic value
425   for the design of future oenological studies and investigations focusing on beverages more
426   generally.
427   In this application, which focused upon the oenological and statistical significance of wine
428   Tasters' binary judgments, the following occurred: Results that were oenologically significant
429   (had practical meaning) were uniformally statistically significant, although the reverse was not
430   always true, that is to say, a number of results were statistically significant, but not oenologically
431   important. A method developed more than six decades ago [19] was shown to simplify the
432   understanding of chance corrected agreement coefficients in any given oenological or other type of
433   scientific investigation. Finally, the correlation between the reliability and validity or accuracy of
434   binary judgments was shown to be exceedingly high whether on an overall omnibus level; or on a
435   case by case basis. With appropriate adjustments, this methodology would apply to ordinal and
436   interval variables, as well.

## References

1. Cicchetti DV. Opinions versus facts: A bio-statistical paradigm shift in oenological research. Proc J Wine Res. 2017; 1: 1-8.

2. Cohen J. A coefficient of agreement for nominal scales. Educ Psychol Meas. 1960; 23: 37-40

3. Cohen J. Weighted kappa: nominal scale agreement with provision for scaled disagreement or partial credit. Psychol Bull. 1968; 70: 195-201.

4. Bartko JJ. The intraclass correlation coefficient as a measure of reliability. Psychol Rep. 1966; 19:3-11.

5. Bartko JJ. Corrective note to "The intraclass correlation coefficient as a measure of reliability." Psychol Rep. 1974; 34: 1-11.

6. Shrout P E, Fleiss J. Intraclass correlations: Uses in assessing rater reliability. Psychol Bull. 1979; 86:420–428.

7. Cicchetti DV, Klin A, Volkmar FR. Assessing binary diagnoses of bio-behavioral disorders: The clinical relevance of Cohen's Kappa. JNMD. 2017; 205: 58-65.

8. Landis JR, Koch GG.The measurement of observer agreement for categorical data. Biom. 1977; 3:159-174.

9. Fleiss J. Statistical methods for rates and proportions. 1981; New York: Wiley (2nd ed).

10. Fleiss J, Levin B, Paik, MC. Statistical methods for rates and proportions. 2003; New York: Wiley (3rd ed).

11. Cicchetti DV, Sparrow SS. Developing criteria for establishing interrater reliability of specific items: Applications to assessment of adaptive behavior. Am J Men Defic. 1981; 86: 127-137.

12. Cicchetti DV. Guidelines. criteria, and rules of thumb for evaluating normed and standardized assessment instruments in psychology. Psychol Assess. 1994; 6:284-290.

13. Cohen J. Statistical power analysis for the behavioral sciences. 1988; Hillsdale, NJ: Erlbaum (2nd ed.).

14. Cicchetti DV, Volkmar FR, Klin   A, Showalter D. Diagnosing autism using ICD-10 criteria: A comparison of neural networks and standard multivariate procedures. Child Neuropsychol. 1995; 1: 26-37.

15. Fleiss J. Measuring agreement between two judges on the resence or absence of a trait. Biom. 1975; 31: 651-659.

16. Fleiss JL, Cohen J. The equivalence of weighted kappa and the intraclass correlation coefficient as measures of reliability. Educ Psychol Meas. 1973; 33: 613-619.

17. Kaltenhauser J, Lee Y. Correlation coefficients for binary data. Geogr Anal. 2010; 8: 305-313.

18. Robinson W. The statistical measurement of agreement. Am Sociol Rev. 1957; 22:17-25.

19. Borenstein M. The shift from significance testing to effect size estimation. Res Meth Compr Clin Psychol. 1998; 3: 319-349.

20. Fleiss JL, Cohen J, Everitt BS. Large sample standard errors of kappa and weighted kappa. Psychol Bull.,1969; 72:323-327.

21. Cohen J, Cohen P, West SG, Aiken IS. Appliede multiple regression/correlation for the behavioral sciences. 2003; Mahwah NJ: Lawrence Erlbaum (2nd ed).

22. Cicchetti DV, Cicchetti AF. As wine experts disagree, consumers' taste buds flourish: How two experts rate the 2004 Bordeaux vintage. J Wine Res. 2013; 24:311-317.

23. Cicchetti DV, Cicchetti AF. Two enological titans rate the 2009 Bordeaux wines. Wine Econ Policy. 2014; on line publication, http://dx.doi.org/10.16/j.wep.2014.01.001.

24. DiDonfrancesco B, Guitierrez Guzman N, Chambers E. Comparison of results from cupping and descriptive sensory analysis of Columbian brewed coffee. J Sens Anal. 2014; 29:301-311.