

Testing the accuracy of automated classification systems using only expert ratings that are less accurate than the system

Paul Lehner
The MITRE Corporation

Abstract

A method is presented to estimate the accuracy of automated classification systems using only expert ratings that may be substantially less accurate than the systems being evaluated. The estimation method begins with multiple expert ratings on test cases, uses the level of inter-rater agreement to estimate rater accuracy, uses Bayesian updating based on estimated rater accuracy to estimate a “ground truth” probability for each classification, and then estimates system accuracy by comparing the relative frequency that the system agrees with the most probable classification at different probability levels. A simulation analysis provides evidence that the method is robust and yields reasonable estimates of system accuracy under diverse and predictable conditions.

Acknowledgment: This research was supported by the Intelligence Advanced Research Projects Activity (IARPA). The U.S. Government is authorized to reproduce and distribute reprints for Governmental purposes notwithstanding any copyright annotation thereon.

Disclaimer: The views and conclusions contained herein are those of the author and should not be interpreted as necessarily representing the official policies or endorsements, either expressed or implied, of IARPA or the U.S. Government.

Testing the accuracy of automated classification systems using only expert ratings that are less accurate than the system

INTRODUCTION

Information technology is advancing to develop systems that address problems of increasing sophistication and complexity. New systems are being developed to address complex problems as diverse as automated medical and clinical diagnoses, technology readiness evaluation, detection of emerging technologies, classification of the contents of unstructured video segments, recognition and classification of metaphors used in natural language text and many others.¹

The complexities of the problems that these advanced systems address make it difficult to evaluate the accuracy of such systems. It is usually necessary to resort to using expert raters to establish the ground truth of test cases. However, the fact that these systems are addressing complex problems also presents a challenge to the expert raters. Expert raters often disagree as to the correct answer. Furthermore as future systems address problems of ever increasing sophistication and complexity, it seems likely that the experts will be even more challenged and exhibit even lower levels of agreement. Ground truth data sets based on expert judgment are fallible and are likely to become more so in the future.

Using expert raters to establish the ground truth of test cases is certainly not a new practice. For classification problems, which are the focus of this paper, it is a common scientific practice to measure the level of agreement amongst raters, with a statistic such as Kappa, and to refine the rating process until a satisfactory level of agreement is reached. Once the agreement threshold is reached, then the judgments of individual raters or collaborating teams of raters are treated as the ground truth (See Gwet, 2010 for review).

For several reasons, this common scientific practice does not adequately meet the needs of advanced system evaluation.

First, the level of agreement among raters will rarely meet a satisfactory level. The problems that these systems address are simply too complex. About the only way to increase the level of agreement is to select simple, non-representative test cases.

Second, estimating system accuracy by measuring the level of agreement with expert raters makes the *de facto* assumption that the experts are more accurate than the system. This assumption runs contrary to a substantial body of empirical research where it is very often found that simple algorithms

¹ A quick look at the web sites for the Intelligence Advanced Research Projects Activity Office of Incisive Analysis (http://www.iarpa.gov/office_incisive.html), the Defense Advanced Research Projects Agency (<http://www.darpa.mil/>), the National Science Foundation (<http://www.nsf.gov/>), and the National Institutes of Health (www.nih.gov) will reveal many such development efforts.

outperform human experts in complex judgments (Dawes, 1979; Grove, et. al. 2001, Tetlock 2005). It should not be presumed that the experts are more accurate than the system.

Third, there is considerable evidence to suggest that for a wide variety of judgment tasks collaborative team judgments are not substantially more accurate than the judgments of a randomly selected individual team member (e.g., Surowiecki, 2005; Armstrong, 2006). In judgment tasks, where there is no obvious correct answer, it should not be presumed that collaboration will reliably lead the raters to converge on the correct answer.

Finally, when evaluating a classification system the statistic of greatest interest is the accuracy of the system - the proportion of system assignments that are correct. Unfortunately the relationship between inter-rater reliability statistics such as Kappa, the probability of correct ground truth assignments and the accuracy of any systems tested against error-prone ground truth assignments is unclear.

This paper presents a different approach to using expert ratings to estimate the accuracy of complex systems. Rather than treat expert ratings as a surrogate for ground truth, expert ratings are treated as error prone estimates of ground truth where independent ratings are fused to generate estimated ground truth probabilities, and the ground truth probabilities are then used to estimate system accuracy.

This paper makes several strong claims for this approach to estimating the accuracy of classification systems.

First, under diverse conditions, this approach reliably yields estimates of system accuracy that are approximately correct. If the system is 90% accurate then this approach will yield an estimate of system accuracy that is close to 90%.

Second, the accuracy of the estimate of system accuracy is largely independent of whether the expert raters are more or less accurate than the system. If the system is in fact 90% accurate, and the raters are individually 60% accurate, then the estimate of system accuracy will still be approximately 90%.

Third, reliable estimates of system accuracy can often be obtained with reasonable sample sizes. Some of the simulation runs shown below use a sample size of just 50 test cases with three independent raters. In complex domains it is important to keep sample sizes as small as possible, since it may be time consuming and costly to obtain expert ratings.

Fourth, and importantly, the conditions under which the above claims may break down are predictable. Therefore test data sets can be intentionally constructed to ensure that the conditions are met that are needed for accurate estimation of system accuracy.

The objective of this paper is to present evidence for the above four claims. This is achieved by describing one practical instantiation of this approach, along with sufficient test results to present clear evidence for the first three claims and some evidence for the fourth claim.

A METHOD FOR ESTIMATING THE ACCURACY OF SYSTEM CLASSIFICATIONS

The method for estimating accuracy used in this paper was derived from the following assumptions.

- A1. For each case there is a unique correct classification.
- A2. For each case raters independently assign classifications.
- A3. Expected agreement between raters increases as expected rater accuracy increases.

Assumption A3 refers to *expected* agreement and accuracy. Here “accuracy” refers to the total proportion of correct classifications made by all the raters, irrespective of which raters are making correct and incorrect classifications. And “agreement” refers to the total proportion of pairwise agreement among all of the raters and cases. For any particular set of cases, accuracy may be low yet agreement high (the raters made the same mistakes), but A3 asserts that *in general* there is an expected positive relationship between accuracy and agreement.

Theorem 1: A1 – A3 are ensured if and only if the raters behave as though their selection for each case is determined by a single confusion matrix where the conditional probability of correct assignment is constant and the conditional probability of all incorrect assignments is equal.

That is to say all raters on all problems are selecting from a single confusion matrix with a structure such as shown in Table 1.

Table 1: Implied Structure of Rater Confusion Matrices for Four Category Problem
(A through D are true category and “A” through “D” are selected categories.)

	“A”	“B”	“C”	“D”
A	P_c	$(1-P_c)/3$	$(1-P_c)/3$	$(1-P_c)/3$
B	$(1-P_c)/3$	P_c	$(1-P_c)/3$	$(1-P_c)/3$
C	$(1-P_c)/3$	$(1-P_c)/3$	P_c	$(1-P_c)/3$
D	$(1-P_c)/3$	$(1-P_c)/3$	$(1-P_c)/3$	P_c

The proof of this theorem is found in the Appendix. The general structure of the proof shows that if the raters are assigning classifications using any process other than selecting from a common confusion matrix with the structure illustrated in Table 1, then it is always possible to construct a classification process with lower expected accuracy and higher expected agreement, or higher accuracy and lower agreement; thereby violating the assumed monotonic relationship between expected accuracy and expected agreement.

A1 through A3 also seem to be implicitly assumed in many contexts where the Kappa statistic is applied. Indeed it is A3 that would seem to warrant the common practice of using expert ratings as surrogates for ground truth when high levels of inter-rater agreement are found. Consequently it is reasonable to claim that the estimation method described below is derived from assumptions implicit in the Kappa statistic and how Kappa is often used. Because of this relationship to the Kappa statistic, in the remainder of this paper A1 – A3 will be referred to as *K-assumptions*. Furthermore, the properties of equal rater accuracy, equal error probabilities and equal problem difficulty that are implied by the *K-assumptions* will be referred to as *K-properties*.

Table 2: Sample data of expert ratings and system assignments for 10 test cases

Case #	Rater 1	Rater 2	Rater 3	Rater 4	System
1	"C"	"D"	"C"	"C"	"A"
2	"B"	"D"	"C"	"C"	"C"
3	"C"	"C"	"D"	"C"	"C"
4	"B"	"B"	"D"	"D"	"B"
5	"A"	"B"	"B"	"B"	"B"
6	"C"	"B"	"D"	"A"	"A"
7	"A"	"A"	"A"	"A"	"A"
8	"A"	"D"	"B"	"C"	"C"
9	"D"	"B"	"A"	"A"	"D"
10	"A"	"D"	"A"	"B"	"B"

The estimation method is straightforward to explain in the context of an example. Consider the test data in Table 2. There are 10 test cases, 4 categories, 4 raters and the system's proposed answers. When referring to ground truth label the four categories are labeled A, B, C, D; when referring to rater and system assignments they are labeled "A", "B", "C", "D".

As described below the estimation method is composed of four basic steps.

Estimate rater accuracy

Given that each rater has an identical confusion matrix, with the structure shown in Table 1, the probability that two raters will agree on any one case is

$$P_a = P_c^2 + \frac{(1-P_c)^2}{N-1}$$

Here P_a is the probability of agreement, P_c is the probability that a rater will make the correct assignment, and N is the number of categories.

Solving for P_c yields

$$P_c = \left(\frac{1}{N}\right) + \sqrt{\left(\frac{(N-1)*P_a - \frac{N-1}{N}}{N}\right)} \quad (\text{Eq. 1})$$

Eq. 1 is used to estimate rater accuracy. In the 10 cases in Table 1 there was 33% agreement (20 pairs out of 60). Setting P_a to .33 and solving for P_c yields $P_c = 0.5$; which is the estimate of rater accuracy.

Estimate base rates

The probability that a rater will assert a category, say "A", is as follows:

$$P("A") = P("A"|A) * P(A) + \left(1 - \frac{P("A"|A)}{N-1}\right) * (1 - P(A))$$

Here $P("A")$ is the marginal probability that the rater asserts "A", $P("A"|A)$ is the conditional probability that the rater will assert "A" if A is true, and $P(A)$ is the marginal probability of A.

Solving for $P(A)$ yields

$$P(A) = \frac{(N-1)*P("A") - 1 + P("A"|A)}{N*P("A"|A) - 1}$$

Setting $P("A")$ to be the observed relative frequency of "A", and $P("A"|A)$ to be the estimate of P_c from above, yields

$$P(A) = \frac{(N-1)*P("A") - 1 + P_c}{N*P_c - 1} \quad (\text{Eq. 2})$$

Eq. 2 is used to estimate the base rate for each category by setting P_c to be the estimate from above and $P("X")$ to be the observed relative frequency across all raters and ratings that category X was assigned. In Table 1 there are 11 instances of each of the categories; so the estimated base rate is 0.325 for category A. Applying Eq. 2 to the other categories yields base rates of 0.25, 0.25 and 0.175 for B, C and D respectively.

Estimate Ground Truth Probabilities

Use Bayes rule, assuming conditional independence for each rater, to estimate ground truth probabilities. For example, in case 1 above the raters selected "CDCC". So for each possible ground truth value calculate $P(\dots|"CDCC")$ and normalize.

$$P(A|"CDCC") \sim P(A)*P("C"|A)*P("D"|A)*P("C"|A)*P("C"|A) = .325*.167*.167*.167*.167 = .00025 \rightarrow .041$$

$$P(B|"CDCC") \sim P(B)*P("C"|B)*P("D"|B)*P("C"|B)*P("C"|B) = .25*.167*.167*.167*.167 = .00019 \rightarrow .032$$

$$P(C|CDCC) \sim P(C) \cdot P("C"|C) \cdot P("D"|C) \cdot P("C"|C) \cdot P("C"|C) = .25 \cdot .5 \cdot .167 \cdot .5 \cdot .5 = .00521 \rightarrow .860$$

$$P(D|CDCC) \sim P(D) \cdot P("C"|D) \cdot P("D"|D) \cdot P("C"|D) \cdot P("C"|D) = .175 \cdot .167 \cdot .5 \cdot .167 \cdot .167 = .00041 \rightarrow .067$$

Repeating this step for the other 9 cases yields the estimated probability distributions shown in Table 3.

Table 3: Estimated ground truth probabilities for sample data

Case #	Ground Truth Probability				System Answer
	A	B	C	D	
1	0.041	0.032	0.860	0.067	"A"
2	0.084	0.195	0.584	0.136	"C"
3	0.041	0.032	0.860	0.067	"C"
4	0.074	0.511	0.057	0.358	"B"
5	0.120	0.828	0.031	0.021	"B"
6	0.325	0.250	0.250	0.175	"A"
7	0.975	0.009	0.009	0.006	"A"
8	0.325	0.250	0.250	0.175	"C"
9	0.657	0.169	0.056	0.118	"D"
10	0.657	0.169	0.056	0.118	"B"

Estimate System Accuracy

Assume any probability distribution over the categories for each test case. For any test case, let P_g be the probability of the classification with the highest probability, P_s be the probability that the system will assign the correct answer, P_a be the probability that the system will assign the same classification as the highest ground truth probability. It follows that ²

$$P_a = P_g * P_s + (1 - P_g) * \frac{1 - P_s}{N - 1}$$

Solving for P_s yields

$$P_s = \frac{(N - 1) * P_a - 1 + P_g}{N * P_g - 1} \tag{Eq. 3}$$

Eq. 3 is used to estimate system accuracy as follows. First separate the test cases into bins with approximately the same highest estimated ground truth probability. In this paper the ranges (.9, 1.0], (.8, .9], (.7, .8], etc. are used. For example, in Table 3 there is one case in the (.9, 1.0] range, 3 cases in the (.8, .9] range, 2 cases in the (.6, .7] range, etc. Second for each bin calculate the average ground

² Note that this relationship holds whether or not the classification with the highest probability is correct.

truth probability within the bin; record the proportion of system assignments that agree with the most probable answer; then estimate system accuracy for each bin using equation Eq. 3. Third estimate overall system accuracy by taking the average of the estimated accuracy in each bin weighted by the number of cases in each bin. This is shown in Table 4.

Table 4: Estimation of System Accuracy for Sample Data

Probability Bin	Average Ground Truth Probability	Number in Bin	Proportion of Agreement	Estimated Accuracy
.9 – 1.0	0.975	1	1.000	1.000
.8 - .9	0.849	3	0.667	0.776
.6 - .7	0.657	2	0.000	0.000
.5 - .6	0.548	2	0.333	0.452
.2 - .3	0.325	2	0.500	1.000
Weighted Average =				0.731

The reader may be curious as to why the estimate of system accuracy is not simply the average of the estimated ground truth probabilities for the system answers. The reason is that taking the average will consistently underestimate system accuracy; because the system answer is additional evidence for each category. So, for example, if the system answer is “C” and the estimated ground truth probability for C is 0.6; then a better estimate for C would be somewhat higher than .6. But until system accuracy is estimated it cannot be determined how much more than .6 is appropriate. In the above example, the average estimated ground truth probability of the system answers is .466, but the estimate of system accuracy in Table 4 is 0.731.

Finally note that the value of Kappa for the data in Table 2 is

$$\begin{aligned}
 \text{Kappa} &= (\text{Observed Agreement} - \text{Random Agreement}) / (1.0 - \text{Random Agreement}) \\
 &= (.333 - .25) / (1 - .25) \\
 &= 0.11^3
 \end{aligned}$$

Standard thresholds normally require a level of Kappa = 0.7 before the expert ratings are considered usefully reliable (Landis and Koch, 1977; Gwet 2010). Kappa = 0.11 is considered “slight agreement” and is far too low for the ratings to be considered useful for establishing ground truth.

Overall then, in the sample data provided in Table 2; inter-rater agreement is “slight” (Kappa = 0.11), estimated rater accuracy is 0.50, and estimated system accuracy is 0.731.

³ There is more than one way to calculate Kappa, depending on how one calculates ‘Random Agreement’. In this paper random agreement is always set to 1/N.

PERFORMANCE AND ROBUSTNESS

The above example illustrates how to estimate system accuracy for classification problems even when inter-rater agreement and estimated rater accuracy are very low. This section examines the accuracy of estimates of system accuracy, and the robustness of those estimates, through a series of simulations.

All of the simulations described below use the following procedure to assign the confusion matrix for each rater and the system, based on the values set to four parameters: an initial probability of correct assignment (IPC), a problem difficulty adjustment (PDA), degree of asymmetric dispersion (AD), and a proportional error range (PER).

Each confusion matrix is constructed as follows:

1. Initially assign the conditional probability of a correct classification to be IPC for all categories.
2. Add PDA to the conditional probabilities of correct assignment.
3. For each category distribute the remaining probability ($1 - IPC - PDA$) to the incorrect classifications in a manner that is proportional to the distance from the correct classification, where the probability of a classification that is M steps removed from the correct classifications is AD times more likely than a classification that is M+1 steps removed.
4. For each conditional probability of incorrect assignment (IC) set the range to be $[IC - PER*IC, IC + PER*IC]$, then randomly select a new probability by uniform sampling over this range.
5. Normalize the modified confusion matrix after the random changes in step 4 so that expected accuracy is equal to $IPC + PDA$.

For example, if there are five categories and $(IPC, PDA, AD, PER) = (.6, 0, 1.0, 0)$, then the resulting confusion matrix is shown in Table 5.

Table 5: Confusion matrix where $(IPC, PDA, AD, PER) = (0.6, 0, 1.0, 0)$

Correct Category	Classification				
	"A"	"B"	"C"	"D"	"E"
A	0.6	0.1	0.1	0.1	0.1
B	0.1	0.6	0.1	0.1	0.1
C	0.1	0.1	0.6	0.1	0.1
D	0.1	0.1	0.1	0.6	0.1
E	0.1	0.1	0.1	0.1	0.6

On the other hand, if $(IPC, PDA, AD, PER) = (.6, -.2, 2.0, 1.0)$, then the confusion matrix after the first three steps would be as shown in Table 6.

Table 6: Confusion matrix where (IPC, PDA, AD, PER) = (0.6, -0.2, 2.0, 0)

Correct Category	Classification				
	"A"	"B"	"C"	"D"	"E"
A	0.400	0.320	0.160	0.080	0.040
B	0.218	0.400	0.218	0.109	0.055
C	0.100	0.200	0.400	0.200	0.100
D	0.055	0.109	0.218	0.400	0.218
E	0.040	0.080	0.160	0.320	0.400

Then after adding random variation around the incorrect probability assignments in step 4, and renormalizing in step 5, the resulting confusion matrix would look something like the randomly generated confusion matrix shown in Table 7.

Table 7: Example of randomly generated confusion matrix where (IPC, PDA, AD, PER) = (0.6, -0.2, 2.0, 1.0)

Correct Category	Classification				
	"A"	"B"	"C"	"D"	"E"
A	0.349	0.438	0.106	0.082	0.025
B	0.015	0.439	0.291	0.183	0.073
C	0.034	0.225	0.377	0.301	0.064
D	0.107	0.088	0.085	0.512	0.207
E	0.010	0.008	0.098	0.469	0.415

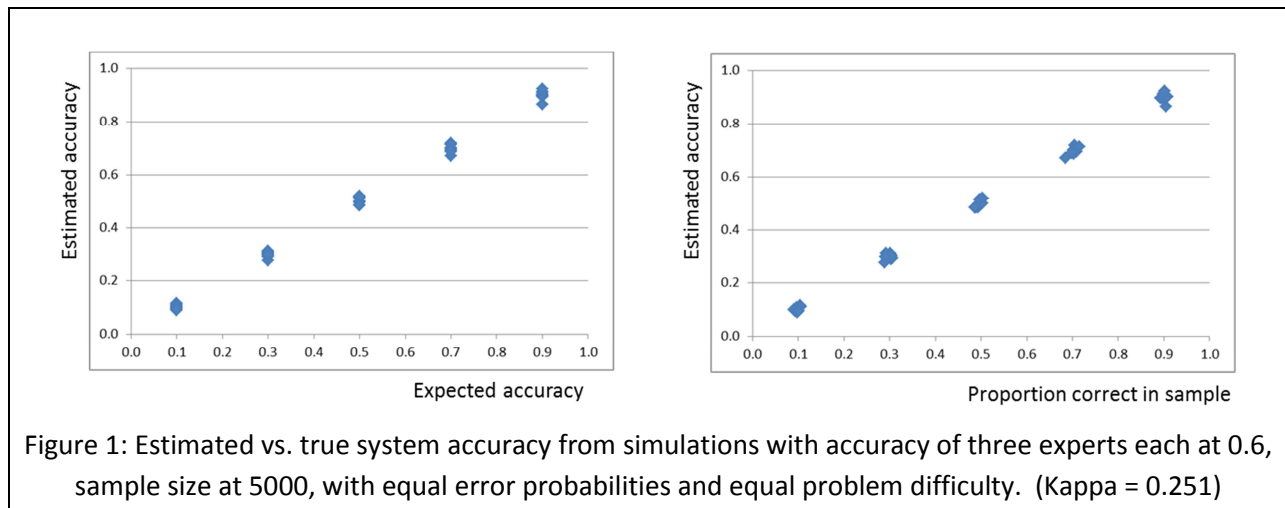
For a selected sample size, N, a "simulation run" executes the following:

1. Randomly select the base rate probability for each classification
2. Generate the confusion matrices for each rater and the system
3. Use the base rate probability and confusions matrices to randomly generate N cases.
4. Estimate system accuracy (using method described above)
5. Compare estimated system accuracy to "true" system accuracy, where there are two measures of true system accuracy
 - a. Expected accuracy (i.e. $P(A)*P("A" | A) + P(B)*P("B" | B) + \dots$)
 - b. Proportion correct in sample

When K-assumptions are satisfied

This section examines circumstances where the assumptions implicit in Kappa are satisfied. That is to say where the raters are selecting from a single confusion matrix of the structure shown in Table 1 and where the system confusion matrix also has the same well-behaved structure.

Figure 1 illustrates the asymptotic behavior of the estimation method. The simulation runs in Figure 1 involved five categories, three experts each with 60% accuracy, 5000 test cases for each run, and where there are 10 runs each with system accuracy set to .1, .3, .5, .7 and .9 respectively.



Clearly the estimation method behaves exactly as it should. Estimates of system accuracy cluster tightly around true system accuracy at all levels. When true system accuracy is 0.1, which is less accurate than random guessing (0.2), estimates of system accuracy cluster tightly around 0.1. When true system accuracy is 0.9, which is far better than the experts' accuracy (0.6), estimates of system accuracy cluster tightly around 0.9. Across all fifty simulation runs the average value of Kappa was just 0.251.

The results below examine what happens when sample size and expert accuracy are varied.

Figures 2-4 show the results of fifty simulations with a sample size of 200 per run and where expert accuracy is set to .4, .6 and .8 respectively.

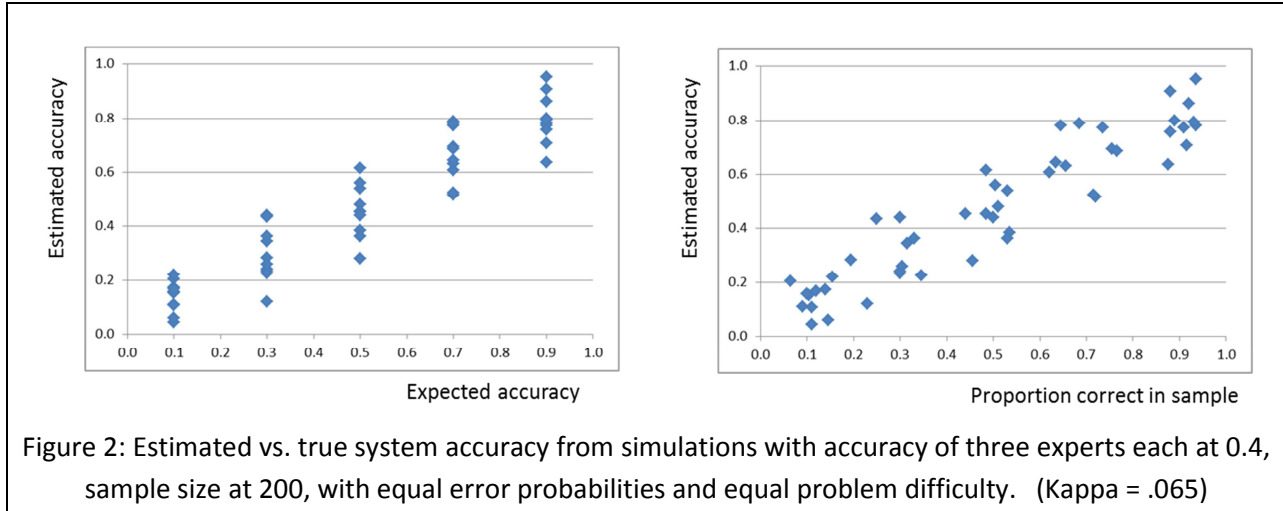


Figure 2: Estimated vs. true system accuracy from simulations with accuracy of three experts each at 0.4, sample size at 200, with equal error probabilities and equal problem difficulty. (Kappa = .065)

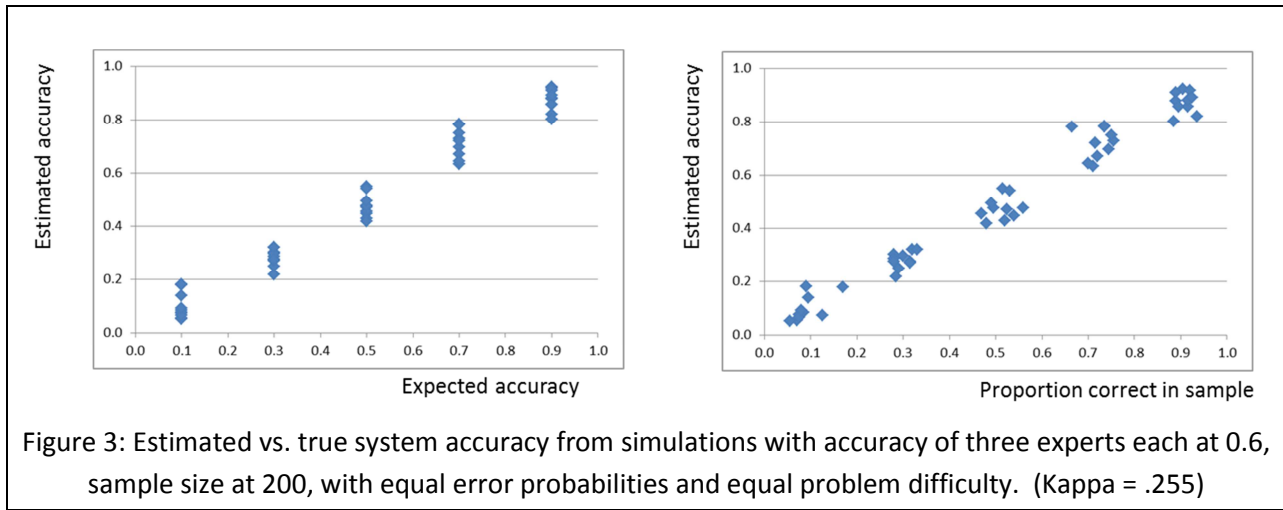


Figure 3: Estimated vs. true system accuracy from simulations with accuracy of three experts each at 0.6, sample size at 200, with equal error probabilities and equal problem difficulty. (Kappa = .255)

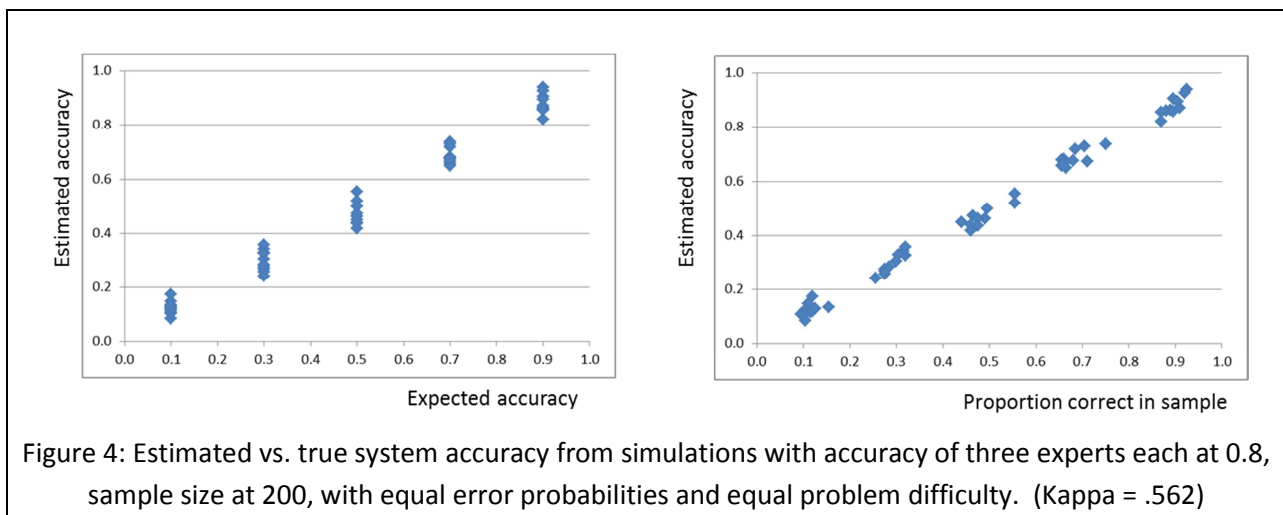
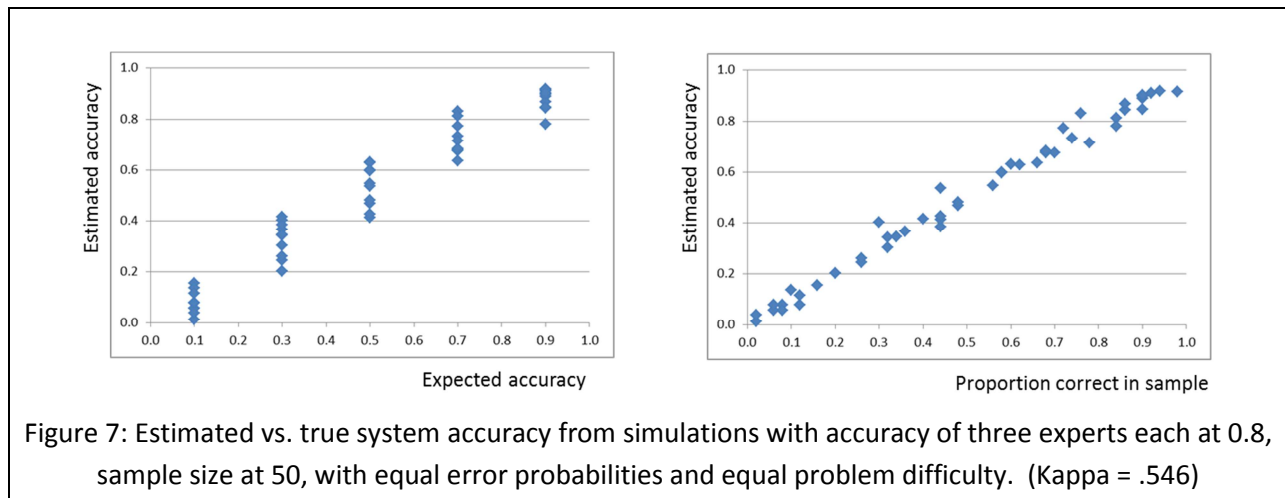
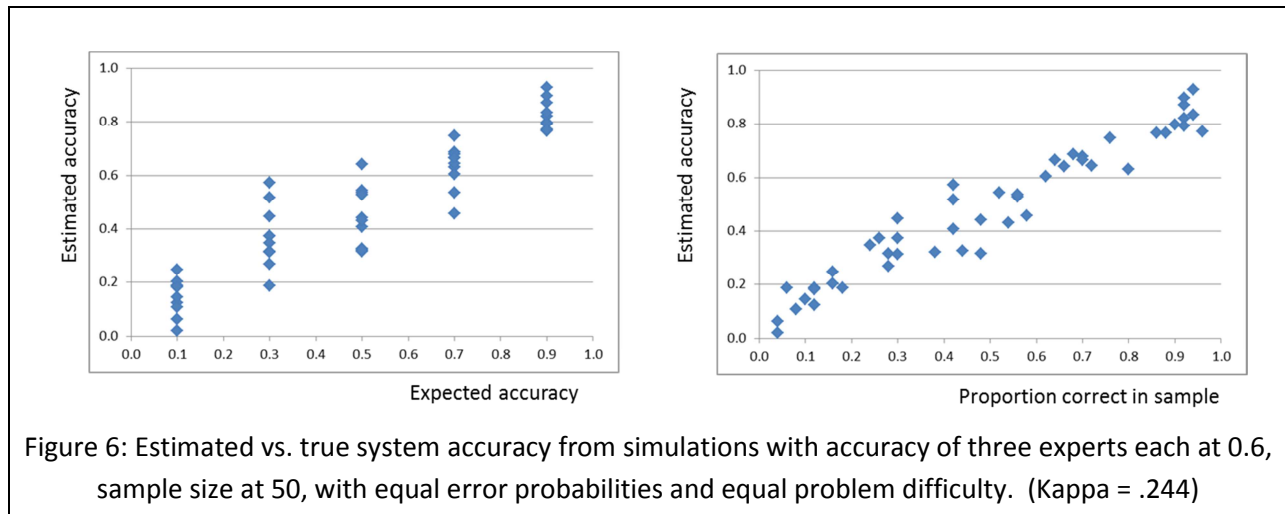
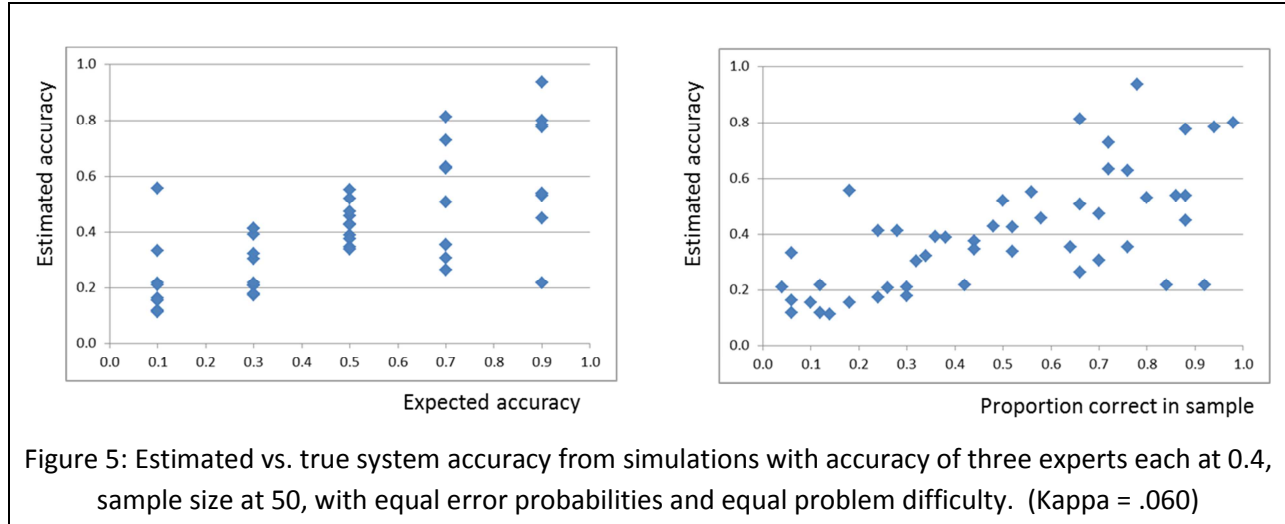


Figure 4: Estimated vs. true system accuracy from simulations with accuracy of three experts each at 0.8, sample size at 200, with equal error probabilities and equal problem difficulty. (Kappa = .562)



It is apparent in Figures 2 - 4 that the correspondence between estimated and true system accuracy improves rapidly as the experts become more accurate. Even when the experts are just 60% accurate, estimates of system accuracy are consistently within ± 0.1 of true system accuracy.

Figures 5 - 7 show what happens when the sample size is reduced further to just 50 cases per run. When rater accuracy is 0.4 there is hardly any correspondence between estimated and true system accuracy. However when rater accuracy is 0.6 and 0.8 this correspondence improves quickly.

Readers may note that in Figures 6 and 7 the two different measures of true accuracy yield slightly different results. Estimated accuracy corresponds more closely to proportion correct in sample than to expected accuracy. This occurs because the true proportion correct in a sample varies according to a binomial distribution defined by system accuracy. So even if there is perfect correspondence between estimated accuracy and proportion correct (as is the case when expert accuracy is set to 1.0), the standard deviation of the estimate around expected accuracy (Exp) would still be equal to $\text{SQRT}(\text{Exp} \cdot (1 - \text{Exp}) / N)$. When $N=50$ and Exp is .7, one would expect around 90% of the estimates to be between .6 and .8. When $N=200$ then one should expect nearly 90% of the samples to have a proportion correct between 0.65 and .75.

In summary, when the K-assumptions are satisfied, the estimation method exhibits a very orderly relationship between estimated and true system accuracy. Estimates of system accuracy are largely unbiased, improve quickly with increasing rater accuracy, and improve nicely with increasing sample size.

When K-assumptions are substantially violated

In practice it is difficult to imagine a circumstance where the K-assumptions and the implied K-properties are satisfied. All raters are not equally accurate; some are typically more experienced and expert than others. All types of errors are not equally probable; this property is certainly false when the categories are naturally ordered or when the raters have some idea of which categories have the highest base rates. And all problems are not equally difficult; unless the test cases are carefully pre-selected and therefore unrepresentative of real world diversity.

This section examines the behavior of the estimation method when the K-properties are violated. In all of the simulation runs summarized in this section the K-properties of equal rater accuracy, equal problem difficulty, and equal error probabilities are substantially violated. Specifically:

Rater accuracy (IPC) was varied by .1. For example, instead of three raters with .6 accuracy, initial accuracy would be set to .5, .6 and .7 respectively.

Problem difficulty (PDA) was varied by .2. For about a third of the test cases rater and system accuracy were reduced by .2 (or set to a minimum of 0.0) and for about another third accuracy was increased by .2 (or set to the maximum probability of 1.0).

Asymmetric dispersion (AD) was set to 2.0. An incorrect answer that is 'next to' the correct answer is twice as likely as one two steps removed and 4 times as likely as one 3 steps removed, etc.

Error probabilities were randomly varied by up to 100% (PER=1.0). For example, if the error probability is initially set to .2 then that error probability would be randomly selected from the range [0, .4]. This random variation is done independently for each error probability.

To appreciate the magnitude of impact of these parameter settings consider again Tables 5 and 7 above. Table 5 is exactly how the confusion matrix looks when initial rater accuracy is set to .6 and the K-properties are satisfied. Table 7 is representative of about 1/3 of the cases when initial rater accuracy is set to .6 but with the above parameter settings. It seems fair to characterize Table 7 as a substantial variation from Table 5.

All of the simulation runs in this section use the above parameter settings to systematically and then randomly vary the rater and system confusion matrices.

Figure 8 illustrates the asymptotic behavior of the estimation method when the K-properties are substantially violated. Please note that when system accuracy is preset to .1 and .9, expected accuracy is .133 and .867 respectively. This occurs because problem difficulty is varied plus and minus 0.2, but accuracy can be no lower than 0.0 or higher than 1.0. So for example when system accuracy is preset to 0.1, one third of the problems have system accuracy reset to 0.3, one third stay at 0.1 and the remaining third are reset to 0.0; then averaged expected system accuracy is then .133.

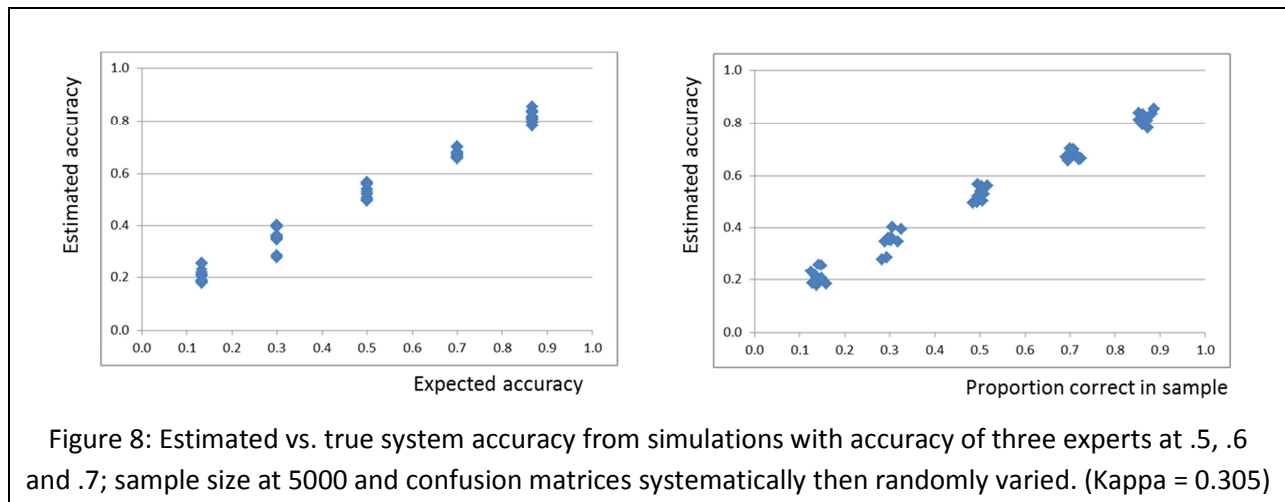


Figure 8: Estimated vs. true system accuracy from simulations with accuracy of three experts at .5, .6 and .7; sample size at 5000 and confusion matrices systematically then randomly varied. (Kappa = 0.305)

Although there is a linear relationship between estimated and true accuracy. There is also some bias in the estimates; overestimating accuracy when the system performs poorly and underestimating accuracy when the system performs very well. Note though that when the system was more accurate than the raters the estimates of system accuracy were still consistently higher than the raters' accuracy. The estimate of system accuracy may be conservative, but it is not bounded by the experts' accuracy.

The explanation for this estimation bias is straightforward. Violations of the K-properties will generally inflate inter-rater agreement. Indeed this disconnect between rater accuracy and inter-rater agreement was the basis for the proof of Theorem 1.⁴ Because inter-rater agreement is used to estimate rater accuracy, as per Eq. 1, this leads to a slightly inflated estimate of rater accuracy. Inflated estimates of rater accuracy in turn lead to overestimates of the ground truth probabilities for the categories with the highest estimated ground truth probabilities. Finally given the equation for deriving system accuracy from the ground truth probabilities (Eq. 3) this leads to the estimation bias. In comparing Figures 1 and 8, note that Kappa was .251 and .305 respectively; and the average estimated accuracy for the runs in Figure 1 was exactly 0.60 while the average estimated rater accuracy for the runs in Figure 8 was 0.64.

In Figures 9 – 11 the sample size is 200 cases per run and expected rater accuracy is set to .4, .6 and .8 respectively. In Figures 12 – 14 sample size is reduced to 50 cases per run.

Except for the bias toward underestimating high system accuracy (and overestimating low system accuracy) these results are similar to the results with the matrices that satisfy the K-properties. Increasing rater accuracy and sample size both serve to decrease the variance of the estimate. The estimation bias is pronounced when rater accuracy is very low (0.4), noticeable when rater accuracy is moderate (.6), and is negligible when rater accuracy is high (0.8).

In practice, most efforts to evaluate system accuracy address systems that are hypothesized to perform well. For such evaluations the estimates derived from this method become increasingly conservative as the ground truth judgments of the experts are increasingly suspect. This, it can be argued, is a positive feature of the estimation method.

⁴ To be more precise, all violations of the K-properties will inflate expected inter-rater agreement with one exception. Differences between rater accuracy decreases rather than increases expected inter-rater agreement. This effect however is small and is overwhelmed by the larger effect of the other violations. For example, if expert accuracy is set to .4, .6, .8 and true system accuracy is 0.9 then the estimate will be approximately 0.924 – a 0.024 overestimate. But when expert accuracy is fixed at .6, but problem difficulty is varied to be .4, .6, .8 then there is a 0.043 underestimate, 100% dispersion around the error probabilities leads to a 0.036 underestimate, and asymmetric dispersion set at 2.0 leads to a 0.068 underestimate.

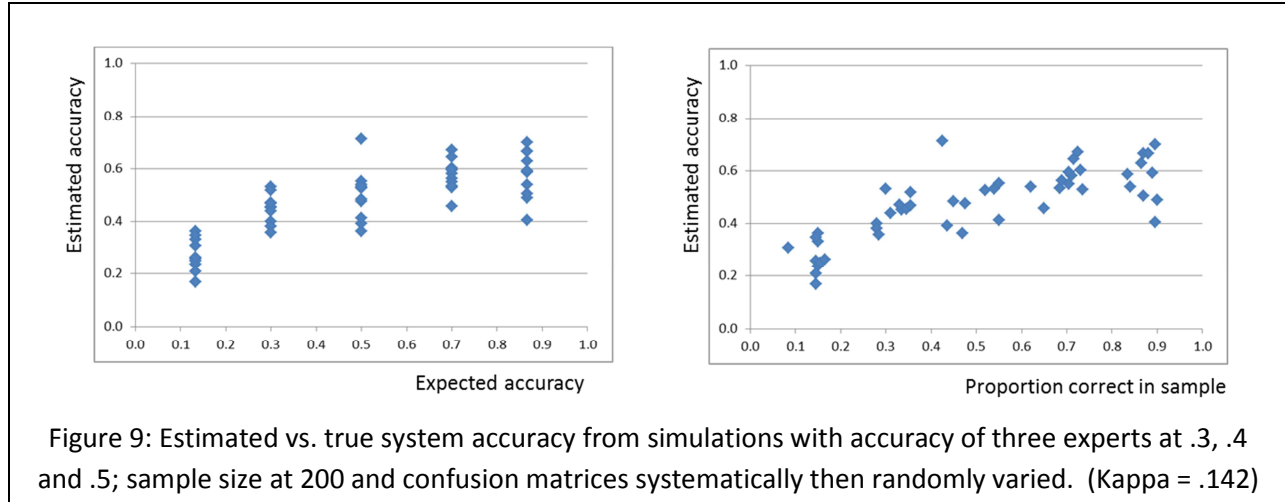


Figure 9: Estimated vs. true system accuracy from simulations with accuracy of three experts at .3, .4 and .5; sample size at 200 and confusion matrices systematically then randomly varied. (Kappa = .142)

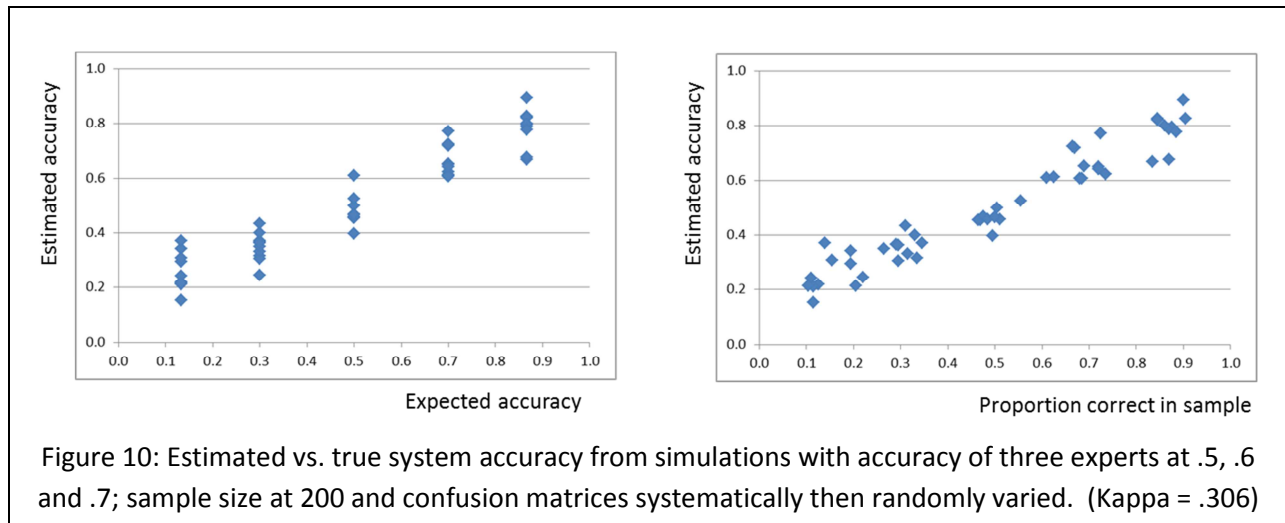


Figure 10: Estimated vs. true system accuracy from simulations with accuracy of three experts at .5, .6 and .7; sample size at 200 and confusion matrices systematically then randomly varied. (Kappa = .306)

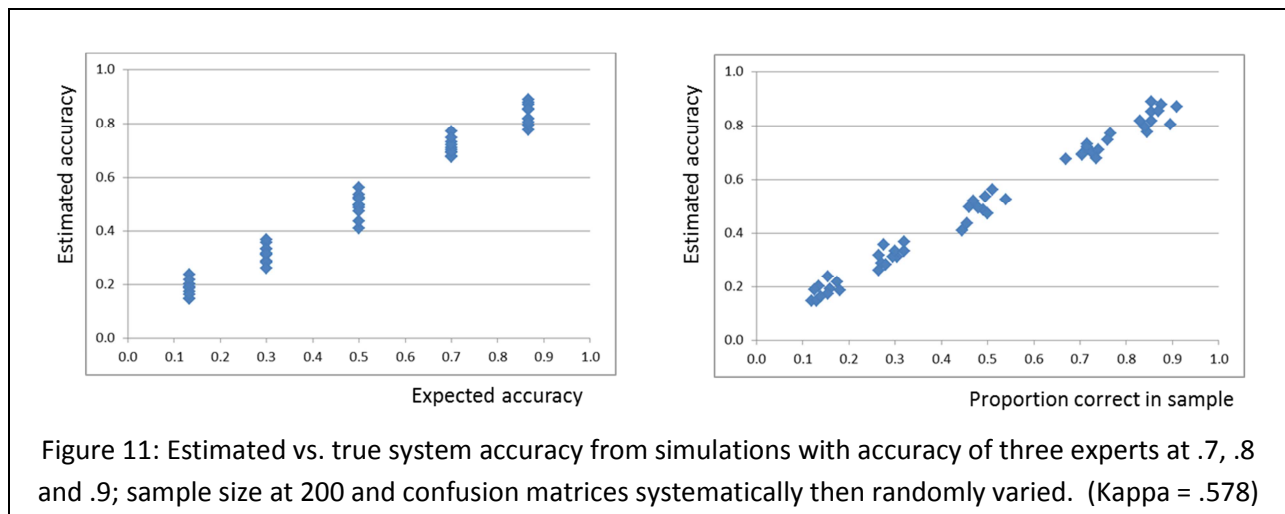
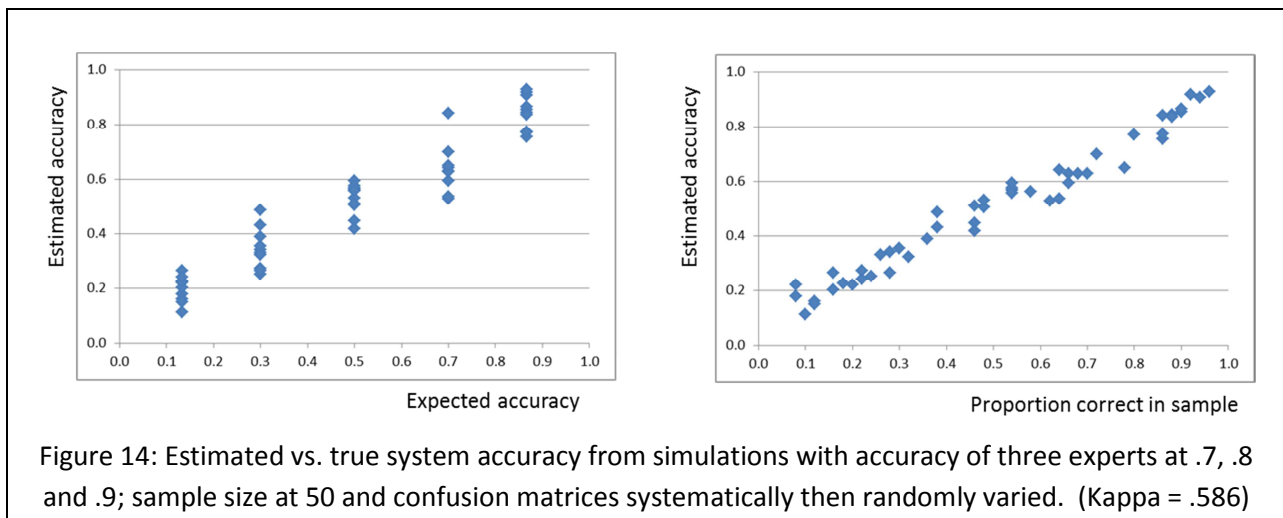
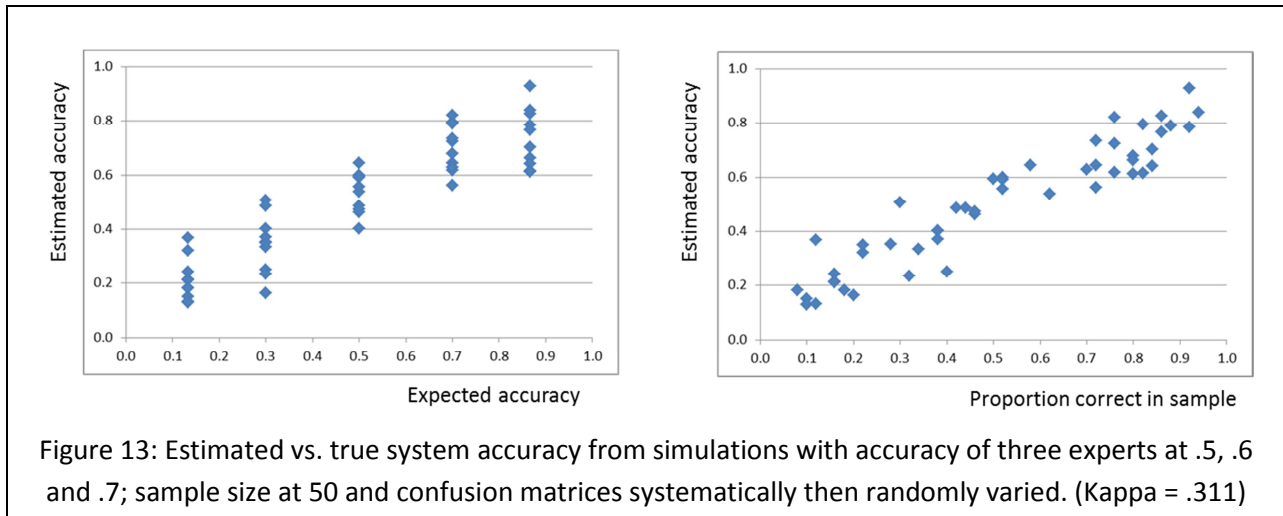
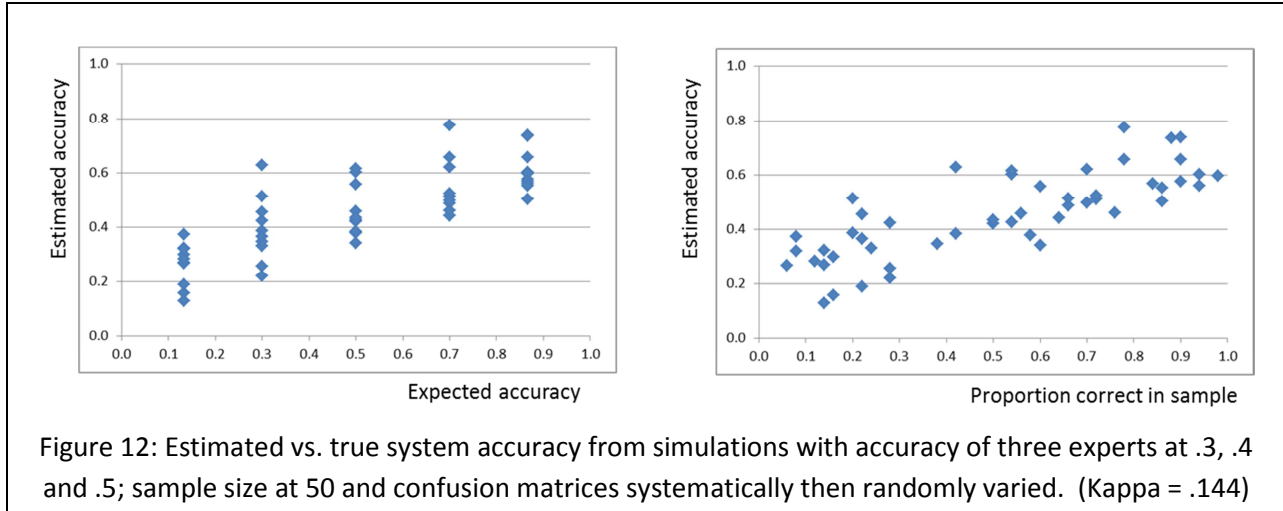


Figure 11: Estimated vs. true system accuracy from simulations with accuracy of three experts at .7, .8 and .9; sample size at 200 and confusion matrices systematically then randomly varied. (Kappa = .578)



DISCUSSION

The objective of this paper was to demonstrate that it is feasible to reliably estimate the accuracy of system classifications when ground truth can only be estimated with fallible expert ratings. The simulations described above provide clear evidence for the claims made in the introduction, namely that reliable estimates of system accuracy can be obtained from fallible expert ratings under a diverse conditions, that the reliability of these estimates is approximately the same whether the system is more or less accurate than the expert raters, and that the conditions under which these accuracy estimates become unreliable are predictable (e.g. when inter-rater agreement is low and sample size is small).

The estimation method uses the level of inter-rater agreement to estimate the overall accuracy of the expert ratings, uses Bayesian updating based on the estimated expert accuracy to estimate a “ground truth” probability for each classification, and then estimates system accuracy by comparing the relative frequency that the system agrees with the most probable classification at different probability levels.

Although the estimation method was derived from assumptions that are implicit in the Kappa statistic (and how it is often used), a simulation analysis shows that the accuracy of the estimates of system accuracy are robust against substantial variations from the rater behavior implied by those assumptions. The accuracy of the estimates of system accuracy is driven primarily by overall rater accuracy (which can be estimated from inter-rater agreement) and sample size.

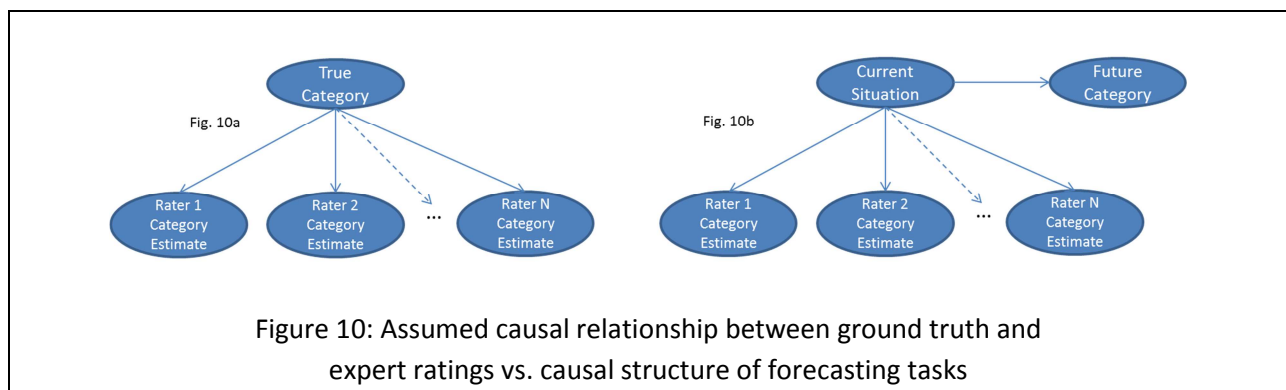
Recommended use and uses to avoid

The above simulation results suggest an overall data collection and estimation approach where measured inter-rater agreement is used to determine the number of test cases needed to obtain high confidence in system accuracy estimates. So for example, for five category problems with three raters if Kappa is around .3, then to obtain 90% “confidence” that estimated accuracy is within .1 of true accuracy data collection should continue for at least 200 cases. On the other hand, if Kappa is around .55 then our additional simulation runs show that a sample size of 100 cases is more than enough to ensure the same “confidence interval” or better. As the number of raters and categories varies, so does the parametric relationship between sample size and confidence in system accuracy estimates; so additional simulation runs such as those shown here would be needed to determine sample size requirements.

In this approach all test cases are useable, even the ones where the raters disagree. This makes it feasible to randomly select test cases from the population of problems from which the system is likely to be applied, and this in turn greatly facilitates the ability generalize test results to practice. This seems superior to a process that limits test cases to ones where there is a high level of expert agreement – which almost certainly would have the effect of limiting testing to a non-representative set of relatively easy cases.

As noted above, when system accuracy is in fact high, variations in test problem difficulty and unequal error probabilities will bias the estimates toward underestimating system accuracy, while variations in rater accuracy will bias the estimates toward overestimating accuracy, with a net effect clearly tilted toward underestimation. The magnitude of these biases drops quickly as overall rater accuracy increases. But when inter-rater agreement is very low (e.g. Kappa = .15 for five category problem), suggesting low overall rater accuracy (e.g. around 0.4), then it might be helpful to examine the inter-rater agreement data for patterns of violations of the K-properties. For example, the K-property of equal error probabilities implies that all pairwise disagreements are equally likely (e.g. “AB” as likely as “AE”) and a statistical test can be performed to help determine if this pattern is violated. If it is, then the estimate of system accuracy can be adjusted upwards. There is much work to be done to determine how and when such adjustments should be made, but doing so appears quite feasible.

The estimation method described herein is specifically intended for cases where each rater is an independent measure of ground truth classifications. The procedure assumes the causal structure shown in Figure 10a.



There are many applications that involve aggregation of independent estimates from multiple individuals but do not have the causal structure shown in Figure 10a. For many such applications use of the estimation method described here would be an egregious error. For example, it is becoming common practice in forecasting to systematically combine the ratings of multiple independent forecasters (e.g. Surowieki, 2005). Although the estimation method presented here could be mechanically applied to such forecasting tasks, such an application would yield spurious results. Forecasting tasks do not have the causal structure shown in Figure 10a, but tend to have a causal structure closer to the one shown in Figure 10b where expert ratings are not in any sense direct measures of the future outcomes. On the other hand the estimation method can and has been used to retrospectively assess whether a forecasted outcome actually occurred. For example Lehner et al. (2012) examined the accuracy of the imprecise forecasts typically found in published forecasts by using multiple raters to retrospectively assess whether the forecasted outcome occurred and then using an estimation method similar to the one presented here to estimate the accuracy of a collection of forecasts.

The distinction between Figures 10a and 10b is essentially the distinction between medical diagnosis and medical prognosis. It would be very appropriate to apply the method to estimate the accuracy of a new diagnostic system by comparing system diagnoses to those of medical professionals, but it would be inappropriate to use it to estimate the accuracy of a new system's prognoses by comparing them to the prognoses of medical professionals.

In general it is important that the causal structure relating the rater and system selections to ground truth match the structure assumed by the estimation method. The process of collecting ratings from the experts should be engineered to ensure this causal structure; such as by ensuring that the expert ratings are independent and to the extent possible having available the same data for each rater for each test case.

The estimation method presented here was developed to address test and evaluation of an automated classification system after development. However it does seem feasible to also employ this approach during system development. Specifically if one is building a system from training data sets then it should be possible to use the method described here to develop training sets with a probability distribution of correct classifications for each training case. Such an approach would get around the need to have domain experts assign a definitive correct answer for each training case. This would be particularly useful in domains where there is concern that the expert ratings are error prone.

Related and future research

The research presented in this paper had the very specific goal of demonstrating that it was feasible to reasonably estimate system accuracy using fallible expert ratings even when the system was substantially more accurate than the experts. The results presented here support a claim that this objective was achieved.

On the other hand nothing in this paper would support a claim that the estimation method presented here is in any sense optimal. There are many opportunities for improvement. Below are two directions of possible improvements.

First, the estimation method was designed for use with classification problems for which there is no natural ordering to the categories. The simulation study showed that the method was robust even when there is a natural ordering, but it seems likely that the accuracy of estimates of system accuracy will be improved if the method specifically accounts for the fact that certain types of errors are much more likely than others. For example, if the natural ordering is A, B, C, D, E, then a rating of "A" should be more evidence for category B than for category E. The method presented here treats B and E equally.

Second, as noted above, it should be feasible to develop statistical procedures to estimate whether and by how much K-properties are violated. From these estimates it should be also feasible to adjust the system accuracy estimates to correct for bias. This area is unexplored.

Third, the estimation method presented here is entirely algebraic. Everything is derived directly from some percent-of-agreement statistics. No effort was made to estimate base rates and confusion matrices that represent a “best fit” to the inter-rater agreement data. But there are best fit methods that could be used for this purpose. For example, the non-linear optimization methods in Latent Class Analysis (McCutcheon, 1987) could be used to find maximum likelihood estimates for the base rate and confusion matrix probabilities. Both Uebersax (1988) and Carpenter (2008) applied this approach to binary classification problems; and Carpenter also used Bayes inference to aggregate ratings and estimate classification probabilities. Similarly one could use non-linear optimization to find base rates and confusion matrix probabilities that minimize the difference between expected and observed relative frequency of each inter-rater pair (relative frequency of “AA”, “AB”, “AC” ...).⁵ It remains an open and interesting question as to whether use of such optimization methods yields better results.

⁵ In fact this was exactly the procedure initially used on this effort, where the optimizer found the base rate and confusion matrices that minimized the squared error between observed and expected pattern of relative frequencies of inter-rater agreement and disagreement. It worked, but non-linear optimization was found to be unnecessarily cumbersome for the needs of this study.

References

Armstrong, J.S. (2006) How to Make Better Forecasts and Decisions: Avoid Face-to-face Meetings. *The International Journal of Applied Forecasting*, 5, 3-15.

Bishop, M.A. & Trout, J.D. (2002). 50 Years of Successful Prediction Modeling Should be Enough. *Philosophy of Science*, 69(S3), S197-S208.

Carpenter, B. (2008). *Multi-level Bayesian Models of Categorical Data Annotation*. Manuscript found at <http://citeseerx.ist.psu.edu/viewdoc/download?doi=10.1.1.174.1374&rep=rep1&type=pdf>.

Dawes, R.M. (1979). The robust beauty of improper linear models. *American Psychologist*, 34, 571- 582.

Grove, W.H., Zald, D.H., Lebow, B.S., Snitz, B.E. & Nelson, C. (2001). Clinical versus mechanical prediction: A meta-analysis. *Psychological Assessment*, 12, 19-30.

Gwet, K. L. (2010). *Handbook of Inter-Rater Reliability (2nd Edition)*. Advanced Analytics, LLC.

Lehner, P., Michelson, A., Adelman, L. & Goodman, A. (2012). Using Inferred Probabilities to Measure the Accuracy of Imprecise Forecasts. *Judgment and Decision Making*, 7(6), 728-740.

McCutcheon, A.L. (1987). *Latent Class Analysis*. Thousand Oaks, California: Sage Publications.

Surowiecki, J. (2005). *The Wisdom of Crowds*. Anchor Books.

Tetlock, P.E. (2005). *Expert political judgment: How good is it? How can we know?* Princeton, N.J.: Princeton University Press.

Uebersax, J.S. (1988). Validity inferences from Inter Observer Agreement. *Psychological Bulletin*, 104(3), 405-416.

APPENDIX
Proof of Theorem

Restating the assumptions:

- A1. For each case there is a unique correct classification
- A2. For each case raters independently assign classifications
- A3. Expected agreement between raters increases as expected rater accuracy increases.

Begin with a few definitions.

Definition of *correct classification* in A1: For each case there is a vector $\langle c_1, c_2 \dots c_n \rangle$ where for some index i , $c_i=1$ and the remaining values are 0.

Definition of *independent assignment* in A2: For each case, the probability that a rater will select a class is conditionally independent of the other raters' selections.

Independent assignments allow the description of each rater's selection behavior as a probability vector. That is to say, for each case each rater has a selection probability for each category. These will be called *selection vectors*.

Definition of *rater accuracy* in A3: For M raters and N cases, rater accuracy is defined as the total proportion of correct selections.

For example, if there are 10 cases and three raters who make correct assignments in 7, 5 and 9 of the cases respectively, then rater accuracy = 0.7.

The three lemmas below all use the same proof strategy. Begin with any two selection vectors that are not identical. Construct a selection vector that is the average of the two. The average vector will necessarily have the same expected accuracy but a different level of expected agreement than the original two vectors. If the average vector has higher/lower expected agreement, then create a new vector by slightly reducing/increasing the probability of correct assignment in the average vector. When the change is sufficiently small the new vector will have higher/lower expected accuracy and lower/higher expected agreement than the original two vectors. Most of the algebraic complexity in these proofs is the result of showing one way to calculate a change that is always "sufficiently small".

Lemma 1: To ensure A1 – A3 within each case all raters must behave as though they are selecting a category using the same selection vector.

Proof:

Let $\langle p_{11}, p_{12} \dots p_{1n} \rangle$ and $\langle p_{21}, p_{22} \dots p_{2n} \rangle$ be the selection vectors of 2 raters for a specific case; where some probabilities do not agree (e.g. $p_{11} \neq p_{21}$). For purposes of the proofs below, assume that category 1 is the correct category. (The arguments below apply no matter which category is correct.)

Below it is shown how to construct from two different selection vectors a common selection vector for both raters where expected accuracy is lower but expected agreement higher. Consequently unless the two raters have the same selection vector, there will always be another pair of vectors with lower expected accuracy and higher expected agreement – violating A3.

Set $p_i = (p_{1i} + p_{2i})/2$, $e_i = (p_{1i} - p_{2i})/2$, $d = (e_1^2 / (2 * (p_2 - p_1)))$ if $p_1 < p_2$, $d = -((e_1^2 / (2 * (p_2 - p_1))))$ if $p_1 > p_2$, and $d=0$ if $p_1 = p_2$.

For selection vectors $\langle p_{11}, p_{12} \dots p_{1n} \rangle$ and $\langle p_{21}, p_{22} \dots p_{2n} \rangle$ (1)

Expected accuracy = p_1

Expected agreement

$$\begin{aligned} &= p_{11} * p_{21} + p_{12} * p_{22} + \dots + p_{1n} * p_{2n} \\ &= (p_1 + e_1) * (p_1 - e_1) + (p_2 + e_2) * (p_2 - e_2) + \dots + (p_n + e_n) * (p_n - e_n) \\ &= p_1^2 + p_2^2 + \dots + p_n^2 - e_1^2 - e_2^2 - \dots - e_n^2 \end{aligned}$$

For selection vectors $\langle p_1, p_2 \dots p_n \rangle$ and $\langle p_1, p_2 \dots p_n \rangle$ (2)

Expected accuracy = p_1

Expected agreement

$$= p_1^2 + p_2^2 + \dots + p_n^2$$

For selection vectors $\langle p_1 - d, p_2 + d, p_3 \dots p_n \rangle$ and $\langle p_1 - d, p_2 + d, p_3 \dots p_n \rangle$ (3)

Expected accuracy = $p_1 - d$

Expected agreement, if $p_1 \neq p_2$,

$$\begin{aligned} &= (p_1 - d)^2 + (p_2 + d)^2 + p_3^2 + \dots + p_n^2 \\ &= p_1^2 + p_2^2 + \dots + p_n^2 + 2 * d * (p_2 - p_1) + 2d^2 \\ &= p_1^2 + p_2^2 + \dots + p_n^2 + 2 * (e_1 / (p_2 - p_1))^2 * (p_2 - p_1) + 2(e_1 / (p_2 - p_1))^4 \\ &= p_1^2 + p_2^2 + \dots + p_n^2 - e_1^2 + e_1^4 / 2 * (p_2 - p_1)^2 \end{aligned}$$

Expected agreement, if $p_1 = p_2$

$$= p_1^2 + p_2^2 + \dots + p_n^2$$

Expected accuracy in (1) is higher than in (3), but expected agreement is lower; where the common selection vector in (3) was constructed from a difference between the vectors in (1). Consequently,

whenever there is a difference between the selection vectors of two raters a selection probability vector for the two raters can be constructed with lower expected accuracy but high expected agreement.

Within each case if the selection vectors of the raters differ A3 is not guaranteed. ***

Lemma 2: To ensure A1 – A3 within each case the error probability is the same for all incorrect categories.

Proof:

From Lemma 1 it is known that A1-A3 imply that for each case all raters have the same selection vector. Let that vector be $\langle p_1, p_2 \dots p_n \rangle$. Assume category 1 is the correct assignment and that the remaining probabilities are not all equal.

Below it is shown how to construct selection vector, with equal probability for all incorrect assignments, where expected accuracy is higher but expected agreement lower. Consequently the error probabilities are unequal, there will always be a vector with higher expected accuracy and lower expected agreement – violating A3.

Set $p_e = (p_2 + \dots + p_n)/(n-1)$, $e_i = (p_i - p_e)$ for all $i > 1$, set $e_{\min} = \min(|e_2| \dots |e_n|)$ and $d = e_{\min}^2/2$
 Note that $(e_2 + \dots + e_n) = 0$ and that there are at least 2 e_i that are not zero.

For the vector $\langle p_1, p_2 \dots p_n \rangle$, (1)

Expected accuracy = p_1

Expected agreement

$$\begin{aligned} &= p_1^2 + p_2^2 + \dots + p_n^2 \\ &= p_1^2 + (p_e + e_2)^2 + \dots + (p_e + e_n)^2 \\ &= p_1^2 + p_e^2 + \dots + p_e^2 + 2p_e(e_2 + e_3 + \dots + e_n) + e_2^2 + e_3^2 + \dots + e_n^2 \\ &= p_1^2 + (n-1)p_e^2 + e_2^2 + e_3^2 + \dots + e_n^2 \end{aligned}$$

For the vector $\langle p_1, p_e \dots p_e \rangle$ (2)

Expected accuracy = p_1

Expected agreement

$$= p_1^2 + (n-1)p_e^2$$

For the vector $\langle p_1+d, p_e-d, p_e \dots p_e \rangle$ (3)

Expected accuracy = $p_1 + d = p_1 + e_{\min}^2/2$

Expected agreement

$$\begin{aligned} &= (p_1+d)^2 + (p_e-d)^2 + p_e^2 + \dots + p_e^2 \\ &= p_1^2 + (n-1)p_e^2 + 2p_1d - 2p_e d + 2d^2 \end{aligned}$$

$$\begin{aligned}
 &= p_1^2 + (n-1)p_e^2 + 2d(p_1 - p_e) + 2d^2 \\
 &= p_1^2 + (n-1)p_e^2 + e_{\min}^2 * (p_1 - p_e) + e_{\min}^4 / 2
 \end{aligned}$$

Since $e_{\min}^2 * (p_1 - p_e) + e_{\min}^4 / 2 < e_{\min}^2 + e_{\min}^2 \leq e_2^2 + e_3^2 + \dots e_n^2$, expected agreement in (1) is higher than expected agreement in (3) even though expected accuracy is lower.

Consequently, whenever the probability of incorrect assignment is unequal, there will always be a selection vector with higher expected accuracy and lower expected agreement, violating A3.

Within each case and selection vector if the error probabilities are unequal A3 is not guaranteed. ***

Lemma 3: To ensure A1 – A3 the selection vector must be the same across all cases.

Proof:

Lemmas 1 and 2 show that A1-A3 imply that for each case the raters have identical selection vectors of the form $\langle p_e \dots p_c \dots p_e \rangle$ where p_c is the probability of assigning the correct category and $p_e = (1-p_c)/(n-1)$ where n is the number of categories.

Below it is shown that across different cases the selection vectors must have the same values for p_c (and therefore p_e) else a violation of A3 can be constructed.

Let p_{c1} and p_{c2} be the probability of correct assignment on two different cases, and p_{e1} and p_{e2} the corresponding error probabilities. For each case, order the cases such that the correct assignment is first. So for all raters the probability vector is $\langle p_{c1}, p_{e1}, \dots p_{e1} \rangle$ for case 1 and $\langle p_{c2}, p_{e2}, \dots, p_{e2} \rangle$ for case 2, but the categories may be in a different order. The proof below makes no reference to matching categories across cases so this ordering does not affect the proof.

Set $p_c = (p_{c1} + p_{c2})/2$, $p_e = (p_{e1} + p_{e2})/2$, $e_c = (p_{c1} - p_{c2})$, $e_e = (p_{e1} - p_{e2})$, $e_{\min} = \min(|e_c|, |e_e|)$, $d = e_{\min}^2/2$

For two cases with accuracy $p_{c1} \neq p_{c2}$, then

$$\text{Expected accuracy} = p_c \tag{1}$$

Expected agreement

$$\begin{aligned}
 &= (p_{c1}^2 + (n-1)p_{e1}^2 + p_{c2}^2 + (n-1)p_{e2}^2)/2 \\
 &= ((p_c + e_c)^2 + (n-1)(p_e + e_e)^2 + (p_c - e_c)^2 + (n-1)(p_e - e_e)^2)/2 \\
 &= (2p_c^2 + 2(n-1)p_e^2 + 2e_c^2 + 2(n-1)e_e^2)/2 \\
 &= p_c^2 + (n-1)p_e^2 + e_c^2 + (n-1)e_e^2
 \end{aligned}$$

For two cases with accuracy $p_{c1} = p_{c2}$ (2)

Expected accuracy = p_c

Expected agreement

$$= p_c^2 + (n-1)p_e^2$$

For two cases with accuracy vectors $\langle p_c+d, p_e-d, p_e \dots p_e \rangle$ (3)

Expected accuracy = $p_c + d$

Expected agreement

$$\begin{aligned} &= (p_c+d)^2 + (p_e-d)^2 + (n-2)p_e^2 \\ &= p_c^2 + (n-1)p_e^2 + 2p_c d + d^2 - 2p_e d + d^2 \\ &= p_c^2 + (n-1)p_e^2 + 2d(p_c-p_e) + 2d^2 \\ &= p_c^2 + (n-1)p_e^2 + 2(e_{\min}^2/2)(p_c-p_e) + 2(e_{\min}^2/2)^2 \\ &= p_c^2 + (n-1)p_e^2 + (p_c-p_e)e_{\min}^2 + e_{\min}^4/2 \end{aligned}$$

Since $e_{\min}^2(p_c - p_e) + e_{\min}^4/2 < e_{\min}^2 + e_{\min}^2 \leq e_c^2 + e_e^2$, expected agreement in (1) is higher than expected agreement in (3) even though expected accuracy is lower.

Consequently, whenever the probability of correct assignment across cases is unequal, there will always be a probability vector that is the same across cases with higher expected accuracy and lower expected agreement, violating A3.

Across cases, if the selection vectors differ then A3 is not guaranteed.

Theorem 1: A1 – A3 are ensured if and only if the raters behave as though their selection for each case is determined by a single confusion matrix where the conditional probability of correct assignment is constant and the conditional probability of all incorrect assignments is equal.

Proof:

The “only if” necessity portion follows directly from Lemmas 1 - 3. Sufficiency follows the fact that with a constant conditional probability of correct assignment (P_c) and incorrect assignments (P_e), expected accuracy is P_c and expected agreement is $P_c^2 + (n-1)P_e^2 = P_c^2 + (1-P_c)^2/(n-1)$. Clearly expected agreement increases monotonically with P_c .
