

Running head: COMMENT ON MICKES ET AL.

Latent Mnemonic Strengths are Latent: A Comment on Mickes, Wixted, and Wais (2007)

Jeffrey N. Rouder and Michael S. Pratte

University of Missouri

Richard D. Morey

University of Groningen

Jeff Rouder

210 McAlester Hall

University of Missouri

Columbia, MO 65211

rouderj@missouri.edu

3992 Words

Abstract

Mickes, Wixted, and Wais (2007) propose a simple test of latent strength variability in recognition memory. They asked participants to rate their confidence using either a 20-point or 99-point strength scale and plotted distributions of the resulting ratings. They find 25% more variability in ratings for studied than for new items, which they interpret as providing evidence that latent mnemonic strength distributions are 25% more variable for studied than for new items. We show here that this conclusion, as well as those from ROC analysis, are critically dependent on assumptions—so much so that these assumptions determine the conclusions. In fact, opposite conclusions, such that study does not affect the variability of latent strength, may be reached by making different though equally plausible assumptions. Because all measurements of mnemonic strength variability are critically dependent on untestable assumptions, all are arbitrary. Hence, there is no principled method to assess the relative variability of latent mnemonic strength distributions.

KEYWORDS: Recognition Memory; ROC Curves; Signal Detection Models.

Latent Mnemonic Strengths are Latent: A Comment on Mickes,
Wixted, and Wais (2007)

It has long been debated whether memory is served by a single mnemonic process or several distinct ones (Schacter & Tulving, 1994; Wixted, 2007; Yonelinas & Parks, 2007). A current focus in this debate is the role of these processes in recognition memory. In a recognition memory paradigm, participants decide at test if items were previously studied or are new. The results are often summarized with receiver operating characteristic (ROC) curves. Empirically observed ROC curves tend to be asymmetric around the negative diagonal (see the dashed line in Figure 1B). These asymmetries, first popularized by Ratcliff, Sheu, & Gronlund (1992) and replicated repeatedly (see Glanzer, Kim, Hilford, & Adams, 1999; Yonelinas & Parks, 2007, for a review) have served as a first-order phenomenon to be explained by theories of mnemonic process.

One single-process model of ROC asymmetry is the unequal-variance signal detection model. This model posits that participants evaluate the mnemonic strength of items against criteria. The strengths for studied and new items are distributed as normal random variables. The effect of study is in two parameters: it both displaces the mean and increases the variability of the strength distribution. If the distributions for studied and new items have the same variance, then the model predicts symmetric ROC plots (Figure 1B, solid line), which is not characteristic of observed data. Asymmetric ROCs result if the variances are unequal (Figure 1A, dashed line), and the unequal-variance model agrees with observed data (Wixted, 2007). For the unequal-variance normal model, the degree of asymmetry in ROCs is a function of the ratio of standard deviations of the underlying strength distribution. Let σ_n and σ_s denote the standard deviations of the strength distributions for new and studied items, respectively. Glanzer et al. (1999) and Yonelinas & Parks (2007) performed large meta-analyses with the unequal variance

normal model and found that estimates of σ_n/σ_s tend to be around .8; that is, the studied-item distributions tend to be 25% larger in standard deviation than new-item distributions. Figure 1C shows z -ROC curves; if strengths are distributed as normals, then the resultant curves are straight-lines with slopes given by σ_n/σ_s .

Mickes, Wixted, and Wais (2007, hereafter referred to as MWW) advocate a new and seemingly direct test of the ratio of variability for latent strength distributions without recourse to ROC analysis. Participants rate confidence on a unidimensional scale with a large number of options. MWW plotted the distributions of ratings for studied and new words; Figure 2 shows their results. The distribution for studied items (solid bars) is more dispersed than that for new items (hatched bars). For the purposes of this paper, we term these distributions as *response-category distributions* and the ratio of the standard deviations of these distributions as the *response-category standard-deviation ratio*. The key finding of MWW is that the response-category standard-deviation ratio is .83, which is quite close to the .8 value observed from ROC analysis. Moreover, MWW find that this near equivalence holds on a participant-by-participant basis. Based on these equivalences, MWW conclude that latent strengths are more variable for studied than for new items. They use this conclusion to bolster support for the unequal-variance signal detection model, a single-process account of recognition memory performance.

We argue here that MWW's conclusion about the variance of latent distributions is unjustified. The problem is that this conclusion is exceedingly dependent on initial assumptions, so much so that the assumptions determine the conclusion. In fact, an opposing conclusion may be reached by relying on different but equally plausible assumptions. Moreover, MWW's assumptions are not testable—as a matter of mathematical logic, there is no way of gathering evidence for or against them. As is shown here, there are unavoidable limits on what may be learned about latent strength distributions.

Latent Distributions and ROC Analysis

Researchers assume that latent strengths are normally distributed when interpreting ROC asymmetries as an index of the ratio of variance (see Figure 3). Normality is a common assumption for analyzing observables in experimental psychology and underlies t -tests, ANOVA, and regression. The assumption of normality is benign in most applications because these tests are robust to moderately large violations (Young & Veldman, 1965). Restated, as long as the underlying distributions are not extremely different than normals, these tests have real Type I error rates not much inflated over the nominal values. Moreover, concerned researchers can always check whether their data are normally distributed and use nonparametric tests if gross deviations are detected.

The situation is different for the analysis of ROC data. We follow Egan (1975) and show here that distributional assumptions, such as normality, completely determine the ratio of standard deviations. Figure 3A-C show three examples in which ROCs are exactly identical, yet the standard-deviation ratio varies with distributional assumptions. The top row shows the case in which strengths are distributed as normals with greater mean and variance for studied than new items. The standard-deviation ratio, σ_n/σ_s in the figure is .80. The resulting asymmetric ROC and straight-lined z -ROC are shown. A criterion is drawn (vertical dashed line) for illustrative purposes and it corresponds to a hit and false alarm rate of .745 and .250, respectively. These values are shown as a point on the ROC and z -ROC plots.

Figure 3B shows the case for a different set of distributions: *log-normals*. These particular distributions were obtained by exponentiating the normals in the top row. The standard deviation ratio for these distributions is .114, which is more extreme than the .8 ratio in the top row. Surprisingly, the single criteria at a false-alarm rate of .250 corresponds to a hit rate of .745, the same values as with the normals. In fact, even though the distributions in Figure 3B are not normal, the resulting ROC and z -ROC plots

are *exactly identical* to the normal-distribution ROC and z -ROC plots in Figure 3A, respectively. Because the normals and log-normals in 3A and 3B produce exactly identical ROC and z -ROC plots, they can never be differentiated with ROC data, no matter how large the sample size. Importantly, any statistical goodness-of-fit metric from ROC data will yield numerically identical support for the distributions in 3A and 3B.

Figure 3C shows a different example of the same equivalence. The distributions come about by passing the normals through an inverse-probit transform as follows: Let X be a normal. The distributions shown, denoted Y , are produced from $Y = \Phi(2X/3)$, where Φ is the cumulative distribution function of the standard normal. The standard-deviation ratio for these distributions is 1.1; that is, the noise distribution is more dispersed than the signal distribution. Though the distributions may seem strange, they may be used to model variables with finite upper and lower bounds, such as the firing rates of nerve cells. Even though the new-item distribution is more dispersed than the studied-item distribution, the ROC and z -ROC plots are identical to the previous examples. The three models in Figure 3A-C make exactly identical ROC and z -ROC predictions and can never be distinguished. This equivalence of ROCs demonstrates a key fact—models with different standard deviation ratios may produce identical ROC predictions. Therefore, estimation of standard deviation ratio from ROC data is impossible; any numeric result reflects arbitrary and untestable assumptions.

It is worthwhile to consider how these ROC equivalences occur. Let X_1 and X_2 denote the normal distributions in Figure 3A. In this case $X_1 \sim \text{Normal}(0, 1)$ and $X_2 \sim \text{Normal}(1.5, 1.25)$, where the arguments of the normal are the mean and standard deviation, respectively. Let Y_1 and Y_2 denote the log-normal distributions in Figure 3B. We constructed these as $Y_1 = \exp(X_1)$ and $Y_2 = \exp(X_2)$. Although the log-normal and normal distributions are different, there is an important invariance concerning the areas under the curves. The criterion in Figure 1A is at .67 and divides the bottom 75% of the

noise strength distribution from the top 25%. When these upper 25% are exponentiated, they are all above $\exp(.67)$, likewise, the bottom 75% of Y_1 are below $\exp(.67)$. The same holds true for the signal distribution. The criterion at .67 divides the bottom 25.5% from the top 74.5%. Consequently, 25.5% and 74.5% of Y_2 are above and below, respectively, the criteria of $\exp(.67)$. The hit and false alarm rates are the percentages above criteria, or .25 and .745, respectively, for X_1 and X_2 and for Y_1 and Y_2 . Hence, the point (.25, .745) is on both ROC curves in Figure 3A and 3B. In fact, this equivalency can be shown for all criteria, and consequently, the ROC curves in Figure 3A and 3B must be exactly identical. The same holds for z -ROC curves.

The equivalence of ROC curves holds as follows: Let $Y_1 = g(X_1)$ and $Y_2 = g(X_2)$. If g is a strictly monotonic function, such as an exponential or logarithm, then the ROC of Y_1 vs. Y_2 is identical to that of X_1 vs. X_2 . In Figure 3C, we constructed the distributions through the function $\Phi(3x/2)$, which is strictly monotonic. There are uncountably many different strictly monotonic functions g , hence there are uncountably many non-normal distribution pairs that give rise to identical ROCs and exactly straight-line z -ROCs. In these non-normal pairs, the slope of the straight-line is unrelated to the standard-deviation ratio.

Not all distribution pairs produce the same ROCs. Figure 3D shows the case of uniforms: $X_1 \sim \text{Uniform}(0, 1)$, $X_2 \sim \text{Uniform}(.5, 1.5)$. These distributions produce ROCs different than normals. We can transform X_1 to a normal by taking $g = \Phi^{-1}$, e.g., $Y_1 = \Phi^{-1}(X_1) \sim \text{Normal}(0, 1)$. We cannot, however, use the same transform on X_2 as $\Phi^{-1}(X_2)$ is undefined when $X_2 > 1$. Hence, there is no common function g that can map two different uniforms each into normals, and, consequently, the ROCs are different than normal pairs.

We discriminate the above argument from Egan (1975) about the identity of ROCs from Lockhart and Murdock's (1970) well-known argument about ROC mimicry. Lockhart

and Murdock noted that many distribution pairs mimic straight line z -ROCs, for instance, the gamma distributions shown in Figure 3E produce nearly straight-line z -ROCs. Most recognition models, such as the dual-process model (Yonelinas, 1994), the extreme-value distribution model (DeCarlo, 1998), or the mixture-of-normals model (Decarlo, 2002) mimic straight-line z -ROCs rather than predict linearity exactly. These distribution families that mimic straight line z -ROCs may be theoretically distinguished from normals with exceptionally large sample sizes. Egan's argument is stronger than Lockhart and Murdock's because ROCs from transformed distributions are exactly identical to the original distributions, as shown in Figure 3A-C. The fact that standard deviation ratios vary across distribution pairs that yield identical ROCs and z -ROCs (Figure 3A-C) shows these ratios reflect nothing more than arbitrary and untestable distributional assumptions.

Response-Category Distributions

MWW note that their inferences about standard deviation ratios from ROCs are dependent on the assumption of normality. This dependence, in fact, is part of the rationale for measuring standard deviation from response category distributions. Response categories form an *ordinal scale*; that is, higher ratings indicate greater strength. Measurements of central tendency and dispersion, however, are predicated on interval scales, which are stronger than ordinal scales. For response categories to form an interval scale, differences in ratings must be linearly related to differences in latent strength. For instance, the difference in strength between Categories "14" and "15" must be the same as that between Categories "4" and "5." Clearly, the interval scale assumption is too strong as there is no reason to believe that the difference between Categories "14" and "15" is the same as that between Categories "4" and "5." The interval-scale assumption is equivalent to assuming that the criteria on latent strength are equally spaced and cover the support of the distributions. As shown by MWW, if this assumption does not hold true, then the

variance of the response categories will not reflect that of the underlying latent distributions. MWW provide the example in Figure 4 in which an equal variance normal model may give rise to larger response-category standard deviations for studied than for new items, or vice-versa. Once again, the inference about standard deviation ratios is critically dependent on an assumption; furthermore, this assumption is untestable.

MWW use far more response categories than usual (20 and 100 in Experiments 1 and 2, respectively). Unfortunately, adding more response categories does not ameliorate measurement difficulties. As an extreme case, consider an infinite number of response categories, which may be implemented by having participants turn a dial to indicate the strength of the test item. It is reasonable to ask whether this approach can provide a principled assessment of latent distribution variability, because if so, then a large number of response categories may be seen as an approximation to setting a dial. The situation seems promising because unlike response category data, the values of the dial may be measured on a ratio scale (such as angular displacement) and judgments may be made directly without recourse to criteria. One problem, however, is that these judgments are made on a physical scale of angular displacement rather than on a psychological scale of mental strength. Psychologists have long known that the transformation between the two is not trivial (e.g., Fechner, 1966; Stevens, 1957). If this transformation is linear then the variance of the distributions on the dial settings reflect the relative variance of latent strengths. Yet, this assumption of linearity is problematic as the resulting conclusions about variance are completely dependent upon it. If the transformation is logarithmic or exponential, then the wrong conclusion will be reached. Moreover, there is no way of testing the linearity assumption in MWW's paradigm.

The Numerical Equivalence of the Standard Deviation Measures

MWW show that standard deviation ratios, whether computed by ROC or by response category, are about the same. Moreover, this near equivalence holds, more or less, on a participant-by-participant level; in fact, the measures correlate .83 and .61 in two experiments. As we have shown, the normality assumption is critical in interpreting the ROC standard-deviation ratio; the equal-spacing assumption is critical in interpreting the response-category standard deviation ratio. Perhaps this numerical equivalence indicates that both assumptions are likely and that both standard-deviation ratio values are valid.

MWW's data, however, show that both assumptions cannot hold simultaneously. If both held, then the distributions in Figure 2 would be normal. Instead, the distributions have substantial and opposing skewness. The implication is that at least one of the assumptions is wrong. Given that at least one of the assumptions is wrong, and hence at least one of the ratios does not measure the intended construct, it is difficult to take seriously the near equality of the two measures.

Two Equivalent Models

We demonstrate the arbitrariness of MWW's conclusion that studied-item latent strengths are more variable than new-item latent strengths by constructing two completely equivalent models of their data from Experiment 1. The first account, shown in Figure 5A, is an unequal-variance account, which is in line with MWW's conclusions. The standard deviation of the studied item distribution was fixed to 1.25, hence the true slope of z -ROC curves is .8. Free parameters were d' and the criteria (shown as dotted vertical lines); these free parameters were estimated by minimizing the mean-squared error between the predicted response-category proportions and observed proportions. The resulting predicted response-category distributions are shown in Figure 5C. These distributions correlate .98 with the empirical distributions obtained by MWW (see Figure 2). The

second account is shown in Figure 5B. The distributions in Figure 5B have about equal variances (the standard deviation ratio is .99). We constructed these distributions by monotonically transforming the distributions in Figure 5A; in this case the transform was $\Phi(2X/3)$, where X denotes the normal distribution in Figure 5A. Because this transform is strictly monotonic, the resulting z -ROC curves are exactly straight lines with a slope of .8. With the drawn criteria¹, the model yields exactly identical predictions as the model in Figure 5A, namely the predictions in Figure 5C. Hence, a model with equal variances explains important aspects of MWW’s data as well as one with unequal variances.

As mentioned previously, MWW report a positive correlation between ROC and response category standard deviation ratios across participants. To explore whether these correlations are expected under the equal-variance model of Figure 5B, we simulated data. In our first simulation, we assumed that all participants shared the same underlying mnemonic strength distributions and used the same criteria. All of the variability, therefore, was due to sampling noise. For each hypothetical experiment, we generated data for 14 participants, each tested on 150 studied and 150 new items. These sample sizes are those from MWW, Experiment 1. From these data, we computed both standard deviation ratios for each participant (see MWW for details). Over 1000 such hypothetical experiments, the average correlation between these two ratios was .77, with 95% of the values between .46 and .93. In a second simulation, we set a high degree of participant variability in sensitivity, criteria, and true standard deviation ratio.² The averaged correlation between the ratios was .60, with a 95% of the values between .14 to .88. These simulation results show that correlations in standard-deviation ratios are expected. They are unsurprising as both measures are conditioned on the same raw data. In sum, there are no aspects of MWW’s data that are incompatible with an equal-variance model.

Conclusion

Latency and Variability

The key finding of the above analyses is that it is logically impossible to measure the ratio of variability across latent distributions. Any conclusion about the ratio simply reflects *a priori* assumptions that are not testable. MWW used the measurement of latent variability to bolster support for unequal-variance normal signal detection, their preferred single-process explanation of recognition memory phenomena. We show here that MWW's data offer no support for unequal variances in particular or for single-process accounts in general.

We suspect that ROC analysis will prove important and helpful in adjudicating between single and multiple process accounts. We are not, however, convinced that the current approach of specifying parametric models is best. Popular models (such as the unequal-variance normal model or Yonelinas' dual-process model) are perhaps specified too finely as ROC data provide only ordinal constraints on latent strengths. As an alternative, researchers may wish to focus on ordinal rather than parametric properties of latent strength distributions.

Dominance of ROC Curves

One promising ordinal property is dominance of ROC curves, which is illustrated in Figure 6 and defined as follows: Consider an experiment with two levels of a factor manipulated at study; for example, two levels of study duration. Figure 6A provide an example of strength distributions for new items (solid) as well as for studied items from Condition 1 (dashed) and Condition 2 (dash-dotted). Studied-item strength in Condition 1 is unambiguously larger than that in Condition 2. More formally, the studied-item strength distribution for Condition 1 *stochastically dominates* that for Condition 2. Stochastic dominance implies that CDFs of the studied-item distributions order as in

Figure 6B. This stochastic dominance implies that the ROC curves order as well and never cross (see Figure 6C). The right column of Figure 6D-F provides an example where stochastic dominance is violated. In Figure 6D, for example, the strongest strengths for Condition 1 are stronger than the strongest strengths for Condition 2, yet the opposite holds for the weakest strengths. Consequently, the ROCs do not order (Figure 6F).

A single-process model would have to be quite complex to handle ROC crossings such as that in Figure 6F. If these cases can be documented, such phenomena may be more parsimoniously accounted for by two-process models. Conversely, a systematic failure to find ROC crossings lends support to the applicability of one-process, strength-based accounts. Because ROC dominance is an ordinal property, it is defined without recourse to untestable parametric assumptions.

References

- DeCarlo, L. M. (1998). Signal detection theory and generalized linear models. *Psychological Methods, 3*, 186-205.
- DeCarlo, L. T. (2002). Single detection theory with finite mixture distribution: Theoretical developments with applications to recognition memory. *Psychological Review, 109*, 710-721.
- Egan, J. P. (1975). *Signal detection theory and ROC analysis*. New York: Academic Press.
- Fechner, G. T. (1966). *Elements of psychophysics*. New York: Holt, Rinehart and Winston.
- Glanzer, M., Kim, K., Hilford, A., & Adams, J. K. (1999). Slope of the receiver-operating characteristic in recognition memory. *Journal of Experimental Psychology: Learning, Memory, and Cognition, 25*, 500-513.
- Lockhart, R. S., & Murdock, B. B. (1970). Memory and the theory of signal detection. *Psychological Bulletin, 74*, 100-109.
- Mickes, L., Wixted, J. T., & Wais, P. E. (2007). A direct test of the unequal-variance signal detection model of recognition memory. *Psychonomic Bulletin and Review, 14*, 858-865.
- Ratcliff, R., Sheu, C. F., & Gronlund, S. D. (1992). Testing global memory models using ROC curves. *Psychological Review, 99*, 518-535.
- Rouder, J. N., & Lu, J. (2005). An introduction to Bayesian hierarchical models with an application in the theory of signal detection. *Psychonomic Bulletin and Review, 12*, 573-604.
- Schacter, D., & Tulving, E. (1994). What are the memory systems of 1994? In D. Schacter & E. Tulving (Eds.), *Memory systems 1994*. Cambridge, MA: MIT Press.

- Stevens, S. S. (1957). On the psychophysical law. *Psychological Review*, *64*, 153-181.
- Wixted, J. (2007). Dual-process theory and signal-detection theory of recognition memory. *Psychological Review*, *114*, 152-176.
- Yonelinas, A. P. (1994). Receiver-operating characteristics in recognition memory: Evidence for a dual-process model. *Journal of Experimental Psychology: Learning, Memory, and Cognition*, *20*, 1341-1354.
- Yonelinas, A. P., & Parks, C. M. (2007). Receiver operating characteristics (ROCs) in recognition memory: A review. *Psychological Bulletin*, *133*, 800-832.
- Young, R. K., & Veldman, D. J. (1965). *Introductory statistics for the behavioral sciences*. New York: Holt, Rinehart, and Winston.

Author Note

Address correspondence to Jeff Rouder, Department of Psychological Sciences, 210 McAlester Hall, University of Missouri, Columbia, MO 65211 or to rouder@missouri.edu. We thank Geoff Iverson and Paul Speckman for helpful conversations. This research is supported by NSF grant SES-0720229 and NIMH grant R01-MH071418.

Footnotes

¹Let c_1, \dots, c_{19} denote the 19 criteria in Figure 5A and let c_1^*, \dots, c_{19}^* denote the same in Figure 5B. Response category predictions for the two models are preserved when $c_i^* = \Phi(2c_i/3)$.

²Individual d' (before transform) was sampled from $\text{Normal}(1.145, .33)$; individual's σ (before transform) was sampled from $\exp(\text{Normal}(.2, .2))$; individual criteria (after transform) were the order statistics from 19 draws of a $\text{beta}(1.5, 1.1)$ (see Rouder & Lu, 2005, for a discussion and parameterization of the beta distribution). These settings correspond to a large amount of individual differences on all parameters.

Figure Captions

Figure 1. The effect of variance in the unequal-variance signal detection model with normal distributions. **A.** The solid line shows a studied-item distribution with equal variance to the new-item distribution; the dashed line shows unequal variances. **B.** Corresponding ROC curves. The solid line (equal variance) is symmetric around the negative diagonal while the dashed line (unequal variance) is asymmetric. The negative diagonal is shown as a dotted line. **C.** Corresponding z -ROCs curves are straight lines with slopes reflecting the ratio of standard deviations.

Figure 2. Response-category distributions from Mickes, Wixted and Wais (2007), Experiment 1. The standard deviation for studied items (solid bars) is 20% greater than that for new items (hatched bars). Figure adapted from Mickes et al., p. 860; permission pending.

Figure 3. Latent strength distributions and corresponding ROC and z -ROC plots. The first three panels (A-C) yield exactly identical ROC and z -ROC curves even though the relationship of standard deviation differs. **A.** Normal distributions with unequal variance. **B.** Log-normal distributions. **C.** Inverse probit transforms of unequal-variance normals (see text for details). **D.** Uniform distributions. The ROC curves are discriminable from normal distribution ROC curves. **E.** Gamma distributions. The ROC curves closely mimic but are not exactly identical to normal-distribution ROC curves.

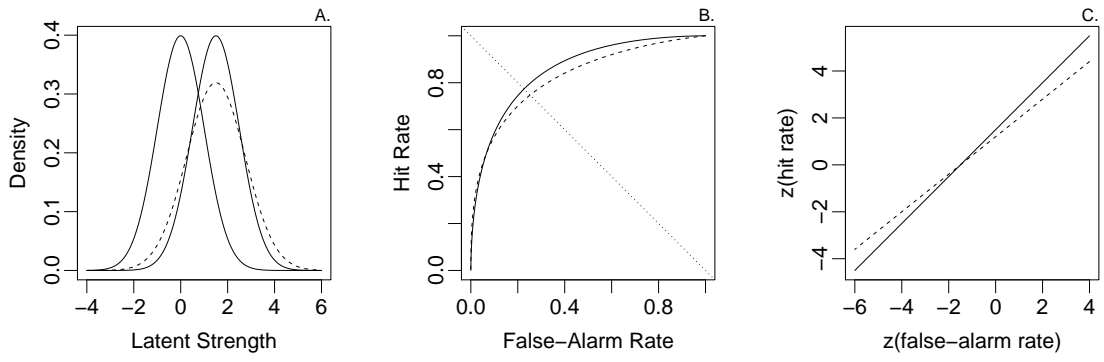
Figure 4. Criteria placement determines standard deviation ratio. **A.** The response category mass is more concentrated in a few categories for new items relative to studied ones. The resulting response-category standard deviations are larger for studied than for new items. **B.** The response category mass is more concentrated for studied items relative to new ones resulting in the reciprocal standard deviation ratio. Figure is adapted from

Mickes, Wixted, and Wais, 2007, p. 862; permission pending.

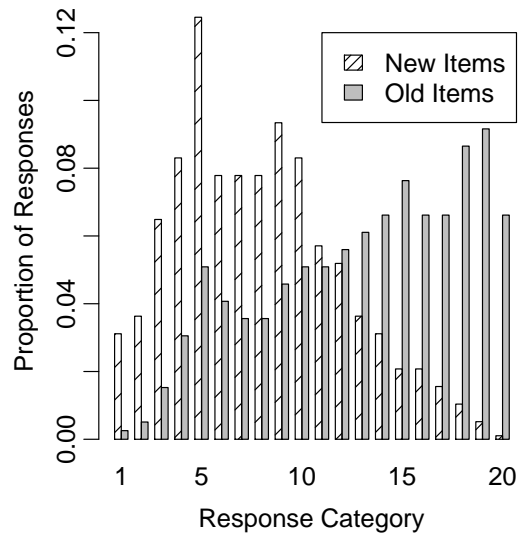
Figure 5. Equivalent accounts of MWW's data. **A.** An unequal-variance account. **B.** An equal-variance account. **B.** Resulting response-category distributions are exactly identical for both accounts.

Figure 6. ROC dominance (left column) and a violation thereof (right column). **A, D.** Strength distributions for new items (solid) and for studied items in Condition 1 (dashed) and Condition 2 (dash-dotted). **B, E.** Cumulative probability distribution functions of studied-item strength distributions. **C, F.** Resulting ROC curves.

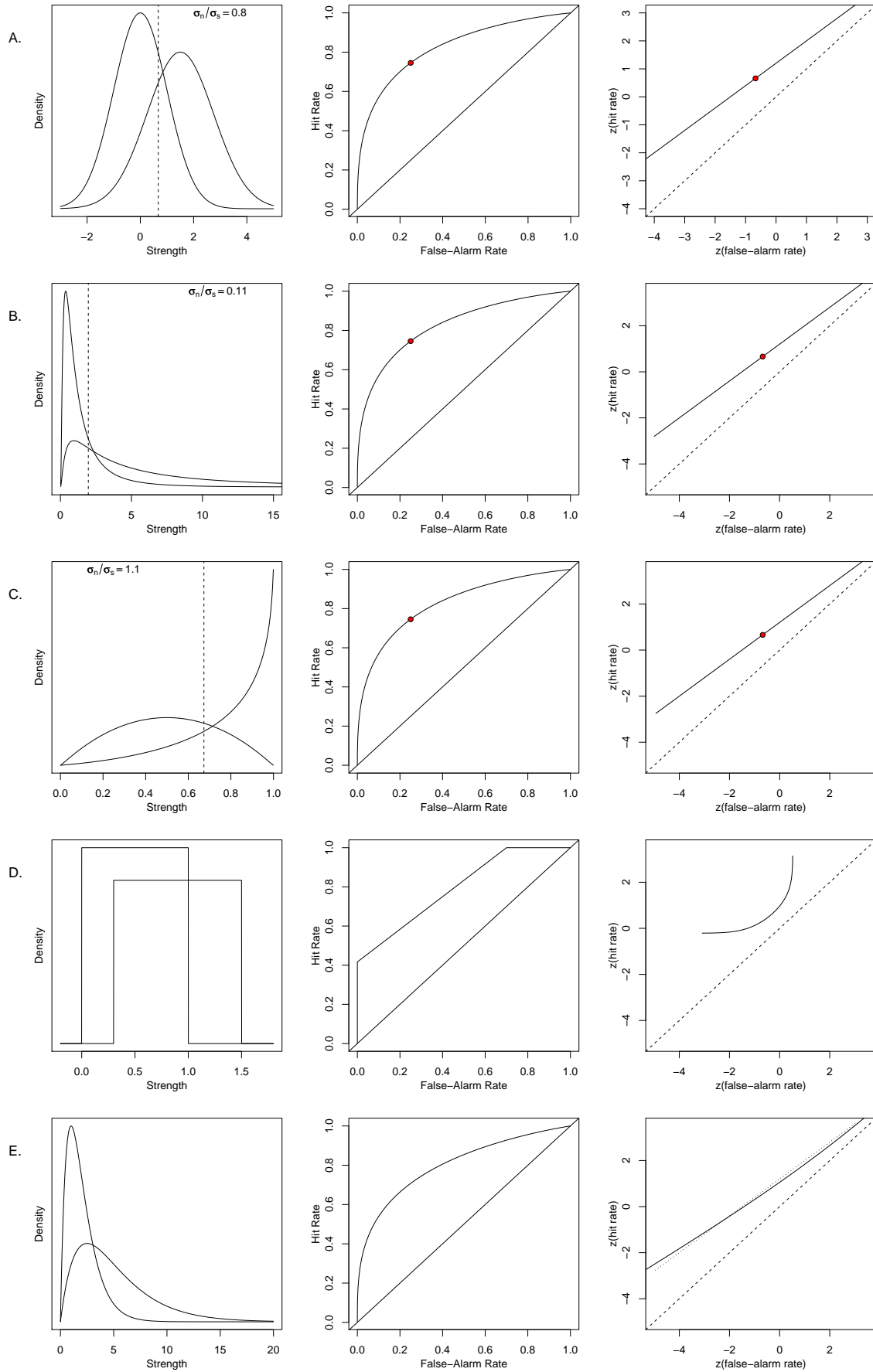
Comment on Mickes et al., Figure 1



Comment on Mickes et al., Figure 2



Comment on Mickes et al., Figure 3



Comment on Mickes et al., Figure 4

