

1

2

3

4 Easy-to-Compute Response Times Based Statistics for Detecting Aberrant Behaviors of

5

Test-takers

6

7

8

Zhen Li

9

eMetric, LLC

10

Nathan L. Wall

11

eMetric, LLC

12

13

August 8, 2022

14

15 Address all correspondence to: Zhen Li, Research Scientist, eMetric, LLC. [211 N Loop](#)  
16 [1604 E, Suite 170](#), San Antonio, Texas 78232. Email: [zli@emetric.net](mailto:zli@emetric.net)

17

18

1 Abstract

2 In this study, we develop an easy-to-compute statistic for detecting examinees'  
3 aberrant response times in large-scale computer-based assessments. Using this statistic,  
4 a response time is flagged as aberrant when it is longer or shorter than expected. The  
5 flagged response times are summarized to indicate examinees' abnormal test behaviors,  
6 such as pre-knowledge, rapid guessing and item memorization. A simulation study was  
7 conducted to evaluate this method's performance in various conditions. Results showed  
8 that the proposed statistic approximately followed a normal distribution in the null  
9 condition, performed equivalently well to van der Linden & Guo's (2008) Bayesian  
10 procedure in detecting aberrant response times, and reduced computational burden  
11 monumentally. A high-stake educational assessment was used to illustrate its  
12 application.

13 Key words: test security breach, response times, abnormal test behavior

14

1 Easy-to-Compute Response Times Based Statistics for Detecting Aberrant Behaviors of  
2 Test-takers

3 Introduction

4 A breach to test security may have serious implications for the psychometric integrity of  
5 the reported test scores and on the interpretations and consequences of those scores  
6 (Standards for Educational and Psychological Testing, 2014; p. 225). Statistical methods  
7 developed for test security purposes have become increasingly popular (Drasgow,  
8 Levine & Williams, 1985; Reise & Due, 1991; Impara et al., 2005; van der Linden & Guo,  
9 2008; Marianti & et al., 2014; Li & Smith, 2015; Fox & Marianti, 2017). So far, statistical  
10 methods based on examinees' responses are the most studied, such as erasure analysis  
11 indices, answer copying indices and person fit indices. Applying statistical methods for  
12 test security purposes provides methods that are highly efficient but low cost, as they  
13 serve as a screening tool before more expensive investigations are conducted.

14 Compared to item responses, the study of response times is relatively new. According  
15 to recent research, response times can provide valuable information for improving test  
16 development, test security, as well as score use and interpretability. For example, Wise  
17 & Kong (2015) studied examinees' engagement during tests based on response time  
18 patterns. Fox & Marianti (2016) explore the relationship between response speed and  
19 accuracy. It is well known that traditional person fit indices examine the congruence  
20 between an item response pattern and a specified item response theory model (Reise &

1 Due, 1991). Similarly, aberrant test behaviors can be identified by comparing expected  
2 response times and observed response times. Recently, using latent variable modeling  
3 of response times, several methods have been developed for detecting aberrance in  
4 response time patterns (van der Linden & Guo, 2008; van der Linden, 2009). However,  
5 these methods have the shortcomings of being highly complex to implement and  
6 computationally ineffective. The purpose of this study is to develop easy-to-compute  
7 statistics for detecting aberrant response time patterns in large-scale assessments.  
8 The practical value of the three proposed statistics are illustrated and expanded in the  
9 paper. So far, no adequate approach has been developed to detect aberrant response  
10 time patterns during testing. Most of the current methods remain in the psychometric  
11 lab and are not widely used in practice. Either the field needs to develop new software  
12 to implement the more complex procedures based on response time models, or develop  
13 an easy-to-compute index, so that existing testing software can employ these methods  
14 without additional effort.

#### 15 Data Forensic Methods and Response Time Modeling

16 Post-hoc data analysis for test security purposes has long existed. Early studies focused  
17 on the similarities between paired examinees' answers, indices of answer copying, and  
18 the likelihood of test-takers' responses with well-known test theory models.  
19 Nonetheless, test security breaches remain a significant issue in large-scale assessments  
20 since many existing statistical methods are only based on examinees' responses and a

1 large number are highly complex. Recently, with the popularity of computer-based  
2 assessments, more information can be captured during the testing process. Therefore,  
3 new statistics have been proposed for this purpose.

4 Early studies focused on developing statistics used for detecting aberrant item response  
5 patterns, which are commonly known as person-fit indices (Reise & Due, 1991). The  
6 shift to computer-based testing allows examinees' response times to be easily recorded,  
7 however, and the response time data can be used for test security analyses. van der  
8 Linden & Guo (2008) proposed a Bayesian procedure to identify aberrant response time  
9 patterns, specifically those indicating pre-knowledge and item memorization in  
10 adaptive testing. Their method was based on a hierarchical latent response time model.  
11 In their first step, latent variables and item parameters are estimated based on a set of  
12 real data using Monte-Carlo Markov-Chain (MCMC). In the second step, the posterior  
13 distribution of each test-taker's response time is calculated. In the third step, the  
14 observed response time for each test-taker on each item is compared with the posterior  
15 distribution. A  $p$  value is computed, assuming that the posterior distribution is log-  
16 normal.

17 Qian et al. (2016) applied the Bayesian procedure to detect item pre-knowledge and  
18 potentially compromised items in two computer-based large-scale licensure  
19 examinations. The results indicated this procedure was helpful in monitoring aberrant

1 examinee behaviors, as well as enhancing future item writing. However, this is the only  
2 publication that could be found applying the Bayesian procedure.

3 Marianti & et al. (2014) developed another set of statistics for detecting aberrant test  
4 behaviors based on the lognormal response time model. The statistics were derived  
5 from the well-known person fit statistic  $l_z$  (Reise & Due, 1991). The log likelihoods of  
6 the response time patterns were used to evaluate the fit of a response time pattern to a  
7 specified model. Furthermore, Fox & Marianti (2017) developed a person fit index with  
8 a hierarchical response time model, taking into consideration both response time and  
9 response accuracy. Additionally, Wang & Gong (2015) proposed a hierarchical mixture  
10 response time model for detecting examinees' engagement during testing, however this  
11 may be harmful to test validity. All these methods have their own merits in different  
12 ways. Nonetheless, the complexity of the methods based on latent variable modeling  
13 constrain them from being widely used in practice. Additionally, these methods are  
14 often computationally demanding.

### 15 Response Time Models

16 Using response time to explore and interpret test-takers' testing behaviors and  
17 outcomes is not necessarily a new approach. Thissen (1983) developed an extended item  
18 response theory model, taking into account person speed and ability simultaneously.

19 Wang & Hansen (2005) developed a four-parameter logistic response time model,  
20 incorporating response time information to predict the probability of answering an item

1 correctly. With these models, however, the aberrance of response time patterns is not  
 2 known.

3 One of the popular response time models is van der Linden's (2006) lognormal model.  
 4 The model assumes that test-takers' speed remains constant during testing, and that  
 5 examinees answer each item independently. It contains three latent variables:  $\tau_j$   
 6 represents the speed of test taker  $j$ ,  $\beta_i$  is the time intensity of item  $i$ , and  $\alpha_i$  is a  
 7 discrimination parameter. This model is often presented as a lognormal density for the  
 8 distribution of test-takers' response times ( $RT_{ij}$ ):

$$f(RT_{ij}; \tau_j, \alpha_i, \beta_i) = \frac{\alpha_i}{RT_{ij}\sqrt{2\pi}} e^{\{-\frac{1}{2}[\alpha_i(\log RT_{ij} - (\beta_i - \tau_j))]^2\}}. \quad (1)$$

9 The advantages of this model include that, when it is combined with an item response  
 10 model, a hierarchical response time model is formed. Not only could we study the  
 11 relationship among latent variables, the information from response times might also be  
 12 used for improving item calibration and test scoring.

13 van der Linden and other researchers have applied the above-mentioned response time  
 14 model in many aspects. For example, van der Linden & Guo (2008) applied the  
 15 hierarchical lognormal response time model for detecting aberrant test-takers' aberrant  
 16 test behaviors in computer adaptive testing (CAT); van der Linden (2009) proposed  
 17 another bivariate lognormal response time model for the detection of collusion between  
 18 test-takers. van der Linden (2008) used this model to improve the accuracy of item  
 19 selection in a CAT design.

1 As mentioned above, the idea of detecting aberrant test behaviors is not new. Any  
2 differences between the expected response pattern and the observed response pattern  
3 or the expected response time pattern and the observed response time pattern could  
4 imply an aberrant test behavior. The difficulty comes from how we detect the  
5 aberrances efficiently and accurately. Although van der Linden and Guo's (2008)  
6 Bayesian estimation can provide relative accurate latent variable estimates, the  
7 estimation process is computationally demanding. In the following section a simplified  
8 approach is proposed for detecting test-takers' aberrant behavior. Not only is this  
9 approach easy to implement in any testing software, it requires much less computation.

### 10 The Proposed Statistics

#### 11 *Basic Summative Statistics and the Indicator of Aberrance*

12 Suppose there are  $N$  examinees and  $n$  items. The response time of examinee  $j$  to item  $i$  is  
13 indicated by  $RT_{ij}$ . First, the average response time to item  $i$  is indicated by  $\overline{RT}_i$ , which is

$$\overline{RT}_i = \frac{\sum_1^N RT_{ij}}{N}. \quad (2)$$

14 Furthermore, the total test time of examinee  $j$  is indicated by  $RT_j = \sum_1^n RT_{ij}$ . The average  
15 total test time among  $N$  examinees is:

$$\overline{RT}_j = \frac{\sum_1^N RT_j}{N} = \sum_1^n \overline{RT}_i, \quad (3)$$



1 which is also equal to the sum of the average response time on each item. Thereby, the  
 2 expected response time for examinee  $j$  on item  $i$  could simply be the average response  
 3 time across all persons and all items:  $\widehat{RT}_{ij} = \overline{RT}_j/n$ .

4 An aberrance indicator is defined as the difference between the observed response time  
 5  $RT_{ij}$  and the expected value  $\widehat{RT}_{ij}$ :

$$Aberrance = (RT_{ij} - \widehat{RT}_{ij}). \quad (4)$$

6 *Statistic I: Standardized Aberrance*

7 Assuming that examinees' response time to each item follows a normal distribution, the  
 8 simplest statistic for detecting aberrant response times can be obtained by standardizing  
 9 the index of aberrance in equation 4:

$$Z_0 = \frac{RT_{ij} - \widehat{RT}_{ij}}{S.D. (RT_{i.})}. \quad (5)$$

10  $Z_0$  will be used as a baseline for index comparison. It is conjectured that it won't work  
 11 well, as response times don't often follow a normal distribution.

12 *Statistic II: Standardized the Aberrance using Logarithm of Response Times*

13 It is commonly shown that students' response times follow a log-normal distribution.

14 Therefore, we take the logarithm of the response time matrix, and compute a

15 standardized residual index like this:

$$Z_l = \frac{\log (RT_{ij}) - \log (\widehat{RT}_{ij})}{S.D. (\log (RT_{i.}))}. \quad (6)$$

1 Notice that both  $Z_0$  and  $Z_l$  have a strong assumption that all test-takers have an equal  
 2 speed during testing. This assumption can often be violated in the real world. Following  
 3 the idea of a standardized response time, a relative standardized response time is  
 4 proposed. This new statistic considers person speed in the computation.

5 *Statistic III: Standardized Logarithm Aberrance with Person Speed*

6 Van der Linden (2006, 2009, 2011) defined person speed as latent variables in his  
 7 lognormal RT model. As we know, person speed is defined as the item's time loading  
 8 divided by response times on the item:  $\tau_j = \frac{\beta_i}{\log(RT_{ij})}$ , where  $\tau_j$  is assumed to be constant  
 9 across all items.

10 The definition of person speed used here follows the same structure. However, instead  
 11 of using the latent variables, simple summative statistics are used. Assuming that  
 12 person speed is a constant parameter across all items. Examinees' speed is quantified as  
 13 the average total test time divided by the total test time of examinee  $j$ :

$$Speed_j = \frac{\overline{RT}_j}{RT_j}. \quad (7)$$

14 The longer the time an examinee takes for a test, the lower the speed is. With speed  
 15 computed, the expected response time of examinee  $j$  to item  $i$  is obtained, which is the  
 16 average response time on item  $i$  divided by person speed:

$$\widehat{RT}_{ij} = \frac{\overline{RT}_i}{Speed_j}. \quad (8)$$

1 Similarly, when the response times are transferred onto the logarithmic scale, person  
 2 speed can be calculated by  $\overline{\log(RT_j)} - \log(RT_j)$ . Furthermore, considering random  
 3 errors reflected in the variances of response times on each item, person speed is  
 4 calculated by the weighted average response time residual:

$$speed_j = \frac{\sum_{i=1}^n \frac{(\overline{\log(RT_i)} - \log(RT_{ij}))}{var(\log RT_i)}}{\sum_{i=1}^n \frac{1}{var(\log RT_i)}}. \quad (9)$$

5 The expected value of  $\log(RT_{ij})$  after adjusting for person speed is  $(\overline{\log(RT_i)} - speed_j)$ .  
 6 The variance of  $\log(RT_{ij})$  is the total variance minus the group (person) variance:

$$var(\log RT_i) - var(speed_j). \quad (10)$$

7 Thereby, the standardized value of aberrance taking into consideration person speed  
 8 variance is as follows:

$$Z_s = \frac{\log(RT_{ij}) - (\overline{\log(RT_i)} - speed_j)}{\sqrt{var(\log RT_i) - var(speed_j)}}. \quad (11)$$

9 In the following two sections, we use a simulation study and an empirical study to  
 10 examine the performance of these three proposed statistics.

### 11 Simulation Study

12 A simulation study was conducted to examine the performance of the proposed  
 13 statistics in various conditions. In the simulation study, the proposed three statistics  
 14 based on Z-score method were compared with two statistics based on van der Linden &  
 15 Guo's latent variable modeling method. Furthermore, four factors, including proportion

1 of compromised items, sample size, test length, as well as the population aberrance  
2 rates, might have an influence on the performance of these statistics. The effects of the  
3 four factors were examined in this study.

#### 4 *Simulation Design*

5 Four factors were considered in the simulation study. Sample size and number of items  
6 in the test were considered because both are critical features of any assessment. A pilot  
7 study showed that the proposed statistic performed equivalently well in various  
8 conditions of sample sizes (500, 1000, 10000). Therefore, in the current simulation study,  
9 only the small sample size ( $N=500$ ) condition was considered. Additionally, the  
10 aberrance rates of items (proportion of compromised items), and the aberrance rates in  
11 the sample (proportion of test-takers who have aberrant test behaviors) were  
12 considered based on previous research (Marianti & et al., 2014). The details of these  
13 factors are listed in Table 1.

#### 14 *Data Generation*

15 1) Generating observed response times

16 Response times were generated based on a lognormal model (van der Linden, 2006):

$$\log(RT_{ij}) = \beta_i - \tau_j + \varepsilon_{ij} , \quad \varepsilon_{ij} \sim N(0, 1/\alpha_i^2) \quad (12)$$

17 where  $\beta_i$ ,  $\tau_j$ , and  $\alpha_i$  were obtained from empirical data analysis. The longest response  
18 time for an item is constrained to be 20 minutes, and the shortest response time on one  
19 item is constrained to be 1 second, according to students' observed response times for

1 multiple choice items in state assessments. Outliers, e.g., an extremely long response  
2 time or short response time, were not a factor of interest, therefore were not simulated  
3 in the current study.

#### 4 2) Generating aberrant response times

5 Aberrant RTs were generated by three steps. First, a proportion of items were chosen  
6 according to the factor ARI in the simulation design. Second, a proportion of examinees  
7 were randomly selected from the sample according to the factor ARS in the simulation  
8 design. Third, with true person parameters and item parameters, the lognormal  
9 distribution of response times for the selected person on the chosen item was known.  
10 An aberrant RT is generated by taking a value of the cut points which are located at  
11 three standard deviations from the mean (half negative and half positive).

#### 12 *Statistics of Interest*

13 The proposed three easy-to-compute statistics ( $Z_0, Z_1, Z_s$ ) are of major interest in this  
14 study. In addition, two statistics based on van der Linden's (2006) lognormal response  
15 time model are computed for comparison.

16 After fitting the generated response time matrix with the lognormal model, parameter  
17 estimates  $\hat{\beta}_i$ ,  $\hat{\alpha}_i$  and  $\hat{\tau}_j$  can be obtained. Assuming that these parameter estimates were  
18 the true parameters, the expected response time for examinee  $j$  on item  $i$  is  $(\hat{\beta}_i - \hat{\tau}_j)$ , and  
19 the variance of response times on this item as  $(1/(\hat{\alpha}_i)^2)$ . A standardized residual index  
20 for detecting aberrance will be the fourth statistical index:

$$Z_m = \frac{\log(RT_{ij}) - (\hat{\beta}_i - \hat{\tau}_j)}{1/\hat{\alpha}_i}. \quad (13)$$

1 The reason for introducing  $Z_m$  is that it has the same structure as the proposed statistic  
 2  $Z_s$ . The only difference is that  $Z_m$  uses parameters estimated from latent variable  
 3 modeling for persons' speed and items' time loadings.

4 The other statistic is a posterior predictive checking method, taking into consideration  
 5 the posterior distribution of the person speed parameters. This method was first  
 6 introduced by van der Linden & Guo (2008), known as a Bayesian procedure. Based on  
 7 the posterior distribution of response time of person  $j$  on item  $i$ , the  $p$  values and  
 8 standardized residuals were calculated by comparing the observed response time to  
 9 this posterior distribution.

10 These two indices are comparable to the proposed easy-to-compute statistics above, as  
 11 both of them follow a standard normal distribution. The difference lies in whether the  
 12 estimated parameters from latent variable modeling or the observed response times are  
 13 used directly to approximate the expected response times.

#### 14 *Evaluation Criteria*

15 To determine whether a response time pattern is aberrant, the five statistics are  
 16 calculated for each observed response time. RTs with extreme values will be flagged if  
 17 the statistic is higher than 1.96 or lower than -1.96. After this process, each examinee  
 18 will have or not have several flagged RTs. The proportion of flagged RTs of an

1 examinee is computed. This number is then compared with a designated cut (varies  
 2 with total number of items, see attachment). The power of the statistics is compared  
 3 with respect to the following criteria:  
 4 Probability of detection: proportion of examinees flagged with aberrant RT patterns,  
 5 among the actual number of examinees with true aberrant RT patterns (as designed in  
 6 data generation). It is often regarded as a sensitivity index in statistics:

$$Detection\ rate = \frac{N * P_{true\ positive}}{N * (P_{true\ positive} + P_{false\ negative})}, \quad (14)$$

7 Precision of detection: Proportion of examinees with true aberrant RT patterns, among  
 8 the number of examinees flagged as having an aberrant RT pattern. The equation is as  
 9 following:

$$Precision\ of\ Detection = \frac{N * P_{true\ positive}}{N * (P_{true\ positive} + P_{false\ positive})}. \quad (15)$$

## 10 *Results of Simulation Study*

### 11 1) Type I error

12 In the null conditions, the asymptotic distributions of the proposed statistics are  
 13 examined. As the statistics follow a standard normal distribution, when the  $\alpha$  level was  
 14 set to be 0.025 on each side, the empirical rejection rates at both sides should be close to  
 15 0.025.

16 In Table 2, the empirical rejection rates for  $Z_l$  and  $Z_s$  approximate 0.025 on both the left  
 17 side ( $\alpha = -0.025$ ) and the right side ( $\alpha = 0.025$ ). However, the empirical rejection rates

1 for  $Z_0$  were not close to 0.025 on both sides. The large value of empirical rejection rates  
2 on the right side shows that the distribution of  $Z_0$  was positively skewed.

3 In addition, flagging rules for individuals were used to check how likely an individual  
4 with no aberrance will be falsely flagged by the statistics. Results in Table 2 show that  
5  $Z_s$  has zero probability of randomly flagging any examinee when no aberrant response  
6 time pattern exists. On the contrary, both  $Z_0$  and  $Z_l$  flagged 6%-11% examinees  
7 incorrectly by chance.

## 8 2) Power

9 The power of the proposed statistics was tested in various alternative conditions, where  
10 aberrant response times were generated. The probability of detection (P.d.) and  
11 precision of detection (P) of statistics  $Z_0$ ,  $Z_l$  and  $Z_s$  were compared with van der Linden  
12 & Guo's (2008) Bayesian procedure and  $Z_m$ .

13 From Table 3, it is obvious that  $Z_s$  performed better than  $Z_0$  and  $Z_l$ . It had a higher  
14 precision of detection in all the conditions, and had a higher probability of detection in  
15 most conditions. When the number of items were small ( $n=20$ ), and the proportion of  
16 compromised items was low (10%), which meant that only 2 items were compromised  
17 in the test, all statistics had very low ( $\leq 0.08$ ) probability of detection. However, the  
18 precision of detection for  $Z_s$  was much higher than those of  $Z_0$  and  $Z_l$ . As the number of  
19 items increased and the proportion of compromised items increased, both the detection  
20 rates and precision of detection for  $Z_s$  grew. For example, when the number of items



1 was 70, aberrance rates in the sample were 0.05, and the proportion of compromised  
2 items was 25%, the detection rate grew to 97%, and the detection precision grew to 98%.  
3 Furthermore, compared to the two approaches based on latent variable modeling (van  
4 der Linden and Guo's Bayesian procedures and  $Z_m$ ),  $Z_s$  performed equivalently well in  
5 all the conditions we tested. The detection rates and precision of detection of the  
6 proposed new statistic were close to those of van der Linden and Guo's (2008) Bayesian  
7 procedure. In the following section, some factors are discussed that exert an influence  
8 on the performance of the proposed statistic.

### 9 3) The Influence of Factors

10 To better illustrate how each of the statistics' performance were influenced by the three  
11 factors: the proportion of compromised items, test length and the proportion of  
12 examinees with aberrant response times, the following figures are displayed.  
13 Figure 1 shows the influence of aberrance rates of items on the probability of detection  
14 and precision when the number of items was increased from 20 to 70. The aberrance  
15 rate in the sample was fixed to 10%. When the percentage of compromised items was  
16 10%, all five indices had low probability of detection. However, when the percentage of  
17 compromised items increased from 10% to 25%, three of the indices' probability of  
18 detecting aberrant response time patterns improved adequately for the short test ( $n=20$ ),  
19 and improved rapidly for longer tests ( $n=40$  and  $n=70$ ).  $Z_m$  had the highest probability  
20 of detection in the tested conditions, while  $Z_s$  and the Bayesian procedure followed very

1 closely. As a comparison,  $Z_0$  and  $Z_l$  didn't improve significantly when the percentage of  
2 compromised items increased.

3 The bottom three plots in Figure 1 show that,  $Z_s$ , the Bayesian procedure and  $Z_m$  also  
4 performed similarly with respect to precision of detection in various conditions. When  
5 the aberrance rates of items increased to 25%, the precision of detection improved  
6 rapidly for the short test ( $n=20$ ). For longer tests, the precision of these three detection  
7 indices was high even if the aberrance rate of items was low (10%), therefore the  
8 increased rates were small. When the aberrance rate of items reached 25%, no matter  
9 how many items there were, their precision of detection was close to 1. On the contrary,  
10  $Z_0$  and  $Z_l$  had very low precision of detection across all condition levels, with slight  
11 improvement when the aberrance rates of items increased to 25%.

12 Figure 2 shows the influence of aberrance rates of sample on the probability and  
13 precision of detection. The top plots in Figure 2 show that when the aberrance rates of  
14 sample increased from 5% to 10%, the three statistics with high detection rates ( $Z_s$ ,  
15 Bayesian procedure, and  $Z_m$ ) had lower probabilities of detecting aberrant response  
16 time patterns. The other two statistics ( $Z_0$  and  $Z_l$ ) had equivalently low probability of  
17 detection in both conditions. On the contrary, the bottom three plots of Figure 2 show  
18 that with the increased aberrance rate in the sample, all statistics' precision of detection  
19 improved slightly.

1 Both Figure 1 and Figure 2 show the influence of test length on the performance of the  
2 statistics: with the increased total number of items, the probability and precision of  
3 detection improved for all statistics. In particular, the precision of detection increased to  
4 almost 1 with a test length of 70, even if the percentages of compromised items and the  
5 percentages of aberrant examinees were relatively small.

#### 6 4) Summary

7 A simulation study was carried out to evaluate the performance of the proposed  
8 statistics compared to an existing Bayesian procedure (van der Linden & Guo, 2008). In  
9 the null condition, two of the three new statistics,  $Z_l$  and  $Z_s$ , approximately follow a  
10 standard normal distribution. The simplest statistic  $Z_0$  turns out to be extremely skewed  
11 to the right. Specifically, the left-tail p values were close to 0, while the right-tail p  
12 values approximated 0.05.

13 A most interesting finding was that, compared with the two procedures based on latent  
14 variable modeling,  $Z_s$  performed equivalently well in detecting response time  
15 aberrance. Among the proposed easy-to-compute statistics,  $Z_s$  performed much better  
16 than  $Z_0$  and  $Z_l$ . Not only did it have higher precision of detection in all the simulation  
17 conditions, it also had a higher probability of detection in most conditions.

18 Moreover, several factors, especially the aberrance rates in the sample and the  
19 proportion of compromised items, exerted a significant influence on the performance of  
20 the proposed statistics. The more compromised items the test had, the less the

1 examinees with aberrant response time patterns, the more likely that an aberrant  
2 response time pattern could be detected (high probability of detection). However, the  
3 precision of detection increased with both the number of compromised items and the  
4 proportion of examinees with aberrant response time patterns.

5 It was also found that the compromised items will always have more aberrant RTs than  
6 the other non-compromised items. In other words, based on the aberrance flagging of  
7 RTs, the compromised items were flagged correctly.

#### 8 Empirical Study

9 Examinees' response times in a computer-based Math test consisting of 58 items were  
10 analyzed to illustrate the application of the proposed statistics. This test was part of a K-  
11 12 state testing program from the 2016-17 school year. The sample included 6,827 Grade  
12 6 students. Demographically, the sample was diverse. Results from the real data  
13 analysis are discussed below.

14 We found that examinees' response times to most items approximately follow a  
15 lognormal distribution. One item was removed from analysis as its response time did  
16 not fit the lognormal distribution. The QQ-plot in Figure 3 shows that the total response  
17 time on the remaining 57 items was well approximated by a lognormal distribution,  
18 even though the slightly uplifted right tail indicated this distribution was a little heavy  
19 tailed.

1 Next, we computed the proposed statistics, as well as statistics based on Guo and van  
2 der Linden (2008)'s Bayesian procedure. Only a small proportion of aberrant response  
3 time patterns were detected (Table 4). With a close look at these aberrant response time  
4 patterns, possible reasons of the detected aberrance include speededness towards the  
5 end, time management strategies, and rapid guessing. For example, one examinee spent  
6 1-5 seconds through all 57 items, which indicates a rapid guessing behavior. It appeared  
7 that no test-taker was engaged in any potential cheating behavior during testing.

8 Results in Table 4 show that about 5% of the examinees' response times were flagged by  
9 our statistic of interest ( $Z_s$ ). According to the property of standard normal distribution,  
10 this probability is very close to the significance level 0.05. This finding indicates that no  
11 significant aberrance was found in the empirical data set. In addition, it was noticed  
12 that more positive response times were flagged than the negative response times by  $Z_0$   
13 statistic, which further proved a positively skewed distribution of response times. For  
14  $Z_1$ ,  $Z_s$ , and the Bayesian procedure, the difference between positive and negative  
15 flagging decreased to a small amount. Moreover, the extent to which two statistics flag  
16 the same examinee's response time pattern was checked.

17 In Table 5, we found that, among the flagged individuals, more than 88% were flagged  
18 simultaneously by  $Z_s$  and the Bayesian procedure, even if none of the examinees were  
19 identified with cheating behaviors. It is very likely that, when real aberrant response

1 times exist,  $Z_s$  could be used to detect the aberrance as accurately as the Bayesian  
2 procedure, which was also proved in the simulation study.

3 Additionally, the percentages of flagged examinees by each statistic in each school were  
4 computed. Schools were flagged when the percentage of aberrant examinees was higher  
5 than the state-level percentage (bottom row in Table 4). Results show that 23%, 29%,  
6 26%, 27% and 26% schools were flagged by Index  $Z_0$ ,  $Z_l$ ,  $Z_s$ , Bayesian procedure, and  $Z_m$   
7 respectively. Specifically, 96% of 596 schools in the state were simultaneously flagged  
8 by both  $Z_s$  and the Bayesian procedure. Meanwhile, only 63% of schools in the studied  
9 state were flagged by both  $Z_0$  and the Bayesian procedure at the same time.

10 Furthermore, the criterion of flagging a school should be higher than the state-level  
11 average rate in practice. When the flagging cut increased to 0.2, 99% of schools were  
12 flagged by both  $Z_s$  and the Bayesian procedure at the same time.

### 13 Discussion

14 Results from the simulation study demonstrated that the proposed easy-to-compute  $Z_s$   
15 statistic performed well in null and alternative conditions. Specifically, in the conditions  
16 tested,  $Z_s$  had similar probability and precision of detecting aberrant response time  
17 patterns as van der Linden & Guo's (2008) Bayesian procedure. On the contrary, the  
18 other two simpler statistics proposed in this paper ( $Z_0$  and  $Z_l$ ) do not have adequate  
19 power in detecting aberrant response times. Furthermore, it was shown that test

1 lengths, proportion of aberrant examinees in the sample, as well as the number of  
2 compromised items all exert an influence on the performance of the indices.

3 van der Linden & Guo (2008) argued that the Bayesian procedure accounts for the  
4 presence of estimation error in any of the parameters of the psychometric model (e.g.,  
5 ability parameters). Our method, on the other hand, eliminated any estimation error  
6 because no latent variable needs to be estimated. Moreover, we noticed that the  
7 Bayesian procedure had a slightly higher probability of detection than  $Z_m$ , which is also  
8 based on the lognormal response time model. However, the probability and precision of  
9 detection of  $Z_s$  are always close to van der Linden & Guo's (2008) Bayesian procedure.

10 Another important finding was that, when the test was short (20 items) and the  
11 proportion of compromised items was low, all the statistics had a low probability of  
12 detecting aberrant response time patterns. The precision of detection was lower or equal  
13 to 54%. Therefore, it is recommended that no individual-level detection should be  
14 carried out in this situation. Only when the test length is adequate (40 items), are the  
15 detection results by the proposed  $Z_s$  statistic sufficiently reliable.

16 Additionally, empirical data analysis showed that  $Z_s$  could flag a large proportion of  
17 examinees flagged as aberrance by the Bayesian procedure. When the statistics were  
18 aggregated at the school level,  $Z_s$  and the Bayesian procedure almost flagged the same  
19 schools even if the aberrance rates were very low in the real data. This finding provided  
20 further evidence that  $Z_s$  would be a useful and powerful statistic in practice.

1 This study has its limitations. First, the data for the simulation study was generated  
2 with a lognormal response time model. This doesn't take into consideration all the  
3 features of a real response time data set. For example, test-takers' response time on one  
4 item might not follow a lognormal distribution and their response speed might vary  
5 during testing. Secondly, only one set of empirical data set was used to illustrate the  
6 new statistics. This limits the number of types of aberrant test behaviors we detected in  
7 this study. The performance of the proposed statistics need to be examined in more  
8 situations. Furthermore, future study should consider improving the current data  
9 forensic methods by incorporating more information sources, detecting aberrant  
10 examinee behaviors using item response, response time and answer changing patterns.

#### 11 References

- 12 American Educational Research Association, American Psychological Association,  
13 National Council on Measurement in Education, Joint Committee on Standards for  
14 Educational and Psychological Testing (U.S.). (2014). *Standards for educational and*  
15 *psychological testing*. Washington, DC: AERA.
- 16 Drasgow, F., Levine, M. V., & Williams, E.A. (1985). Appropriateness measurement  
17 with polychotomous item response models and standardized indices. *British Journal*  
18 *of Mathematical and Statistical Psychology*, 38, 67-86.
- 19 Fox, J. P. & Marianti, S. (2016). Joint modeling of ability and differential speed using  
20 responses and response times. *Multivariate Behavioral Research*, 51(4), 540-553.



- 1 Fox, J. P. & Marianti, S. (2017). Person-fit statistics for joint models for accuracy and  
2 speed. *Journal of Educational Measurement*, 54(2), 243-262.
- 3 Impara, J. C., Kingsbury, G., Maynes, D., & Fitzgerald, C. (2005). *Detecting cheating in*  
4 *Computer Adaptive Tests using data forensics*. Paper presented at 2005 Annual Meeting  
5 of the National Council on Measurement in Education, Montreal, Canada.
- 6 Marianti, S., Fox, J., Avetisyan, M., Veldkamp, B.P. & Tijmstra, J. (2014). Testing for  
7 aberrant behavior in response time modeling. *Journal of Educational and Behavioral*  
8 *Statistics*, 39(6), 426-451.
- 9 Li, Z., & Smith, J. (2015). *Detecting aberrant behaviors with a hierarchical lognormal response*  
10 *time model*. Paper presented at the annual meeting of the National Council on  
11 Measurement in Education (NCME), Chicago, IL.
- 12 Qian, H., Staniewska, D., Reckase, M. & Woo, A. (2016). Using response time to detect  
13 item preknowledge in computer-based licensure examinations. *Educational*  
14 *Measurement: Issues and Practice*, 35(1), 38-47.
- 15 Reise, S. P. & Due, A. M. (1991). The influence of test characteristics on the detection of  
16 aberrant response patterns. *Applied Psychological Measurement*, 15(3), 217-226.
- 17 Thissen, D. (1983). Timed testing: An approach using item response theory. In D.J.  
18 Weiss(Ed.) *New horizons in testing: Latent trait test theory and computerized adaptive*  
19 *testing*. New York: Academic Press.

- 1 van der Linden, W. J. (2006). A lognormal model for response times on test items.  
2 *Journal of Educational and Behavioral Statistics, 31*, 181-204.
- 3 van der Linden, W. J. (2008). Using response times for item selection in adaptive testing.  
4 *Journal of Educational and Behavioral Statistics, 33(1)*, 5-20.
- 5 van der Linden, W. J., & Guo, F. (2008). Bayesian procedures for identifying aberrant  
6 response-time patterns in adaptive testing. *Psychometrika, 73(3)*, 365-384.
- 7 van der Linden, W. J. (2009). A bivariate lognormal response-time model for the  
8 detection of collusion between test takers. *Journal of Educational and Behavioral*  
9 *Statistics, 34(3)*, 378-394.
- 10 van der Linden, W.J. (2011). Modeling response times with latent variables: principles  
11 and applications. *Psychological Test and Assessment Modeling, 53(3)*, 334-358.
- 12 Wang, T., & Hanson, B. A. (2005). Development and calibration of an item response  
13 model that incorporates response time. *Applied Psychological Measurement, 29(5)*, 323-  
14 339.
- 15 Wang, C. & Xu, G. (2015). A mixture hierarchical model for response times and  
16 response accuracy. *British Journal of Mathematical and Statistical Psychology, 68*, 456-  
17 477.
- 18 Wise, S. & Kong, X. (2005). Response time effort: A new measure of examinee  
19 motivation in computer-based tests. *Applied Measurement in Education, 18*, 163-183

1 Appendix

2 Table 1

3 *Simulation Study Design*

---

Factors	Conditions
Sample size (N)	500
Number of items (n)	20, 40, 70
Aberrance rates of items (ARI)	5%, 10%
Aberrance rates in sample (ARS)	10%, 25%

---

4 \*The factors are fully crossed in a  $1 \times 3 \times 2 \times 2$  design with 500 replications in each  
 5 condition.

6 Table 2

7 *Type I error rates and the probability of false alarm*

---

n	$Z_0$			$Z_l$			$Z_s$		
	$\alpha_{-0.025}$	$\alpha_{0.025}$	$R_{flagged}$	$\alpha_{-0.025}$	$\alpha_{0.025}$	$R_{flagged}$	$\alpha_{-0.025}$	$\alpha_{0.025}$	$R_{flagged}$
20	.000	.053	.08	.023	.023	.06	.026	.024	.00
40	.000	.052	.09	.029	.021	.09	.026	.024	.00
70	.000	.052	.11	.027	.022	.10	.026	.024	.00

---

8  
 9  
 10  
 11

1 Table 3

2 *Statistical power at individual level*

n	A.R. in Sample	A.R. of Items	Proposed Simple Statistics						Latent Model Statistics			
			$Z_0$		$Z_l$		$Z_s$		Bayes		$Z_m$	
			P.d.	P	P.d.	P	P.d.	P	P.d.	P	P.d.	P
20	.05	.10	.08	.06	.08	.06	<b>.05</b>	<b>.49</b>	.03	.54	.04	.50
		.25	.12	.09	.13	.10	<b>.52</b>	<b>.94</b>	.46	.97	.54	.95
	.10	.10	.08	.11	.07	.12	<b>.03</b>	<b>.62</b>	.02	.63	.03	.60
		.25	.11	.17	.12	.18	<b>.26</b>	<b>.95</b>	.20	.97	.26	.96
40	.05	.10	.13	.06	.11	.07	<b>.17</b>	<b>.75</b>	.13	.80	.16	.77
		.25	.19	.09	.21	.12	<b>.86</b>	<b>.95</b>	.86	.97	.88	.96
	.10	.10	.12	.12	.11	.13	<b>.09</b>	<b>.79</b>	.07	.81	.08	.79
		.25	.17	.17	.20	.22	<b>.58</b>	<b>.98</b>	.55	.99	.62	.98
70	.05	.10	.11	.06	.13	.06	<b>.34</b>	<b>.90</b>	.30	.92	.33	.91
		.25	.18	.10	.25	.12	<b>.97</b>	<b>.98</b>	.97	.99	.97	.98
	.10	.10	.12	.12	.13	.13	<b>.15</b>	<b>.91</b>	.12	.91	.14	.91
		.25	.16	.17	.23	.20	<b>.84</b>	<b>.99</b>	.83	.99	.85	.99

3

4 Table 4

5 *Percentages of aberrance flagged by different statistics in empirical data analysis*

Aberrance Type		$Z_0$	$Z_l$	$Z_s$	Bayesian procedure	$Z_m$
	All	3.6%	5.0%	5.0%	5.0%	4.9%
Flagged response times	Positive	3.6%	2.8%	2.8%	2.8%	2.8%
	Negative	0%	2.2%	2.2%	2.2%	2.1%
Flagged Examinees	All	3.2%	5.8%	4.5%	4.5%	4.3%

6

1 Table 5

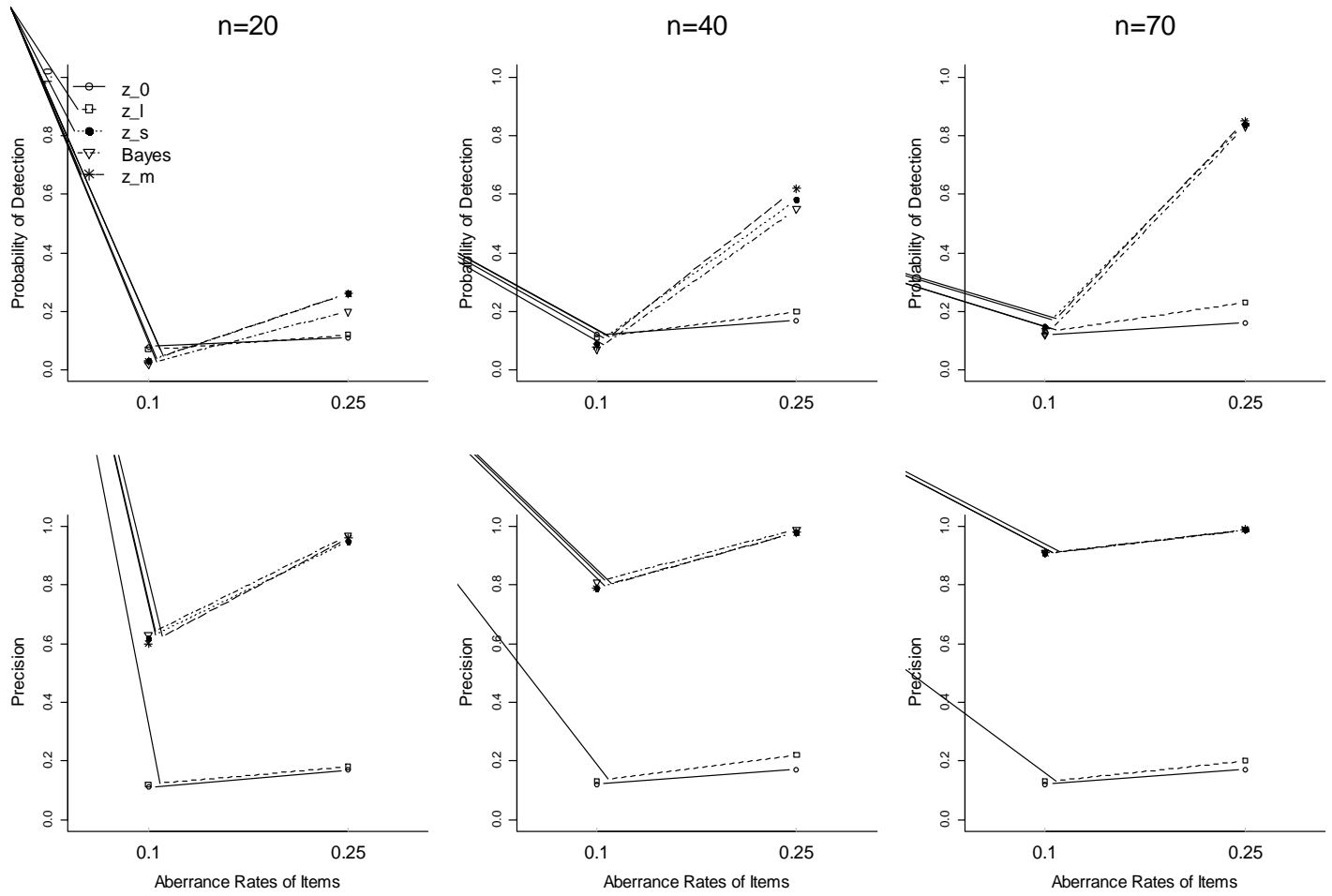
2 *Similarity of examinees' flagging with different statistics*

	$Z_0$	$Z_l$	$Z_s$	Bayesian procedure
$Z_l$	34.8%	--	--	--
$Z_s$	11.8%	45.5%	--	--
Bayesian procedure	11.8%	46.9%	88.4%	--
$Z_m$	11.8%	47.5%	89.6%	93.9%

3

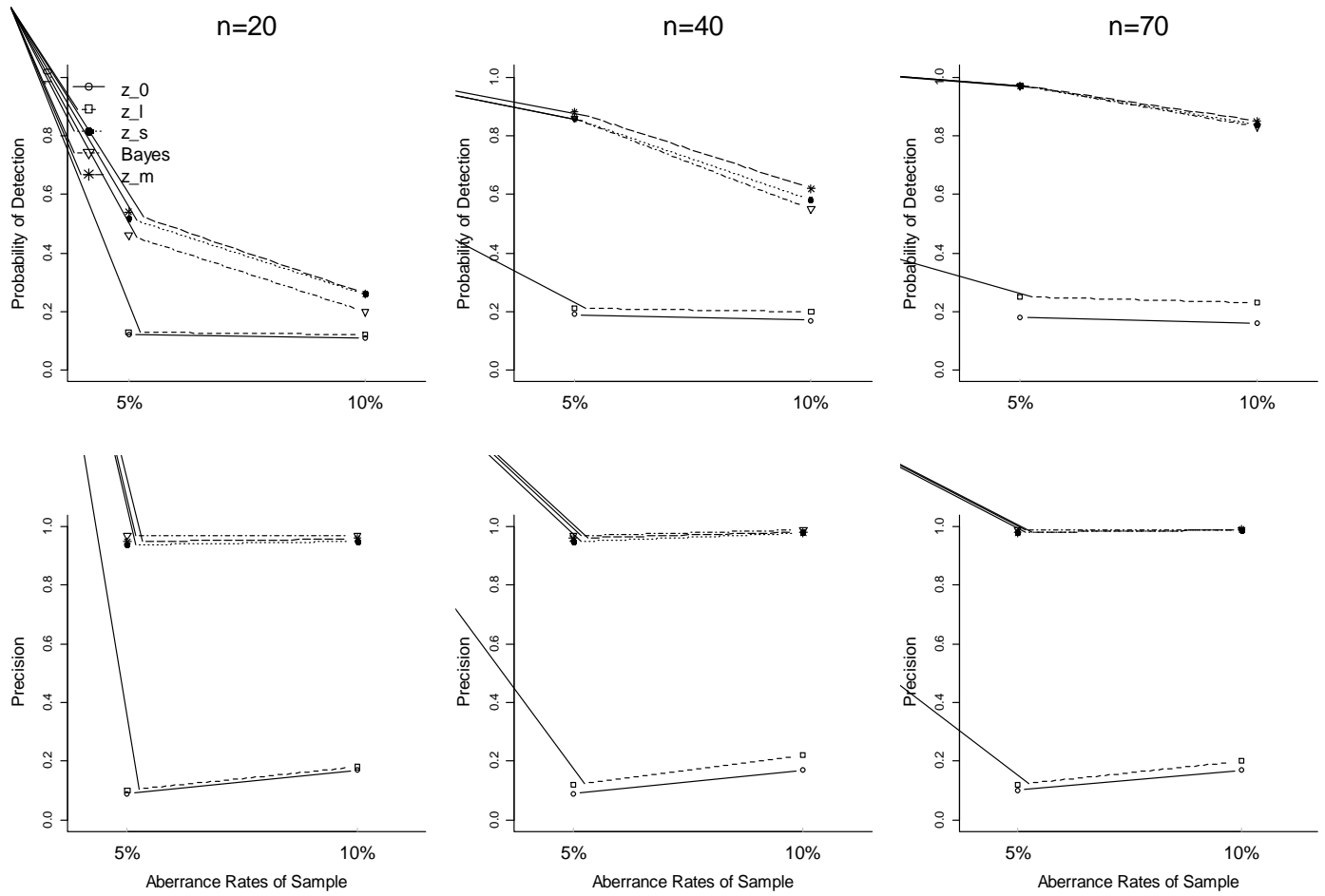
4

# Easy-to-Compute Statistics for Detecting Aberrant Test Behaviors



1

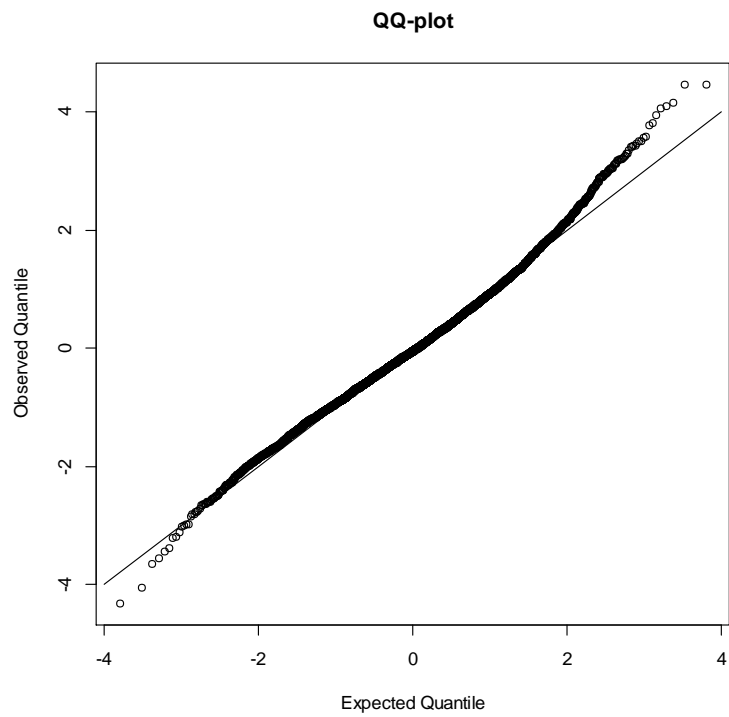
2 *Figure 1. Probability of detection (top) and precision of detection (bottom) of five*  
 3 *statistics across test length and aberrance rates of items*



1

2 *Figure 2. Probability of detection (top) and precision of detection (bottom) of all statistics*

3 *across test length and aberrance rates of sample*



1

2 *Figure 3.* Quantile-quantile plot for total response times on the test in Empirical data

3 analysis

4

5



## 1 Appendix

## 2 Table A1

3 *The likelihoods of randomly flagging a student at or above different cut levels (fixed number of*  
 4 *items)*

<b>Number of Items as the Cut Level</b>										
n	1	2	3	4	5	6	7	10	15	20
<b>20</b>	64.15%	26.42%	7.55%	1.59%	0.26%	0.03%	0.00%	0.00%	0.00%	0.00%
<b>30</b>	78.54%	44.65%	18.78%	6.08%	1.56%	0.33%	0.06%	0.00%	0.00%	0.00%
<b>40</b>	87.15%	60.09%	32.33%	13.81%	4.80%	1.39%	0.34%	0.00%	0.00%	0.00%
<b>50</b>	92.31%	72.06%	45.95%	23.96%	10.36%	3.78%	1.18%	0.02%	0.00%	0.00%
<b>60</b>	95.39%	80.84%	58.26%	35.27%	18.03%	7.87%	2.97%	0.07%	0.00%	0.00%
<b>70</b>	97.24%	87.08%	68.63%	46.61%	27.21%	13.72%	6.04%	0.25%	0.00%	0.00%
<b>80</b>	98.35%	91.39%	76.94%	57.16%	37.11%	21.08%	10.53%	0.65%	0.00%	0.00%
<b>90</b>	99.01%	94.33%	83.36%	66.42%	47.03%	29.48%	16.39%	1.45%	0.00%	0.00%
<b>100</b>	99.41%	96.29%	88.17%	74.22%	56.40%	38.40%	23.40%	2.82%	0.01%	0.00%
<b>120</b>	99.79%	98.45%	94.25%	85.56%	72.18%	55.85%	39.37%	7.86%	0.10%	0.00%

5

6 Table A1 shows the probability of randomly flagging an individual when the cut

7 increases from 1 to 20 aberrant RTs, for test lengths of 20-120 items. For example, when

8 the test only contains 20 items, an individual with 6 aberrant RTs will have a low false

9 positive rate (0.03%, only 3 students will be randomly flagged in 10,000 students).

10 However, when the test has 70 items, a cut of 6 will have a high false positive rate

11 (13.72%, 1,372 students will be randomly flagged in 10,000 students). By this table, a cut

12 for flagging individuals can be determined after the total number of aberrant RTs for

13 each individual is computed.