

# Effects of Item-Selection Criteria on Classification Testing with the Sequential Probability Ratio Test

Chuan-Ju Lin

Judith Spray

**Effects of Item-Selection Criteria on Classification Testing  
with the Sequential Probability Ratio Test**

**Chuan-Ju Lin  
Judith Spray**

## ABSTRACT

This paper presents comparisons among three item-selection criteria for the sequential probability ratio test. The criteria were compared in terms of their efficiency in selecting items, as indicated by average test length (ATL) and the percentage of correct decisions (PCD). The item-selection criteria applied in this study were the Fisher information function, the Kullback-Leibler information function, and a weighted log-odds ratio. We also examined the effects of the cutoff scores, the width of the indifference region, the item pool size, and the item exposure rate under the different item-selection criteria. The results of the computer simulations showed that the three criteria yielded very small differences in the outcome measures, regardless of the conditions imposed.

## **EFFECTS OF ITEM-SELECTION CRITERIA ON CLASSIFICATION TESTING WITH THE SEQUENTIAL PROBABILITY RATIO TEST**

### **Introduction**

Computerized adaptive testing (CAT) is receiving more attention and has been applied more commonly over the last few years. Adaptive testing can yield more efficient tests by saving testing time (i.e., shorter tests) and increasing measurement precision. If the purpose of a test is to classify examinees into one of two or more mutually exclusive categories rather than estimating ability levels, the CAT procedure can be applied to make efficient decisions of classification by selecting and administering optimal items with algorithms based on statistical hypothesis testing, such as the sequential probability ratio test or SPRT (Spray & Reckase, 1994, 1996). The main purpose of this study was to compare three item-selection criteria in terms of average test length (ATL) and percentages of correct decisions (PCD) in the context of item selection with the SPRT. Variables hypothesized to affect ATL and PCD included the choice of the item-selection criteria, position of cutting points on the ability metric, the width of the indifference region, item pool size, and item exposure rate. Three types of selection criteria, three different cutting points, 11 indifference regions, two different item pool sizes, and three item exposure rates were examined.

### **The SPRT**

Wald's (1947) SPRT has been applied for classifying examinees into two mutually exclusive categories using a computerized adaptive test (Eggen, 1999; Spray & Reckase, 1996). In order to distinguish the computerized SPRT from conventional CAT, the SPRT is usually regarded as a computerized classification test or CCT (Spray, Abdel-fattah, Huang, & Lau,

the upper and lower bounds of the likelihood ratio test are defined as functions of  $\alpha$  and  $\beta$ . The actual observed error rates,  $\alpha^*$  and  $\beta^*$ , may be different from those predetermined, where usually  $\alpha^* \leq \alpha / (1 - \beta)$  and  $\beta^* \leq \beta / (1 - \alpha)$ . With the specified nominal error rates, the decision (or stopping) rules used can be defined as follows (Wald, 1947):

Continue selecting another item when:	$\beta / (1 - \alpha) < LR(\underline{x}) < (1 - \beta) / \alpha;$
Accept $H_1$ when:	$LR(\underline{x}) \leq \beta / (1 - \alpha);$
Accept $H_2$ when:	$LR(\underline{x}) \geq (1 - \beta) / \alpha.$

Any test administered using SPRT is adaptive in terms of test length. The items are administered, one by one, to an examinee until a classification decision is made, so that examinees with different ability levels obtain different average lengths of tests. Examinees with ability  $\theta_1 < \theta < \theta_2$  are expected to have longer tests than those with ability  $\theta \leq \theta_1$  or  $\theta \geq \theta_2$ , because it is more difficult to make decisions about those examinees with ability levels in the indifference region, especially those near the cutting score.

In practice, a minimum and maximum test length are usually specified. Even though a decision may not be achieved after the specified maximum number of items have been administered from the item pool, a forced classification can be made: reject  $H_1$  if  $LR(\underline{x})$  is greater than the midpoint of the interval  $[\beta / (1 - \alpha), (1 - \beta) / \alpha]$ ; otherwise accept  $H_1$ .

### **Item-selection criteria**

#### *(Fisher) Item Information*

In computer-based classification tests, the items in the item pool are usually ranked from maximum to minimum in terms of some item-selection criteria at the specified cutting point.

where  $K_i(\theta_2||\theta_1)$  denotes an item information index for item  $i$  for any two  $\theta$  values ( $\theta_2$  and  $\theta_1$ ), and  $E$  is the expected value operator, taken relative to  $\theta_2$ . The K-L test information function (i.e.,  $K(\theta_2||\theta_1)$ ) is the sum of the K-L information functions over all  $k$  items in the test, which equals

$$K(\theta_2||\theta_1) = \sum_{i=1}^k K_i(\theta_2||\theta_1). \quad (6)$$

The items with maximum K-L information are selected sequentially. The discrepancy between the likelihood function under the null and alternative hypotheses is a maximum when the K-L information is maximized. Therefore, testing is expected to be quite efficient because K-L information is, itself, a likelihood ratio; thus, the number of items needed to make decisions is expected to be minimized. With the dichotomously-scored IRT model, K-L item information can be computed as:

$$K_i(\theta_2||\theta_1) = p_i(\theta_2) \log \frac{p_i(\theta_2)}{p_i(\theta_1)} + q_i(\theta_2) \log \frac{q_i(\theta_2)}{q_i(\theta_1)}, \quad (7)$$

where  $p_i(\theta_2)$  and  $p_i(\theta_1)$  are the probabilities of a correct response to item  $i$  at  $\theta_2$  and  $\theta_1$ , respectively, and  $q_i(\theta_2)$  and  $q_i(\theta_1)$  are the complement probabilities.

#### *Weighted Log-odds Ratio*

An alternative measure on which to rank items for selection using the SPRT procedure is a weighted log-odds ratio criterion. This value is based on the following premise:

The likelihood ratio,  $LR(\underline{x})$ , is equal to 1.0 at the beginning of the testing session. The value,  $p_i(\theta_2)/p_i(\theta_1)$ , is multiplied to the likelihood ratio if the item is answered correctly or when  $x = 1$ . Likewise,  $LR(\underline{x})$  is multiplied by  $q_i(\theta_2)/q_i(\theta_1)$  when  $x = 0$ , or when the item is answered incorrectly. As testing continues,  $LR(\underline{x})$  is compared to the two boundaries,  $\beta / (1 - \alpha)$  and

The rationale for using this value to select items within the SPRT framework is that we are searching for items that will cause the SPRT likelihood ratio to cross the decision boundaries,  $(1-\beta)/\alpha$  and  $\beta/(1-\alpha)$ , or  $\log[(1-\beta)/\alpha]$  and  $\log[\beta/(1-\alpha)]$ , most quickly. Therefore, it makes sense to find the value of (10) for all items in the item pool. Thus, in theory, those items with greater weighted log-odds ratios should be selected earlier so that a decision will be made as soon as possible with the fewest number of items.

### **Item Exposure Control**

With computerized adaptive testing, the *best* items will be frequently selected, which is undesirable for test security reasons. Therefore, in order to protect the item pool, many item-exposure control strategies have been developed (e.g., Davey & Parshall, 1995; McBride & Martin, 1983; Sympson & Hetter, 1985). Item-exposure control is not only an important issue in CAT but also in CCT. Within the context of the current study, the *best* or optimal items refer to those with the best criterion values (e.g., highest Fisher information) at the cutting point. Without item-exposure control, the item-overlap rate between two CCT examinations would be very high because optimal items would be selected first in the test administration sequence and would eventually lead to overexposure.

A randomization scheme is a typical approach to controlling item exposure for CCT examinations (Spray et al., 1997; Way, Zara, & Leahy, 1996), especially in simulation studies. This approach for CCT is similar to the 5-4-3-2-1 randomization procedure used in CAT for ability estimation (McBride & Martin, 1983). The randomization methods indirectly control item exposure by randomly selecting an item from a group of a particular number (e.g.,  $m$ ) of

purpose of the current study was to investigate the efficiency of these two item-selection criteria more thoroughly by including several manipulations hypothesized to maximize possible differences in the criteria, as well as to include the weighted log-odds ratio criterion in the comparison.

## Method

### *Item Pools*

This study utilized two sizes of item pools – a *whole* pool and a *half* pool. The *whole item pool* used in this study was the ACT Assessment Mathematics Usage Test containing six equivalent (i.e., previously administered, intact) test forms. Each form was composed of 60 items, and thus, 360 items comprised the pool. Although two dimensions have been identified for each form based on previous multidimensional studies, the unidimensional SPRT procedure can be used with this item pool because it is robust to the violation of the unidimensionality assumption (Spray et al., 1997). The items were calibrated with the 3-PL IRT model.

In addition to using the whole pool, the item pool was split into two similar pools, each of which included three equivalent test forms and, thus, 180 items. One of these smaller pools was subsequently used for this study and was labeled as the *half pool*.

### *Item-selection Criteria*

Three item-selection criteria or functions were used for item selection:

1. Fisher information function.
2. Kullback-Leibler information function.
3. Weighted log-odds ratio.

### *Design*

In this study, the randomization scheme was used to control item-exposure rate, and different *stratum depths* were used. A stratum depth referred to the number of items grouped

compared on the outcome variables, average test length (ATL) and the percent of correct decisions made (PCD). Therefore, there were 99 possible conditions (i.e., 3 information criteria times 3  $\theta_0$  values times 11  $\delta$  values) under each of five combinations of pool size and exposure-control conditions listed previously.

### Results and Discussion

The ATL and PCD for three item-selection criteria with various indifference regions under five conditions are presented in Tables 1-5. It appeared that, for a particular cutting score with a given  $\delta$ , there were almost no differences in either ATL or PCD among the three item-selection criteria. This was especially surprising when  $\delta$  was largest around the cutting point,  $\theta_0$  (i.e., when  $\theta_1$  and  $\theta_2$  were farthest apart). See Table 1 vs. 2 vs. 3 and Table 4 vs. 5.

These tables also showed several expected results, namely that (1) as  $\theta_0$  moved farther away from the mean of the  $\theta$  distribution, the PCD increased; (2) ATL increased when  $\delta$  decreased and when item-exposure control increased (i.e., when a larger stratum depth was used); and (3) when a smaller pool was used, the ATL and PCD decreased. The latter finding resulted from more optimal items being administered more frequently under the half-pool condition (and, thus, the test lengths were shorter for all simulees). However, those simulees near the cutting point were missclassified at slightly higher rates **because** of the shortened test lengths. Thus, a decrement in classification accuracy occurred.

Further evidence of the similar behavior of the three item-selection criteria was exhibited by the rank correlations of the items at the cutting point. Table 6 provides the rank-order correlations among the three criteria and for three values of  $\delta$  (representing small, medium, and large indifference regions) at the three different cut-off scores. All of the correlation coefficients

## References

- Davey, T., & Parshall, C. G. (1995, April). *New algorithms for item selection and exposure control with computerized adaptive testing*. Paper presented at the annual meeting of the American Educational Research Association, San Francisco.
- Eggen, T. (1999). Item selection in adaptive testing with the Sequential Probability Ratio Test. *Applied Psychological Measurement*, 23(3), 249-261.
- Lord, F. M. (1980). *Applications of item response theory to practical problems*. Hillsdale, NJ: Lawrence Erlbaum Associates.
- McBride, J. R., & Martin, J. T. (1983). Reliability and validity of adaptive ability tests in a military setting. In D. J. Weiss (Ed.), *New horizons in testing: Latent trait test theory and computerized adaptive testing* (pp. 223-236). New York, Academic Press.
- Spray, J. A., Abdel-fattah, Abdel-fattah A., Huang, C. Y., & Lau, C. A. (1997). *Unidimensional approximations for a computerized test when the item pool and latent space are multidimensional*. (Research Report 97-5). Iowa City, IA: American College Testing.
- Spray, J. A. & Reckase, M. D. (1994, April). *The selection of test items for decision making with a computer adaptive test*. Paper presented at the Annual Meeting of the American Educational Research Association, New Orleans.
- Spray, J. A. & Reckase, M. D. (1996). Comparison of SPRT and sequential Bayes procedures for classifying examinees into two categories using a computerized test. *Journal of Educational and Behavioral Statistics*, 21(4), 405-414.
- Sympson, J. B., & Hetter, R. D. (1985, October). *Controlling item-exposure rates in computerized adaptive testing*. Paper presented at the 17<sup>th</sup> annual meeting of the Military Testing Association. San Diego, CA: Navy Personnel Research and Development Center.
- Wald, A. (1947). *Sequential analysis*. New York: Wiley.
- Way, W. D., Zara, A. R., & Leahy, J. (1996, April). *Modifying the NCLEX CAT item selection algorithm to improve item exposure*. Paper presented at the Annual Meeting of the American Educational Research Association, New York.

TABLE 1

Average Test Length and Percentage of Correct Decisions for All Possible Item Selection Procedures with Whole Pool and Stratum Depth = 1.

$\delta$	Procedure	$\theta_0 = -.32$		$\theta_0 = .81$		$\theta_0 = 1.79$	
		ATL	PCD	ATL	PCD	ATL	PCD
.20	Fisher	64.66	.945	16.64	.969	6.28	.991
	LR	65.74	.946	16.70	.969	5.86	.991
	K-L	64.72	.946	16.48	.968	6.37	.991
.21	Fisher	60.49	.945	15.38	.967	5.99	.991
	LR	60.42	.945	15.04	.968	5.48	.991
	K-L	61.39	.944	15.40	.967	6.02	.991
.22	Fisher	56.47	.945	14.10	.967	5.69	.991
	LR	56.59	.944	13.98	.967	5.23	.991
	K-L	56.48	.943	13.98	.966	5.69	.991
.23	Fisher	52.49	.944	13.08	.965	4.54	.990
	LR	53.14	.943	12.48	.965	4.81	.990
	K-L	52.61	.943	12.96	.965	4.49	.991
.24	Fisher	49.48	.942	11.83	.963	4.26	.991
	LR	49.93	.941	11.72	.964	4.52	.990
	K-L	49.73	.941	12.06	.963	4.30	.990
.25	Fisher	46.21	.941	10.64	.962	3.85	.990
	LR	45.99	.941	11.00	.963	4.36	.990
	K-L	46.43	.940	11.18	.963	3.89	.990
.26	Fisher	43.72	.939	10.10	.962	3.70	.989
	LR	43.26	.940	10.41	.961	4.09	.990
	K-L	44.01	.939	10.49	.962	3.71	.989
.27	Fisher	40.14	.938	9.73	.959	3.49	.989
	LR	40.43	.937	9.59	.960	3.92	.990
	K-L	40.23	.936	9.26	.961	3.51	.988
.28	Fisher	37.67	.936	8.96	.959	3.29	.989
	LR	37.60	.936	8.88	.959	3.03	.989
	K-L	37.52	.937	8.86	.958	3.28	.989
.29	Fisher	35.34	.935	8.42	.958	3.14	.988
	LR	35.23	.935	8.58	.958	2.87	.989
	K-L	35.36	.936	8.55	.958	3.19	.988
.30	Fisher	33.15	.935	8.03	.956	2.99	.988
	LR	33.28	.934	8.05	.955	2.76	.988
	K-L	33.47	.933	8.18	.957	3.05	.988

Note: Fisher: Fisher Information  
 LR: Weighted Log-Odds Ratio  
 K-L: Kullback-Leibler Information  
 ATL: Average Test Length  
 PCD: Percentage of Correct Decisions

TABLE 3

Average Test Length and Percentage of Correct Decisions for All Possible Item Selection Procedures with Whole Pool and Stratum Depth = 10.

$\delta$		$\theta_0 = -.32$		$\theta_0 = .81$		$\theta_0 = 1.79$	
		ATL	PCD	ATL	PCD	ATL	PCD
$\delta = .20$	Fisher	108.27	.945	34.67	.968	13.33	.991
	LR	108.47	.946	34.57	.968	12.43	.991
	K-L	108.28	.946	34.93	.968	13.50	.991
$\delta = .21$	Fisher	103.50	.945	32.34	.967	12.48	.990
	LR	103.08	.944	32.04	.966	11.65	.991
	K-L	102.64	.944	32.28	.967	12.73	.991
$\delta = .22$	Fisher	98.10	.943	30.01	.966	11.75	.991
	LR	97.79	.943	30.17	.965	10.90	.990
	K-L	97.39	.943	30.36	.967	11.67	.990
$\delta = .23$	Fisher	93.06	.942	28.37	.964	11.02	.991
	LR	92.72	.943	28.08	.965	10.33	.990
	K-L	93.67	.941	28.28	.965	11.21	.990
$\delta = .24$	Fisher	88.88	.941	26.45	.962	10.37	.989
	LR	88.95	.941	26.30	.964	9.60	.990
	K-L	88.63	.943	26.38	.963	10.52	.990
$\delta = .25$	Fisher	84.15	.941	24.79	.961	9.87	.989
	LR	84.82	.941	24.52	.962	9.21	.990
	K-L	84.78	.939	24.88	.962	9.81	.990
$\delta = .26$	Fisher	80.86	.940	23.36	.960	9.12	.989
	LR	80.37	.938	23.17	.960	8.52	.990
	K-L	80.59	.939	23.32	.961	9.31	.989
$\delta = .27$	Fisher	76.62	.937	21.77	.958	8.65	.988
	LR	76.76	.938	21.72	.958	8.10	.989
	K-L	76.49	.938	21.92	.959	8.80	.989
$\delta = .28$	Fisher	72.83	.937	20.69	.957	8.27	.988
	LR	73.23	.937	20.51	.957	7.58	.989
	K-L	72.93	.937	20.75	.958	8.21	.989
$\delta = .29$	Fisher	69.98	.935	19.35	.958	7.79	.988
	LR	69.93	.935	19.34	.956	7.08	.988
	K-L	69.85	.936	19.56	.957	7.96	.988
$\delta = .30$	Fisher	66.81	.933	18.49	.955	7.45	.988
	LR	66.69	.934	18.37	.955	6.90	.988
	K-L	66.11	.933	18.72	.954	7.66	.987

Note: Fisher: Fisher Information  
 LR: Weighted Log-Odds Ratio  
 K-L: Kullback-Leibler Information  
 ATL: Average Test Length  
 PCD: Percentage of Correct Decisions

TABLE 5

Average Test Length and Percentage of Correct Decisions for All Possible Item Selection Procedures with Half Pool and Stratum Depth = 5.

$\delta$	Method	$\theta_0 = -.32$		$\theta_0 = .81$		$\theta_0 = 1.79$	
		ATL	PCD	ATL	PCD	ATL	PCD
.20	Fisher	86.57	.926	34.86	.963	14.71	.990
	LR	86.88	.926	34.31	.963	13.78	.989
	K-L	86.51	.928	34.87	.964	14.78	.990
.21	Fisher	83.00	.927	33.14	.963	14.23	.989
	LR	83.05	.927	32.49	.962	12.88	.989
	K-L	83.30	.928	32.92	.963	14.24	.990
.22	Fisher	79.89	.927	31.06	.963	13.44	.989
	LR	79.94	.927	30.91	.961	12.04	.990
	K-L	80.00	.927	31.30	.961	13.47	.989
.23	Fisher	76.89	.928	29.50	.962	12.60	.990
	LR	76.98	.927	29.09	.962	11.60	.989
	K-L	76.60	.927	29.43	.961	12.58	.990
.24	Fisher	73.57	.927	27.95	.960	11.94	.989
	LR	73.93	.926	27.35	.960	11.02	.989
	K-L	73.62	.926	27.83	.960	12.21	.988
.25	Fisher	71.03	.926	26.42	.959	11.42	.988
	LR	70.98	.927	26.03	.959	10.55	.988
	K-L	71.20	.925	26.56	.959	11.49	.988
.26	Fisher	68.27	.925	24.93	.958	10.88	.988
	LR	68.08	.925	24.92	.958	9.90	.989
	K-L	68.26	.926	25.10	.958	10.95	.989
.27	Fisher	65.87	.924	23.57	.957	10.34	.988
	LR	65.59	.925	23.48	.957	9.48	.988
	K-L	65.56	.924	23.77	.958	10.41	.989
.28	Fisher	63.29	.924	22.40	.956	9.86	.988
	LR	63.55	.924	21.94	.956	9.04	.988
	K-L	63.61	.923	22.58	.957	9.94	.988
.29	Fisher	60.98	.924	21.17	.956	9.55	.987
	LR	60.86	.923	21.24	.955	8.54	.988
	K-L	60.84	.923	21.39	.955	9.54	.987
.30	Fisher	58.46	.922	20.24	.953	8.97	.988
	LR	58.30	.923	20.15	.954	8.08	.987
	K-L	58.79	.922	20.40	.955	9.23	.987

Note: Fisher: Fisher Information  
 LR: Weighted Log-Odds Ratio  
 K-L: Kullback-Leibler Information  
 ATL: Average Test Length  
 PCD: Percentage of Correct Decisions