5000 System Panels' MCQs – what have we learnt?

> Curriculum Retreat 2009 Faculty of Medicine, CUHK

Background and goals of the present study

Background:

There appears to be more failures in PHOM panel-end MCQ tests when compared with those of other system panels.

Goals:

- Whether there are differences in the psychometric properties of MCQ used in SA of different system panels?
- What is/are the indicator(s) of a MCQ test (with a pass mark of 50) that would predict the presence of failures?

Scope of the analyses

- System panel summative assessments (Year 1 to 3) since the implementation of the new medical curriculum (2001 to 2009) were included in the analysis. Data from skill panel SA (COSK, LLSK. PCLM) and PHES were not included.
 - There were 154 SA involving 6475 items.

Year	Panels	No. of SA	No. of items
1 <i>(2001-2009)</i>	PCAR1, PFOS1, PGIN1, PHOM1, PHUS1, PMUS1	50	2377
2 (2002-2009)	PCAR2, PGIN2, PHAE2, PHOM2, PHUS2, PMUS2, PMDT2, PNEU2	56	2426
3 <i>(2003-2009)</i>	PCAR3, PGIN3, PHAE3, PHOM3, PMUS3, PMDT3, PNEU2, PREP3	48	1672

Psychometric properties of MCQ (i.e. difficulty level, discriminating power), number of options picked and mean-2SD of students' score in each test.

Item analysis report provided by OES

		HISTOGRAM 1						
****TEST STATISTICS**** NUMBER OF EXAMINEES 137 TOTAL SCORE: MEAN 21.04 VARIANCE 10.05 STANDARD DEVIATION 3.18 KR-20 RELIABILITY 0.5247 S.E. OF MEASUREMENT 2.1895	FREQUENCY 22 21 20 19 18 17 16 15 14 13 12 11 10 9 8 7 6 5 4	HISTOGRAM 1 1 0 0 1 0 4 1 6 4 6 12 20 16 22 15 12 8 6 2 1 Bottom 27% 0.46 $\cdot \cdot $						
	3 2 1 ·							
	CLASS INTERVAL	9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28						
ITEM2: $DIF=0.620$, RPB= 0.380 , CRPB= 0.228 (95% CON= 0.062 , 0.381)DP is defined as the difference between the proportions of the HIGHGROUPNINVNFOMITABC*DETOTAL137000.070.170.620.020.12HIGH4400.020.140.820.000.02proportions of the HIGHMID5800.090.140.570.030.17(top 27%) and LOWLOW3500.110.260.460.030.14(top 27%) and LOWTESTSCOREMEAN %:6266736366(bottom 27%) groupsDISCRIMINATING POWER-0.09-0.120.36-0.03-0.12giving this response. $0.82 - 0.46 = 0.36$								
ITEM 4: DIF=1.000, RPB= 0.0 RBIS= 0.0 GROUP N INV NF OMIT TOTAL 137 0 0 0 HIGH 44 0 MID 58 0 LOW 35 0 TEST SCORE MEAN %: DISCRIMINATING POWER STANDARD ERROR OF D.P.	00, CRPB= 0 00, CRBIS= 0 C* 1.00 1.00 1.00 81 0.00 0.00	$\begin{array}{c} 0.000 & (95\% \ \text{CON} = \ -0.174, \ 0.174) \\ 0.000, \ \text{IRI} = 0.000 \\ \end{array}$ $\begin{array}{c} \text{ITEM 16: } \text{DIF} = 0.934, \ \text{RPB} = \ 0.356, \ \text{CRPB} = \ 0.278 \ (95\% \ \text{CON} = \ 0.116, \ 0.426) \\ \text{RBIS} = \ 0.691, \ \text{CRBIS} = \ 0.540, \ \text{IRI} = 0.088 \\ \end{array}$ $\begin{array}{c} \text{GROUP N INV NF OMIT C D E*} \\ \text{TOTAL 137 0 0 0 0.02 0.04 0.93} \\ \text{HIGH 44 0 0 0.00 0.00 1.00} \\ \text{MID 58 0 0.02 0.03 0.95} \\ \text{LOW 35 0 0.06 0.11 0.83} \\ \text{TEST SCORE MEAN $$: 60 54 71 \\ \text{DISCRIMINATING POWER -0.06 -0.11 0.17} \\ \text{STANDARD ERROR OF D.P. 0.04 0.05 0.06} \end{array}$						

Medical Year 1

- Mean score is given by the average of the difficulty levels of MCQs used in the test.
- There are significant differences in the mean scores of SA (i.e. difficulty levels of items) offered by different system panels.
- Differences among panels become more evident when mean-2SD are compared, and more of these values are close to the 50% mark.
 - Those SA with higher mean scores are made up of MCQs with fewer numbers of options picked by students.

Open circle indicates the mean score of individual SA, and the bar and error bar show their mean \pm SD (n=8).



Percentage distribution of Year 1 Summative MCQ based on the number of options picked by students



Number of usable options is related to difficulty level and discriminating power



Data represents mean<u>+</u>SD with the number of items shown within brackets. Differences between groups are statistically significant (P<0.001), by Kruskal-Wallis one-way ANOVA on ranks.

Average number of usable options in MCQ is inversely related to the mean score of the test



What determines SD of the mean score (i.e. spread of marks in a MCQ test)?

SD of the mean score shows a strong positive correlation to the discriminating power of MCQ.



How is the discriminating power of an item related to its difficulty level?



For items that few students can get the "right" answer, are they really difficult?

Difficult items could have their stems/options poorly written, and therefore they are confusing to students. Or it is possible that a wrong key was used.

ITEM	21:	DIF=0	.180,	RPB= RBIS=	0.035,	CRPB= CRBIS	-0.010 =-0.015,	(95% CON IRI=0.0	= -0.176,	0.157)
	GROUP	N	INV	NF	OMIT	A	В	C*	D	E
	TOTAL	139	0	0	0	0.39	0.16	0.18	0.23	0.04
	HIGH	36	0			0.53	0.08	0.14	0.22	0.03
	MID	66	0			0.35	0.20	0.20	0.24	0.02
	LOW	37	0			0.32	0.16	0.19	0.22	0.11
	TEST S	SCORE	MEAN	8:		71	66	70	69	61
	DISCRI	MINAT	ING PO	OWER		0.20	-0.08	-0.05	0.01	-0.08
	STANDA	ARD ER	ROR O	F D.P.		0.12	0.08	0.09	0.10	0.06

ITEM	9:	DIF=0	.000,	RPB= RBIS:	0.000,	CRPB= CRBIS=	0.000	(95% CON= -0.171, IRI=0.000	0.171)
G	RÓUP	N	INV	NF	OMIT	D			
г	OTAL	132	0	0	0	1.00			
	HIGH	27	0			1.00			
	MID	67	0			1.00			
	LOW	38	0			1.00			
Т	EST S	SCORE	MEAN	8:		78			
D	ISCR	IMINAT	ING P	OWER		0.00			
S	TAND	ARD EF	RROR O	F D.P	•	0.00			

How to determine pass/fail?

- Professor C. B. Hazlett has addressed the issue of standard setting in Curriculum Retreats held in 2004 and 2007.
- What we are using for our SA is an "outdated" approach of setting an absolute and fixed pass/fail point: "a particular score or a % that has been determined prior to administering the test is set as the pass mark (e.g. 50%)".
- Another approach would be norm-referencing of setting a relative and not fixed pass/fail point: "students are compared with each other and those who fail are "X" SDs below the mean performance of all candidates".
- More "modern" practice would be to use test-centered or examinee-centered approach (or a combination of both) along with the judgment of "subject matter experts" factored into the method of choice.

What have I learnt?



More likely to identify outliers who might score below 50, therefore requiring to take the supplementary exam.

How is the number of failures related to mean-2SD score?



The mean-2SD score approaching or falling below 50 represents the indicator that predicts the presence of failures in a MCQ test.

What we could learn?

Keep track of the psychometric properties of each item, and the number of options picked by students, when being used over the years in different SA.

Year used	Difficulty Level	Discrimination Index					Selection percentage				
		Α	В	С	D	E	Α	В	С	D	E
03-04	0.682	-0.11	0.46	-0.09	-0.10	-0.17	0.08	0.68	0.06	0.06	0.11
05-06	0.655	-0.13	0.52	-0.18	-0.09	-0.12	0.10	0.65	0.14	0.06	0.05
07-08	0.766	0.02	0.45	-0.15	-0.18	-0.15	0.02	0.77	0.09	0.08	0.05
08-09	0.611	-0.15	0.50	-0.23	-0.07	-0.06	0.13	0.61	0.15	0.06	0.06

- Weed out those items that all students can pick the right answer (Diff level = 1).
- Improve on the writing of MCQ to eliminate confusion and to improve on the quality of the distracters (i.e. maximize the number of options that may be picked by students).
- Keep those items that test on important concepts/knowledge that students should know, even though they are not too difficult.
- Avoid over-testing students with difficult items that are beyond the level required or cover unimportant concepts/knowledge.

Results of PHOM Panel SA in 2008-09



Trends in the number of failures in Panel SA



Over the years, there appears to be a downward trend in the number of panel SA having failures, and the number of students failing in SA could also be on the decline.

Acknowledgement

- Thanks to Prof CB Hazlett and Dr SPY Yip for allowing me to access the item analysis data of Panel SA.
- Thanks to Prof CB Hazlett for his comments.
- Thanks to Prof SM Kumta for allowing me to present the study here.
- Thanks to Miss Diana Kwan for providing well organized records of the item analysis.