



THE CLERICAL TEST BATTERY

the
technical
manual

7

Test Batteries

P	=	α
S		β
	Y)
	T	σ
+		E
√	(C
	Σ	H

C

ONTENTS

1	THEORETICAL OVERVIEW.....	5
2	THE CONSTRUCTION OF THE CLERICAL TEST BATTERY.....	9
3	PSYCHOMETRIC PROPERTIES OF THE CLERICAL TEST BATTERY.....	15
4	REFERENCES.....	25

②

LIST OF TABLES

- 1** CTB2 STANDARDISATION SAMPLES
- 2** GENDER DIFFERENCES ON CTB2
- 3** RELATIONSHIP WITH CTB2 SUB-SCALES AND AGE
- 4** COEFFICIENT ALPHA AS A MEASURE OF CTB2 INTERNAL CONSISTENCY
- 5** INTER-CORRELATIONS CTB2 SUB-SCALES
- 6** CORRELATIONS CTB2 WITH OTHER MEASURES.
- 7** CORRELATIONS CTB2 WITH PERFORMANCE FACTORS

1

THE ROLE OF PSYCHOMETRIC TESTS IN PERSONNEL SELECTION AND ASSESSMENT

THE ROLE OF PSYCHOMETRIC TESTS IN PERSONNEL SELECTION AND ASSESSMENT

One of the main reasons for using reasoning tests to aid selection decisions is that they can provide information which cannot be obtained in other ways. If such tests are not used then our knowledge of the applicant is limited to the information that can be gleaned from an application form or CV, an interview and references. If we wish to gain information about a person's specific aptitudes and abilities then we have little option but to use psychometric tests. But such tests can do more than simply provide additional information about the applicant. They can add a degree of reliability and validity to the selection procedure that is impossible to achieve in any other way. How they do this is best addressed by examining the limitations of the information obtained through interviews, application forms and references and exploring how some of these limitations can be overcome by the use of reasoning tests.

While much useful information can be gained from the interview, which clearly has an important role in any selection procedure, it does nonetheless suffer from a variety of weaknesses. Perhaps the most important of these is that it is not a reliable way to judge a person's level of reasoning ability. While the interview enables us to probe each applicant in depth and discover individual strengths and weaknesses it will not enable us to objectively

assess an applicant's aptitudes and abilities.

There are similar limitations on the range and usefulness of the information that can be gained from application forms or CV's. While work experience and qualifications may be prerequisites for certain occupations, in and of themselves they do not determine whether a person is likely to perform well or badly. Past experience and academic achievement is not always a good predictor of ability or future success. While such information is important it may not be sufficient on its own to enable us to confidently choose between applicants. Thus aptitude and ability tests are likely to play a significant role in the selection process as they provide information on a person's potential and not just their achievements to date.

While past performance can be a good predictor of future performance references may not be good predictors of past performance. If the applicant supplies the referee then it is likely that they have chosen someone whom they expect to speak highly of them and have probably avoided supplying the names of those who may have a less positive view of their abilities. Aptitude and ability tests, on the other hand, give an indication of the applicant's likely performance, obtained under exam conditions and are thus likely to be an objective, true reflection of the person's ability.

So what advantages does the use of reasoning tests have over these other forms of assessment? The first advantage they have is that they are standardised. That is to say the same test is given to all the applicants under the same conditions and a standard method is used for scoring and interpreting the test results. Thus the test should produce the same results no matter who administers and interprets it. Moreover, the test results can be represented numerically making it easy both to compare applicants with each other, and with pre-defined groups (e.g. successful and unsuccessful job incumbents). In addition, as noted above, they provide a range of information which is not easily and reliably assessed in other ways. Such information can fill in important gaps which have not been assessed by application forms, interviews and references and can also raise questions which can later be directly addressed in the interview. It is for this reason that psychometric tests are being increasingly used in personnel selection. Their use adds a degree of objectivity, reliability and breadth to assessment decisions which can not be achieved any other way.

2

THE CONSTRUCTION OF THE CLERICAL TEST BATTERY

- 1** THE COMPONENTS OF THE CLERICAL BATTERY

The growing awareness during the 1970's (cf. Schmidt and Hunter, 1977) that psychometric tests which were valid for one job could also be shown to be valid for other similar jobs led to interest in the development of job specific test batteries. From this was developed the idea of job families; groups of jobs which share similar skill requirements, where performance on each job could be predicted from the same test battery. Until the pioneering work of Schmidt and colleagues it had been thought necessary to demonstrate the validity of a test for each specific job, as validity was thought to vary substantially from job to job, due to job specific factors (Ghiselli, 1966). The demonstration that many test batteries were valid for a whole job family encouraged the development of test batteries for specific occupational areas. The CTB2 is one such test battery, which has been developed to assess the aptitudes which are predictive of performance in a wide range of clerical jobs. Thus the CTB2 assess Verbal Reasoning, Numerical Ability, Perceptual Speed & Accuracy (clerical checking) and Spelling Ability.

Research has clearly demonstrated that in order to accurately assess aptitudes and reasoning ability it is necessary to develop tests which have been specifically

designed to measure that ability in the population under consideration. That is to say, we need to be sure that the test has been developed for use on the particular group being tested, and thus that it is appropriate for that group. There are a number of ways in which this is important. Firstly, it is important that the test has been developed in the country in which it is intended to be used. This ensures that the items in the test are drawn from a common, shared cultural experience, giving each candidate an equal opportunity to understand the logic which underlies each item. Secondly, it is important that the test is designed for the particular ability range on which it is to be used. A test designed for those of average ability will not accurately distinguish between people of high ability as all their scores will cluster towards the top end of the scale. Similarly, a test designed for people of high ability will be of little use if given to people of average ability. Not only will it not discriminate between applicants, as all the scores will cluster towards the bottom of the scale, but also as the questions will be too difficult for most of the applicants they are likely to be demotivated, producing artificially low scores. Consequently the CTB2 was developed specifically for clerical staff of general, as opposed to gradu-

ate level ability, who have received a basic level of education.

In constructing the items in the CTB2 a number of guidelines were borne in mind. Firstly, and perhaps most importantly, each item was constructed so that only a general educational level was needed in order to be able to correctly solve that item. Thus we tried to ensure that each item was a measure of aptitude, rather than being a measure of specific knowledge or experience. For example, in the case of the numerical items the calculations involved in solving each item are relatively simple, with the difficulty of the item being due to the logic which underlies that question rather than being due to the need to use complex mathematical operations to solve that item. Secondly, a number of different item types (e.g. odd one out, word meanings etc.) were used to measure each aspect of aptitude. This was done in order to ensure that each subscale measures a broad aspect of aptitude (e.g. Verbal Reasoning Ability), rather than measuring a very specific aptitude (e.g. vocabulary). In addition, the use of different item types ensures that the test is measuring different components of each aptitude. For example the ability to understand analogies, inclusion/exclusion criteria for class membership etc.

THE COMPONENTS OF THE CLERICAL BATTERY

NUMERICAL REASONING (NA2)

This is a measure of numerical ability which has been isolated as a component of general reasoning. Additionally, routine computational ability is clearly relevant as a measure of capacity to perform certain types of numerically based clerical work.

VERBAL REASONING (VA2)

This is a standard verbal reasoning test aimed at the general population level of reasoning. While it assesses reasoning ability it also gives a guide to the relative level of literacy of a candidate. Thus, it will be of relevance in determining if the requisite verbal skill is present.

FILING (FL2)

This test is composed of a verbal and numerical filing operation. It assesses the ability to file quickly and accurately. Either an alphabetical or a numerical test of filing can be used depending on the job in question. The Filing test is only available for on-screen administration using the GeneSys software.

SPELLING (SP2)

This test gives the respondent a choice of a number of different spellings for a common English word. The ability to spell correctly is still an important asset in a number of clerical occupations.

CLERICAL CHECKING (CC2)

The Clerical Checking test assesses the ability to check words and numbers for accuracy. This is an important clerical skill representing as it does the ability to quickly and accurately check information for accuracy.

TYPING TEST (TY2)

This is a GeneSys-based on-screen test of typing speed and accuracy. There are several methods of assessing a candidate's typing ability other than pure speed. For example one can require that errors be corrected while copy typing or that a table be reproduced. Not all clerical positions require typing ability and for this reason the typing test is available without the Clerical Battery and can stand alone with its own report being generated.

3

PSYCHOMETRIC PROPERTIES OF THE CLERICAL TEST BATTERY

- 1** STANDARDISATION
- 2** RELIABILITY
- 3** THE RELIABILITY OF THE CTB2
- 4** VALIDITY
- 5** THE VALIDITY OF THE CTB2

STANDARDISATION

For each of the sub-tests of the Clerical Test Battery, information is provided on age and gender differences where they apply. All normative data is available from within the GeneSys system which computes for any given raw score, the appropriate standardised scores for the selected reference group. In addition the GeneSys™ software allows users to establish their own in-house norms to allow more focused comparison with profiles of specific groups.

The total norm base of the CTB2 is based on a three different groups. These include clerical trainees, Customer Service Staff and Clerical Staff as detailed in **Table 1**.

CTB2 GENDER AND AGE DIFFERENCES

Gender differences on CTB2 were examined by comparing results of males and female trainees matched as far as possible for educational and socio-economic status. **Table 2** provides mean scores for males and females on each of the CTB2 scales as well as the t-value for mean score differences.

As a test of gender differences, scores differences between males and females from a sample of administrative staff (N=125) were compared.

With the exception of the Numerical test (NA2), the CTB2 sub-scales showed no significant gender differences.

The effect of age on CTB2 scores was examined using the sample of respondents in **Table 3**. Only the Spelling test just reaches statistical significance, suggesting that the CTB2 is not strongly subject to age effects.

	Total N	Males N	Females N	Mean Age	Age Range	SD Age
Clerical Trainees	101	23	78	18	16-22	2.3
Customer Service Staff	115	28	87	26	19-50	7.61
Clerical Staff	54	14	40	25	18-45	7.19

Table 1: GRT2 Standardisation Samples

	Mean F	Mean M	t-value	p	Std.Dev.F	Std.Dev.M	F-ratio	p
VER	22.25	24.64	-1.29	.203	5.44	7.36	1.84	.142
NA2	9.475	14.21	-3.72	.000	3.99	4.44	1.24	.579
CC2 Verbal	11.77	11.64	.164	.871	2.29	3.37	2.15	.065
CC2 Numerical	12.03	11.93	.107	.916	2.67	3.54	1.75	.176
CC2 Total	23.80	23.57	.137	.892	4.80	6.86	2.04	.087
SP2	17.68	16.57	.927	.358	3.85	3.78	1.04	.992

Table 2: Gender Differences on CTB2

CTB2	AGE
Verbal	-.01
Numerical	-.04
Checking	.02
Spelling	.24

Table 3: Relationship with CTB2 sub-scales and Age

RELIABILITY

If an aptitude test is to be used for selection and assessment purposes the test needs to measure each of the aptitude or ability dimensions it is attempting to measure reliably, for the given population (e.g. clerk/typists, personal assistants etc.). That is to say, the test needs to be consistently measuring each aptitude so that if the test were to be used repeatedly on the same candidate it would produce similar results.

TEST-RETEST RELIABILITY

Test-retest reliability statistics estimate the reliability of an aptitude test by administering it repeatedly to the same applicants. If a test is reliable then we would expect it to produce consistent results when repeatedly administered over short periods of time. Thus we would not expect a reliable test to classify someone as having a high level of aptitude on one occasion and a low level of aptitude on another. Thus repeated test administration can provide an estimate of a test's reliability.

INTERNAL CONSISTENCY RELIABILITY

Internal consistency statistics estimate the reliability of a test by exploring whether each of the items which measure one ability or aptitude combine to produce a consistent scale. That is to say, we would expect people of superior aptitude to do well on all the items which measure that aptitude, and not simply on a subset of these items. If the latter were the case then we might suspect that those items which they did not perform well on were in fact not good measures of the underlying reasoning ability. These statistics are the most commonly used ways to estimate a test's reliability.

It is generally recognised that aptitude tests are more reliable than personality tests and for this reason high standards of reliability are usually expected from such tests. While many personality tests are considered to have acceptable levels of reliability if they have reliability coefficients in excess of .7, aptitude tests are not usually considered to have acceptable levels of reliability unless they have reliability coefficients in excess of .8.

THE RELIABILITY OF THE CTB2

Table 4 presents alpha coefficients for four of the sub-tests which form the CTB2. The data was collected on a sample of 126 clerical staff attending an in-service training course. Inspection of the above table reveals that each of these four sub-tests have high levels of reliability.

CTB2 test	Chronbach's Alpha
VA2	0.83
NA2	0.83
CC2	0.90
SP2	0.81

Table 4 Coefficient Alpha as a measure of CTB2 internal consistency

VALIDITY

Whereas reliability assesses the degree of measurement error of an aptitude test, that is to say the extent to which the test is consistently measuring one underlying ability or aptitude, validity addresses the question of whether or not the scale is measuring the characteristic it was developed to measure. This is clearly of key importance when using an aptitude test for assessment and selection purposes. In order for the test to be a useful aid to selection we need to know that the results are reliable and that the test is measuring the aptitude it is supposed to be measuring. Thus, after we have examined a test's reliability we need to address the issue of validity. We traditionally examine the reliability of a test before we explore its validity as reliability sets the lower bound of a scale's validity. That is to say a test cannot be more valid than it is reliable.

There are two main ways in which we can say that a test is valid. We call these Construct Validity and Predictive Validity. When tests are used for individual assessment Construct Validity is the more important and when tests are used to predict performance the test's Predictive Validity is the more important.

Construct Validity

Construct Validity assesses whether the characteristic which the test is actually measuring is psychologically meaningful and is consistent with the scale's definition.

Criterion-related Validity

This assesses whether the test is capable of predicting some real-world criterion; for example job performance. Thus while a test may have criterion-related validity it may not have construct validity. That is, it may predict a given criterion but may not be measuring the psychological construct it purports to measure.

ASSESSING CONSTRUCT VALIDITY

Unlike reliability which can be easily measured, Construct Validity is a much more difficult characteristic to assess. Rather than there being one generally agreed way to assess a test's construct validity the validity of a test is usually established by presenting a variety of evidence which converges to demonstrate the test's validity. For example, we will want to know that the aptitudes or abilities which the test measures are stable over time and have intuitive, consensual meaning. Moreover we will want to show that a variety of statistical properties hold for the test. These concern the test's:

Internal Structure

Specifically we are concerned that the test's sub-scales are correlated with each other in a meaningful way. For example, we would expect the different sub-scales of an aptitude test battery to be moderately correlated as each will be measuring a different facet of the general aptitude. Thus if such sub-scales are not correlated with each other we might wonder whether each is a good measure of the aptitude. Moreover, we would expect different facets of verbal reasoning ability (e.g. vocabulary, similarities etc.) to be more highly correlated with each other than they are with a measure of numerical reasoning ability. Consequently, the first way in which we might assess the validity of an aptitude test battery is by exploring the relationship between the test's sub-scales.

Concurrent Validity

Here we are concerned to demonstrate that the test produces results which are consistent with those produced by other widely used, recognised, validated tests. To explore the concurrent validity of a test we might correlate the candidates' scores on the test which is being validated with their scores on a test which is already known to be valid.

Criterion Validity

Here we are concerned to demonstrate that the test discriminates between criterion groups which we would predict to obtain different scores on the test's sub-scales. For example, we might validate a test measuring verbal and numerical ability by showing that graduates perform better on the test than non-graduates, and that science students perform better on the numerical ability test than arts students.

THE VALIDITY OF THE CTB2

As a demonstration of validity, a number of studies have been undertaken. These include inter-correlating the CTB2 sub-scales and correlating CTB2 sub-scales with other recognised measures of the same or similar characteristics and performance ratings in a work context.

THE STRUCTURE OF THE CTB2

The inter-correlations of CTB2 sub-scales are presented in **Table 5** below. These demonstrate that the sub-scales are fairly independent of each other. The highest inter-correlation of 0.55 between the verbal and numerical reasoning tests, is in line with expectations. The correlation between Spelling scores and Checking is 0.48 appears to indicate that spelling ability may be an important contributor to checking ability.

CONSTRUCT VALIDITY

As a measure of construct validity, sub-scales of the CTB2 were correlated with other measures purporting to assess the same or similar characteristics.

The verbal and numerical sub-scales were correlated with the AH series counter-parts, whereas the Checking tests was correlated with the IPI series Perception measure.

Table 6 show good to excellent correlations with similar measures.

CRITERION-RELATED VALIDITY

A sample of 125 customer service operators with a High Street Bank completed the CTB2 Verbal, Numerical and Clerical Checking tests and were concurrently rated on 14 performance indicators (competencies). These were factored and four performance factors were extracted and subsequently correlated with the CTB2 sub-scales.

The results in **Table 7** show that the CTB2 sub-scales used were more effective in predicting the skill areas of Software and Numerical Competency areas than the first three performance indicators.

Subscale	Verbal	Numerical	Checking	Spelling
VA2	1.00	.51	.37	.34
NA2	.51	1.00	.34	.31
CC2	.37	.34	1.00	.48
SP2	.34	.31	.48	1.00

Table 5 Inter-correlations CTB2 sub-scales

VER	.63	AH2 Verbal
NA2	.76	AH2 Numerical
CC2	.68	IPI Perception

Table 6 Correlations CTB2 with other measures.

CTB2 subscale	Factor I General	Factor II Interpersonal skills	Factor III Planning & Organising	Factor IV Software & Numerical
VA2	-.06	-.09	.08	.14*
NA2	-.16	-.08	.03	.31*
CC2 P1	.03	.13	-.05	.28*
CC2 P2	.01	.11	-.03	.34*
CC2 Total	.02	.13	-.04	.32*

Table 7 Correlations CTB2 with performance factors

*p<0.05

4

REFERENCES

Ghiselli E. E. (1966) *The Validity of Occupational Aptitude Tests*, New York: Wiley

Gould, S.J. (1981). *The Mismeasure of Man*. Harmondsworth, Middlesex: Pelican.

Heim, A.H. (1968). *AH5 Group Test of High-grade Intelligence*; Manual. Windsor: NFER-Nelson

Heim, A.H. (1970). *Intelligence and Personality*. Harmondsworth, Middlesex: Penguin.

Heim, A.H., Watt, K.P. and Simmonds, V. (1974). *AH2/AH3 Group Tests of General Reasoning*; Manual. Windsor: NFER Nelson.

Industrial Psychologists Inc. (1981) *Perception. Test Examiner's Manual*, Test Agency, Bucks.

Schmidt F & Hunter J (1977) *Development of a General Solution to the Problem of Validity Generalisation*. Journal of Applied Psychology, 1977, 62, 529-540