

第二章 文獻探討

第一節 常用之試題反應理論模式簡介

在社會科學領域裡，我們常藉由考生作答的反應來估計其潛在的能力，而考生的能力與其在題目反應的關係可藉由試題反應理論來建立(Lord, 1980)，此理論的特色有(1)試題的參數估計值不變性(invariance)(2)能力的參數估計值不變，和(3)能力估計值的測量誤差大小，隨能力不同而異。

現代測驗理論發展至今，已有許多試題反應理論(IRT)模式被發展出來，因為國中學力測驗的內容皆是四選一的單選選擇題，在此，我們僅介紹大多適用於在大型學業成就評量的模式。

(一) 1-PL 模式：

1-PL 模式相通於 Rasch 模式，因此又可稱為 Rasch 模式，Rasch 認為受試者的潛在能力 θ 與受試者對試題 i 的反應可以用以下之試題特徵函數表示：

$$P_i(\theta) = \frac{\exp(\theta - b_i)}{1 + \exp(\theta - b_i)}$$

θ 表示考生的能力， b 代表題目的難度，若在 Rasch 模式中加入常數 D ，通常 D 值為 1.7，則原來的 Rasch 模式則稱為 1-PL 模式，其試題特徵函數表示為：

$$P_i(\theta) = \frac{\exp D(\theta - b_i)}{1 + \exp D(\theta - b_i)}$$

(二) 2-PL 模式：

2-PL 模式比 1-PL 模式多考慮一個試題參數 a ，即鑑別度參數，其試題特徵函數表示為：

$$P_i(\theta) = \frac{\exp Da_i(\theta - b_i)}{1 + \exp Da_i(\theta - b_i)}$$

(三) 3-PL 模式：

此模式又比 2-PL 模式多使用了一個猜測參數 c ，來描述試題，其試題特徵函數表示為：

$$P_i(\theta) = c_i + (1 - c_i) \frac{\exp Da_i(\theta - b_i)}{1 + \exp Da_i(\theta - b_i)}$$

c 參數為能力極低時仍有答對該題的機率。

第二節 測驗等化

由於測驗等化牽涉到不同測驗或不同受試者之間的關係，為了使得不同的估計數值可轉變成相同，或容易解釋與應用，測驗學家發展各種測驗等化的理論與技術。

一、建立評量量尺的重要性

若要申請美國研究所(除了商學院外)，大部份的學校都會要求 GRE 的成績，作為他(她)們申請入學之必要條件。假如甲生在今年 1 月份在計量部份考 700 分，語文考 400 分。乙生在今年 3 月份在計量部份考 600 分，語文考 500 分。雖然兩者的考試時間與試卷皆不同，但對於以上兩位考生的成績，我們將會有

一致的看法：甲生的計量部份比乙生好；乙生的語文部份比甲生好。此外，從甲生的計量得分，我們可以推論甲生的計量成績是屬於上，但語文成績部份則為中下，而甲生的數學能力可能比語文好。從乙生的得分，我們可以推論其數學能力是中上，語文能力是中等，同樣地，乙生的數學能力可能比語文好。且間隔一段期間後(如二個月)，甲生與乙生再重考一次，他們的得分應不會有太大的改變，除非在間隔的這一段時間，他們非常用功地準備與復習，或在兩次考試時之身心狀況有極大的差異。

二、等化的設計(王寶壙, 民 84; 余民寧, 民 82; Hambleton & Swaminathan, 1985; Kolen & Brennan, 1995)

1. 單團體設計(single group design)：讓同一組受試者，接受兩個或兩個以上的測驗，是最簡單的一種設計。但易受重覆練習和疲倦因素的影響。
2. 等團體設計(equivalent group design)：以隨機選取或配合抽取能力分配相等的不同受試團體，接受兩個不同的測驗而進行等化。
3. 等測驗設計(equivalent test design)：隨機分派題庫或很多要校準的題目成不同的測驗組合，但有如傳統測驗的平行複本試題，將此平行測驗施測於受試能力分配不必相同的團體，但嚴格的平行測驗實際上不易達成。
4. 校準團體設計(anchor group design)：此設計是以不同測驗施測於不同團體，但校準團體接受每個測驗的施測，而其他團體則只施測一種測驗。此設計的關鍵在校準團體的樣本數決定估計數的穩定性。

5. 校準測驗設計(anchor test design)：此設計是將校準試題放入不同的測驗中，而對不同的團體施測。但將校準試題放入不同的測驗，有下列幾種方式：
- 甲、內在校準測驗(internal anchor test)：將校準試題放入不同的測驗內。
 - 乙、外在校準測驗(external anchor test)：除了已有的測驗外，另外加測校準試題。
 - 丙、分段校準測驗(cascading anchor test)：不同校準試題放入不同的測驗內。
 - 丁、一致校準測驗(uniform anchor test)：相同校準試題放入不同的測驗內。
6. 校準試題矩陣抽樣設計(matrix sampling design with overlapping items)：此設計主要有兩種型態，一是區塊共同試題多重矩陣抽樣設計，另一種是分散型共同試題多重矩陣抽樣設計。兩種將校準試題、非共同試題及不同的受試團體，靈活交互應用，是建立題庫時有效的設計。
7. 同時校準法：藉由測驗資料的收集設計與 IRT 電腦軟體之應用結合，將所收集之數筆測驗資料同時執行試題校準，經過校準之後，所有题目的參數估計值，皆有相同的量尺單位，此乃同時校準法之涵意。
- 將所收集到的試題，經專家學審題後，組成數份不同之試卷，且試卷裡之

部份題目(定錨試題)是相同的(比率約 15%)(Vale, 1986)，將其施測於預先選好其能力分佈與學測相似的不同學校中，在蒐集並整理測驗的資料中，將幾組樣本對幾份試卷反應資料做調整，並整理成一筆測驗資料矩陣，如下圖。然後以 IRT 電腦軟體，如 BILOG-MG(Zimowski, Muraki, Mislevy & Bock, 1995)，對此一資料矩陣進行試題校準工作。假設此資料矩陣含有 300 道題目，則此 300 題之參數估計值將對應於相同量尺上。經過校準後的試題，必須能夠通過適合度的考驗者，方可被保留在題庫裡，因為它們可以被適當的反應模式所解釋。

KA		AB					
		AB		BC			
			BC		CD		
.....							
					IJ	JK	
						JK	KA

灰色部份通常假設其為考生未完成部份。

一般測驗機構所採用的量尺分數根據其產生的方式，大概可以分為兩種：

- (1) 將原始分數常態轉換後所得 (normalizing raw scores) 之量尺分數。例

如美國教育測驗服務社 (ETS) 的 GRE 或 TOEFL 測驗分數，就是將原始分數常態化轉換後的量尺分數。智力測驗中的比西量表分數 (Stanford-Binet Intelligence Scale) 及魏氏兒童智力量表分數 (Wechsler Intelligence Scale for Children-Revised Form; WISC-R) 等也是一種將原始分數常態化轉換後的量尺分數。(2) 均等測量標準誤 (equalizing measurement error variability) 之量尺分數。這是 E. L. Lindquist 在發展「愛荷華教育發展測驗」(Iowa Tests of Educational Development, ITED) 時所提出來建立量尺分數的方法，也是美國知名測驗機構 ACT 公司的 ACT Assessment Test (ACT, 1997) 所採用的量尺分數型態。基本上，這是在原始分數轉換成量尺分數的同時，利用數學的方法將每一個量尺分數點上的測量標準誤 (或稱測量誤差) 調整成相等或是非常接近 (Kolen & Hanson, 1989; Kolen, Hanson, & Brennan, 1992)。

測驗等化的方法可分為古典測驗理論與現代測驗 IRT 的方法。古典測驗理論等化法，是利用原始總分進行等化，又可分為線性等化(linear equating)和等百分位數等化(equipercen-tile equating)。現代測驗 IRT 的等化法常用的有迴歸法(regression method)、平均數與標準差法(mean and sigma method)、強韌平均數與標準差法(mean and sigma method)、真分數等化法(true score equating)及特徵曲線法(characteristic curve method)(王寶壙, 民 84; Angoff, 1971 ; Kolen, 1988 ; Hambleton & Swaminathan, 1985)。

一、古典測驗理論等化的方法

1. 線性等化法：

兩個測驗的分數分配相同，於各測驗中找出對應相同標準分數之觀察分數，即兩測驗之觀察分數可置於相同量尺比較，亦即除了平均數和標準差不同外，X 和 Y 分數的分佈是相同的。如果這個假設能夠成立，那就能從 X 和 Y 找出具有相同的 Z 分數配對，這些分數即為等化分數。亦即如果

$$(X - \mu_X) / \sigma_X = (Y - \mu_Y) / \sigma_Y$$

則 X 與 Y 等值。其中 X 與 Y 分別代表測驗 X 和 Y 的原始分數， μ_X 、 μ_Y 分別代表兩測驗的平均數， σ_X 與 σ_Y 分別代表兩測驗的標準差。其直線轉換為：

$$Y = aX + b$$

$$\mu_Y = a\mu_X + b$$

$$\sigma_Y = a\sigma_X$$

$$Y = (\sigma_Y / \sigma_X)X + [\mu_Y - (\sigma_Y / \sigma_X)\mu_X]$$

2. 等百分位數等化法

等百分位數等化法是在決定兩個測驗的哪些分數具有相同的百分等級。第一步是決定兩個測驗分數分佈上的百分等級。然後將兩個測驗的百分等級對原始分數的散佈圖畫出來。

在連結各資料點時必須利用直線內插法。如假設在 A 測驗卷的原始分數是 a，其對應的百分等級為 P，而在 B 測驗卷，從百分等級 P 所對應的原始分數是

b。則我們可以說一個人在 A 測驗的得分 a 與在 B 測驗的得分 b 是具有相同能力的。

Lord(1980)認為測驗分數要能公平的等化，需滿足下列需求：

- A. 測量不同特質或能力的測驗不能等化。
- B. 信度不相等的測驗之原始分數不能等化。
- C. 難度不相等的測驗之原始分數不能等化。
- D. 除非兩測驗是完全複本測驗，否則兩測驗不能等化。
- E. 具有完全信度的測驗分數才可以等化。

此外，對稱性(symmetry)與不變性(invariance)等二條件亦是進行測驗分數等化所必備的。對稱性是指無論是 X 轉化為 Y 或 Y 轉化為 X，等化結果必須相同。不變性則指等化的程序是獨立的，不受所選用樣本的影響。

由於古典測驗理論是依據弱勢假設而來，理論淺顯易懂、便於計算，但卻有下列的缺失與限制：

- A. 難度、鑑別度與樣本相依，即難度、鑑別度隨樣本的不同而變化。
- B. 不同測驗不易比較不同受試者的能力。
- C. 理論中假設所有受試者等測量標準誤與實際不符。

亦即以原始分數來進行等化有其限制存在，如公平性、對稱性與不變性等要求常無法獲得滿足。

二、試題反應理論的等化方法

常用的試題反應理論的等化方法有下列幾種：迴歸法(regression method)、平均數與標準差法(mean and sigma method)、強韌平均數與標準差法(robust mean and sigma method)、真分數等化法(true score equating)及特徵曲線法(characteristic curve method) (王寶壙, 民 84; Angoff, 1971; Kolen, 1988; Hambleton & Swaminathan, 1985)。

1. 迴歸法(regression method)

兩測驗分數的直線關係，可由下列迴歸式來表示：

$$y = \alpha x + \beta + e, \quad \alpha = r_{xy} \frac{s_y}{s_x}, \quad \beta = \bar{y} - \alpha \bar{x}$$

y ， x 可以是能力，即 $y = \theta_y$ ， $x = \theta_x$ ，也可以是難度，如 $y = b_y$ ， $x = b_x$ 。誤差 e 則視為獨立而相等的隨機變項，迴歸係數 α 與 β 是斜率與截距。 r_{xy} 是兩測驗相關係數， s_x, s_y 是標準差， \bar{x} 與 \bar{y} 是平均數。迴歸法的缺點是不對稱，只有 Rasch 模式例外。

2. 平均數與標準差法(mean and sigma method)

平均數與標準差法是假如 $y = \alpha x + \beta$ ，則 $\bar{y} = \alpha \bar{x} + \beta$ 且 $s_y = \alpha s_x$ ，

所以得到 $\alpha = s_y / s_x$ ， $\beta = \bar{y} - \alpha \bar{x}$ ，此關係是對稱的。

3. 強韌平均數與標準差法(mean and sigma method)

如果連結是包含兩個以上的參數，使用平均數與標準差法的估計誤差變化很

大，而且離異值(outlier)會影響係數的計算。因此，Linn, Levine, Hastings, & Wardrop(1981)提出強韌平均數與標準差法。

強韌平均數與標準差法認為應給每一對(x, y)的能力或參數加權，此加權得自兩個難度參數估計值的變異數中較大一個的倒數，其步驟為：

a. 決定每個配對 (x_i, y_i) 的加權植 w_i ，

$$w_i = \max\{Var(x_i), Var(y_i)\}, \quad j=1, 2, \dots, k,$$

b. 加權值量尺化，新加權值： $w_i' = w_i / (\sum_{i=1}^k w_i)$

c. 計算新的 x 與 y 值

$$x'_i = w'_i x_i, \quad y'_i = w'_i y_i \quad j=1, 2, \dots, k,$$

d. 求出加權後的平均數與標準差：

$$\bar{x}', \bar{y}', s'_x, s'_y$$

e. 以加權後的平均數與標準差求出 α, β ：

$$\alpha = s'_y / s'_x \quad \beta = \bar{y}' - \alpha \bar{x}'$$

4. 特徵曲線法(characteristic curve method)

上述所提的各種等化方法，都只考慮能力與難度指數，其餘之鑑別度及猜測度都未列入考慮，因此，Haebara(1980)及 Stocking & Lord(1983)提出特徵曲線法，方法如下：

具有 θ_a 能力的受試者在測驗 X 的真分數為 ξ_{xa}

$$\xi_{xa} = \sum_{i=1}^n P(\theta_a, a_{xi}, b_{xi}, c_{xi})$$

具有 θ_a 能力的受試者在測驗 Y 的真分數為 ξ_{ya}

$$\xi_{ya} = \sum_{i=1}^n P(\theta_a, a_{yi}, b_{yi}, c_{yi})$$

試題參數是：

$$b_{yi} = \alpha b_{xi} + \beta,$$

$$a_{yi} = a_{xi} / \alpha,$$

$$c_{yi} = c_{xi}$$

為求兩真分數間差異的最小值，其函數為

$$F = \frac{1}{N} \sum_{a=1}^N (\xi_{xa} - \xi_{ya})$$

用 Newton-Raphson 解出以下偏導數方程式，即可得兩真分數的最小值：

$$\frac{\partial F}{\partial \alpha} = \frac{\partial F}{\partial \beta} = 0$$

因特徵曲線法涉及三個試題參數，等化的效果最佳。

5. 真分數等化法(true score equating)

IRT 能力分數之評量單位不易為一般人所接受，而且長久以來大眾皆習慣採用原始分數來計分。有鑑於此，IRT 原始分數(IRT observer score equating)與真實分數等化法(IRT true score equating)乃被發展運用來克服此問題。IRT 真實分數為考生在 n 題題目答對機率的總和，公式如下(Lord, 1980)：

$$T | \theta = \sum_{i=1}^n P_i(\theta)$$

真實分數等化方法包括兩個主要步驟(Kolen & Brennan, 1995)：

- A. 將兩份試卷的試題參數估計值對應在相同之量尺上。
- B. 使用 IRT 分數為橋樑，製做兩份試卷真實分數的對應表。

第二步驟之原理為：假定兩組(分為基礎組與等化組)樣本之能力參數之單位相同，且兩份試卷的試題參數值之單位亦同，將基礎組每一位受試者的 IRT 能力分數轉化為真實分數；同樣地，也等化組每一位受試者的 IRT 能力分數轉化為真實分數。然後以 IRT 分數為橋樑，製造出兩份試卷真實分數的對應表。IRT 真實分數等化方法之最大優點為：當我們從題庫裡選定試題樣本後，即可開始製造試卷間之分數轉化表，而不需要等到考完試取得考生作答資料再來做。但先決條件為題庫裡試題的參數估計值要準確。

6. 觀察分數等化法

真實分數 ξ 和觀察原始分數 r 有相同的量尺，而

$$r = \sum_{i=1}^n U_i$$

若項目反應理論是有效的，則

$$E(r) = \xi$$

一個測驗理論上的觀察分數分佈 $f(r | \theta)$ 可以由下列等式中獲得

$$\sum_{r=0}^n f(r | \theta) t^r = \prod_{i=1}^n [Q_i(\theta) + tP_i(\theta)]$$

若 θ 的值已知，如 N 個考生的能力為 $\theta_1, \theta_2, \dots, \theta_a, \dots, \theta_N$ ，則 $f(r)$ 的邊際分配為

$$f(r) = \sum_{a=1}^N f(r | \theta_a)$$

我們可以利用下列例子來做說明：

假設有一個 2 參數的 logistic 模式，有三個題目，其

難度 $b = [b_1, b_2, b_3] = [1.0, 0.0, -1.0]$

鑑別度 $a = [a_1, a_2, a_3] = [1.5, 1.0, 0.5]$

給定下列 5 個能力 $\theta = [\theta_1, \theta_2, \theta_3, \theta_4, \theta_5] = [-2, -1, 0, 1, 2]$

根據上面的條件，我們可以獲得在不同能力的理論上觀察分數分佈如下：

條件相對次數 $f(r \theta)$					
原始分數	$\theta = -2$	$\theta = -1$	$\theta = 0$	$\theta = 1$	$\theta = 2$
0	.678	.420	.139	.012	.000
1	.313	.500	.475	.143	.009
2	.010	.080	.361	.488	.158
3	.000	.000	.025	.357	.833

假設總共的考生有 100 人，而且具有能力 $\theta = -2, -1, 0, 1, 2$ 的人分別有 5, 15, 30, 40, 及 10 人，則原始分數的理論邊際次數分配可計算得如下表：

原始分數	條件相對次數 $f(r \theta)$					邊際次數 分配 $f(r)$
	$\theta = -2$	$\theta = -1$	$\theta = 0$	$\theta = 1$	$\theta = 2$	
0	3	6	4	0	0	13
1	2	8	14	6	0	30
2	0	1	11	20	2	34
3	0	0	1	14	8	23
考生人數	5	15	30	40	10	100

根據以上，觀察分數等化法的步驟如下：

- a. 根據 $\sum_{r=0}^n f(r | \theta) t^r = \prod_{i=1}^n [Q_i(\theta) + tP_i(\theta)]$ 及考生的能力及測驗 X 的項目參數，可以得到條件次數分配
- b. 再從 $f(r) = \sum_{a=1}^N f(r | \theta_a)$ 得到邊際次數分配
- c. 對於測驗 Y 重覆 a, b 的步驟
- d. 用等百分位數來等化此兩個測驗的原始分數

與觀察分數等化法比較，IRT 真實分數等化方法之計算過程較為簡易，此方法無須依賴考生能力之分佈狀態，而且經過 Lord & Wingersky (1983) 的研究發現，此兩種方法的等化結果相似，況且 ETS 測驗公司亦採用 IRT 真實分數等化方法。因此，本研究擬採用 IRT 真實分數法來對兩次學測的數學科做等化，然後與國中基本學力測驗中心的結果比較做探討。

第三節 國中基本學力測驗與聯考實務作法的差異

基本學力測驗(以下簡稱基本學測)與傳統聯招考試在作法上是有實質上的不同。這些差異主要是圍繞在四個相關聯的議題：(一)、基本學測的內容(基本學測到底是什麼?)；(二)、入闈組題；(三)、量尺分數；(四)、一年多次(目前是一年兩次)測驗等化的問題(林世華，民91)。

基本學測與傳統聯考最大的差異是：測驗內容與測驗工具發展方法的差異。基本學測主要是設計來評量國中學生在修完國中三年後，所學習到的基本知識與所培養的基本能力。所謂「基本」意指核心的、重要的、統整的。國中三年老師在學校所教的，學生所學習的，包羅萬象，其中主要的當然是基本的範疇，當然也不乏複雜專業方面的涉略。基本學測是排除複雜專業的部分。基本學測的發展方法是以標準化心理測驗發展的科學方法為依據。科學方法步驟中，包括：清楚界定基本學力；測驗試題編寫，審查與試驗研究；測驗分數如何使用等。基本學測研究團隊從八十七年起便著手建造題庫，其中每一試題均需經歷嚴密審查並預試，以便清楚掌握試題的品質。題庫中的試題是不停地在增加與淘汰。這與聯考在考前入闈後才開始命題形成測驗是截然不同的。

其次是入闈組題的議題。基本學測是以持續發展的題庫入闈，並在闈場以自行開發的組題電腦系統來選題，組成當次的基本學測。理論上，系統會組成

哪一種測驗試題，事先大家並不清楚。所能清楚的是組題的原則。由於基本學測的界定相當清楚，所以組題原則也隨之具體。幾個主要原則如：試題評量符合基本學測目的；試題取材均勻分佈（到國三下第二次段考前的範圍），生活化、統整化素材優先；試題難度以 50-75% 通過率為主。上述組題原則的目的是為了保證所組出來的測驗能吻合基本學測的原有目的。這也等同於是說，基本學測有一定的設計規格，基本上來說，基本學測是朝穩定性的方向來發展。傳統聯考中易出現「去年簡單今年會難」、「去年難今年會簡單」的現象，在基本學測中這樣的情況會比較少出現，但預測難度及真實難度若有差別的話則是在所難免。

至於量尺分數則是基本學測分數使用上的考量。測驗分數的使用有兩大重要考量；第一是，測驗相對機制的建立（專業術語稱建立常模），它是希望能建立一個類似於「IQ100 是中等」，「IQ100 以上是中等資質以上」的具體意義。這套作法，在理論上還可以將量尺分數作跨年的比較，然而由於教育部所公佈的命題範圍會因不同學年而異等因素，以測驗的角度來看量尺分數作跨年的比較，暫不可完全實行，但有其參考的價值。第二是測驗分數的採計與運算。基本學測是以答對題數作為原始分數。原始分數在測量性質上是屬於名次性質，它會有在不同區段中名次間的差距看似等距離，其實不是等距離（如第 1 與第 2 名；第 15 與 16 名，間距均為 1，但實質並不是相等的 1），因此將其轉換成一

個可以採計與運算的量尺分數（基本學測量尺分數是1~60）是一個測驗標準化方法中的重要步驟。

最後基本學測是一年多次的（目前是一年兩次），國人也不易接受。來自各方的質疑有兩方面，第一、「學生會因考到簡單題而佔便宜，也會因考到難題而吃虧」。其實基本學測採用量尺分數的目的之一，也是因應此一問題。不過處理這問題的基本作法有二：第一、不同次的基本學測內容不會一樣，但是組題分佈、難易程度均要控制成一樣（再一次保證了基本學測穩定性的必要）；第二、如果在控制上出現小瑕疵，則另一個專業機制：等化，便可以派上用場。簡單的說測驗等化設計可以讓學生在考到稍難試題時，只要答對較少題時，就可以得到考簡單題時要答對較多題時的量尺分數，舉例來說：第一次數學科測驗較簡單，第二次數學科測驗較難，那是有可能透過等化，第一次對10題與第二次對9題的量尺分數是一樣的。另外一個質疑則是兩次基本學測考生不同（有第二次考生資質較優，亦有較低之說），所建立量尺分數在相對機制上有公平性的問題。確實，基本學測如果每次測驗均獨立建造常模，那會是問題。但基本學測其實不是這樣做。以90年為例：只有第一次基本學力測驗真正建立量尺分數對照表，第二次其實是以第二次答對題數等化到第一次答對題數，再由第一次答對題數對照量尺分數。

基本學力測驗的分數等化過程(涂柏原，民92)：

- (1) 由於幾乎所有的國三學生都參加了第一次國中基本學力測驗，第一次測驗的結果非常適合用來建立量尺分數。因此，根據前面所提到的原始分數轉換成量尺分數的方法，我們先將原始分數與量尺分數對照表算出來。如此，我們可以得到每一位考生的原始分數與量尺分數等資料。因為要用 IRT 的方法來進行等化，所以我們得進一步用 IRT 的方法來估算每一位考生的能力參數，根據全部考生的資料，經過綜合整理之後，我們可以得到下表的第二、三及四欄的資料，也就是答對題數、能力值與量尺分數的對照表。
- (2) 根據參加第二次學力測驗的考生在第二次測驗上面的作答資料，用 IRT 的方法來估計考生的能力值。如此，每一位考生的答對題數原始分數與其能力估計值都可以被收集到，而將這些資料與從第一個步驟所得到的資料合併，我們可以得到如下的對照表範例。
- (3) 第二次國中基本學力測驗不再建立新的量尺分數，而是使用我們在上一個步驟所得到的表，來將考生在第二次學力測驗所得到的原始分數，對應到在第一次學力測驗時所建立的量尺分數。舉例而言，如果某一考生在第二次基本學力測驗某學科的原始分數為 58 分，從下表我們可以得知他所對應的量尺分數為 51，相當於在第一次基本學力測驗該科的原始分數為 62。換言之，對該科來說，在第一次基本學力測驗原始分數為 62 分的考生，與在第二次基本學力測驗原始分數得到 58 分的考生，有相同的量尺分數 51 分。

某科量尺分數對照表範例：

第二次答對題數	能力值	第一次答對題數	量尺分數
63——→	4.93	66	60
	4.4		59
62——→	4.1		58
	3.71	65	57
61——→	3.35		56
	2.99	64	55
60——→	2.77		54
	2.56	63	53
59——→	2.4		52
58——→	2.25	62	51
:	:	:	:
:	:	:	:

IRT 最為有力的地方，是同一個考生若在同一個時段考兩次試，即使作答不同的試題，所估計出來的兩個能力參數的值，理論上是一樣大小的。如果考完第一次之後隔了一些時日再考第二次，以至於兩次考試之間有成長（或進步）的現象存在，那麼根據第二次測驗的結果所估計得到的能力參數值，將會大於

第一次測驗後所估計得到的。因為 IRT 具備這樣的特性，因此透過用 IRT 方法所進行的等化之後，我們有信心考生在第二次測驗所得到的量尺分數絕對是合理、公平、公正的；而且從兩次測驗所得到的量尺分數之間也是可以相互比較的。

影響整個基本學力測驗等化程序的關鍵點，其實是在於各個題目的試題參數 (item parameter) 是否被精確地估計出來。因為這裡所說明的等化方法在某個程度之下，其原理與目前常見的電腦化適性測驗 (computerized adaptive testing, CAT, 如 TOEFL、GRE 以及美國的護士執照考試等) 原理是一樣的，是在試題參數已知的情形下估計考生的能力參數 (ability parameter) (註：在 IRT 中，每一個試題會有幾個試題參數來描述該題的特性，而也有一個能力參數來描述考生的能力)。目前國中基本學力測驗每一道試題皆經過至少 240 至 320 位不同地區的國三學生「預試」過，截至目前為止，參與過預試的國三學生已超過數萬人，涵蓋全國各縣市的國民中學。以此大規模的預試工作所得的答題反應資料，可以用 IRT 來估計出每一個題目的試題參數 (難度、鑑別度) 等。經過實際驗證，九十年預試所得的題目難度與該年度第一次測驗三十萬人的資料，算出來的題目難度相當接近。即使如此，學力小組仍然依據從三十萬人的資料所得到的試題參數，來將題庫中每一個題目的試題參數加以校正。