Mesure et évaluation en éducation

MESURE ÉVALUATION en Éducation

La capacité discriminante d'un instrument de mesure

Louis Laurencelle

Volume 20, numéro 1, 1997

URI: https://id.erudit.org/iderudit/1091385ar DOI: https://doi.org/10.7202/1091385ar

Aller au sommaire du numéro

Éditeur(s)

ADMEE-Canada - Université Laval

ISSN

0823-3993 (imprimé) 2368-2000 (numérique)

Découvrir la revue

Citer cet article

Laurencelle, L. (1997). La capacité discriminante d'un instrument de mesure. Mesure et évaluation en éducation, 20(1), 25–39. https://doi.org/10.7202/1091385 ar

Résumé de l'article

La « capacité discriminante » est cette propriété d'un test, instrument de mesure ou examen scolaire, grâce à laquelle on peut séparer les personnes ou objets mesurés selon leurs valeurs et les discriminer les uns des autres. Le concept, partiellement inspiré de l'indice de pouvoir classificatoire de Ferguson, est élaboré sur trois bases : la probabilité d'une catégorisation correcte, la largeur suffisante des intervalles de mesure et le nombre d'intervalles efficaces d'une distribution de mesures empirique ou théorique. La capacité discriminante, qui s'exprime comme une fonction du coefficient de fidélité de l'instrument de mesure, s'ajoute donc à l'arsenal conceptuel de la théorie classique des tests.

Tous droits réservés © ADMEE-Canada - Université Laval, 1997

Ce document est protégé par la loi sur le droit d'auteur. L'utilisation des services d'Érudit (y compris la reproduction) est assujettie à sa politique d'utilisation que vous pouvez consulter en ligne.

https://apropos.erudit.org/fr/usagers/politique-dutilisation/



Cet article est diffusé et préservé par Érudit.

Érudit est un consortium interuniversitaire sans but lucratif composé de l'Université de Montréal, l'Université Laval et l'Université du Québec à Montréal. Il a pour mission la promotion et la valorisation de la recherche.

La capacité discriminante d'un instrument de mesure

Louis Laurencelle Université du Québec à Trois-Rivières

La « capacité discriminante » est cette propriété d'un test, instrument de mesure ou examen scolaire, grâce à laquelle on peut séparer les personnes ou objets mesurés selon leurs valeurs et les discriminer les uns des autres. Le concept, partiellement inspiré de l'indice de pouvoir classificatoire de Ferguson, est élaboré sur trois bases : la probabilité d'une catégorisation correcte, la largeur suffisante des intervalles de mesure et le nombre d'intervalles efficaces d'une distribution de mesures empirique ou théorique. La capacité discriminante, qui s'exprime comme une fonction du coefficient de fidélité de l'instrument de mesure, s'ajoute donc à l'arsenal conceptuel de la théorie classique des tests.

(capacité discriminante - erreur de mesure - catégorisation correcte - indice de Ferguson)

« Discriminative capacity » is defined as a property of a test, measuring device or scholastic exam, which enables us to segregate and categorize objects or people according to their measured values. The concept, partially derived from Ferguson's index of classificatory power, is developed upon three bases: the probability of categorizing an object (or person) in its proper measuring interval; the sufficient length of measuring intervals; the number of efficacious intervals in an empirical or theoretical distribution of measures. Expressed as a function of the reliability coefficient of a measuring device, discriminative capacit appears as a brand new tool in the conceptual apparatus of classical test theory.

(discriminative capacity - measurement error - rightful categorization - Ferguson's index)

Une règle à mesurer de 1 mètre, graduée au 0,01 m, nous permet de catégoriser tous les objets selon leurs grandeurs en 101 groupes distincts: objets de longueurs inférieures à 0,01 m, objets de 0,01 m à moins de 0,02 m, et cœtera jusqu'aux objets de 1 m et plus. De même, un pèse-personne en kilogrammes, gradué au ½ kg et plafonnant à 150 kg, peut séparer les personnes et les objets selon leurs masses en 301 catégories distinctes. En combien de catégories distinctes une épreuve d'aptitude en mathématiques sépare-t-elle les élèves? Combien de degrés d'intensité vraiment distincts

nous permettent d'obtenir une échelle de tendances suicidaires? En combien de catégories un test de quotient intellectuel (QI) va-t-il répartir une population?

La « capacité discriminante » d'un instrument de mesure est le nombre de catégories de mesure qu'exploite cet instrument et parmi lesquelles il peut classer les objets. Pour un objet de mesure physique, tels la règle à mesurer et le pèse-personne, la capacité discriminante, qu'on peut dénoter D', est facilement obtenue par

$$D' = E / u,$$

c'est-à-dire en divisant l'étendue (E) de l'instrument par son unité de mesure $(u)^1$. Cette opérationnalisation se heurte toutefois à deux difficultés. La première est que, pour la plupart des instruments, voire pour la règle à mesurer, l'étendue effective (E) reste indéfinie et ambiguë:

- indéfinie parce que, pour la règle par exemple, elle peut être appliquée bout à bout et que l'étendue virtuelle serait alors infinie et parce que, en général et pour les mesures biologiques, psychologiques et cognitives, il n'y a pas de vrai minimum ni de vrai maximum et qu'il est impossible de déterminer l'étendue;
- ambiguë parce que, pour une valeur donnée de E, la répartition des éléments d'un univers (d'objets ou de personnes) varie d'une catégorie à l'autre et est rarement uniforme, la zone centrale étant ordinairement plus densément occupée que les zones extrêmes, et que la capacité discriminante devrait refléter le nombre de catégories effectives (ou effectivement occupées) du système de mesure.

La seconde raison, très prégnante pour les édumétriciens et les psychométriciens, vient de ce que l'unité de mesure (u), dans un examen scolaire ou un test, n'est aucunement définie² et que la catégorisation d'une personne dans une mesure X d'unité u (ou X_u) est fortement incertaine, en raison d'une erreur de mesure importante caractérisant la plupart de nos instruments de mesure.

C'est pourquoi nous proposons d'aborder différemment le concept de capacité discriminante en tenant compte à la fois de l'erreur de mesure, en fait de la *fidélité* (ρ_{XX}) de l'instrument de mesure, et en invoquant une répartition normale (ou gaussienne) du caractère mesuré dans la population. Par « capacité discriminante » d'un instrument de mesure, nous entendons alors

le nombre d'intervalles de mesure efficaces, ou catégories efficaces, parmi lesquels l'instrument de mesure à fidélité ρ_{XX} peut répartir une population statistique normale de telle façon qu'un élément mesuré ait une probabilité d'au moins ½ d'être classé dans son intervalle propre.

La capacité discriminante, que nous dénoterons D, a pour formule générale:

(2)
$$D = 2,654 / \sqrt{1 - \rho_{XX}}.$$

Dans les sections qui suivent, nous développons et justifions les bases mathématiques du concept de capacité discriminante, en payant d'abord tribut aux auteurs qui ont précédemment abordé la question.

Notes historiques

Les publications en psychométrie ne font pas mention d'un concept semblable à celui de la capacité discriminante. Nous retrouvons cependant quelques idées apparentées à notre concept et, notamment, l'indice de pouvoir classificatoire de Ferguson, dont nous avons tiré profit.

Bassière et Gaignebet, dans leur traité de Métrologie générale (1966), font seulement mention de la « capacité d'information » d'un instrument, définie comme le nombre d'états distincts qu'il peut prendre. « La capacité d'information de l'instrument dépend à la fois de son pouvoir de résolution et de son temps de réponse qui limite le nombre de mesure [sic] par unité de temps » (op. cit., p. 140-141). L'implication du temps, ou d'un taux d'information par unité de temps, évoque en même temps le concept de « capacité de canal » en théorie de l'information (voir, par exemple, Ralston & Reilley, 1983) et crée, pour nous, une confusion notionnelle.

D'un autre côté, les auteurs en psychométrie et en docimologie reconnaissent depuis longtemps l'impact des niveaux de réussite des items d'un test sur le bon étalement des scores du test, donc sur sa capacité discriminante. Cette discussion, qui concerne l'analyse et la sélection des items devant constituer un test afin de produire un bon étalement des scores totaux, se retrouve par exemple dans Davis (1951) et Anastasi (1994); on peut la résumer comme suit. S'ils sont peu inter-corrélés, il faut sélectionner des items à niveaux de réussite (ou « indices de facilité ») moyens, c'est-à-dire des niveaux proches de ½ pour des items dichotomiques. Au contraire, si les items sont fortement corrélés et ont tendance à discriminer les mêmes répondants, on doit les choisir avec des niveaux de réussite gradués. Bien que pertinentes pour notre propos, ces considérations sur la construction d'un test

à base d'items se situent en amont du concept de capacité discriminante, une propriété globale de l'instrument de mesure.

Le vocable de « pouvoir discriminant » revient dans les textes de psychométrie classique, mais dans une acception restreinte, et il désigne alors l'efficacité d'un item pour discriminer les meilleurs des moins bons répondants, selon le caractère mesuré (Henrysson, 1971). Cette propriété d'item est quantifiée au moyen de divers indices, tels les soi-disant « indices d'homogénéité », le coefficient de corrélation bisériale entre item et score total, etc. (voir aussi Guilford, 1954). Elle se retrouve aussi dans la théorie des réponses aux items, précisement dans le paramètre de pente (« a ») de la courbe caractéristique d'item (Hambleton, Swaminathan & Rogers, 1991).

Ferguson publie, en 1949, un article intitulé « On the theory of test discrimination », dans lequel il propose de considérer et de quantifier, d'une certaine manière, l'étalement des scores d'un test ou des données d'un instrument de mesure; il propose aussi un « coefficient de discrimination » (voir aussi Guilford, 1954). Ferguson argue que, si n répondants se répartissent en k scores différents à raison de f_1 , f_2 , …, f_k répondants par valeurs de score, les discriminations obtenues sont égales au nombre de distinctions deux à deux que le test permet d'opérer parmi les n individus. Ce nombre, calculé par

(3)
$$N_d = f_1 \times f_2 + f_1 \times f_3 + \dots + f_{k,1} \times f_k,$$

reflète donc l'efficacité de l'instrument de mesure pour discriminer les uns des autres individus ou objets. Le nombre de distinctions (N_d) s'obtient aussi par la formule équivalente:

(4)
$$N_{d} = {n \choose 2} - \sum_{j=1}^{k} {f_{j} \choose 2}.$$

Ce nombre N_d dépend fondamentalement de trois facteurs: le nombre total de mesures (n), le nombre de valeurs (ou catégories de mesure) disponibles (k), et la distribution des mesures d'une valeur à l'autre (f_1, f_2, \cdots) . Pour n constant, le pouvoir discriminatif pourra donc croître avec k, mais il dépend de façon essentielle de la distribution des fréquences f_j . En effet, la seule adjonction de nouvelles catégories de mesure, au-delà de k, n'entraînera de distinctions supplémentaires que si ces catégories sont occupées. En outre, le nombre maximal de distinctions qu'un système à k catégories permet d'obtenir advient lorsque ces catégories sont occupées également, c'est-à-dire si la distribution des fréquences f_j est uniforme. Ce nombre maximum est alors, approximativement:

(5)
$$\max N_d \approx \binom{n}{2} - k \binom{n/k}{2}$$
$$= \frac{n(n-1)}{2} - \frac{n(n-k)}{2k}$$

Le « coefficient de discrimination » 8, proposé par Ferguson, est alors le quotient du nombre réel de distinctions opérées sur le nombre maximal, soit

$$\delta = N_d / \max N_d.$$

Thurlow (1950), dans un long article publié l'année suivante, réclame la copaternité du coefficient de discrimination tel que présenté par Ferguson. Fait intéressant, Thurlow rajoute la notion de « discriminations stables », c'est-à-dire des discriminations qui persistent du test au retest par contraste aux discriminations non stables, et il relie cette différence à la fidélité de l'instrument de mesure. Cependant, l'auteur n'opérationnalise pas davantage son idée.

L'élaboration du concept de « capacité discriminante »

Notre concept de capacité discriminante est essentiellement basé sur celui de Ferguson. Cependant, nous avons généralisé le concept emprunté à Ferguson, qui s'appuyait sur une distribution de fréquences, pour l'appliquer à toute distribution de probabilités d'intervalles, puis nous l'avons inversé afin de déterminer le nombre équivalent k^* d'intervalles (ou catégories de mesure) efficaces. La capacité discriminante, en outre, est tout aussi essentiellement basée sur la fidélité de l'instrument de mesure, dénotée par le coefficient habituel ρ_{XX} , ceci en faisant intervenir explicitement une probabilité (dénotée γ) de catégorisation de chaque mesure dans son intervalle propre. Les paragraphes suivants présentent les étapes de développement du concept.

Catégorisation correcte d'une personne (ou d'un objet)

Plaçons-nous d'emblée dans le contexte de la mesure d'un caractère dans un objet, ou dans une personne, au moyen d'un instrument de mesure. L'instrument, appliqué à l'objet, fournit une mesure X; l'unité de mesure n'est pas spécifiée. Le coefficient de fidélité ρ_{XX} est connu.

Pour un objet donné (ou une personne donnée) i, la valeur précise, ou valeur vraie V_i , de cet objet existe³ et chaque mesure en constitue une approximation. En fait, la mesure X_i va s'écarter plus ou moins de V_i selon que la précision, la fidélité du procédé de mesure est petite ou grande. En fait, plus grand est l'intervalle de mesure formé autour de V_i , plus grande est

la probabilité que la mesure prise X_i s'y retrouve. Soit L la largeur de l'intervalle considéré et γ la probabilité que la mesure X_i de l'objet considéré s'y retrouve, nous avons

(7)
$$\gamma(L) = \Pr\{X_i \in (V_i - \frac{1}{2}L, V_i + \frac{1}{2}L)\}.$$

La probabilité γ qu'un objet soit catégorisé dans son intervalle propre, c'est-à-dire dans le voisinage immédiat de sa valeur vraie V_i , varie directement avec la largeur d'intervalle considérée. À la limite, posant L=0, on voit qu'il est (mathématiquement) impossible que X_i égale V_i : ainsi, il serait illusoire de croire, pour un instrument de mesure à échelle infiniment divisible, que la mesure obtenue X_i soit, jusqu'à la dernière décimale, la valeur exacte de l'objet mesuré. Pour catégoriser numériquement un objet avec une certaine plausibilité, il faut que la probabilité de catégorisation correcte soit déterminée et suffisante, ce qui implique en retour la détermination d'un intervalle de grandeur suffisante.

Détermination d'un intervalle de mesure suffisant

L'écart de la mesure X_i par son rapport à V_i est généralement nommé erreur de mesure et dénotée ε (Lord & Novick, 1968; Laurencelle, sous presse). La valeur attendue (ou espérance) de ε est 0, sa variance (ou espérance de ε^2) σ_{ε}^2 et sa distribution, symétrique. C'est expressément pour cet écart ε , entre mesure et valeur vraie, que le modèle de la loi normale fut élaboré au XVIII siècle (Stigler, 1986). Nous pouvons donc, en toute légitimité, faire se reporter la variable ε au modèle normal et en indiquer la distribution par

$$\varepsilon \sim N(0, \sigma_{\varepsilon}^2).$$

Dans le contexte indiqué, la mesure $X_{i, o}$ de l'objet i à l'occasion o s'écrira donc

(8)
$$X_{i,o} = V_i + \varepsilon_o.$$

De plus, dénotant par Z une variable normale de distribution N(0,1), nous pouvons ré-exprimer (7) plus explicitement, comme

(9)
$$\gamma(L) = \Pr\{ V_i + \varepsilon_o \in (V_i - \frac{1}{2}L, V_i + \frac{1}{2}L) \}$$

$$= \Pr\{ \varepsilon_o \in (-\frac{1}{2}L, +\frac{1}{2}L) \}$$

$$= \Pr\{ Z \in (-\frac{1}{2}L/\sigma_e, +\frac{1}{2}L/\sigma_e) \}$$

$$= 1 - 2 \times \Pr\{ Z > \frac{1}{2}L/\sigma_e \} .$$

Fixant la probabilité γ à une valeur déterminée, nous pouvons inverser (9) afin de trouver la grandeur d'intervalle L suffisante pour qu'une mesure X_i y

soit catégorisée avec une probabilité (d'au moins) γ. L'équation résultante est simplement

(10)
$$L(\gamma) = 2\sigma_{\varepsilon} z_{\nu_{0}(1+\gamma)};$$

dans cette expression, $z_{\nu(1+\gamma)}$ dénote le centile $100 \times \frac{1}{2}(1+\gamma)$ de la loi normale standard.

Nous pouvons enfin ramener la grandeur de l'intervalle suffisant L en valeur standard, c'est-à-dire sur une échelle standardisée de moyenne 0 et variance 1, en divisant les deux parties de l'équation (10) par l'écart-type de l'échelle de mesure, $\sigma_{\rm X}$. La fidélité $\rho_{\rm XX}$ étant également définie par

(11)
$$\rho_{XX} = \sigma_V^2 / \sigma_X^2 = 1 - \sigma_\epsilon^2 / \sigma_X^2,$$

nous pouvons écrire de façon équivalente $\sigma_{\epsilon} = \sigma_{x}\sqrt{(1-\rho)}$. Ainsi, la grandeur d'intervalle suffisante standardisée, $\lambda(\gamma)$, devient

(12)
$$\lambda(\gamma) = 2\sqrt{(1-\rho)}z_{\nu_{i}(1+\gamma)}.$$

Dans une échelle de mesure X catégorisée en intervalles de grandeur commune $\lambda(\gamma)\sigma_X$, la mesure X_i d'un objet serait placée dans son intervalle propre avec une probabilité γ . La question qui reste à résoudre est la suivante: combien une telle échelle comporte-t-elle d'intervalles suffisants, de catégories?

Le nombre d'intervalles d'une distribution de probabilité et les intervalles efficaces

Adoptons d'emblée un procédé de mesure idéal, relatif à une variable X de fidélité parfaite $(\rho_{XX} = 1)$ et ayant dans une population hypothétique une distribution quelconque, de densité de probabilité f(X) et de répartition F(X). La probabilité qu'une mesure X soit égale à une valeur précise x est nulle, par définition. Cependant, la probabilité p que X tombe dans un intervalle, par exemple $p_{a,b} = \Pr[X \in (a,b)], b > a$, est facilement calculable par

(13)
$$p_{ab} = \Pr[X \in (a, b)] = F(b) - F(a);$$

ce calcul peut être effectué pour tous les intervalles de X.

Maintenant, l'indice de discrimination de Ferguson, en particulier le nombre de distinctions N_d de l'équation (4), est basé sur les fréquences d'occupation des intervalles plutôt que sur leurs probabilités. Imaginons alors un échantillon de taille n, les n observations étant réparties à raison de la

fréquence f_j dans l'intervalle j. En valeurs attendues f_j ou pour la moyenne de tous les échantillons de taille n tirés de la population, nous avons

$$\hat{f}_i = \mathbf{E}(f_i) = np_i.$$

Admettant un nombre quelconque (k) d'intervalles de mesure, l'expression de Ferguson (4) devient, en espérance et pour un échantillon hypothétique de taille n:

(15)
$$E(N_d) = \binom{n}{2} - \sum_{j}^{k} \binom{np_j}{2}$$
$$= \frac{n^2}{2} \left(1 - \sum_{j}^{k} p_j^2 \right).$$

La taille d'échantillon n étant arbitraire, de même que la constante « 2 » dans la formule, nous simplifions la formule en divisant (15) par $n^2/2$, obtenant enfin

(16)
$$C_k = 1 - \sum_{i=1}^{k} p_i^2,$$

cette quantité étant proportionnelle (plutôt qu'égale) au nombre de distinctions opérées en vertu des fréquences des k intervalles considérés. À l'instar de Ferguson (1949) et Thurlow (1950), il est aisé de montrer que le nombre maximum de distinctions, ou la valeur maximale de (16), survient lorsque toutes les probabilités sont égales⁴, c'est-à-dire lorsque $p_j = 1/k$ pour tout j, ce maximum étant

(17)
$$\max C_k = 1 - 1/k$$
.

Le nombre maximum de distinctions survient donc quand toutes les fréquences, ou de façon équivalente, toutes les probabilités, sont égales: nous désignerons alors les intervalles de mesure ayant des probabilités égales intervalles efficaces, et nous désignerons provisoirement par le symbole k^* le nombre d'intervalles efficaces d'un système de mesure. La règle suivante permet de déterminer k^* dans un cas donné. Soit C_k défini en (16), le nombre de distinctions opérées par un système de mesure en k intervalles de fréquences $\{f_j\}$. Alors, en égalisant $C_k = \max C_k$, et en inversant (17) pour obtenir k^* , nous avons

(18)
$$k^* = (1 - C_k)^{-1};$$

l'indice k^* indique le nombre d'intervalles efficaces, c'est-à-dire des intervalles virtuels à probabilités égales tels qu'ils déterminent le nombre de distinctions C_k constaté dans le système de mesure X. En d'autres mots, si notre système de mesure X était découpé en k^* intervalles à probabilités (ou fréquences d'occupation) égales, il serait caractérisé par un indice de discrimination C_k . égal à l'indice C_k observé. La quantité k^* constitue ainsi une mesure standardisée du nombre d'intervalles d'un système de mesure, quelle que soit la forme de distribution de ses données.

Il est enfin important de noter que la détermination de k^* ne dépend pas du nombre k d'intervalles originaux du système de mesure; ce nombre k pourrait même rester indéterminé. Nous pouvons alors réécrire les définitions (16) et (18) en les généralisant pour qu'elles s'appliquent à des systèmes de mesure non bornés, ayant un nombre d'intervalles indéterminé: la distribution normale en est un exemple illustre. Les définitions généralisées sont simplement

(16')
$$C = 1 - \sum_{-\infty}^{\infty} p_j^2,$$

et

(18')
$$k^* = (1 - C)^{-1}.$$

Le nombre d'intervalles efficaces à catégorisation correcte

La discussion de la section précédente, portant sur le nombre k^* d'intervalles efficaces, concernait un système de mesure idéal, à fidélité parfaite: dans un tel système, les intervalles de mesure, ou l'unité de mesure, peuvent être définis *ad libitum*, d'une grandeur quelconque, et les objets mesurés y seront toujours catégorisés correctement. Toutefois, les mesures sont, en pratique, rarement parfaites, et le coefficient de fidélité ρ_{XX} qui les caractérise est habituellement inférieur à 1.

Ainsi, pour que les mesures découlant d'un procédé de mesure à fidélité ρ_{XX} donnent lieu à une catégorisation « fiable », il faut tenir compte de l'erreur de mesure comprise dans chaque évaluation. Nous avons montré plus haut qu'il est possible de déterminer une grandeur d'intervalle, $L_X(\gamma) = \lambda(\gamma)\sigma_X$, telle que la probabilité qu'un objet soit catégorisé dans son intervalle propre soit d'au moins γ . Nous pouvons alors découper d'une manière ou d'une autre l'axe des X en une séquence d'intervalles limitrophes de grandeur $L_X(\gamma)$. Reprenant la distribution hypothétique des X déjà mentionnée, nous pouvons trouver par (13) la probabilité d'occurrence de chaque intervalle j, puis calculer (16') et enfin (18'); ce dernier calcul fournit l'indice k^* - le nombre de catégories efficaces caractérisant notre procédé de mesure X -,

lequel est explicitement une fonction de γ (la probabilité de catégorisation correcte) et de ρ_{XX} (la fidélité du procédé de mesure).

L'invocation du modèle normal et la capacité discriminante

Les idées et formules élaborées plus haut permettent de déterminer k^* pour chaque procédé de mesure particulier, selon sa fidélité ρ_{XX} , en fonction de la distribution, empirique ou théorique, des mesures X_i dans la population, et pour une probabilité γ donnée de catégorisation correcte. Cependant, dans le but de parvenir à un concept mieux cadré et plus net, nous proposons deux simplifications importantes, l'une sur la distribution des X, l'autre sur la probabilité γ .

En premier lieu, il apparaît commode d'adopter un modèle de référence pour la distribution des mesures X dans la population, modèle qui tient lieu de la distribution réelle, obvie à son estimation et, en quelque sorte, l'idéalise. Nous faisons allusion, bien sûr, au modèle de la distribution normale, qui sert ordinairement à représenter la répartition des caractères comme l'intelligence (ou QI), les grandeurs en anthropométrie et bien d'autres caractéristiques dans la population⁵.

En second lieu, le choix d'une valeur particulière de la probabilité de catégorisation correcte, γ , reste arbitraire. Des valeurs telles que $\gamma \geq \frac{1}{2}$ paraissant souhaitables, nous retenons $\gamma = \frac{1}{2}$ comme probabilité de référence. Toute autre valeur dans l'intervalle ($\frac{1}{2} \leq \gamma < 1$) serait acceptable.

Adoptant donc le modèle de la distribution normale et la probabilité $\gamma = \frac{1}{2}$ pour la catégorisation correcte, l'indice k^* devient une fonction de la fidélité ρ_{XX} seulement, et nous définissons ainsi la capacité discriminante, soit

(19)
$$D_0 = k^*[f(X) \text{ normale, } \gamma = \frac{1}{2}, \rho_{XX}].$$

Pour une valeur donnée de $\rho = \rho_{XX}$, donc de $\sigma_e = \sqrt{(1-\rho)}$ en valeurs standardisées, la grandeur d'intervalle standardisée $\lambda(\frac{1}{2})$ s'obtient par l'équation (12). Nous appliquons alors cet intervalle en segmentant l'axe standardisé X' en intervalles contigus centrés, produisant ainsi le système d'intervalles suivant:

(20)
$$-\infty$$
 ... $\left(-\frac{5}{2}\lambda, -\frac{3}{2}\lambda\right), \left(-\frac{3}{2}\lambda, -\frac{1}{2}\lambda\right), \left(-\frac{1}{2}\lambda, \frac{1}{2}\lambda\right), \left(\frac{1}{2}\lambda, \frac{3}{2}\lambda\right), \left(\frac{3}{2}\lambda, \frac{5}{2}\lambda\right), \dots \infty$.

Les probabilités p_i associées à chaque intervalle j dans la distribution normale peuvent être calculées et leurs carrés sommés en vue de l'indice (16'). Partant de l'intervalle central, qui a la probabilité p la plus grande, et allant vers les intervalles plus extrêmes, à probabilités fortement décroissantes, la

sommation se poursuit jusqu'aux deux infinis en ce sens que la valeur de $\sum p_j^2$ converge et se stabilise rapidement près des ailes de la distribution.

L'inversion par l'équation (18') fournit enfin $D_{\rho} = k^*$ pour la valeur de ρ considérée. Le graphique de la figure 1 montre la relation entre k^* et ρ telle qu'établie à partir du modèle normal et avec la probabilité de catégorisation correcte $\gamma = \frac{1}{2}$. La relation entre k^* et ρ devient presque parfaitement linéaire quand nous utilisons, en abscisses, la fonction $1/\sqrt{(1-\rho)}$, comme en fait foi la figure 2. Remaniant l'équation (11), nous voyons que cette fonction $1/\sqrt{(1-\rho)}$ est égale au quotient de l'écart-type des mesures X sur leur erreurtype, σ_X/σ_E .

L'étude de l'ensemble de données (k^*, ρ) , pour $0 < \rho < 1$, tel que représenté à la figure 2, permet d'établir une équation de régression linéaire approximative, soit

(20)
$$D_{\rho} = k^* [\gamma = \frac{1}{2}] = 2.654 / \sqrt{(1 - \rho_{XX})};$$

le coefficient de détermination (r^2) obtenu pour cette régression est d'environ 0,99996. La relation est globalement linéaire, sauf aux extrémités. Ainsi, à l'extrémité basse, quand $\rho \rightarrow 0$, le nombre minimal d'intervalles à probabilité $\gamma = \frac{1}{2}$ de catégorisation correcte est évidemment 2, violant ainsi la linéarité stricte.

Nous avons étudié la relation (k^*, ρ) pour d'autres valeurs de probabilité γ . À toutes fins utiles, les régressions linéaires obtenues sont:

(21a)
$$k^*[\gamma = 0.75] = 1.576/\sqrt{(1-\rho_{XX})}$$
 $(r^2 = 0.99994);$

(21b)
$$k^*[\gamma = 0.90] = 1.113/\sqrt{(1-\rho_{XX})}$$
 $(r^2 = 0.9998);$

(21c)
$$k^*[\gamma = 0.95] = \max\{1, 0.939/\sqrt{(1-\rho_{xx})}\}\$$
 $(r^2 = 0.9998)$.

Enfin, d'après les équations (11), la forme générale obtenue pour notre indice du pouvoir discriminant d'un instrument est

(22)
$$D_{\rho} \approx C_{\gamma} / \sqrt{(1 - \rho_{XX})}$$
$$= C_{\gamma} \times \frac{\sigma_{\chi}}{\sigma_{\alpha}} ,$$

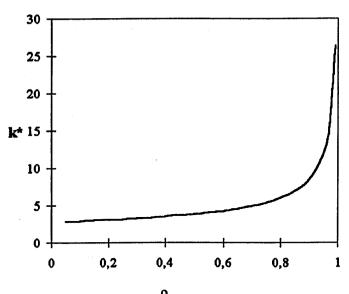


Figure 1 - Relation entre k^* et ρ selon un modèle $X \sim N(\mu_x \sigma^2_x)$

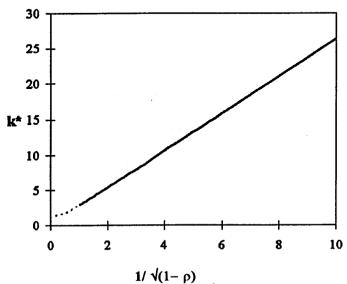


Figure 2 - Relation entre k^* et $1/\sqrt{(1-\rho_{xx})}$ (= σ_x/σ_s) selon un modèle X ~ N (μ_x , σ_x^2)

 C_{γ} étant la constante de pente, fonction de γ et caractérisant l'équation de régression. Cette forme rappelle avec évidence l'indice naïf du départ, D' = E/u, avec au numérateur une quantité relative à la dispersion des mesures sur l'axe X et, au dénominateur, une quantité reflétant la précision, structurelle ou statistique, de l'instrument employé.

La forme (22) ci-haut suggère aussi une analogie plus profonde⁶, cette fois avec la *fonction d'information* d'un test, $I(\theta)$ (Hambleton *et al.*, 1991, p. 94-95). On peut exprimer cette dernière par:

(23)
$$\sqrt{\mathbf{I}(\boldsymbol{\theta})} = \frac{1}{\sigma_{\boldsymbol{e}}(\boldsymbol{\theta})};$$

la valeur moyenne de cette fonction (dans le domaine de θ) est certainement corrélée avec la capacité (22), ne serait-ce qu'en raison de la parenté mathématique des concepts.

Un exemple d'interprétation

Prenons, comme exemple fictif, un test d'intelligence donnant une mesure de QI; le test, de moyenne 100 et écart-type 15, aurait une fidélité $\rho_{XX} = 0.85$. Réécrivant (10) autrement, la grandeur d'intervalle suffisante pour catégoriser une personne dans son intervalle propre avec probabilité γ est

(24)
$$L(\gamma) = 2\sigma_{X}\sqrt{(1-\rho_{XX})}z_{\nu(1+\gamma)}.$$

Appliquant nos valeurs de σ_X , ρ_{XX} et $\gamma = \frac{1}{2}$, nous avons

(25)
$$L(\frac{1}{2}) = 2 \times 15 \times \sqrt{(1 - 0.85)} \times z_{0.75}$$

en utilisant $z_{0.75} \approx 0,6745$. L'unité, conventionnelle et artificielle, du QI est de 1 point; en redéfinissant cette unité (à travers les tableaux de conversion normative) pour qu'elle englobe une suite de 7,84 points originaux, la nouvelle échelle de QI serait telle que le score attribué à une personne serait son score vrai avec probabilité $\gamma = \frac{1}{2}$. De plus, le nombre de catégories efficaces et à probabilité de $\frac{1}{2}$ dont est capable ce procédé de mesure serait de D = $2,654/\sqrt{(1-0.85)} \approx 6,85$; il permet donc de classer la population aussi efficacement que s'il la disposait dans environ 7 catégories de tailles comparables.

En conclusion

Le pouvoir de séparer et de catégoriser les objets selon leurs valeurs est l'une des propriétés fondamentales de la mesure. La « capacité discriminante » que nous proposons ici est une opérationnalisation directe de cette propriété, qui tient compte aussi de l'erreur de mesure telle qu'elle est conçue dans la théorie des tests.

Le développement du concept de capacité discriminante nous a conduit à définir une autre quantité intéressante, caractéristique d'un système de mesure: la largeur d'intervalle suffisante $L(\gamma)$. Dans les mesures en sciences humaines, pour lesquelles l'unité de mesure n'a de réalité que symbolique, la quantité $L(\gamma)$, qui est fonction de la probabilité γ de catégorisation correcte, constitue peut-être un palliatif.

La capacité discriminante, telle que nous la présentons, dénote le pouvoir qu'a un instrument de mesure de répartir les objets mesurés en catégories vraiment distinctes et de capacités (ou tailles) comparables. Le concept, qui exprime une propriété d'un système de mesure X à variation continue, pourrait aussi être généralisé à un système de catégories fermées tel qu'on retrouve dans les systèmes d'observation dits qualitatifs. Si une telle généralisation du concept de capacité discriminante était faite, nous serions enfin en possession d'un concept métrologique commun aux mesures exactes des sciences physiques, aux mesures relatives des sciences humaines et aux mesures dites qualitatives.

NOTES

- La grandeur précise serait plutôt donnée par D = 1 + LE / u J, où Lx J dénote la partie entière de x. Ainsi, pour la règle à mesurer de 1 mètre graduée aux unités u de 0,01 m, D = 1 + L1/u J = 101.
- 2. Elle est stipulée implicitement à la valeur 1 (i.e. le nombre de points, le nombre de réponses correctes), comme dans les échelles de QI, les examens scolaires, etc.
- 3. Une définition constructive de V_i est $[\Sigma X_{i,o}]/n_o \rightarrow V_i$ si $n_o \rightarrow \infty$, l'objet i étant mesuré dans une infinité d'occasions o équivalentes.
- 4. Une démonstration simple de ce résultat est la suivante. Soit $var(p_j)$, la variance entre les p_j , et $k \times var(p_j) = \sum p_j^2 [\sum p_j]^2$. Puisque $\sum p_j = 1$, on a $k \times var(p_j) + 1 = \sum p_j^2$. Or, par définition de variance, $var(p_j) \ge 0$. La valeur minimale de $\sum p_j^2$ correspond à $var(p_j) = 0$; toutes les valeurs p_j étant égales et leur somme égale à 1, on obtient $p_j = 1/k$ pour tout j.

- 5. Nous sommes bien conscient que, à l'opposé, plusieurs caractères mesurables ont une distribution dans la population qui s'écarte sensiblement du modèle normal, tels le revenu des particuliers, le temps d'exécution d'une tâche standardisée, le diamètre des axones.
- 6. L'auteur tient à remercier l'évaluateur anonyme pour cette intéressante suggestion.

RÉFÉRENCES

- Anastasi, A. (1994). <u>Introduction à la psychométrie</u> (traduction par F. Gagné de <u>Psychological testing</u>, 6^e éd., Macmillan). Montréal : Guérin Universitaire.
- Bassière, M. & Gaignebet, E. (1966). <u>Métrologie générale : théorie de la mesure, les instruments et leur emploi</u>. Paris : Dunod.
- Davis, F.B. (1951). Item selection techniques. In E. F. Lindquist (éd.), <u>Educational</u> measurement (pp. 266-328). Washington (D.C.): American Council on Education.
- Ferguson, G.A. (1949). On the theory of test discrimination. <u>Psychometrika</u>, <u>14</u>, 61-68.
- Guilford, J.P. (1954). <u>Psychometric methods</u> (2^e édition). New York: McGraw-Hill.
- Hambleton, R.K., Swaminathan, H. & Rogers, H.J. (1991). <u>Fundamentals of item response theory</u>. London: Sage.
- Henrysson, S. (1971). Gathering, analyzing, and using data on test items. In R. L. Thorndike (éd.), <u>Educational measurement</u> (2^e édition) (pp. 130-159). Washington (D.C.): American Council on Education.
- Laurencelle, L. (sous presse). <u>Théorie et techniques de la mesure instrumentale</u>. Sainte-Foy: Presses de l'Université du Québec.
- Lord, F.M. & Novick, M.R. (1968). <u>Statistical theories of mental test scores</u>. Reading (Mass.): Addison-Wesley.
- Ralston, A. & Reilley, E.D. Jr., éd. (1983). <u>Encyclopedia of computer science and engineering</u> (2^e édition). New York: Van Nostrand.
- Stigler, S.M. (1986). The history of statistics. The measurement of uncertainty before 1900. Cambridge (Mass.): Harvard University Press.
- Thurlow, W.R. (1950). Direct measures of discriminations among individuals performed by psychological tests. <u>Journal of Psychology</u>, 29, 281-314.