





































































































































































































que le score de l'individu à chaque temps ne peut être calculé si l'individu n'a pas répondu à assez d'items. Dans le cas du dropout, les données manquantes ont le même impact sur les deux méthodes car il n'y a pas d'information partielle pour le patient : soit l'information est complète aux temps où le patient a été mesuré soit il n'y a aucune information pour le patient lorsque celui-ci est sorti de l'étude. Dans ce cadre, la propriété d'objectivité spécifique du modèle de Rasch ne permet pas de prendre en compte plus d'information dans l'analyse avec la méthode LRM qu'avec la méthode SM.

## Chapitre 6

# Comparaison de méthodes d'analyse des effets temps et groupe pour données subjectives longitudinales sujettes au dropout

Dans le traitement et le suivi des patients, les études évaluant des Patient-Reported Outcomes ont souvent pour but d'étudier l'évolution d'un PRO au cours du temps dans une population donnée. Il est également fréquent que les études visent à comparer l'évolution d'un PRO au cours du temps dans deux ou plusieurs groupes distincts par une caractéristique clinique comme le traitement reçu ou par une caractéristique socio-démographique. C'est le cas, par exemple, dans l'étude de Durna et al. [27] qui compare la qualité de vie de survivantes de cancer du sein ayant reçu ou non une hormonothérapie substitutive. Dans ce type d'études, outre l'intérêt dans l'évolution du critère au cours du temps, la détection de la présence ou non d'un effet groupe et sa quantification sont également importantes. Les performances des méthodes d'analyse de PRO longitudinales ont également été évaluées dans ce cadre à travers une étude de simulation.

L'étude de simulation suivante vise à évaluer les performances des méthodes SM et LRM dans le cadre de données subjectives longitudinales sujettes au dropout recueillies auprès de deux groupes de patients. Elle a été menée de la même manière que les études de simulation précédentes. Lorsqu'un effet groupe a été simulé, il vaut  $\Delta_\theta = 0,5$  pour la variable latente

et a été approximé à  $\Delta_S = 0,38$  lorsque  $J = 4$  et  $\Delta_S = 0,63$  lorsque  $J = 7$  pour le score (cf. equations 3.8 et 3.9 p.60). L'ensemble des paramètres de l'étude sont rappelés dans le tableau 6.1. La combinaison des différents paramètres nous a amené à considérer 624 cas différents.

TABLE 6.1 – Paramètres utilisés pour la simulation des données et l'analyse

	Nombre de cas	624
	Nombre de jeux simulés/cas	500
	Taille d'échantillon (N)	100 ; 200
	Nombre d'items (J)	4 ; 7
Variable latente	Effet temps ( $d_\theta$ )	0 ; 0,2
	Variance ( $\sigma_\theta^2$ )	1
	Corrélation ( $\rho_\theta$ )	0,4 ; 0,7 ; 0,9
	Effet groupe ( $\Delta_\theta$ )	0 ; 0,5
Score	Effet temps ( $d_S$ )	0 ; 0,15 (J=4) ; 0,25 (J=7)
	Effet groupe ( $\Delta_S$ )	0 ; 0,38 (J=4) ; 0,63 (J=7)
Dropout	Proportion ( $\pi^{(t)}$ )	0% ; 5% ; 10% ; 20%
	Corrélation variable latente	MCAR : 0
	et propension au dropout ( $\rho_{\theta\chi}$ )	MNAR : -0,4 ; -0,7 ; -0,9
Analyse	Méthodes comparées	Score and Mixed Models (SM) Longitudinal Rasch Mixed model (LRM)
	Structures de matrice	sans contraintes (UN) auto-régressive d'ordre 1 (AR(1)) "compound symmetry" hétérogène (CSH)

## 6.1 Résultats

Les structures de matrice de variance covariance UN, AR(1) et CSH ont été étudiées pour la méthode SM. Pour la méthode SM, l'AIC était minimisé plus souvent dans les modèles avec une structure de variance covariance AR(1) que dans les modèles avec une autre structure pour des valeurs simulées du coefficient de corrélation  $\rho_\theta = 0,4$  et  $\rho_\theta = 0,7$ . Lorsque  $\rho_\theta = 0,9$ , l'AIC était plus souvent minimisé par les modèles ayant une structure de variance covariance CSH que par les modèles ayant une autre structure. Les résultats présentés ci-dessous ont été obtenus avec une structure AR(1). La méthode LRM estime une matrice de structure UN.





Au contraire, la puissance décroît lorsque le coefficient de corrélation de la variable latente  $\rho_\theta$  augmente.

La puissance pour les données complètes semble être généralement supérieure à la puissance des données sujettes au dropout quelle que soit la méthode d'analyse utilisée. Comme attendu, le dropout semble entraîner une perte de puissance. Les résultats sont comparables qu'un effet temps ait été simulé ou non.

### 6.1.2 Estimation de l'effet groupe

Le tableau 6.4 présente l'estimation de l'effet groupe et de son erreur standard pour chaque méthode étudiée et pour certaines valeurs de la taille d'échantillon, du nombre d'items, de la corrélation de la variable latente et du type de dropout. Les résultats ont été obtenus à partir de données simulées sans ( $d_\theta = 0$ ) ou avec un effet temps ( $d_\theta = 0.2$ ) et sans ( $\Delta_\theta = 0$ ) ou avec ( $\Delta_\theta = 0,5$ ) effet groupe. Les résultats complets pour l'ensemble des valeurs de paramètres simulées sont présentés en annexe dans les tableaux D.5 p.176, D.6 p.178, D.7 p.180 et D.8 p.182. La plupart des estimations de l'effet groupe sont non-biaisées. Lorsque l'effet temps simulé est nul, la plupart des cas biaisés l'étaient pour les deux méthodes.

En revanche, l'estimation de l'effet groupe a tendance à être plus souvent biaisée pour la méthode SM que pour la méthode LRM lorsque l'effet temps simulé vaut  $d_\theta = 0.2$ . Sous cette hypothèse, la méthode SM semble sous-estimer l'effet groupe alors que la méthode LRM semble le sur-estimer. L'augmentation de la taille d'échantillon, du nombre d'items et du coefficient de corrélation de la variable latente ne semble pas avoir d'impact sur les estimations de l'effet groupe.













pouvoir être utilisées sur des données sujettes à du dropout non-ignorable.

Les deux méthodes ont également montré des résultats comparables en terme de risque de première espèce de l'effet groupe. Le risque est maintenu proche de 5% mais a tendance à augmenter en présence de dropout. Le dropout entraîne une perte de puissance quel que soit le type et la quantité de dropout. Comme attendu, la puissance de l'effet groupe augmente avec la taille d'échantillon et le nombre d'items. Au contraire, la puissance diminue avec l'augmentation de la corrélation entre les mesures de la variable latente. Une augmentation de l'erreur standard de l'estimation de l'effet groupe avec la corrélation de la variable latente a été observée alors que la corrélation ne semble pas avoir d'impact sur les estimations de l'effet groupe. Cette augmentation de l'erreur standard peut expliquer la diminution de la puissance en raison de l'expression du test de l'effet groupe. En effet, le test de Wald utilisé s'exprime comme le ratio de l'estimation de l'effet groupe sur son erreur standard.

L'estimation de l'effet groupe est non-biaisée dans la plupart des cas. Lorsqu'un biais a été observé, celui-ci est faible, de l'ordre de 0,02. En l'absence d'effet groupe, les deux méthodes montrent autant de cas biaisés l'une que l'autre. En revanche, lorsqu'un effet groupe a été simulé et lorsque les données sont sujettes à du dropout MNAR, il semble que la méthode SM présente plus souvent des cas biaisés que la méthode LRM. L'effet groupe est alors sous-estimé pour la méthode SM. Pour la simulation du dropout MNAR, la corrélation entre la variable latente  $\theta$  et la propension à être en dropout  $\chi$  est négative. Ainsi, les patients dont le niveau du PRO étudié est le plus bas ont plus de risque de sortir de l'étude que les autres. On peut donc s'attendre à une sur-estimation de plus en plus importante de la moyenne du PRO étudié au fur et à mesure des temps d'évaluation. En présence d'un effet groupe, il semble que la sur-estimation de la moyenne est plus importante pour le groupe de patients dont le niveau moyen du PRO était plus bas (groupe 0 de la figure 3.1 p.54). Ceci pourrait expliquer la sous-estimation de l'effet groupe car il est estimé par la différence de moyennes entre les deux groupes.

Seul un faible impact du dropout a été observé sur l'effet groupe quel que soit le type et la quantité de données manquantes alors que les résultats en terme d'effet temps sont fortement impactés par la présence de dropout, surtout celui de type MNAR. L'effet groupe simulé valait  $\Delta_\theta = 0,5$  alors que l'effet temps simulé entre deux temps consécutifs valait  $d_\theta = 0,2$ . On peut donc s'interroger sur l'impact de la taille d'effet sur les résultats. Les

bonnes performances des méthodes SM et LRM en terme d'effet groupe sont peut-être dues à la taille de l'effet groupe et pourraient être moins bonnes avec une taille d'effet plus petite et donc plus difficile à détecter.

Les deux groupes de patients ont été simulés en faisant l'hypothèse que les deux groupes avaient des niveaux moyens de PRO différents au premier temps d'étude. Cette hypothèse est réaliste dans des études où les patients ne sont pas randomisés. En revanche, dans les essais cliniques, il est généralement requis que les groupes de patients soient comparables à l'inclusion. Il est alors souvent attendu que l'évolution du PRO dans le temps soit différente d'un groupe à l'autre. La présence d'une interaction est alors possible et, dans ce cas, les deux méthodes présentées doivent être adaptées en ajoutant un terme d'interaction entre le temps et le groupe pour pouvoir analyser de type d'étude.



# Chapitre 7

## Application à deux études de qualité de vie

Les méthodes comparées dans les études de simulation ont été appliquées sur des données réelles. Ces applications visent à mettre en œuvre et à comparer ces méthodes dans deux cas de figure. Dans le premier cas, la dimension étudiée est issue d'un questionnaire de qualité de vie générique et est formée d'items dichotomiques. Cette application se place dans un cadre proche de celui des simulations. Au contraire, la deuxième application se place dans un cadre plus étendu. Les dimensions étudiées sont issues d'un questionnaire de qualité de vie spécifique aux patients atteints de pathologie cancéreuse et composées d'items polytomiques. La méthode LRM, basée sur l'IRT, a donc dû être adaptée par l'utilisation du rating-scale model, extension aux données polytomiques du modèle de Rasch. Les données de ces deux applications sont sujettes à du dropout.

### 7.1 Hyperparathyroïdie et SF-36

#### 7.1.1 Symptômes non spécifiques et qualité de vie dans l'hyperparathyroïdie primaire modérée

L'hyperparathyroïdie primaire est due à une hypersécrétion de parathormone qui entraîne un excès de calcium dans le sang. Il n'existe pas de traitement médical de l'hyperparathyroïdie primaire. Le seul traitement définitif est un traitement chirurgical. La maladie est souvent découverte avant même l'apparition des signes cliniques dits "classiques". Des symptômes non















Le modèle complet comprend des effets temps, des effets groupes ainsi que les interactions d'ordre 1 temps x groupe comme effets fixes, une constante et une pente aléatoires et une structure de variance covariance pour les effets aléatoires ou variance covariance résiduelle de type UN, AR(1), ARH(1), CS, CSH ou identité. Parmi l'ensemble des modèles, le modèle avec le plus faible AIC a été retenu pour chaque dimension. Les effets fixes du modèle retenu ont ensuite été réduits lorsque le test de l'interaction temps x groupe ne rejetait pas l'hypothèse de la nullité du paramètre associé à l'interaction au seuil de 5%.

**Estimation et test des effets** Les effet temps et groupe ont été estimés comme précédemment par la différence des moyennes estimées. Les tailles d'effet (effect size - ES) de l'effet temps sont définies comme le ratio de l'effet temps entre deux temps donnés  $t$  et  $t'$  sur la racine de la moyenne de la variance de ces deux temps :  $\widehat{ES}_{tt'} = \frac{\hat{d}_{tt'}}{\sqrt{\frac{n_t \hat{\sigma}_t^2 + n_{t'} \hat{\sigma}_{t'}^2}{n_t + n_{t'}}}}$  avec  $n_t$  le nombre d'individus au temps  $t$ .

La nullité des effets et des tailles d'effet a été testée avec un test de Student. Le test global de l'effet temps ou de l'effet groupe utilise le test de Fisher.

### 7.2.3 Dimension "fonctionnement physique" du QLQ-C30

Le tableau 7.2 présente les résultats de l'analyse de la dimension "fonctionnement physique" avec les méthodes SM et LRSM. La structure de variance covariance retenue pour la méthode LRSM est de type CS. Le modèle retenu pour la méthode SM comprend pour les effets fixes : les effets temps et les effets groupes et une structure de variance covariance de type CSH. Il ne comprend pas d'interaction temps x groupe ni d'effets aléatoires. Les méthodes SM et LRSM montrent des résultats comparables. Elles concluent toutes deux à la présence d'un effet temps global. L'effet temps entre les temps 1 et 2 est négatif et l'effet temps entre le temps 2 et 3 est positif. Le "fonctionnement physique" s'est dégradé entre le diagnostic ( $t_1$ ) et la fin des traitements ( $t_2$ ). Il s'est ensuite amélioré entre la fin des traitements ( $t_2$ ) et le sixième mois après les traitements ( $t_3$ ). Néanmoins, l'augmentation sur la deuxième période est de taille plus faible que la diminution de la première période entraînant une diminution globale du niveau de "fonctionnement physique" sur l'ensemble de la période de l'étude. On retrouve les mêmes tendances pour les tailles d'effet estimées. Les tailles d'effet estimées sur le score (méthode SM) sont plus petites que pour la variable latente (méthode LRSM). Les deux méthodes concluent à l'absence d'un effet groupe global. De faibles différences entre







des symptômes sur la dimension “fatigue” sur l’ensemble de la période de l’étude. On retrouve les mêmes tendances pour les tailles d’effet estimées. Les tailles d’effet estimées par les deux méthodes sont proches. Les deux méthodes concluent à l’absence d’un effet groupe global. Il existe un effet groupe significatif entre le groupe 1 et le groupe 2 au seuil de 5%. La différence entre le groupe 2 et le groupe 3 est également importante mais n’est pas significative. Les groupes 1 et 3 sont proches en terme de symptômes sur la dimension “fatigue”. Le groupe 2 présente plus de symptômes que les groupes 1 et 3 sur la dimension “fatigue”.

### 7.3 Discussion

Ces deux applications ont permis d’utiliser les différentes méthodes d’analyse comparées dans les études de simulation sur des données réelles. Dans l’étude sur l’hyperparathyroïdie, la dimension étudiée était composée d’items dichotomiques et les méthodes ont pu être utilisées sans modifications. Dans l’étude sur le cancer du sein, les dimensions du QLQ C-30 sont composées d’items polytomiques et la méthode LRM basée sur le modèle de Rasch a du être adaptée en développant la méthode LRSM basée sur le rating-scale model, modèle de la famille de Rasch pour les items polytomiques. Cette extension du modèle se trouve confrontée au problème de l’augmentation du nombre de paramètres à estimer. Avec la méthode LRSM, l’introduction d’une interaction temps x groupe ne menait pas à la convergence du modèle quelle que soit la structure de variance covariance choisie. La pertinence de l’ajout d’une interaction dans le modèle n’a donc pas pu être testée. La méthode SM a également été confrontée à des problèmes de convergence des modèles lorsqu’ils contenaient beaucoup d’effets différents. Ce problème risque fort d’être fréquent en santé où les tailles d’échantillon sont plus petites qu’en sciences de l’éducation, discipline dans laquelle les modèles IRT ont initialement été développés. Il convient donc de privilégier la parcimonie dans la construction des modèles ce qui est confirmé par les modèles finalement retenus dans l’étude sur le cancer du sein.

Toute étude de simulation est confrontée au choix des paramètres de simulation. Ces paramètres doivent être le plus proche possible des valeurs fréquemment rencontrées en pratique. Dans les deux applications, les valeurs estimées du coefficient de corrélation pour la variable latente sont proches des valeurs simulées. En revanche, les estimations des variances sont plus élevées que les valeurs simulées et semblent varier dans le temps. De plus, les

structures de variance covariance retenues sont éloignées de la structure AR(1) utilisée pour la simulation des données. Il est envisageable d'étendre l'étude de simulation pour étudier l'impact de la valeur de la variance sur les résultats des différentes méthodes.

Par ailleurs, l'effet temps simulé augmente linéairement dans le temps. Ce n'est pas le cas dans les deux applications présentées. L'hypothèse faite pour les simulations peut se révéler inadéquate dans beaucoup d'études. Il est fréquent que l'évolution d'un PRO soit évaluée avant la phase de traitement des patients puis plusieurs fois après les traitements. Dans ce cas, on peut s'attendre à une dégradation du niveau du PRO sur la première période puis à une stabilité ou une amélioration sur les périodes post-traitement. Toutefois, les méthodes étudiées estiment l'effet temps sans faire l'hypothèse de linéarité de l'effet temps et peuvent donc être utilisées quelle que soit sa forme.

Les tailles d'effet estimées dans les applications étaient généralement plus élevées que la taille d'effet simulée. On peut donc s'attendre à une meilleure puissance des méthodes dans les applications pour des effectifs similaires. Enfin, les modèles utilisées ne prennent pas en compte l'espacement des temps de mesure car le temps est considéré comme une variable discrète. Or, dans les deux applications, les patients ne sont pas évalués à des temps régulièrement espacés. De plus, les temps d'évaluation déterminés dans les protocoles sont approximatifs et les patients peuvent, par exemple, être vus entre le cinquième et le huitième mois après traitement pour une évaluation fixée dans le protocole au sixième mois après traitement. Il semble important de pouvoir prendre en compte cette information en traitant le temps comme une variable continue.

# Discussion générale

Dans ce travail, nous avons comparé différentes méthodes basées sur les deux approches existantes pour l'analyse de Patient-Reported Outcomes : la théorie classique des tests et la théorie de réponse aux items. Le principal objectif était de déterminer la méthode la plus adéquate pour analyser des PRO recueillis de manière longitudinale et issus d'un questionnaire validé avec un modèle de Rasch.

Plusieurs méthodes d'analyse ont été comparées dans des situations différentes à travers diverses études de simulation. La première étude a permis de comparer quatre méthodes d'analyse, une basée sur la CTT et trois sur l'IRT, dans le cadre de données complètes et un groupe de patients. Dans les deuxième et troisième études, une méthode basée sur la CTT et une méthode basée sur l'IRT ont été comparées dans le cadre de données complètes ou sujettes à du dropout de type MCAR ou MNAR et portant sur un ou deux groupes de patients.

## Méthodes d'analyse à privilégier

Parmi les méthodes basées sur l'IRT, au vu des résultats de la première étude de simulation, seule la méthode utilisant un modèle IRT longitudinal (Longitudinal Rasch Mixed model) a montré des résultats satisfaisants. Les méthodes estimant l'effet temps à partir d'estimations des valeurs individuelles du trait latent (Rasch and Mixed models et Plausible Values) se sont révélées inappropriées pour l'analyse des données dans le cadre de cette étude. En effet, ces méthodes présentaient des puissances peu élevées ainsi que des effets temps sous-estimés au contraire des méthodes SM et LRM.

Dans l'ensemble des études de simulation, les méthodes Longitudinal Rasch Mixed model, basée sur l'IRT, et Score and Mixed models, basée sur la CTT, ont montré des résultats similaires. Il semble donc que l'une ou l'autre approche puisse être utilisée à travers ces

méthodes pour l'analyse de PRO longitudinaux. En pratique, il est probable que la CTT sera privilégiée. La méthode SM est, en effet, simple à mettre en œuvre puisqu'elle repose sur un score observé analysé par un modèle linéaire mixte. De plus, l'interprétation des résultats en terme de score est plus simple et plus parlante pour les cliniciens que l'interprétation en terme de variable latente. Des travaux restent encore à faire dans ce domaine pour rendre plus accessible l'IRT dans le domaine de la santé tant au niveau des modèles utilisés que celui de l'interprétation en terme notamment de taille d'effet mesurée sur une variable latente. Le choix de l'utilisation d'une méthode basée sur un modèle de la famille de Rasch peut être motivé par leurs propriétés psychométriques : l'exhaustivité du score sur le trait latent, l'objectivité spécifique et surtout l'obtention d'une mesure d'intervalle.

La similitude des performances des deux approches CTT et IRT est retrouvée dans les travaux, toutefois encore peu nombreux, de la littérature. Lawson [55] fut un des premiers à comparer de manière empirique les deux approches. Il a observé sur 3 jeux de données différents que les paramètres d'items et d'individus obtenus avec la CTT et le modèle de Rasch étaient très comparables. Fan [31] a ensuite confirmé les résultats de cette première étude par la comparaison de la CTT avec plusieurs modèles IRT (modèle de Rasch, 2-PLM et 3-PLM) dans une analyse transversale de données réelles de tests de mathématiques et de lecture. Il a conclu que les estimations des paramètres d'items et des paramètres des individus étaient très comparables pour les deux approches, en particulier si un modèle de Rasch était utilisé. De plus, il s'est intéressé à l'invariance des paramètres. Dans le contexte de son étude, les paramètres d'items sont invariants de l'échantillon utilisé que ce soit en IRT ou en CTT. L'étude de Fan portait sur 40 échantillons sélectionnés aléatoirement, 80 échantillons sélectionnés par sexe (40 pour chaque sexe) et 80 échantillons sélectionnés par niveau de capacité (40 de bas niveau et 40 de haut niveau), tous composés de 1000 individus, dans une base de 193 000 individus. Cependant, bien qu'intéressantes, ces études comparent les deux approches uniquement sur l'observation de leur comportement sur des données réelles. MacDonald et Paunonen [65] ont mené une étude de simulation dans laquelle sont étudiées la comparabilité, l'invariance et la précision des paramètres selon l'approche utilisée (CTT, modèle de Rasch ou 2-PLM). Cette étude de simulation confirme les résultats des études de Lawson et Fan. Les paramètres de difficulté d'items et d'individus sont comparables et correctement estimés et les paramètres de difficulté d'items sont également invariants pour les deux approches.









s'attendre à ce que la méthode LRM gère mieux les données manquantes intermittentes que la méthode SM. En effet, lorsque le nombre d'items manquants pour chaque questionnaire est important, le score risque de ne pouvoir être calculé et l'information même partielle sera perdue pour la méthode SM. La méthode LRM pourrait prendre en compte plus d'information que la méthode SM, sans imputation, puisque le modèle de Rasch utilise l'information disponible au niveau des items. De plus, la propriété d'objectivité spécifique pourrait assurer une estimation correcte du trait latent car cette estimation ne dépend pas de l'ensemble d'items utilisé pour la mesure. On peut également noter que, même en cas d'imputation des valeurs manquantes pour le calcul du score, il se peut que plus d'information ne soit perdue pour le score que pour la variable latente. En général, les recommandations pour l'imputation dans les questionnaires de qualité de vie sont d'imputer avec la méthode PMS (imputation de la moyenne des items remplis de l'individu) si au moins la moitié des items ont été remplis. Cette méthode peut ne pas être la plus adéquate. De plus, l'imputation peut biaiser les résultats si la méthode d'imputation est mal choisie.

L'impact de données manquantes intermittentes sur les deux méthodes pourrait donc être différent. Le choix d'imputer ou non et la méthode d'imputation utilisée pourrait aussi avoir un impact sur les performances de la méthode SM.

L'impact de la survenue de données manquantes intermittentes sur les performances des méthodes SM et LRM sera étudié dans la thèse d'Elodie Dumas au sein de l'EA 4275.

## Items polytomiques

L'analyse de plusieurs dimensions du QLQ-C30 a confronté les méthodes SM et LRM à l'analyse d'items polytomiques. En pratique, les échelles couramment employées en santé sont principalement composées d'items polytomiques. Le SF-36 comprend 2 dimensions à items dichotomiques et 7 dimensions à items polytomiques. Toutes les dimensions de la QLQ-C30 et du WHOQOL-BREF sont composées d'items polytomiques. Si l'adaptation des méthodes à l'analyse d'items polytomiques est plutôt aisée, leur utilisation pourrait en revanche se révéler problématique.

Une adaptation naturelle de la méthode LRM, basée sur le modèle de Rasch, est de remplacer ce modèle par une de ses extensions aux items polytomiques, le rating-scale model (RSM) ou le partial credit-model (PCM). Ces modèles comprennent un nombre de paramètres plus





shift. Même si cette méthode est couramment utilisée, plusieurs limites à son utilisation sont connues. Les capacités de mémoire des patients doivent être relativement bonnes pour qu'ils puissent réaliser le "then test". De plus, cette méthode mesure uniquement la composante de réétalonnage du response-shift et n'est adaptée qu'aux études comprenant deux temps de mesure. En dépit de ces inconvénients, cette méthode reste utilisée car elle permet d'évaluer le response-shift au niveau de l'individu et la mesure du "then test" serait plus valide que celle du pretest [11] pour mesurer le changement. Dans le cadre de notre travail, l'utilisation du "then-test" n'est pas envisageable car il n'est adapté qu'à des études à deux temps de mesure.

L'utilisation de méthodes statistiques permettant de détecter le response-shift au moment de l'analyse a également été proposée avec notamment l'analyse factorielle, la détection de DIF avec le modèle de Rasch ou l'analyse de courbes de croissance. Mais, ces méthodes présentent l'inconvénient d'évaluer le response-shift au niveau du groupe et non au niveau de l'individu. L'ensemble des méthodes, qu'elles se situent au niveau de la planification ou de l'analyse, ne permettent de détecter que certaines composantes du response-shift. La question du traitement du response-shift reste ouverte et beaucoup de pistes sont à explorer [11]. Il semble que les méthodes SM et LRM ne puissent pas encore intégrer l'étude du response-shift. Cependant, si cette problématique est nouvelle, il est certain qu'elle sera de plus en plus traitée à l'avenir au vu des nombreux travaux de ces dernières années. Les méthodes présentées dans ce travail devront alors être adaptées en conséquence.









## Annexe A

Article : Comparison of CTT and Rasch-based approaches for the analysis of longitudinal Patient Reported Outcomes

# Comparison of CTT and Rasch-based approaches for the analysis of longitudinal Patient Reported Outcomes

Myriam Blanchin,<sup>a,\*†</sup> Jean-Benoit Hardouin,<sup>a</sup> Tanguy Le Neel,<sup>a</sup>  
Gildas Kubis,<sup>a</sup> Claire Blanchard,<sup>b</sup> Eric Mirallié<sup>b</sup>  
and Véronique Sébille<sup>a</sup>

Health sciences frequently deal with Patient Reported Outcomes (PRO) data for the evaluation of concepts, in particular health-related quality of life, which cannot be directly measured and are often called latent variables. Two approaches are commonly used for the analysis of such data: Classical Test Theory (CTT) and Item Response Theory (IRT). Longitudinal data are often collected to analyze the evolution of an outcome over time. The most adequate strategy to analyze longitudinal latent variables, which can be either based on CTT or IRT models, remains to be identified. This strategy must take into account the latent characteristic of what PROs are intended to measure as well as the specificity of longitudinal designs. A simple and widely used IRT model is the Rasch model. The purpose of our study was to compare CTT and Rasch-based approaches to analyze longitudinal PRO data regarding type I error, power, and time effect estimation bias. Four methods were compared: the Score and Mixed models (SM) method based on the CTT approach, the Rasch and Mixed models (RM), the Plausible Values (PV), and the Longitudinal Rasch model (LRM) methods all based on the Rasch model. All methods have shown comparable results in terms of type I error, all close to 5 per cent. LRM and SM methods presented comparable power and unbiased time effect estimations, whereas RM and PV methods showed low power and biased time effect estimations. This suggests that RM and PV methods should be avoided to analyze longitudinal latent variables. Copyright © 2010 John Wiley & Sons, Ltd.

**Keywords:** Item Response Theory; Classical Test Theory; Patient Reported Outcomes; longitudinal data; simulation study

## 1. Introduction

Patient Reported Outcomes (PRO) data are widely used in health sciences to evaluate concepts, such as health-related quality of life (HRQoL), pain, fatigue, or anxiety [1], which are often referred to as latent variables because they cannot be directly observed from patients. PRO data are evaluated using the answers of patients to items often grouped into several dimensions in a questionnaire. Two approaches are commonly used for the analysis of such data: Classical Test Theory (CTT) and Item Response Theory (IRT). The CTT is an approach based on the computation of a score usually computed as the sum of the item responses. This score is an estimation of a 'true' score assumed to represent the evaluated outcome (e.g. HRQoL). The observed and true scores are assumed to be linked by a linear relation. In IRT, item responses have a central role. The probability to answer to an item is a function (not necessary linear) of the latent variable which represents the evaluated outcome. IRT models are a large family

<sup>a</sup>EA 4275 'Biostatistics, Clinical Research and Subjective Measures in Health Sciences', Faculty of Pharmaceutical Sciences, University of Nantes, Nantes, France

<sup>b</sup>Department of Digestive and Endocrine Surgery/Institut des Maladies de l'Appareil Digestif, CHU Nantes, Faculty of Medicine, University of Nantes, France

\*Correspondence to: Myriam Blanchin, EA 4275 'Biostatistics, Clinical Research and Subjective Measures in Health Sciences', Faculté de Pharmacie—Université de Nantes, 1, rue Gaston Veil—44035 Nantes Cedex 1, France.

†E-mail: myriam.blanchin@univ-nantes.fr























and/or concepts of HRQoL. Three components of response-shift are identified: recalibration (a change in the respondent's internal standard of measurement), reprioritization (a change in the importance of component domains constituting the target construct), and reconceptualization (a redefinition of the target construct) [25]. Specific designs like the then-test have been developed to detect response-shift. They have been first used in the area of educational training interventions and then in the area of quality of life, in particular for cancer patients. Treatments of cancer patients can be harmful for quality of life and it has been shown that these patients succeed in adapting to the adverse effects of the disease and its treatments [26]. Since then, the impact of health state changes on an individual's quality of life has gained increased attention in social and medical clinical research. The response-shift that may occur is now considered in studies on evolution of quality of life but the debate on which method to use to detect the response-shift still continues. Some methods are addressing the problem of the response-shift at the design stage of the study, such as the then-test and the individualized methods. Other methods are statistical methods to address response shift, such as factor analysis, growth curve analysis, and Rasch analysis. Among all these possibilities, the then-test is the most commonly used method to measure response-shift. Different components of the response-shift are detected from a method to another. The major point of development remains the quantification of the response-shift. Whereas each method allows to detect it, only the then-test and the factor analysis give a value of change of quality of life adjusted for response-shift effect. Although this simulation study assumed no response-shift, this subject has gained major concern in longitudinal PRO studies and it will be of interest to study the behavior of CTT and Rasch-based methods when response-shift is present.

Finally, many longitudinal studies are faced with the problem of missing observations. In this case, different approaches are often adopted: complete-case analysis, available-data analysis, and imputation. Because the SM method is based on the score computed by summing item responses, in the presence of missing data, the analysis can only be performed through complete-case analysis or imputation approach. In the LRM method, based on item responses, the analysis can also be performed with available-data approach.

Little and Rubin [27] made distinctions between missing value processes. A missing data process is said to be missing completely at random (MCAR) if the missingness is independent of both unobserved and observed data. Data are missing at random (MAR) if the missingness is independent of the unobserved measurements, conditional on the observed data. Otherwise, the missing data process is missing not at random (MNAR). Likelihood-based analyses that ignore the missing data mechanism lead to valid analyses when the missingness is ignorable (MCAR or MAR) [28]. Selection models and pattern-mixture models [29] were proposed to model nonignorable nonresponse. They are an interesting way of dealing with MNAR missing data process by modeling explicitly the missing data mechanism. These models have to be used with caution because untestable assumptions have to be made on the missing data process for selection models and untestable identifying restrictions are used in pattern-mixture models.

Each approach for handling missing data leads to different results and possible bias. It seems important to study the impact of missing data on the performances of methods to analyze longitudinal latent variables. We suspect that there will be a more important loss of information using a CTT-based method than a Rasch-based method because of the necessity to impute for missing data or to use only complete cases in SM method. In the presence of missing data, we expect that the LRM method will present better results than the SM method as it has been shown in the context of sequential analysis of latent variables [30].

This simulation study is based on the assumption that the data follow a Rasch model. The different results on the performance of the methods will probably be affected to different extents if the Rasch model does not correctly fit the data. In this case, we can expect that the CTT approach will perform better than methods based on the Rasch model.

In conclusion, it has been shown that using either the SM or LRM method give comparable and satisfying results. These two methods are adequate for the analysis of longitudinal PRO data following a Rasch model without missing data.

## Acknowledgements

This work was supported by the Ligue Nationale Contre le Cancer.

## References

1. Gotay CC, Kawamoto CT, Bottomley A, Efficace F. The prognostic significance of patient-reported outcomes in cancer clinical trials. *Journal of Clinical Oncology* 2008; **26**(8):1355–1363.
2. Rasch G. *Probabilistic Models for Some Intelligence and Attainment Tests* (expanded edn). University of Chicago Press: Chicago, 1980.
3. Verbeke G, Molenberghs G. *Linear Mixed Models for Longitudinal Data*. Springer: New York, 2001.
4. Fischer GH, Molenaar IW. *Rasch Models, Foundations, Recent Developments, and Applications*. Springer: New York, 1997.
5. Bjorner JB, Petersen MA, Groenvold M, Aaronson N, Ahlner-Elmqvist M, Arraras JI, Brédart A, Fayers P, Jordhoy M, Sprangers M, Watson M, Young T. Use of item response theory to develop a shortened version of the EORTC QLQ-C30 emotional functioning scale. *Quality of Life Research* 2004; **13**(10):1683–1697.
6. Garcia SF, Cella D, Clauser SB, Flynn KE, Lad T, Lai JS, Reeve BB, Smith AS, Stone AA, Weinfurt K. Standardizing patient-reported outcomes assessment in cancer clinical trials: a patient-reported outcomes measurement information system initiative. *Journal of Clinical Oncology* 2007; **25**(32):5106–5112.
7. Rijmen F, Tuerlinckx F, De Boeck P, Kuppens P. A nonlinear mixed model framework for Item Response Theory. *Psychological Methods* 2003; **8**(2):185–205.
8. Embretson S. A multidimensional latent trait model for measuring learning and change. *Psychometrika* 1991; **56**(3):495–515.
9. Fitzmaurice G, Laird N, Ware SJ. *Applied Longitudinal Analysis*. Wiley: Hoboken, 2004.
10. Fitzmaurice G, Davidian M, Verbeke G, Molenberghs G. *Longitudinal Data Analysis: A Handbook of Modern Statistical Methods*. Chapman & Hall/CRC: London, 2008.
11. Hoijsink H, Boomsma A. On person parameter estimation in the dichotomous Rasch model. In *Rasch Models*, Fischer GH, Molenaar IW (eds). Springer: New York, 1997; 53–68.
12. Rubin DB. *Multiple Imputation for Nonresponse in Surveys*. Wiley-IEEE: New York, 2004.
13. Wu M, Adams RJ. *PISA 2000 Technical Report*. OECD Publications: Paris, 2002.
14. Thomas N. Assessing model sensitivity of the imputation methods used in the national assessment of educational progress. *Journal of Educational and Behavioral Statistics* 2000; **25**(2):351–371.
15. Glas CAW, Geerlings H, van de Laar MAFJ, Taal E. Analysis of longitudinal randomized clinical trials using item response models. *Contemporary Clinical Trials* 2009; **30**(2):158–170.
16. Wu M. The role of plausible values in large-scale surveys. *Studies in Educational Evaluation* 2005; **31**(2–3):114–128.
17. Littell RC, Milliken GA, Stroup WW, Wolfinger R. *SAS System for Mixed Models*. SAS Institute Inc: Cary, NC, 1996.
18. Hardouin J. Rasch analysis: estimation and tests with Rasch test. *Stata Journal* 2007; **7**(1):22–44.
19. Zheng X, Rabe-Hesketh S. Estimating parameters of dichotomous and ordinal item response models with gllamm. *Stata Journal* 2007; **7**(3):313–333.
20. Caillaud C, Sebag F, Mathonnet M, Gibelin H, Brunaud L, Loudot C, Kraimps JL, Hamy A, Bresler L, Charbonel B, Leborgne J, Henry JF, Nguyen JM, Mirallié E. Prospective evaluation of quality of life (SF-36v2) and nonspecific symptoms before and after cure of primary hyperparathyroidism (1-year follow-up). *Surgery* 2007; **141**(2):153–160.
21. Mislavy RJ, Beaton AE, Kaplan B, Sheehan KM. Estimating population characteristics from sparse matrix samples of item responses. *Journal of Educational Measurement* 1992; **29**:133–161.
22. Adams RJ. Reliability as a measurement design effect. *Studies in Educational Evaluation* 2005; **31**(2–3):162–172.
23. Wainer H, Thissen D. Estimating ability with the wrong model. *Journal of Educational and Behavioral Statistics* 1987; **12**(4):339–368.
24. Barclay-Goddard R, Epstein JD, Mayo NE. Response shift: a brief overview and proposed research priorities. *Quality of Life Research* 2009; **18**(3):335–346.
25. Schwartz CE, Sprangers MAG. Methodological approaches for assessing response shift in longitudinal health-related quality-of-life research. *Social Science and Medicine* 1999; **48**(11):1531–1548.
26. Sprangers MAG. Response-shift bias: a challenge to the assessment of patients' quality of life in cancer clinical trials. *Cancer Treatment Reviews* 1996; **22**(Supplement 1):55–62.
27. Little RJA, Rubin DB. *Statistical Analysis with Missing Data*. Wiley: New York, 1987.
28. Troxel AB, Fairclough DL, Curran D, Hahn EA. Statistical analysis of quality of life with missing data in cancer clinical trials. *Statistics in Medicine* 1998; **17**:653–666.
29. Little RJA. Modeling the drop-out mechanism in repeated-measures studies. *Journal of the American Statistical Association* 1995; **90**(431):1112–1121.
30. Sébille V, Hardouin J, Mesbah M. Sequential analysis of latent variables using mixed-effect latent variable models: impact of non-informative and informative missing data. *Statistics in Medicine* 2007; **26**(27):4889–4904.



## Annexe B

Article : Analysis of longitudinal Patient-Reported Outcomes with informative and non-informative dropout : Comparison of CTT and Rasch-based methods.

# Analysis of longitudinal Patient-Reported Outcomes with informative and non-informative dropout: Comparison of CTT and Rasch-based methods.

Myriam Blanchin<sup>1</sup>, Jean-Benoit Hardouin<sup>1</sup>, Tanguy Le Neel<sup>1</sup>, Gildas Kubis<sup>1</sup>  
and Véronique Sébille<sup>1</sup>

<sup>1</sup>EA 4275 “Biostatistics, Clinical Research  
and Subjective Measures in Health Sciences”  
Faculty of Pharmaceutical Sciences, University of Nantes, France  
myriam.blanchin@etu.univ-nantes.fr

## ABSTRACT

*Patient-reported outcomes (PRO) are more and more used in health sciences to evaluate concepts such as health-related quality of life. These outcomes cannot be directly observed and are often referred to a latent variable. Two psychometric theories exist for the analysis of PRO: the classical test theory, the most common used in practice and the item response theory with its most used model, the Rasch model. In many studies, PRO are collected longitudinally in order to study the evolution of the outcome through time. Missing data are frequently encountered in longitudinal studies and can be potentially informative. This study aimed at comparing Classical Test Theory (CTT) and Rasch-based approaches to analyze longitudinal PRO collected from a scale validated with a Rasch model and studying the impact of dropout, informative or not, on both approaches. Data with informative dropout have shown estimation bias and have to be analyzed with more appropriate methods. For complete data and data with non-informative dropout, a method of analysis based on the Rasch model may be preferred for the analysis of longitudinal PRO collected from a scale validated with a Rasch model due to the generally observed slight gain of power and the psychometric properties of the model.*

**Keywords:** Classical Test Theory, Rasch model, longitudinal data, dropout, informative missing data.

**2010 Mathematics Subject Classification:** 92B15, 62P10 .

## 1 Introduction

Patient-reported Outcomes, PRO, have gained major concern in the past years. This type of measures, reported by patients based on their perceptions, includes health-related quality of life (HRQoL), functional well-being, patient satisfaction with treatment,... PRO are broadly used in clinical trials as secondary outcome, especially in chronic diseases like cancer and are sometimes used as primary outcome in some contexts. The evaluated outcome (e.g. HRQoL) is often referred to a latent variable and is measured through the responses of patients to items. The special nature of PRO, which are not directly observable, beg the question of the

analysis of the data. Two psychometric theories exist for the analysis of PRO. The most common used approach is the Classical Test Theory (CTT). In this theory, the observed score, usually computed by summing item responses, is used to estimate the 'true' value of the evaluated outcome. In the second theory, Item Response Theory (IRT), the probability to answer to an item is a function of a latent variable, which represents the evaluated outcome, and item parameters. Among the wide family of IRT models, the Rasch model (Rasch, 1980; Fischer and Molenaar, 1995) is the most commonly used for dichotomous items due to its properties: parameters invariance, specific objectivity and exhaustivity of the score on the latent trait. The Rasch model is now widely used in development and validation of scales in health sciences (Lai et al., 2007; Cella et al., 1996).

Many studies including PRO are studies on chronic illness or follow-up studies after treatment or surgery. Thus, patients are evaluated at different time points to allow the analysis of the evolution of the evaluated PRO. In these longitudinal studies, the measures of each patient are therefore unlikely to be independent such as in cross-sectional studies and the correlation between measurements has to be taken into account in the analysis. One way to deal with correlated data is the use of linear mixed models (Verbeke and Molenberghs, 2000b; Fitzmaurice et al., 2009).

Missing data are frequently encountered in longitudinal studies. For instance, a patient can drop out from the study at a certain time point and so answers to questionnaire are missing for this patient after this time. Intermittent missing data can also occur when some items are not answered in a questionnaire. When the reason for missingness may be related to the evaluated outcome level of the patient, the missing data are said to be informative. Otherwise, they are called ignorable. Missing data, depending on their amount and informativity, may have an impact on the analysis and interpretation of the data. There may be a reduction of the statistical power of the analysis and a bias may be introduced leading to incorrect conclusions.

In practice, when longitudinal data coming from a scale validated with a Rasch model have to be analyzed, many methods can be considered. The researchers tend to use more often the CTT approach, probably more from habit than evidence of suitability. The purpose of this paper is to compare methods either based on CTT or Rasch model to analyze longitudinal latent variables through a simulation study. The impact of missing data in this context has also been studied.

## **2 Methods**

### **2.1 Longitudinal data analysis**

When measures are repeated on the same patients through time, linear mixed models are widely used for the analysis of the data. These models allow to deal with the correlation between measures of longitudinal data by specifying fixed effects (population characteristics), random effects (subject-specific effects) and structure of the variance-covariance matrix

(Verbeke and Molenberghs, 2000b). A general linear mixed model can be written as follows:

$$\begin{aligned}
 Y_i &= X_i\beta + Z_i b_i + e_i, \\
 b_i &\sim N_q(0, D), \\
 e_i &\sim N_{n_i}(0, \Sigma_i), \\
 b_1, \dots, b_q, e_1, \dots, e_N &\text{ independent}, \\
 Y_i &\sim N_{n_i}(X_i\beta, Z_i D Z_i' + \Sigma_i),
 \end{aligned}
 \tag{2.1}$$

where  $Y_i$  is the response vector for patient  $i$ ,  $i = 1, \dots, N$ ,  $N$  is the number of patients,  $n_i$  is the number of observations on patient  $i$ ,  $p$  the number of fixed parameters,  $q$  the number of random parameters,  $\beta$  is a  $(p \times 1)$  vector of fixed effects parameters,  $X_i$  is the  $(n_i \times p)$  design matrix for fixed effects,  $b_i$  is a  $(q \times 1)$  vector of random effects parameters,  $Z_i$  is the  $(n_i \times q)$  design matrix for random effects,  $e_i$  is a  $(n_i \times 1)$  vector of residual components,  $D$  is the  $(q \times q)$  among-unit covariance matrix and  $\Sigma_i$  is the  $(n_i \times n_i)$  within-unit covariance matrix.

Two methods of estimation based on the likelihood can be used to estimate the mean parameters  $\beta$  and variance components  $\omega$  (that contains all variances and covariance parameters found in  $V_i = Z_i D Z_i' + \Sigma_i$ ) of the model: Maximum Likelihood estimation (ML) and Restricted maximum Likelihood estimation (REML). The estimations of variance components obtained with ML are known to be biased for finite samples. The REML estimation is used to correct the bias on ML estimates of variance components (Laird and Ware, 1982). The REML estimators of  $\beta$  and  $\omega$  are found by maximizing the so-called REML likelihood function.

$$\begin{aligned}
 L_{REML}(\beta, \omega) &= \left| \sum_{i=1}^N X_i' V_i^{-1}(\omega) X_i \right|^{-1/2} \\
 &\times \prod_{i=1}^N (2\pi)^{-n_i/2} |V_i(\omega)|^{-1/2} \exp \left\{ -\frac{1}{2} (Y_i - X_i\beta)' V_i^{-1}(\omega) (Y_i - X_i\beta) \right\}
 \end{aligned}
 \tag{2.2}$$

## 2.2 Patient Reported Outcomes analysis

Assume that a questionnaire has  $J$  dichotomous items and that measures are repeated  $T$  times on the  $N$  patients of the study. The response of patient  $i$  ( $i = 1, \dots, N$ ) to an item  $j$  ( $j = 1, \dots, J$ ) at time  $t$  ( $t = 1, \dots, T$ ) is denoted by  $Y_{ij}^{(t)}$ .

### 2.2.1 Classical Test Theory approach

The Classical Test Theory (CTT) approach is based on a score usually computed by summing item responses. This approach was the first one developed in psychometrics and has widely spread in PRO analysis because of its simplicity to use and to interpret. The basic model assumes that the observed score is a linear function of the true score and an error term.

A method called **Score and Mixed Models (SM)** and based on the CTT approach was used for the analysis. In the first step, the score of each patient at each time was computed by summing the  $J$  item responses of the patient  $i$  at time  $t$ . A simple linear mixed model was then applied on the scores to investigate whether a time effect was plausible. The SM method was

carried out as follows:

$$\begin{aligned} S_i^{(t)} &= \sum_{j=1}^J Y_{ij}^{(t)}, \\ S_i &= (S_i^{(1)}, \dots, S_i^{(T)})' = X_i \beta + e_i, \\ e_i &\sim N(0, \Sigma_{S,i}), \\ S_i &\sim N_T(\mu_S, \Sigma_{S,i}), \end{aligned} \tag{2.3}$$

with  $\mu_S = (\mu_S^{(1)} \dots \mu_S^{(T)})' = X_i \beta$ . The mean parameters  $\beta$  and variance components  $\omega$  of the model were estimated using REML estimation in SAS Proc MIXED (Littell et al., 1996).

Three covariance structures  $\Sigma_{S,i}$  are often used with longitudinal data : unstructured, first-order autoregressive and heterogeneous compound symmetry denoted by UN, AR(1) and CSH respectively. The unstructured matrix is the most general possible structure. It is used when no hypothesis can be made on the structure of the covariance matrix but leads to estimate an important number of parameters. The AR(1) structure assumed that variances are constant over time and that correlation decreases when measures get further apart from each other in time. On the contrary, the choice of a CSH structure assumes that the variances are not equal but that the correlation is constant over time. For  $T = 3$ , the three structures of covariance can be written as follows with  $\rho$  denoting the correlation coefficient,  $\sigma_{ij}$  the covariance between latent variables at time  $i$  and time  $j$  ( $i \neq j$ ) and  $\sigma_i^2$  the variance of the latent variable at time  $i$ .

$$\begin{aligned} \Sigma_{S,i} &= \begin{pmatrix} \sigma_1^2 & \sigma_{12} & \sigma_{13} \\ \sigma_{12} & \sigma_2^2 & \sigma_{23} \\ \sigma_{13} & \sigma_{23} & \sigma_3^2 \end{pmatrix} \text{ for UN} & \Sigma_{S,i} &= \begin{pmatrix} \sigma^2 & \sigma^2 \rho & \sigma^2 \rho^2 \\ \sigma^2 \rho & \sigma^2 & \sigma^2 \rho \\ \sigma^2 \rho^2 & \sigma^2 \rho & \sigma^2 \end{pmatrix} \text{ for AR(1)} \\ \Sigma_{S,i} &= \begin{pmatrix} \sigma_1^2 & \sigma_1 \sigma_2 \rho & \sigma_1 \sigma_3 \rho \\ \sigma_1 \sigma_2 \rho & \sigma_2^2 & \sigma_2 \sigma_3 \rho \\ \sigma_1 \sigma_3 \rho & \sigma_2 \sigma_3 \rho & \sigma_3^2 \end{pmatrix} \text{ for CSH} \end{aligned}$$

To compare non-nested models for the covariance, one of the information criteria that have been proposed is the Akaike Information Criteria also noted AIC (Akaike, 1974). This tool for model selection aims at comparing models based on their maximized log-likelihood value, ensuring that the retained models show a good fit of data. In order to select the most parsimonious model, the AIC penalizes models for the use of too many parameters. When the likelihood is estimated using REML, the AIC can be expressed as:

$$AIC = -2\hat{l} + 2c \tag{2.4}$$

where  $\hat{l}$  is REML maximum log-likelihood and  $c$  is the number of covariance parameters. The most parsimonious correct model will be the model with the smallest AIC amongst the models with the same mean structure.

## 2.2.2 Item Response Theory and the Rasch model

Item Response Theory (IRT) emerged recently in instrument development and data analysis in health outcomes measurements due to its potential advantages over CTT such as parameters invariance (Hambleton, 2000). IRT is a family of models that express the probability of a patient's particular response to an item as a function of characteristics of the patient (latent variable  $\theta$ ) and characteristics of the item. The latent variable is the evaluated outcome and is considered as latent because it is not observable and must be inferred from item responses. Most of IRT models assume the unidimensionality of the construct, that is a unique latent variable explains the item responses. Among the unidimensional IRT model, the most commonly used model is the Rasch model (Rasch, 1980) due to its properties: the exhaustivity of the score on the latent trait and the specific objectivity. The exhaustivity refers to the property that the total score of a person is a sufficient statistic for the unknown latent trait. This means that no additional information is needed to estimate the person parameter  $\theta$ . Each total score is associated with only one trait level in the Rasch model whatever the pattern of responses. The property of specific objectivity ensures that the difference between two latent traits does not depend on the set of items used to evaluate these traits. It allows to construct shorter versions of questionnaires or several versions of a same questionnaire to adapt the version to the patient's latest variable level.

The Rasch model expresses the probability of a response  $y$  ( $y = 0$  for a negative response (the most pejorative response) and  $y = 1$  for a positive response) of an individual  $i$  ( $i = 1, \dots, N$ ) to a dichotomous item  $j$  ( $j = 1, \dots, J$ ) as a logistic function of the individual value of the latent trait  $\theta_i$  and one item parameter, its difficulty  $\delta_j$ .

$$P(Y_{ij} = y | \theta_i, \delta_j) = \frac{\exp(y(\theta_i - \delta_j))}{1 + \exp(\theta_i - \delta_j)} \quad (2.5)$$

The Rasch model has been extended to situations where responses of individuals to items are observed at several points in time (Meiser, 2007). A longitudinal form of the Rasch model was used for the analysis in the method called **Longitudinal Rasch Model (LRM)**.

$$P(Y_{ij}^{(t)} = y^{(t)} | \theta_i^{(t)}, \delta_j) = \frac{\exp(y^{(t)}(\theta_i^{(t)} - \delta_j))}{1 + \exp(\theta_i^{(t)} - \delta_j)} \quad (2.6)$$

$$\boldsymbol{\theta}_i = (\theta_i^{(1)}, \dots, \theta_i^{(T)})' \text{ iid } N_T(\boldsymbol{\mu}, \boldsymbol{\Sigma}_i)$$

This model assumes that the item parameters  $\delta_j$  remain constant over time. The change in the latent ability  $\theta$  may be person-specific, that is the speed or direction of the evolution may be different from a person to another. As  $\theta$  is assumed to have a multinormal distribution, this model is of the family of the mixed-effects logistic models. The mean parameters  $\boldsymbol{\mu}$  and covariance parameters  $\boldsymbol{\Sigma}_i$  of the model are estimated using Marginal Maximum Likelihood (MML) estimation method. The marginal likelihood is expressed as

$$L(\delta_1, \dots, \delta_J, \boldsymbol{\mu}, \boldsymbol{\Sigma} | \mathbf{y}) = \prod_{i=1}^N \int_{\mathbb{R}^T} \prod_{t=1}^T \prod_{j=1}^J \frac{\exp(y_{ij}^{(t)}(\theta_i^{(t)} - \delta_j))}{1 + \exp(\theta_i^{(t)} - \delta_j)} G(\boldsymbol{\theta}_i / \boldsymbol{\mu}, \boldsymbol{\Sigma}_i) d\boldsymbol{\theta}_i \quad (2.7)$$



$d_\theta$  is the value of the time effect between two consecutive times,  $d_{1,2} = d_{2,3} = d_\theta$ . For data simulated without time effect,  $d_\theta = 0$ . For data simulated with a time effect,  $d_\theta = 0.2$ . The first-order autoregressive structure adopted for the covariance matrix  $\Sigma$  means that variances are constant with time and that correlation between measures of a same patient decreases with time.

We can expect that some parameters have an impact on the performance of the two methods: datasets with different values for the sample size (N), number of items (J) and correlation of the latent variable ( $\rho_\theta$ ) were simulated. The data were assumed to come from a 4-item scale or a 7-item scale with dichotomous items. The values of difficulty parameters were  $\delta_1 = -1$ ,  $\delta_2 = -0.5$ ,  $\delta_3 = 0.5$ ,  $\delta_4 = 1$  for a 4-item scale and  $\delta_1 = -1.5$ ,  $\delta_2 = -1$ ,  $\delta_3 = -0.5$ ,  $\delta_4 = 0$ ,  $\delta_5 = 0.5$ ,  $\delta_6 = 1$ ,  $\delta_7 = 1.5$  for a 7-item scale. The sample size could be of 100 or 200 individuals. Three different values for the correlation coefficient of the latent trait between two consecutive times  $\rho_\theta$  were used:  $\rho_\theta = 0.4$  (small correlation),  $\rho_\theta = 0.7$ , and  $\rho_\theta = 0.9$  (high correlation).

To simulate the dropout of patients from the study, a latent variable denoted  $\chi$  was defined as the dropout propensity. The probability that a patient drops out from the study at time  $t$  depends on its dropout propensity. The dropout process was simulated using the following model derived from a 4-parameter logistic model (Sijtsma and Hemker, 2000):

$$P(DO_i^{(t)} = 1 | \chi_i^{(t)}, \pi_{min}^{(t)}, \pi_{max}^{(t)}) = \pi_{min}^{(t)} + (\pi_{max}^{(t)} - \pi_{min}^{(t)}) \frac{\exp(\chi_i^{(t)})}{1 + \exp(\chi_i^{(t)})} \tag{2.9}$$

with  $DO_i^{(t)} = 1$  represents the situation where a patient  $i$  drops out from the study at time  $t$ ,  $\pi_{min}^{(t)}$  the minimum individual probability of dropout at time  $t$  and  $\pi_{max}^{(t)}$  the maximum individual probability of dropout at time  $t$ .  $\pi_{min}^{(t)}$  and  $\pi_{max}^{(t)}$  were defined such as the expected proportion of dropout at time  $t$  was  $\pi^{(t)} = \frac{\pi_{min}^{(t)} + \pi_{max}^{(t)}}{2}$ . We assume that data are complete at the first time of evaluation ( $\pi^{(1)} = 0$ ). The dropout of the patients is then linear and  $\pi$  of the remaining patients drop out from the study at each time ( $t = 2, 3$ ).

The dropout propensity  $\chi_i$  has a multinormal distribution with mean vector  $(0\ 0\ 0)'$  and a vari-

ance covariance matrix equals to 
$$\begin{pmatrix} 1 & \rho_{\theta\chi}^2 \rho_\theta & \rho_{\theta\chi}^2 \rho_\theta^2 \\ \rho_{\theta\chi}^2 \rho_\theta & 1 & \rho_{\theta\chi}^2 \rho_\theta \\ \rho_{\theta\chi}^2 \rho_\theta^2 & \rho_{\theta\chi}^2 \rho_\theta & 1 \end{pmatrix}.$$

The correlation between the value of the latent variable at time  $t$ ,  $\theta^{(t)}$ , and the dropout propensity at time  $t$ ,  $\chi^{(t)}$ , denoted  $corr(\theta^{(t)}, \chi^{(t)}) = \rho_{\theta\chi}$  and assumed constant with time, was used to determine the type of missingness of the dropout process. When the value of the latent variable for the outcome does not depend on the dropout propensity ( $\rho_{\theta\chi} = 0$ ), the simulated dropout is of MCAR type following the definition of Little and Rubin. As the two methods are based on likelihood and ignores the dropout mechanism, we can expect that analyses on data with MCAR dropout will be valid. However, it is reasonable to assume that missing data mechanism may often be MNAR in studies including PRO, HRQoL for instance. The simulated dropout is MNAR when  $\rho_{\theta\chi} \neq 0$ . Furthermore, the patients with worse HRQoL, due to disease progression or increase of side effects, may be most likely to dropout from the study than other patients (Troxel, Fairclough, Curran and Hahn, 1998). So, we assume that  $\rho_{\theta\chi} < 0$  to simulate

MNAR dropout.

To study the behaviour of SM and LRM in case of missing data, the proportion of dropout in the simulated datasets could be  $\pi^{(t)} = 0\%$  (complete data), 5%, 10% or 20% ( $t = 2, 3$ ). The correlation between the value of the latent variable  $\theta$  and the dropout propensity were  $\rho_{\theta\chi} = 0$  (MCAR dropout),  $\rho_{\theta\chi} = -0.4; -0.7; -0.9$  (MNAR dropout with increasing informativity). The different values of the parameters led to consider 312 different cases. Five hundred simulated datasets were generated and analyzed for each case.

To compare the two methods, a test for time effect was defined using an approximate Wald test:

$$H_0 : \mu_1 = \mu_2 = \mu_3 = \mu \Leftrightarrow L\boldsymbol{\mu} = 0$$

$$H_1 : \exists i | \mu_i \neq \mu \Leftrightarrow L\boldsymbol{\mu} \neq 0$$

$$L = \begin{pmatrix} -1 & 1 & 0 \\ -1 & 0 & 1 \end{pmatrix}$$

Under  $H_0$ ,  $T_L = (\mathbf{L}\hat{\boldsymbol{\mu}})'(\mathbf{L}\hat{\mathbf{V}}\mathbf{L}')^{-1}\mathbf{L}\hat{\boldsymbol{\mu}}$  has approximately a  $\chi_r^2$  distribution (Verbeke and Molenberghs, 2000a) where  $\hat{\boldsymbol{\mu}}$  is the estimate of  $\boldsymbol{\mu}$ ,  $r$  is the rank of  $L$  and  $\hat{\mathbf{V}}$  is the estimated covariance matrix.

In order to compare the methods to analyze longitudinal PRO data, three criteria were studied: the type I error, the power and the bias of the time effect estimation. The type I error of the tests were classically computed as the proportion of rejection of  $H_0$  under the null hypothesis. Rejection of  $H_0$  was based on a test of simultaneous equality of mean estimations, i.e. the absence of time effect. This criteria allowed the study of the aptitude of the method to avoid falsely detecting a time effect. The power calculation used the same tests but calculated the proportion of rejection of  $H_0$  under the alternative hypothesis. On the opposite, the power allowed to study the aptitude of the method to correctly detect the presence of time effect. A one-sided McNemar's test for paired data (McNemar, 1947) was used to compare the power observed for each method.

$$H_0 : power_{LRM} = power_{SM}$$

$$H_1 : power_{LRM} > power_{SM}$$

The comparison of the estimated time effect to the simulated 'true' time effect gave the bias of time effect and informed about the quality of the parameters estimation. The comparison held for the LRM method as the known true value of time effect had been fixed for the latent variable. For SM method, based on score, the 'true' time effect was not known and was estimated by  $d_S$ , the difference of the computed expected score at each time. For  $t = 2, 3$ ,

$$d_S = E(S_i^{(t)}) - E(S_i^{(t-1)}) \quad (2.10)$$















ing mechanism is ignorable, such as MCAR missing mechanism, if the separability condition is met (Verbeke and Molenberghs, 2000b). Both methods are likelihood-based analysis and ignores the dropout mechanism, we can expect that LRM and SM will lead to a valid analysis such as for MCAR case, provided that the parameters describing the measurement process are functionally independent of the parameters describing the dropout process.

This article focused on one type of missing data: the dropout. An important further development to this study concerns intermittent missing data. The missing data are intermittent if the patient has not answered all the items of the questionnaire. As dropout, intermittent missing data can lead to a loss of power and possible bias. The causes for intermittent missing data are multiple. For example, the patient may have not seen the question and the value is missing completely at random. The item can bother the patient because its content concerns religion, politics, sexual life. Thus, the patient choose not to answer this particular item and the missing data is informative. The informative dropout seems to be linked with the quality of life level of the patient whereas informative intermittent missing data might be related to the characteristics of the item. The ways to deal with intermittent missing data are complete case analysis, available case analysis, imputation. The complete case analysis only includes measurements that are complete in the analysis. The available case analysis uses as much data as possible to take advantage of all available information. Imputation methods can be simple or multiple. The problem that arises with intermittent missing data in CTT is the computation of the score. The most commonly used questionnaires (SF-36, EORTC QLQ-C30) have general guidelines regarding treatment of missing data. Generally, the score can still be computed if the patient has filled in half or more of the items of the scale. The patients with more missing items can't be used in the analysis. The Rasch model uses all available items in the analysis. In studies with a high proportion of intermittent missing data, we can expect that Rasch-based approach perform better than CTT-based approach because Rasch model may use more information than CTT. Furthermore, the property of specific objectivity of the Rasch model may ensure that the latent variable may be estimated consistently even for patients with missing items. We can expect that the occurrence of ignorable intermittent missing data (MCAR and MAR) will lead to valid analysis because both approaches are based on likelihood. Estimation problems will probably be observed for non-ignorable dropout but maybe in different extent for each approach.

In health sciences, longitudinal studies evaluating PRO often include two or more groups of patients in order to compare the evolution of the outcome between the groups. For example, study of a new strategy of treatment for a cancer may compare the long-term quality of life of patients receiving the new strategy and the long-term quality of life of patients who have received the usual strategy. This study has to be extended to compare both approaches in the context of longitudinal PRO collected in different groups of patients.

The purpose of this study was to compare CTT-based and Rasch-based models to determine which approach is the most adequate to analyze longitudinal latent variables when the

data were actually collected from a scale validated with a Rasch model. Data with informative dropout have shown parameter estimation problems. They have to be analyzed with appropriate methods taking into account of the dropout process such as selection models or pattern-mixture models. For complete data and data with non-informative dropout, the method of analysis based on the Rasch model may be preferred than the method based on CTT due to the generally observed slight gain of power and the psychometric properties of the Rasch model.

## Acknowledgment

This work was supported by the Ligue Nationale Contre le Cancer.

## References

- Akaike, H. 1974. A new look at the statistical model identification, *IEEE Transactions on Automatic Control* **19**: 716–723.
- Cella, D. F., Dineen, K., Arnason, B., Reder, A., Webster, K. A., Karabatsos, G., Chang, C., Lloyd, S., Steward, J. and Stefoski, D. 1996. Validation of the functional assessment of multiple sclerosis quality of life instrument, *Neurology* **47**(1): 129–139.
- Diggle, P. and Kenward, M. G. 1994. Informative Drop-Out in longitudinal data analysis, *Applied Statistics* **43**(1): 49–93.
- Embretson, S. E. 1996. The new rules of measurement, *Psychological Assessment* **8**(4): 341–349.
- Fischer, G. H. and Molenaar, I. W. 1995. *Rasch models*, Springer.
- Fitzmaurice, G. M., Davidian, M., Verbeke, G. and Molenberghs, G. 2009. *Longitudinal data analysis*, Chapman and Hall/CRC.
- Hambleton, R. K. 2000. Emergence of item response modeling in instrument development and data analysis, *Medical Care* **38**(9 Suppl II): 60–65.
- Lai, J., Cella, D., Kupst, M. J., Holm, S., Kelly, M. E., Bode, R. K. and Goldman, S. 2007. Measuring fatigue for children with cancer: Development and validation of the pediatric functional assessment of chronic illness Therapy-Fatigue (pedsFACIT-F), *Journal of Pediatric Hematology/Oncology* **29**(7): 471–479.
- Laird, N. M. and Ware, J. H. 1982. Random-Effects models for longitudinal data, *Biometrics* **38**(4): 963–974.
- Littell, R. C., Milliken, G. A., Stroup, W. W. and Wolfinger, R. D. 1996. *SAS system for mixed models*, SAS Institute, Inc., Cary, NC.
- Little, R. J. A. 1995. Modeling the Drop-Out mechanism in Repeated-Measures studies, *Journal of the American Statistical Association* **90**(431): 1112–1121.

- Little, R. J. A. and Rubin, D. B. 2002. *Statistical analysis with missing data*, Wiley.
- McNemar, Q. 1947. Note on the sampling error of the difference between correlated proportions or percentages, *Psychometrika* **12**(2): 153–157.
- Meiser, T. 2007. Rasch models for longitudinal data, in M. von Davier and C. H. Carstensen (eds), *Multivariate and Mixture Distribution Rasch Models*, Springer, New York, pp. 191–199.
- Rasch, G. 1980. *Probabilistic models for some intelligence and attainment tests*, University of Chicago Press.
- Rubin, D. B. 2004. *Multiple imputation for nonresponse in surveys*, Wiley-IEEE.
- Sijtsma, K. and Hemker, B. T. 2000. A taxonomy of IRT models for ordering persons and items using simple sum scores, *Journal of Educational and Behavioral Statistics* **25**(4): 391–415.
- Troxel, A. B., Fairclough, D. L., Curran, D. and Hahn, E. A. 1998. Statistical analysis of quality of life with missing data in cancer clinical trials, *Statistics in Medicine* **17**(5-7): 653–66.
- Verbeke, G. and Molenberghs, G. 2000a. Inference for the marginal model, *Linear Mixed Models for Longitudinal Data*, Springer, New York, pp. 55–76.
- Verbeke, G. and Molenberghs, G. 2000b. *Linear Mixed Models for Longitudinal Data*, corrected edn, Springer.
- Zheng, X. and Rabe-Hesketh, S. 2007. Estimating parameters of dichotomous and ordinal item response models with gllamm, *Stata Journal* **7**(3): 313–333.



## Annexe C

# Comparaison de méthodes d'analyse de données subjectives longitudinales sujettes au dropout : Résultats complets











## Annexe D

# Comparaison de méthodes d'analyse des effets temps et groupe pour données subjectives longitudinales sujettes au dropout : Résultats complets







































# Bibliographie

- [1] N. K. Aaronson, S. Ahmedzai, B. Bergman, M. Bullinger, A. Cull, N. J. Duez, A. Filiberti, H. Flechtner, S. B. Fleishman, and J. C. de Haes. The european organization for research and treatment of cancer QLQ-C30 : a quality-of-life instrument for use in international clinical trials in oncology. *Journal of the National Cancer Institute*, 85(5) :365–376, 1993.
- [2] C. Acquadro, R. Berzon, D. Dubois, N. K. Leidy, P. Marquis, D. Revicki, M. Rothman, and P. H. Group. Incorporating the patient’s perspective into drug development and communication : An ad hoc task force report of the Patient-Reported outcomes (PRO) harmonization group meeting at the food and drug administration, february 16, 2001. *Value in Health*, 6(5) :522–531, 2003.
- [3] H. Akaike. A new look at the statistical model identification. system identification and time-series analysis. *IEEE Transactions on Automatic Control*, 19(6) :716–723, 1974.
- [4] E. B. Andersen. Asymptotic properties of conditional Maximum-Likelihood estimators. *Journal of the Royal Statistical Society. Series B (Methodological)*, 32(2) :283–301, 1970.
- [5] E. B. Andersen. The numerical solution of a set of conditional estimation equations. *Journal of the Royal Statistical Society. Series B (Methodological)*, 34(1) :42–54, 1972.
- [6] E. B. Andersen. A goodness of fit test for the rasch model. *Psychometrika*, 38(1) :123–140, 1973.
- [7] E. B. Andersen. Sufficient statistics and latent trait models. *Psychometrika*, 42 :69–81, 1977.
- [8] E. B. Andersen. Estimating latent correlations between repeated testings. *Psychometrika*, 50(1) :3–16, 1985.

- [9] D. Andrich. A rating formulation for ordered response categories. *Psychometrika*, 43(4) :561–573, Dec. 1978.
- [10] D. Andrich and G. Luo. Conditional pairwise estimation in the rasch model for ordered response categories using principal components. *Journal of Applied Measurement*, 4(3) :205–21, 2003.
- [11] R. Barclay-Goddard, J. D. Epstein, and N. E. Mayo. Response shift : a brief overview and proposed research priorities. *Quality of Life Research*, 18(3) :335–346, 2009.
- [12] A. Birnbaum. Some latent trait models and their use in inferring an examinee’s ability. In *Statistical theories of mental test scores*. F. M. Lord & M. R. Novick, New York, Addison-Wesley edition, 1968.
- [13] M. Blanchin, J. Hardouin, T. L. Neel, G. Kubis, C. Blanchard, E. Mirallié, and V. Sébille. Comparison of CTT and Rasch-based approaches for the analysis of longitudinal Patient Reported Outcomes. *Statistics in Medicine*, 30(8) :825–838, 2011.
- [14] M. Blanchin, J. Hardouin, T. L. Neel, G. Kubis, and V. Sébille. Analysis of longitudinal Patient-Reported Outcomes with informative and non-informative dropout : Comparison of CTT and Rasch-based methods. *International Journal of Applied Mathematics & Statistics*, 24(SI-11A) :107–124, 2011.
- [15] R. D. Bock and M. Aitkin. Marginal maximum likelihood estimation of item parameters : Application of an EM algorithm. *Psychometrika*, 46(4) :443–459, 1981.
- [16] B. C. Brown, S. P. McKenna, M. Solomon, J. Wilburn, D. A. McGrouther, and A. Bayat. The patient-reported impact of scars measure : development and validation. *Plastic and Reconstructive Surgery*, 125(5) :1439–1449, May 2010.
- [17] W. Brown. Some experimental results in the correlation of mental abilities. *British Journal of Psychology*, 3(3) :296–322, 1910.
- [18] C. Caillard, F. Sebag, M. Mathonnet, H. Gibelin, L. Brunaud, C. Loudot, J. Kraimps, A. Hamy, L. Bresler, B. Charbonnel, J. Leborgne, J. Henry, J. Nguyen, and E. Mirallié. Prospective evaluation of quality of life (SF-36v2) and nonspecific symptoms before and after cure of primary hyperparathyroidism (1-year follow-up). *Surgery*, 141(2) :153–160, Feb. 2007.

- [19] C. C. Chen and R. K. Bode. Psychometric validation of the manual ability measure-36 (MAM-36) in patients with neurologic and musculoskeletal disorders. *Archives of Physical Medicine and Rehabilitation*, 91(3) :414–420, Mar. 2010.
- [20] Committee for medicinal products for human use. Reflection paper on the regulatory guidance for the use of healthrelated quality of life (HRQL) measures in the evaluation of medicinal products. Technical Report EMEA/CHMP/EWP/139391/2004, European Medicines Agency, Londres, July 2005.
- [21] F. Cousson, M. Bruchon-Schweitzer, B. Quintard, J. Nuissier, and N. Rasclé. Analyse multidimensionnelle d’une échelle de coping : validation française de la W.C.C. (ways of coping checklist). *Psychologie française*, 41(2) :155–164, 1996.
- [22] L. J. Cronbach. Coefficient alpha and the internal structure of tests. *Psychometrika*, 16(3) :297–334, 1951.
- [23] D. Curran, M. Bacchi, S. F. Schmitz, G. Molenberghs, and R. J. Sylvester. Identifying the types of missingness in quality of life data from clinical trials. *Statistics in Medicine*, 17(5-7) :739–56, 1998.
- [24] D. Curran, G. Molenberghs, P. M. Fayers, and D. Machin. Incomplete quality of life data in randomized trials : missing forms. *Statistics in Medicine*, 17(5-7) :697–709, 1998.
- [25] A. P. Dempster, N. M. Laird, and D. B. Rubin. Maximum likelihood from incomplete data via the em algorithm. *Journal of the Royal Statistical Society. Series B*, 39(1) :1–38, 1977.
- [26] P. Diggle and M. G. Kenward. Informative Drop-Out in longitudinal data analysis. *Applied Statistics*, 43(1) :49–93, 1994.
- [27] E. M. Durna, S. M. Crowe, L. R. Leader, and J. A. Eden. Quality of life of breast cancer survivors : the impact of hormonal replacement therapy. *Climacteric*, 5(3) :266–276, 2002.
- [28] S. E. Embretson. The new rules of measurement. *Psychological Assessment*, 8(4) :341–349, 1996.

- 
- [29] S. E. Embretson and S. P. Reise. The new rules of measurement. In *Item response theory for psychologists*, Multivariate Applications Series. Lawrence Erlbaum Associates Inc, 2000.
- [30] B. Falissard. *Mesurer la subjectivité en santé : Perspective méthodologique et statistique*. Masson, 2001.
- [31] X. Fan. Item response theory and classical test theory : An empirical comparison of their Item/Person statistics. *Educational and Psychological Measurement*, 58(3) :357–381, June 1998.
- [32] P. Fayers, N. K. Aaronson, K. Bjordal, D. Curran, and M. Groenvold on behalf of the EORTC Quality of Life Study Group. *EORTC QLQ-C30 Scoring Manual (Third edition)*. EORTC Quality of Life Group, Brussels, 2001.
- [33] P. M. Fayers, D. Curran, and D. Machin. Incomplete quality of life data in randomized trials : missing items. *Statistics in Medicine*, 17(5-7) :679–96, 1998.
- [34] L. S. Feldt and R. L. Brennan. Reliability. In R. L. Linn, editor, *Educational Measurement*, pages 105–146. Macmillan USA, New York, 3rd edition, 1989.
- [35] G. Fischer. Derivations of the rasch model. In G. H. Fischer and I. W. Molenaar, editors, *Rasch models : foundations, recent developments, and applications*. Springer, 1995.
- [36] G. Fischer and P. Parzer. An extension of the rating scale model with an application to the measurement of change. *Psychometrika*, 56 :637–651, 1991.
- [37] G. Fischer and I. Ponocny. An extension of the partial credit model with an application to the measurement of change. *Psychometrika*, 59(2) :177–192, 1994.
- [38] G. H. Fischer. Logistic latent trait models with linear constraints. *Psychometrika*, 48(1) :3–26, 1983.
- [39] G. H. Fischer. Linear logistic models for change. In I. W. Molenaar and G. Fischer, editors, *Rasch models : foundations, recent developments, and applications*. Springer, New York, Apr. 1995.
- [40] R. A. Fisher. *Statistical methods for research workers*. Oliver and Boyd, Edinburgh, 1925.

- 
- [41] G. M. Fitzmaurice, M. Davidian, G. Verbeke, and G. Molenberghs. *Longitudinal data analysis*. Chapman and Hall/CRC, 2009.
- [42] G. M. Fitzmaurice, N. M. Laird, and J. H. Ware. Estimation and statistical inference. In *Applied longitudinal analysis*, Wiley Series in Probability and Statistics. Wiley-IEEE, Hoboken, 2004.
- [43] C. A. Glas, H. Geerlings, M. A. van de Laar, and E. Taal. Analysis of longitudinal randomized clinical trials using item response models. *Contemporary Clinical Trials*, 30(2) :158–170, 2009.
- [44] C. A. W. Glas and I. Hendrawan. Testing linear models for ability parameters in item response models. *Multivariate Behavioral Research*, 40(1) :25, 2005.
- [45] R. Glynn, N. Laird, and D. Rubin. Selection modelling versus mixture modelling with nonignorable nonresponse. In H. Wainer, editor, *Drawing Inferences from Self-selected Samples*, pages 115–142. New York : Springer-Verlag, 1986.
- [46] C. C. Gotay, C. T. Kawamoto, A. Bottomley, and F. Efficace. The prognostic significance of Patient-Reported outcomes in cancer clinical trials. *J Clin Oncol*, 26(8) :1355–1363, 2008.
- [47] H. Gulliksen. *Theory of Mental Tests*. John Wiley & Sons Inc, New York, Dec. 1950.
- [48] R. K. Hambleton and R. W. Jones. Comparison of classical test theory and item response theory and their application to test development. *Educational Measurement : Issues and Practice*, 12(3) :38–47, 1993.
- [49] R. K. Hambleton and W. J. van der Linden. Advances in item response theory and applications : An introduction. *Applied Psychological Measurement*, 6(4) :373–378, 1982.
- [50] J.-B. Hardouin. Rasch analysis : Estimation and tests with raschtest. *Stata Journal*, 7(1) :22–44(23), 2007.
- [51] D. A. Harville. Maximum likelihood approaches to variance component estimation and to related problems. *Journal of the American Statistical Association*, 72(358) :320–338, 1977.

- [52] R. Holman and C. A. W. Glas. Modelling non-ignorable missing-data mechanisms with item response theory models. *The British Journal of Mathematical and Statistical Psychology*, 58(Pt 1) :1–17, May 2005. PMID : 15969835.
- [53] M. G. Kenward and J. H. Roger. Small sample inference for fixed effects from restricted maximum likelihood. *Biometrics*, 53(3) :983–997, 1997.
- [54] N. M. Laird and J. H. Ware. Random-Effects models for longitudinal data. *Biometrics*, 38(4) :963–974, 1982.
- [55] S. Lawson. One parameter latent trait measurement : Do the results justify the effort ? In B. Thompson, editor, *Advances in Educational Research : Substantive findings, methodological developments*, volume 1, pages 159–168. Greenwich, CT : JAI, 1991.
- [56] A. Leplège, E. Ecosse, A. Verdier, and T. V. Perneger. The french SF-36 health survey : translation, cultural adaptation and preliminary psychometric evaluation. *Journal of Clinical Epidemiology*, 51(11) :1013–23, Nov. 1998.
- [57] K.-Y. Liang and S. L. Zeger. Longitudinal data analysis using generalized linear models. *Biometrika*, 73(1) :13–22, 1986.
- [58] M. J. Lindstrom and D. M. Bates. Newton-raphson and em algorithms for linear mixed-effects models for repeated-measures data. *Journal of the American Statistical Association*, 83(404) :1014–1022, 1988.
- [59] R. C. Littell, G. A. Milliken, W. W. Stroup, and R. D. Wolfinger. *SAS system for mixed models*. SAS Institute, Inc., Cary, NC, 1996.
- [60] R. J. A. Little. Pattern-Mixture models for multivariate incomplete data. *Journal of the American Statistical Association*, 88(421) :125–134, 1993.
- [61] R. J. A. Little. A class of pattern-mixture models for normal incomplete data. *Biometrika*, 81(3) :471–483, 1994.
- [62] R. J. A. Little. Modeling the Drop-Out mechanism in Repeated-Measures studies. *Journal of the American Statistical Association*, 90(431) :1112–1121, 1995.
- [63] R. J. A. Little and D. B. Rubin. *Statistical analysis with missing data*. Wiley, 2002.

- 
- [64] F. M. Lord and M. R. Novick. *Statistical Theories of Mental Test Scores*. Addison-Wesley Publishing Company, Inc., June 1968.
- [65] P. Macdonald and S. V. Paunonen. A monte carlo comparison of item and person statistics based on item response theory versus classical test theory. *Educational and Psychological Measurement*, 62(6) :921–943, Dec. 2002.
- [66] G. N. Masters. A rasch model for partial credit scoring. *Psychometrika*, 47(2) :149–174, 1982.
- [67] C. A. McHorney, J. E. Ware, and A. E. Raczek. The MOS 36-Item Short-Form health survey (SF-36) : II. psychometric and clinical tests of validity in measuring physical and mental health constructs. *Medical Care*, 31(3) :247–263, 1993.
- [68] T. Meiser. Rasch models for longitudinal data. In M. von Davier and C. H. Carstensen, editors, *Multivariate and Mixture Distribution Rasch Models*, Statistics for Social and Behavioral Sciences, pages 191–199. New York, springer edition, 2007.
- [69] B. Michiels, G. Molenberghs, L. Bijmens, T. Vangeneugden, and H. Thijs. Selection models and pattern-mixture models to analyse longitudinal quality of life data subject to drop-out. *Statistics in Medicine*, 21(8) :1023–41, 2002.
- [70] P. Minini and M. Chavance. Sensitivity analysis of longitudinal normal data with drop-outs. *Statistics in Medicine*, 23(7) :1039–1054, 2004.
- [71] R. J. Mislevy, A. E. Beaton, B. Kaplan, and K. M. Sheehan. Estimating population characteristics from sparse matrix samples of item responses. *Journal of Educational Measurement*, 29(2) :133–161, 1992.
- [72] R. J. Mislevy, E. G. Johnson, and E. Muraki. Scaling procedures in NAEP. *Journal of Educational Statistics*, 17(2) :131–154, 1992.
- [73] I. W. Molenaar. Estimation of item parameter. In G. H. Fischer and I. W. Molenaar, editors, *Rasch models : foundations, recent developments, and applications*. Springer, Apr. 1995.
- [74] G. Molenberghs, C. Beunckens, C. Sotito, and M. G. Kenward. Every missingness not at random model has a missingness at random counterpart with equal fit. *Journal of the Royal Statistical Society : Series B (Statistical Methodology)*, 70(2) :371–388, 2008.

- 
- [75] G. Molenberghs and M. G. Kenward. A perspective on simple methods. In *Missing data in clinical studies*, pages 41–54. John Wiley and Sons, 2007.
- [76] G. Molenberghs, M. G. Kenward, and E. Lesaffre. The analysis of longitudinal ordinal data with nonrandom drop-out. *Biometrika*, 84(1) :33–44, 1997.
- [77] G. Molenberghs, B. Michiels, M. G. Kenward, and P. J. Diggle. Monotone missing data and pattern-mixture models. *Statistica Neerlandica*, 52(2) :153–161, 1998.
- [78] C. Monseur and R. Adams. Plausible values : how to deal with their limitations. *Journal of Applied Measurement*, 10(3) :320–334, 2009.
- [79] E. Muraki. A generalized partial credit model : Application of an EM algorithm. *Applied Psychological Measurement*, 16(2) :159–176, June 1992.
- [80] J. M. Norquist, R. Fitzpatrick, J. Dawson, and C. Jenkinson. Comparing alternative rasch-based methods vs raw scores in measuring change in health. *Medical Care*, 42(1 Suppl) :I25–36, Jan. 2004.
- [81] J. C. Nunnally. *Psychometric Theory*. Mcgraw-Hill College, 2nd edition, 1978.
- [82] W. H. O. Q. of Life Assessment Group. What quality of life ? the WHOQOL group. *World Health Forum*, 17(4) :354–356, 1996.
- [83] H. G. Osburn. Coefficient alpha and related internal consistency reliability coefficients. *Psychological Methods*, 5(3) :343–355, 2000.
- [84] G. R. Parkerson, W. E. Broadhead, and C. K. Tse. The duke health profile. a 17-item measure of health and dysfunction. *Medical Care*, 28(11) :1056–1072, 1990.
- [85] H. D. Patterson and R. Thompson. Recovery of inter-block information when block sizes are unequal. *Biometrika*, 58(3) :545–554, 1971.
- [86] J. Pfanzagl. On item parameter estimation in certain latent trait models. In G. H. Fischer and D. Laming, editors, *Contributions to mathematical psychology, psychometrics, and methodology*, pages 249–263. Springer, New York, 1994.
- [87] L. Prieto, J. Alonso, and R. Lamarca. Classical test theory versus rasch analysis for quality of life questionnaire reduction. 1 :27, 2003.

- 
- [88] S. Rabe-Hesketh, A. Skrondal, and A. Pickles. GLLAMM manual. *U.C. Berkeley Division of Biostatistics Working Paper Series*, (Working Paper 160), 2004.
- [89] G. Rasch. On specific objectivity : An attempt at formalizing the request for generality and validity of scientific statements. *Danish Yearbook of Philosophy*, 14 :58–94, 1977.
- [90] G. Rasch. *Probabilistic models for some intelligence and attainment tests*. University of Chicago Press, 1980.
- [91] D. B. Rubin. Inference and missing data. *Biometrika*, 63(3) :581–592, 1976.
- [92] D. B. Rubin. *Multiple imputation for nonresponse in surveys*. Wiley-IEEE, 2004.
- [93] F. Samejima. Estimation of latent ability using a response pattern of graded scores. *Psychometrika Monograph*, (17), 1969.
- [94] G. Saporta. *Probabilités, analyses des données et statistiques*. Editions Technip, Paris, 1990.
- [95] F. Satterthwaite. Synthesis of variance. *Psychometrika*, 6 :309–316, 1941.
- [96] C. E. Schwartz and M. A. Sprangers. Methodological approaches for assessing response shift in longitudinal health-related quality-of-life research. *Social Science & Medicine* (1982), 48(11) :1531–1548, June 1999. PMID : 10400255.
- [97] G. Schwarz. Estimating the dimension of a model. *The Annals of Statistics*, 6(2) :461–464, 1978.
- [98] K. Sijtsma and B. T. Hemker. A taxonomy of IRT models for ordering persons and items using simple sum scores. *Journal of Educational and Behavioral Statistics*, 25(4) :391–415, 2000.
- [99] J. A. Sloan, M. Y. Halyard, M. H. Frost, A. C. Dueck, B. Teschendorf, M. L. Rothman, and the Mayo/FDA Patient-Reported Outcomes Consensus Meeting Group. The mayo clinic manuscript series relative to the discussion, dissemination, and operationalization of the food and drug administration guidance on Patient-Reported outcomes. *Value in Health*, 10 :S59–S63, 2007.
- [100] A. C. Smidt, J. Lai, D. Cella, S. Patel, A. J. Mancini, and S. L. Chamlin. Development and validation of Skindex-Teen, a quality-of-life instrument for adolescents with skin disease. *Archives of Dermatology*, 146(8) :865–869, Aug. 2010.

- 
- [101] C. Spearman. The proof and measurement of association between two things. *The American Journal of Psychology*, 15(1) :72–101, 1904.
- [102] C. Spearman. Demonstration of formulae for true measurement of correlation. *The American Journal of Psychology*, 18(2) :161–169, Apr. 1907.
- [103] M. A. Sprangers and C. E. Schwartz. Integrating response shift into health-related quality of life research : a theoretical model. *Social Science & Medicine (1982)*, 48(11) :1507–1515, June 1999. PMID : 10400253.
- [104] M. A. G. Sprangers. Response-shift bias : a challenge to the assessment of patients' quality of life in cancer clinical trials. *Cancer Treatment Reviews*, 22(Supplement 1) :55–62, Jan. 1996.
- [105] S. Stevens. Mathematics, measurement and psychophysics. In *Handbook of Experimental Psychology*. New York, 1951.
- [106] S. S. Stevens. On the theory of scales of measurement. *Science*, 103(2684) :677–680, 1946.
- [107] The WHOQOL Group. Development of the world health organization WHOQOL-BREF quality of life assessment. *Psychological Medicine*, 28(3) :551–558, 1998.
- [108] N. Thomas. Assessing model sensitivity of the imputation methods used in the national assessment of educational progress. *Journal of Educational and Behavioral Statistics*, 25(4) :351–371, 2000.
- [109] L. Thurstone. *The Reliability and Validity of Tests*. Edwards Brothers, Ann Arbor, MI, 1931.
- [110] A. B. Troxel, D. L. Fairclough, D. Curran, and E. A. Hahn. Statistical analysis of quality of life with missing data in cancer clinical trials. *Statistics in Medicine*, 17(5-7) :653–666, 1998.
- [111] R. Turner and R. J. Adams. The programme for international student assessment : an overview. *Journal of Applied Measurement*, 8(3) :237–248, 2007.
- [112] U.S. Department of Health and Human Services FDA Center for Drug Evaluation and Research, U.S. Department of Health and Human Services FDA Center for Biologics

- Evaluation and Research, and U.S. Department of Health and Human Services FDA Center for Devices and Radiological Health. Guidance for industry : patient-reported outcome measures : use in medical product development to support labeling claims : draft guidance. *Health and Quality of Life Outcomes*, 4 :79–79, 2006.
- [113] G. Verbeke and E. Lesaffre. The effect of misspecifying the random-effects distribution in linear mixed models for longitudinal data. *Computational Statistics & Data Analysis*, 23(4) :541–556, 1997.
- [114] G. Verbeke and G. Molenberghs. Estimation of the marginal model. In *Linear Mixed Models for Longitudinal Data*, Springer Series in Statistics, pages 41–54. Springer, New York, 2000.
- [115] G. Verbeke and G. Molenberghs. Inference for the marginal model. In *Linear Mixed Models for Longitudinal Data*, Springer Series in Statistics, pages 55–76. Springer, New York, 2000.
- [116] G. Verbeke and G. Molenberghs. Joint modeling of measurements and missingness. In *Linear Mixed Models for Longitudinal Data*, Springer Series in Statistics, pages 209–219. Springer, New York, 2000.
- [117] G. Verbeke and G. Molenberghs. *Linear Mixed Models for Longitudinal Data*. Springer, Berlin, 2000.
- [118] G. Verbeke and G. Molenberghs. Simple missing data methods. In *Linear Mixed Models for Longitudinal Data*, Springer Series in Statistics, pages 221–229. Springer New York, 2000.
- [119] G. Verbeke, G. Molenberghs, H. Thijs, E. Lesaffre, and M. G. Kenward. Sensitivity analysis for nonrandom dropout : A local influence approach. *Biometrics*, 57(1) :7–14, 2001.
- [120] H. Wainer and D. Thissen. Estimating ability with the wrong model. *Journal of Educational Statistics*, 12(4) :339–368, 1987.
- [121] T. Warm. Weighted likelihood estimation of ability in item response theory. *Psychometrika*, 54 :427–450, 1989.

- [122] B. D. Wright and G. N. Masters. *Rating Scale Analysis*. MESA Press, Chicago, 1 edition, Jan. 1982.
- [123] M. Wu. The role of plausible values in Large-Scale surveys. *Studies in Educational Evaluation*, 31(2) :114–128, 2005.
- [124] X. Zheng and S. Rabe-Hesketh. Estimating parameters of dichotomous and ordinal item response models with gllamm. *Stata Journal*, 7(3) :313–333, 2007.
- [125] A. H. Zwinderman. Pairwise parameter estimation in rasch models. *Applied Psychological Measurement*, 19(4) :369–375, 1995.

