

**ANALYSE DE VALIDITE DES TESTS DE CLASSEMENT EN LIGNE POUR LES
LANGUES A L'AIDE DU MODELE DE RASCH**

Cornelia Gick*, Herv   Platteaux **, Sergio Hoein*, Catherine Blons-Pierre****, Patricia Kohler*******

* *Universit   de Fribourg (Suisse), cornelia.gick@unifr.ch*

** *Universit   de Fribourg (Suisse), herve.platteaux@unifr.ch*

*** *Universit   de Fribourg (Suisse), sergio.hoein@unifr.ch*

**** *Universit   de Fribourg (Suisse), catherine.blons-pierre@unifr.ch*

***** *Universit   de Fribourg (Suisse), patricia.kohler@unifr.ch*

Mots-cl  s : *test de classement en ligne, IPT/Analyse d'item, Mod  le/Analyse de Rasch, Apprentissage des langues, Plateforme Moodle*

R  sum  . *L'Universit   bilingue de Fribourg (Suisse) a r  alis   des tests de classement en ligne adaptatifs pour quatre langues. Avec la plateforme Moodle, les   tudiants passent un test, qui   value leurs comp  tences langagi  res, afin de pouvoir s'inscrire    un cours de langue   trang  re adapt      leur niveau. Depuis quatre semestres, les tests de classement en ligne ont   t   utilis  s avec succ  s (tr  s peu d'  tudiants ont d   changer de niveau de cours). Pour augmenter encore la qualit   de cette proc  dure, nous analysons la validit   et de la fid  lit   des items des diff  rents sous-tests con  us pour l'allemand (DaF) et le fran  ais (FLE)    l'aide du mod  le de Rasch (Winsteps). Cela pour deux actions : la constitution des tests et de leurs sous-tests (mode adaptatif) et l'am  lioration des items (qui participe    la formation continue des comp  tences des   quipes de d  veloppement des items).*

1. Introduction

Le Centre de Langues et le Centre Nouvelles Technologies et Enseignement, de l'Universit   bilingue de Fribourg (Suisse), ont r  alis   des tests de classement (placement tests) en ligne adaptatifs pour l'allemand, l'anglais, le fran  ais et l'italien langues   trang  res. Avec la plateforme d'apprentissage en ligne Moodle, les   tudiants passent un test, qui   value leurs comp  tences langagi  res, avant de pouvoir s'inscrire, via un syst  me d'inscription en ligne,    un cours de langue   trang  re adapt      leur niveau. Mis ensemble, notre test de classement et le syst  me d'inscription permettent d'automatiser le processus d'inscription et d'admission aux cours de langues.

Notre test de classement en ligne est un dispositif d'  valuation formative puisque son objectif premier est de d  terminer le niveau de comp  tence de chaque apprenant afin de lui proposer des cours correspondant    son niveau. Pour le Centre de Langues qui g  re le test, le but est aussi de cr  er des classes de niveaux relativement homog  nes.

Le choix d'une plateforme d'apprentissage en ligne pour mettre en   uvre nos tests de classement adaptatifs (Linacre, 2000) provient d'une volont   de mettre    profit plusieurs forces de ces technologies. En premier lieu, elles permettent de r  aliser une   conomie organisationnelle dans le cadre d'un enseignement/apprentissage en grands groupes. Avec une plateforme en ligne, chaque

semestre, 1'500   tudiants passent le test o   et quand ils veulent sans qu'ils aient besoin de faire la pied de grue dans les couloirs du Centre de Langues. Le stress des enseignants, et du secr  tariat, en est diminu   d'autant durant la m  me p  riode de la rentr  e. L'  valuation de comp  tences r  ceptives permet un calcul automatique du niveau langagier de l'  tudiant    partir de ses r  ponses au test en ligne.

Quelques principes r  sument le mod  le de ces tests :

- Ils sont compos  s de trois cat  gories de sous-tests   valuant les comp  tences r  ceptives la compr  hension orale et la compr  hension   crite ainsi que les structures de la langue.
- Il int  gre l'auto  valuation selon les descripteurs du Cadre Europ  en Commun de R  f  rence pour les langues du Conseil de l'Europe (CECR).
- Pour chaque cat  gorie, plusieurs sous-tests sont propos  s aux   tudiants selon une modalit   adaptative.
- Les deux tests utilisent un instrument qui permet    l'  tudiant d'  tre dirig   de mani  re efficace dans un sous-test de son niveau. Pour le test FLE (Fran  ais Langue Etrang  re), c'est l'auto  valuation sur la base du CECR ; pour le test DaF (Deutsch als Fremdsprache), ce sont 8-10 items de grammaire dont le r  sultat dirige vers un autre test de grammaire de 18-20 items.
- La cr  ation des items (questions et propositions de r  ponse) des sous-tests est bas  e sur le CECR (Blons-Pierre, 2010).

Les tests ont   t   utilis  s avec succ  s : tr  s peu d'  tudiants ont d   changer de niveau de cours. Depuis le lancement des tests en ligne, chaque semestre d  s f  vrier 2009, une r  colte syst  matique de donn  es a permis de mener diverses analyses que nous r  sumons dans une premi  re partie de cet article. Sur cette base, pour augmenter encore la qualit   des banques d'items, nous nous attachons    analyser plus en d  tails la validit   et la fid  lit   (Tagliante, 2005) des items des diff  rents sous-tests con  us pour les tests DaF et FLE.

Pour ce faire, nous utilisons la th  orie de r  ponse aux items (Lord, 1980 ; Wright & Stone, 2004), en particulier le mod  le de Rasch (1980), et le logiciel Winsteps, qui nous permettent de d  terminer le pouvoir de discrimination d'un item et son degr   de difficult   par rapport aux autres, ainsi que de formuler des indications concernant leur qualit  . Nous disposons maintenant d'une population suffisante pour effectuer un tel traitement. Pour le test FLE, 400   tudiants environ ont pass   le test entre les sessions de septembre 2010 et de f  vrier 2011. Pour le test DaF, 1000   tudiants environ, entre f  vrier 2009 et septembre 2011, dont 650 environ durant les sessions septembre 2010 et 2011 et f  vrier 2011 avec un test inchang  .

Ces r  sultats sont utilis  s avec deux finalit  s qui forment le c  ur du pr  sent article :

- En premier lieu, l'analyse de Rasch intervient dans la conception d'une structure adaptative des tests. Elle permet de positionner les items dans les sous-tests, en fonction de leurs degr  s de difficult  , de v  rifier si les degr  s de difficult   des items (et des sous-tests) correspondent    la population test  e (Gick, 2011) et de v  rifier la qualit   des items utilis  s.
- En second lieu, l'analyse de Rasch sert au d  veloppement continu de la production des items. Elle permet d'identifier des items probl  matiques dans les diff  rents tests cibl  s sur un niveau. Ces items sont alors soumis    l'  quipe de conception qui propose des am  liorations. Les analyses selon le mod  le de Rasch servent ensuite    v  rifier si ces changements portent leurs fruits lors de la passation suivante.

2. R  sultats des analyses pr  c  dentes

La base des deux tests en ligne DaF et FLE sont les tests papier pr  existants et valid  s. Depuis le lancement des tests en f  vrier 2009, diverses analyses ont   t   men  es,    partir d'une r  colte syst  matique de donn  es, dans le but de v  rifier la validit   de notre mod  le de test adaptatif et son op  rationnalisation dans Moodle.

On a regard   le cheminement de l'  tudiant dans le test (nombre sous-tests faits, dur  e totale, etc.) et la perception qu'il en construisait (Kohler, Platteaux, & Blons-Pierre, 2012). Ces analyses ont permis de r  gler les seuils de r  ussite des sous-tests, leurs limites en temps et le nombre d'items qui les composent. En parall  le une analyse quantitative, faite    partir des r  sultats obtenus aux tests en ligne par les premiers   tudiants selon le mod  le de Rasch, a permis de faire un premier tri sur la qualit   des items, de supprimer des items estim  s probl  matiques, ainsi que d'ajouter et de valider de nouveaux items (Gick, 2012).

Les feedbacks des   tudiants sur les tests en ligne ont aussi permis de conna  tre et d'am  liorer l'acceptance des tests en ligne (Kohler, Platteaux, & Blons-Pierre, 2010). Cette analyse se poursuit tant sur l'acceptance que sur le degr   d'accord entre l'auto  valuation effectu  e par les test  s et le r  sultat de l'  valuation propos   par le test de classement (Kohler, 2012 ; Blons-Pierre, 2011 ; Blons-Pierre, Kohler, Gick, Hoein et Platteaux, 2012).

On a   galement v  rifi  , qu'   partir des r  ponses au test en ligne de l'  tudiant, un calcul automatique pouvait   tre mis en place dans Moodle en utilisant les fonctions de tests de cette plateforme et les   l  ments d'  valuation associ  s    son tableau de notes pour   valuer des comp  tences r  ceptives langagi  res (Platteaux & Hoein, 2010).

3. Analyse de Rasch dans le contexte de nos tests

Dans la th  orie de la r  ponse aux items (IRT), il y a une estimation pour les personnes,    partir des items r  pondus vers les items sans r  ponse, en partant d'une probabilit   de r  solution de 50 % (Rost 2004). Ceci permet de formuler des observations sur les items et la totalit   du test m  me si, pour une partie du test, nous avons un nombre limit   de participants. Cette analyse nous donne des informations sur le degr   de difficult  , l'adaptation au mod  le et la qualit   des items. Le LMS Moodle permet, avec un module suppl  mentaire, d'obtenir les r  ponses d  taill  es de chaque   tudiant    chaque item. Avec le logiciel Winsteps, nous obtenons relativement rapidement les r  sultats que nous cherchons. Si ce logiciel rend les analyses tr  s rapides et permet   galement d'obtenir un grand nombre de formats de sortie et de repr  sentation des r  sultats, il est n  cessaire de pr  parer les donn  es en les codant en chiffres.

En raison de la structure adaptative de nos tests, il nous faut   galement r  unir toutes les r  ponses de tous les   tudiants aux diff  rents sous-tests dans un seul tableau. Nous traitons uniquement les questions de type QCM, qui constituent la grande majorit   des items du test FLE et sont le seul type dans le test DaF. Pour que l'analyse soit assez fiable, nous ne prenons en compte que les tests pour lesquels nous avons un nombre suffisant de r  ponses. A cause du mod  le adaptatif, il arrive que nous n'obtenions pas assez de r  ponses en une seule session. Si cela se produit, nous d  cisons de n'analyser que les sous-tests avec un nombre de r  ponses suffisant, ou d'attendre un semestre ult  rieur avant de proc  der    des changements. A partir de 70 r  ponses les r  sultats semblent   tre significatifs. Pendant le d  veloppement du test DaF, on a utilis   sous r  serve les r  sultats    partir de 30 r  ponses.

M  me si les deux tests FLE et DaF sont semblables, il existe des diff  rences qui se manifestent   galement lors de l'utilisation de l'analyse de Rasch : les niveaux du CECR (Cadre europ  en commun de r  f  rence pour les langues) sont au centre du test FLE ; chaque sous-test se positionnant par rapport    un des niveaux. Cela n  cessite une discussion approfondie sur ces niveaux dans l'  quipe de d  veloppement, ainsi qu'une phase d'exp  rimentation (Kohler, 2012). Ces deux actions ont pu   tre mises en   uvre et il a   t   possible de construire sur des travaux pr  alables. L'auto-positionnement, qui sert d'entr  e dans chaque comp  tence, a un r  le important. Pour les   tudiants, le niveau test   dans un sous-test est toujours clair. L'  tudiant est dirig   vers un sous-test d'un niveau inf  rieur, d'un niveau sup  rieur ou vers une autre comp  tence. Le chemin    travers le test passe d'abord par la compr  hension orale (CO) et   crite (CE), pour arriver ensuite    la composante « lexique et structure de la langue ». Dans les tests de CO et CE, plusieurs questions sont pos  es sur la base d'un texte (  crit ou oral). Pour l'utilisation de l'analyse de Rasch, cela

signifie que si un item est d'un niveau inférieur ou supérieur, il ne pourra pas simplement être déplacé dans un autre sous-test, mais devra être adapté. Cela justifie donc l'utilisation de l'analyse de Rasch dans une étape ultérieure, pour vérifier les résultats de cette adaptation.

Le test DaF est basé sur une progression et un continuum. Il a ainsi été possible de commencer avec un nombre réduit d'items (18 courtes tâches de compréhension écrite) et peu de sous-tests (trois tests de compréhension écrite). L'analyse de Rasch a permis le placement des items sur une échelle de difficulté et l'identification des items problématiques. Le niveau de cours est déterminé par un calcul, invisible aux étudiants, fait à partir de leurs résultats aux sous-tests. La difficulté majeure était de trouver les seuils de césure entre les niveaux. Le test ne pouvant être mis à l'essai avant son utilisation en situation réelle, le développement d'un tel test en si peu de temps impliquait un certain risque. Par conséquent l'attribution des niveaux de cours a été faite, pendant la première semaine de la première passation du test, par le biais d'une confirmation manuelle. Pour cela, les résultats calculés ont été comparés avec l'auto-évaluation des étudiants, en tenant compte des renseignements personnels, et le calcul a été en partie ajusté. C'est seulement après cette première semaine que tout a été automatisé. L'auto-évaluation du niveau de langue basée sur le CECR se fait à la fin du test, et continue à servir uniquement comme une aide à la décision en cas d'hésitation entre deux niveaux.

Au cours des sessions suivantes, de nouvelles tâches ont pu être ajoutées. La répartition sur plusieurs sous-tests et l'intégration de nouveaux items ont montré des problèmes pendant la session de printemps 2010. Depuis le semestre d'automne 2010, nous pouvons utiliser un test fiable, qui peut être analysé plus en détail et affiné. Cela n'aurait pas été possible sans l'utilisation en parallèle de l'analyse de Rasch.

4. Rasch pour la conception des tests

Depuis le début, le développement du test DaF a été accompagné pas à pas par l'analyse de Rasch. Après chaque session, nous nous intéressons en premier lieu aux valeurs suivantes :

- Le degré de difficulté (Mesure) des items par compétence. Les valeurs indiquées ici ont été utiles pour la distribution des items dans les sous-tests.
- La relation entre la difficulté des items et le niveau de la population testée, donc la question: est-ce que les items proposés couvrent suffisamment l'éventail des niveaux ? Cela est identifiable dans le « Item Map ».
- La relation des trois sous-tests entre eux.
- La qualité des items et des distracteurs par rapport à l'« Infit » et au « Point-Mesure ».

Pour développer le test DaF on s'est basé sur les questions à choix multiples de grammaire du test de classement papier qui fournissaient des indicateurs utiles pour différencier les niveaux A2 - B2. Nous avons analysé les données des tâches de grammaire de la dernière passation du test papier, au semestre d'automne 2008. Cela a fourni des points de repère pour la première version électronique des items de grammaire. A partir du semestre de printemps 2009 les données obtenues dans Moodle ont été préparées pour être ensuite analysées avec le logiciel Winsteps. Ce qui nous a intéressés tout particulièrement, à part les valeurs de difficulté (Mesure), a été le spectre de difficultés que ces tâches couvraient. Pour cela nous avons utilisé l'« Item map ». L'impression subjective qu'il manquait des items dans le segment des niveaux supérieurs a été confirmée. On a montré également ainsi que le test était globalement trop bas par rapport au niveau de notre population ; il fallait donc le compléter avec des items à des niveaux supérieurs. Les valeurs de la « Mesure » et de l'« Infit » ont permis de sélectionner huit items de différents niveaux pour le sous-test de distribution qui représente la porte d'entrée du test. L'analyse a également permis d'implémenter une progression du degré de difficulté dans le test. Pour le semestre d'automne cela signifie qu'au moins un sous-test a pu être construit sur la base d'items validés, et qui a été utilisé comme point d'ancrage pour le développement futur.

Les items nouvellement d  velopp  s pour les comp  tences r  ceptives ont   t   utilis  s pour la premi  re fois    la session d'automne 2009. Ils ont   t   regroup  s, selon le degr   de difficult   suppos  , en trois sous-tests de compr  hension   crite (CE) et quatre sous-tests de compr  hension orale (CO). Pour obtenir un nombre suffisant de tentatives, les   tudiants ont   t   envoy  s, selon leur r  sultat    un sous-test, au sous-test de difficult   sup  rieur ou inf  rieur. Dans cette structure, il peut arriver qu'un   tudiant obtienne un r  sultat bas dans le test de grammaire, et qu'il soit donc dirig   vers un test CE de bas niveau, qu'il r  ussisse avec un bon r  sultat et qu'il se trouve devoir passer tout les trois sous-tests CE. En r  ussissant le test CE de plus haut niveau, il est ensuite dirig   vers un test CO de haut niveau. S'il n'a pas alors un bon niveau dans cette comp  tence, et obtient un r  sultat bas, il est redirig   vers le test CO d'un niveau inf  rieur. Si son r  sultat est    nouveau bas, il est encore redirig   vers un niveau inf  rieur. Dans cette configuration, la dur  e du test est rallong  e, mais cela est une contrainte    accepter si l'on veut obtenir des donn  es fiables pour l'analyse de Rasch.

TABLE 13.1 DaF-Einstufungstest-HS09													ZOU368WS.TXT		Mar 11 17:54 2012	
INPUT: 273 Student 114 Item													REPORTED: 255 Student 18 Item		2 CATS WINSTEPS 3.71.0.1	
Student: REAL SEP.: 1.04 REL.: .52 ... Item: REAL SEP.: 6.03 REL.: .97																
Item STATISTICS: MEASURE ORDER																
ENTRY NUMBER	TOTAL SCORE	TOTAL COUNT	MEASURE	MODEL S.E.	MNSQ	INFIT ZSTD	OUTFIT ZSTD	PT-MEASURE CORR.	MEASURE EXP.	EXACT OBS%	MATCH EXP%	Item				
69	14	37	3.23	.41	.80	-1.1	.78	-.8	.68	.57	78.8	74.1	LV3_03			
67	17	38	2.78	.39	.91	-.5	.84	-.7	.60	.54	73.5	70.4	LV3_01			
72	16	33	2.64	.42	.81	-1.2	.74	-1.3	.66	.54	75.9	69.7	LV3_06			
71	24	36	1.66	.40	1.09	.6	1.17	.8	.37	.44	62.5	69.2	LV3_05			
66	36	157	1.58	.22	.93	-.6	.84	-.7	.49	.45	83.2	80.3	LV2_06			
68	28	36	.98	.44	1.12	.7	1.09	.4	.29	.36	78.1	75.3	LV3_02			
70	29	37	.88	.43	1.04	.3	.94	.0	.35	.36	72.7	76.0	LV3_04			
62	62	145	.43	.19	1.13	1.6	1.20	1.4	.40	.49	66.2	71.7	LV2_02			
63	65	148	.33	.19	1.09	1.1	1.12	.9	.42	.49	71.9	71.7	LV2_03			
60	43	90	-.29	.25	.97	-.4	.90	-.7	.52	.49	65.0	67.3	LV1_06			
65	89	135	-.72	.20	1.08	.9	1.06	.4	.39	.44	69.9	72.9	LV2_05			
61	107	163	-.90	.19	.98	-.1	.88	-.7	.51	.49	72.8	74.1	LV2_01			
58	71	117	-1.08	.22	.95	-.6	.91	-.8	.51	.47	69.5	67.7	LV1_04			
59	83	121	-1.52	.22	1.01	.1	1.02	.2	.45	.46	70.4	72.2	LV1_05			
64	107	140	-1.53	.22	1.12	1.1	1.12	.5	.36	.43	75.4	79.2	LV2_04			
55	119	154	-2.36	.22	.84	-1.3	.73	-1.5	.58	.47	81.3	79.4	LV1_01			
56	113	139	-2.57	.24	.87	-.9	.72	-1.2	.51	.42	83.3	81.5	LV1_02			
57	124	136	-3.55	.33	1.06	.3	1.07	.3	.31	.35	91.1	91.2	LV1_03			
MEAN	63.7	103.4	.00	.29	.99	.0	.95	-.2			74.5	74.7				
S.D.	38.1	50.1	1.89	.09	.11	.8	.15	.8			7.0	5.8				

Figure 1 : DaF HS09 – Item Statistics : Measure Order

L'analyse de Rasch pour la compr  hension   crite (cf. Figure 1) donne l'image suivante :

- La relation entre la difficult   attendue (Entry Number: Item 55-72) et celle effectivement mesur  e (Measure)   tait d  j   bonne lors de la premi  re passation du test.
- La distribution des items dans les trois sous-tests (LV1-LV3) montre   galement des bons r  sultats.
- Les items LV2-04, LV2-02 et LV3-02 ont attir   notre attention    cause de leur « Infit MNSQ », qui dans le cas id  al devrait se situer autour de 1.00. Mais nous pouvons continuer    utiliser ces items en constatant : 1) que la valeur z « Infit ZSTD », qui devrait se trouver entre -2.00 et +2.00, est bon pour les trois items ; 2) que les valeurs du « Point Measure » sont sup  rieures    .20 ; 3) et que les valeurs qui   quilibrent les diff  rences de niveau entre les sous-tests (PT-Measure Exp.) sont   galement bonnes.
- Pour l'item n. 63 (LV2-04), la valeur basse du « Measure » (-1.53) est surprenante. Dans ce cas pr  cis, il a paru int  ressant d'examiner de plus pr  s.

Le graphique des « Expected Score ICC » pour la « Entry Number 63 » (LV2-04) montre que les   tudiants avec un niveau bas ont r  ussi    r  soudre la t  che. Cela est surprenant, car le texte de cet item de compr  hension   crite contient beaucoup de mots qui ne sont probablement pas compris    ce niveau (cf. Figure 2). Nonobstant cela, l'item a pu   tre r  solu correctement.

Verwirrung im Vogelreich

<p>Schlecht geschlafen, weil Drossel, Fink oder Star so laut gepfeifen haben? Selber schuld. Die V�gel konnten nichts daf�ur. Sie passen sich bloss ihrer Umgebung an und machen es demzufolge wie die Menschen: Je h�her der Ger�uschpegel, desto energischer wird geschrien. Amseln, Rotkehlchen oder Kohlmeisen mutieren also insbesondere dann zu Schreih�lsen, wenn sie aus dem stillen Wald in die l�rmige Stadt ziehen. Und das tun die Piepser, weil sie in den Strassen und Gassen vielf�ltigere Nahrung finden und weniger Feinde haben. F�r das sorglose Leben im �berfluss m�ssen die V�gel aber einen hohen Preis bezahlen: Die Urbanisierung bedroht ihre Existenz.</p>	<p>Und zwar deshalb, weil der �berm�ssige Stimmeinsatz die Verst�ndigung mit den Artgenossen im Wald erschwert. Zwitschert zum Beispiel die Stadt-Meise lauter, als sie soll, kann die Land-Meise sie f�lschlicherweise f�r eine Amsel halten und ihren Balzgesang ganz einfach ignorieren. Was macht die Stadt-Meise dann? Na klar: noch lauter pfeifen. Damit lockt sie aber nicht die Meise an, sondern das Rotkehlchen. Totales Chaos im Vogelreich. Und schuld daran ist, wie immer, der Mensch. Liebe V�gel, macht es doch ein bisschen besser als wir!</p> <p style="text-align: right; font-size: small;">Quelle: NZZ, 12.4.2009</p>
---	--

Welche Aussage steht im Text?

Antwort w hlen:

	<input type="radio"/> In der Stadt gibt es weniger Nahrung als auf dem Land.
	<input type="radio"/> Stadt- und Landv�gel verstehen einander gut.
	<input type="radio"/> V�gel, die lange in der Stadt leben, singen lauter.
	<input type="radio"/> Je lauter die Umgebung ist, desto kraftvoller pfeifen die V�gel.

Figure 2 : DaF Item LV

Avec les bonnes strat gies, cet item a pu  tre r solu par des  tudiants avec un niveau bas, car il suffit de retrouver les mots phare « Umgebung » et « je ... desto » dans le texte. De plus, les distracteurs 1 et 2 peuvent probablement  tre identifi s comme faux sur la base de connaissances g n rales. Cet item ne peut  tre consid r  comme valide car il devrait tester la compr hension, et pas les strat gies  voqu es. Ce qui est enseign  en tant que strat gie de lecture et compr hension pendant les cours de langue se montre probl matique dans une situation de test de classement. Ici, des r gles centrales du d veloppement d'items n'ont pas  t  respect es. Une reformulation des distracteurs aurait  t  possible. Il a cependant  t  d cid  de ne pas r utiliser cet item pour les tests   venir. L'analyse de Rasch en elle-m me n'a pas reconnu cet item comme probl matique. Par contre, uniquement l' cart entre le niveau  valu  par une didacticienne de langue  trang re et le degr  de difficult  effectivement calcul  nous a rendus attentifs   cet item.

5. Mod le de Rasch et discussion pour le d veloppement des tests

Les tests se d veloppent en proc dant   une analyse continue de la validit  des items des sous-tests. Pour ce faire, nous voulons  tablir une analyse en boucle compos e des  tapes suivantes :

- Une analyse   l'aide du mod le de Rasch sert   identifier des items probl matiques.
- L' quipe de conception des items discute ceux-ci et propose des am liorations.
- Les am liorations sont impl ment es dans le test pour la session suivante.
- Le mod le de Rasch sert   v rifier si ces changements portent leurs fruits.
- L' quipe de conception discute des r sultats obtenus.

D'autres analyses, men es en parall le (satisfaction, changements de classe, ...), servent   garantir que le test continue   fonctionner suite   ces modifications et   identifier les points avec une marge d'am lioration (Blons-Pierre et al., 2012 ; Gick, 2012).

Nous voulons mettre en place cette proc  dure en boucle pour am  liorer les items des tests de fa  on continue. Actuellement nous sommes en pleine phase d'impl  mentation d'une premi  re boucle. Nous avons fait une analyse des items, identifi   les items probl  matiques, et analysons ces items dans l'  quipe de conception, pour en faire ressortir les am  liorations qui seront impl  ment  es dans le test d  s le prochain semestre.

5.1 Identifier les items potentiellement probl  matiques

Le but de notre premi  re   tape est d'identifier quelques items probl  matiques qui seront analys  s et   ventuellement modifi  s. Comme nous l'avons vu plus haut, nos tests de classement remplissent d  j   leur fonction, et nous voulons   videmment que cela reste ainsi. La proc  dure d'am  lioration se faisant sur le test de classement en utilisation (pour g  n  rer des r  ponses en situation r  elle), nous limitons donc la proportion d'items modifi  s entre deux passations. Cela permet au test, qui a fait ses preuves sur le terrain, de continuer    atteindre son but (placer les apprenants dans des classes le plus homog  nes possibles), tout en permettant d'am  liorer les items.

L'analyse de Rasch nous aide    identifier les items potentiellement    modifier. Nous utilisons ici le mod  le de Rasch de fa  on pragmatique : il nous sert    identifier les items probl  matiques mais ne nous sert pas pour l'identification de la source du probl  me. Pour identifier ces items probl  matiques, nous utilisons principalement deux repr  sentations des r  sultats propos  s par le logiciel Winsteps : l' « item : map » et le tableau de sortie selon Item. Nous regardons d'abord toutes les donn  es ensemble, puis nous les s  parons selon les comp  tences test  es.

Les graphiques des « Item : Map » donne un premier aper  u de la position des items les uns par rapport aux autres selon leur degr   de difficult  . Nous regardons si des items se placent en dehors de la « zone » de difficult   pour laquelle ils ont   t   cr  s. Id  alement, avec la structure adaptative de nos tests, nous aimerions trouver, du bas vers le haut, d'abord les items des tests des niveaux A, ensuite, au milieu, ceux des niveaux B et, tout en haut les items des niveaux C. Si la majeure partie des items se place bien dans cette structure, on en trouve aussi qui sont plac  s plus haut ou plus bas que leur niveau id  al. Ces items sont retenus pour la discussion dans l'  quipe (  tape 2).

Avant cela, ces items sont regard  s dans les tableaux d'analyse par item pour savoir s'ils sont   galement identifi  s comme probl  matiques selon d'autres facteurs (p.ex. Pt.Measure, Infit, Outfit) et la courbe de score attendu (Expected Score ICC). Cela est souvent le cas. Ces tables de r  sultat permettent   galement de ressortir des Items n'attirant pas l'attention dans l'item map, mais pour lesquels d'autres indications font penser qu'ils sont probl  matiques : p.ex. un Pt.Measure bas, un Infit ou Outfit   lev   (ou bas). « Measure », « Point Measure » (Wright et Stone, 1979). Nous identifions ainsi un certain nombre d'autres items potentiellement probl  matiques, qui seront soumis    l'  quipe de conception des items.

Par exemple, pour les donn  es du test FLE des deux semestres 2011 (nous avons regroup   ces deux semestres pour obtenir un nombre suffisant de r  ponses), nous identifions une quinzaine d'items sur un total de 210.

5.2 Identification des probl  mes, proposition et impl  mentation des modifications

Dans cette phase, l'  quipe de conception des items va :

- Regarder les r  ponses des   tudiants    chaque item potentiellement probl  matique (r  sultats des tentatives dans Moodle)
- Discuter et poser des hypoth  ses sur : Est-ce que il y a effectivement un probl  me avec cet item ? Si oui : quelles pourraient   tre les causes ? Si non, en quoi l'analyse l'identifie comme probl  matique ?
- Proposer des possibles am  liorations pouvant aller du changement d'une virgule    l'  limination/substitution compl  te de l'item.
- D  cider des items    changer en cherchant le maintien du fonctionnement des tests.

Nous sommes actuellement dans la phase d'identification des probl  mes et de proposition de modifications. Pour le test FLE nous avons par exemple commenc   une discussion sur certains items. Tous les items probl  matiques ont   t   ou seront discut  s de cette fa  on, et l'  quipe d  cidera lesquels modifier. Voici quelques exemples de discussion :

Un exemple d'item de compr  hension   crite :

- Question : Elle s'int  resse    la planche    voile depuis la pr  -adolescence.
- R  ponses : Vrai, Faux, La r  ponse n'est pas dans le texte.
- Analyse de Rasch : Cet item a un pt. mesure tr  s bas, et un infit haut (avec zsdh haut), la courbe ICC montre   galement un probl  me.

Hypoth  se formul  e :

- Probl  me avec formulation/interpr  tation: « LT monte premi  re fois » dans le texte    lire n'est pas   gale    « s'int  resse » dans la question.
- Possible solution: reformuler « int  resse ». Changer avec : « *****    la planche    voile ».

Ce changement va   tre impl  ment  .

Un exemple d'item de lexique et structure :

- Question : Il s'agit d'une phrase dans laquelle on cherche un accord verbal.
- R  ponses : a choisi, choisira, ait choisi, avait choisi.
- Analyse de Rasch : Cette question a un pt. mesure bas, et sa courbe montre   galement un probl  me. La distribution des r  ponses sur les distracteurs donne un indice : la majorit   des personnes choisissent une m  me mauvaise r  ponse (distracteur 3).

Hypoth  se formul  e :

- Item trop difficile, distracteur interpr  t   comme juste.
- Solutions possibles : a) Il s'agit d'un test de niveau   lev  , donc il doit   tre difficile. b) Il est trop difficile, on pourrait changer le distracteur en « choisissait ».

Ce changement demande qu'il soit encore discut   avant d'  tre impl  ment  .

5.3 Analyse avec les items modifi  s

Lorsque les modifications seront impl  ment  es dans les tests des prochaines sessions, nous pourrions collecter une nouvelle s  rie de donn  es avec lesquelles nous effectuerons une nouvelle analyse sur la base du mod  le de Rasch. Nous pourrions ensuite regarder avec l'  quipe si la modification a port   ses fruits. Si oui, on aura appris que l'hypoth  se sur le probl  me de l'item   tait fond  e. Autrement nous pourrions d  cider si une modification additionnelle s'impose ou si l'item doit   tre substitu  . Dans les deux cas on aura un indice de l'efficacit   des modifications propos  es.

6. Conclusions et perspectives

L'analyse de Rasch nous appara  t utile pour la conception et le d  veloppement de tests adaptatifs. Elle contribue    la transformation rapide, pour les tests DaF et FLE, d'un test papier progressif en un test adaptatif en ligne valide en se servant de donn  es r  colt  es en situation r  elle. Dans une phase ult  rieure, cette analyse soutient   galement l'optimisation des tests en visualisant les degr  s de difficult  s des items les uns par rapport aux autres. De plus elle donne des indices sur la validit   de chaque item individuel.

Cette analyse comporte aussi des limites. Notre exp  rience montre notamment le besoin de continuer les analyses quantitatives tout en y associant d'autres m  thodes plus qualitatives pour d  velopper des tests de classement en langues   trang  res adaptatifs. Ainsi, l'analyse de Rasch ne permet pas de d  tecter tous les items « probl  matiques ». Le regard des experts en didactique des langues   trang  res reste incontournable. En effet, les r  flexions de ces experts sur les items

probl  matiques permettent d'explicitier les probl  mes li  s    certains items. Pour tirer avantage des forces des deux m  thodes, nous avons d  cid   de mettre en place une proc  dure en boucle, alliant analyse de Rasch et regard des experts.

Les premi  res exp  riences avec cette proc  dure portent leurs fruits. Nos perspectives de travail sont de la rendre plus syst  matique. Ensuite, il s'agira d'effectuer un bilan de la proc  dure, en termes des choses apprises sur le d  veloppement des items par les   quipes des tests (formation continue au sein du Centre de Langues).

7. Bibliographie

- Blons-Pierre, C. (2010). *Tensions between change and quality assurance : A dynamic necessary for innovation in university language centres/Tensions entre politique du changement et assurance qualit   : une dynamique n  cessaire pour l'  volution des centres de langues universitaires*. Communication pr  sent  e    la 11th International CercleS Conference, Helsinki, 4 September.
- Blons-Pierre, C. (2011). *L'impact du CECR/GeR/CEFR sur le fonctionnement d'un centre de langues dans une universit   bilingue suisse*. Communication pr  sent  e    la 4th International ALTE Conference, 7-9 July.
- Blons-Pierre, C., Kohler, P., Gick, C., Hoein, S. & Platteaux, H. (2012). Analyse qualitative des tests de classement en ligne pour les langues : acceptance et r  le de l'auto  valuation. Dans *Actes du XXIV  me Colloque de l'ADMEE Europe. L'  valuation des comp  tences en milieu scolaire et en milieu professionnel*, Luxembourg 11-13 janvier.
- Conseil de l'Europe. (2001). Cadre europ  en commun de r  f  rence pour les langues : apprendre, enseigner,   valuer. Paris : Didier.
- Gick, C. (2011). *Online-Einstufung und Online-Einschreibung in studienbegleitende Deutschkurse mittels adaptivem Einstufungstest: Entwicklungsschritte und Stand der Dinge*. Communication pr  sent  e au Bremer Symposium 2011 (Fremdsprachenzentrum der Hochschulen im Land Bremen), Bremen 4-5 M  rz.
- Gick, C. (2012). Einstufungstests im Spannungsfeld von Bologna-Reform, Referenzrahmen f  r Sprachen, Fachwissenschaft und Nutzern, In C. Blons-Pierre (dir.) : *Apprendre, enseigner et   valuer les langues dans le contexte de Bologne et du CECR / Sprachen lernen, lehren und beurteilen im Kontext von Bologna und dem GER*. Bern : Peter Lang, 189-217.
- Kohler, P. (2012). Gestion de l'auto-  valuation dans un dispositif de test en ligne: limites et perspectives. In C. Blons-Pierre (dir.), *Apprendre, enseigner et   valuer les langues dans le contexte de Bologne et du CECR / Sprachen lernen, lehren und beurteilen im Kontext von Bologna und dem GER* (pp. 219-240). Bern : Peter Lang.
- Kohler, P., Platteaux, H., & Blons-Pierre, C. (2010). *R  actions des   tudiants    l'usage des tests en ligne pour   valuer les comp  tences langagi  res*. Communication pr  sent  e au XXVI  me Congr  s AIPU. R  formes et changements p  dagogiques dans l'enseignement sup  rieur, Rabat 17-21 mai.
- Kohler, P., Platteaux, H., & Blons-Pierre, C. (2012). Un test de classement en ligne pour   valuer les niveaux de comp  tence et constituer des groupes classes. *Revue RIPES - Num  ro Sp  cial « Innover dans l'  valuation des apprentissages : pourquoi et comment ? »*, 27(2), R  cup  r   du site de la revue : <http://ripes.revues.org/532>.
- Linacre, J. M. (2000). *Computer-Adaptive Testing: A methodology whose time has come*. Chicago: MESA Psychometric Laboratory - University of Chicago.
- Lord, F. M. (1980). *Applications of item response theory to practical testing problems*. Mahwah, NJ: Erlbaum.
- Platteaux, H., & Hoein, S. (2010). *Tests de positionnement en langues   trang  res sur Moodle*. Communication pr  sent  e    la Conf  rence Moodle Moot 2010, Troyes 28-30 juin.
- Rasch, G. (1980). *Probabilistic models for some intelligence and attainment tests*. Danish Institute for Educational Research Copenhagen 1960, expanded edition with foreword and afterword by B.D. Wright. Chicago: The University of Chicago press.

Actes du 24^e colloque de l'Adm  -Europe
L'  valuation des comp  tences en milieu scolaire et en milieu professionnel
C. Gick, H. Platteaux, S. Hoein, C. Blons-Pierre, P. Kohler
Analyse de validit   des tests de classement en ligne pour les langues    l'aide du mod  le de Rasch

Rost, J. (2004). *Lehrbuch Testtheorie – Testkonstruktion*. Zweite,   berarbeitete und erweiterte Auflage, Bern, G  ttingen, Toronto, Seattle: Hans Huber.

Tagliante, C. (2005). *L'  valuation et le Cadre europ  en commun*. Paris: CLE International.

Wright, B. D., & Stone, M. H. (2004). *Making measure*. Chicago: The Phaneron Press.