

BMJ Open Assessing differential item functioning for the Social Appearance Anxiety Scale: a Scleroderma Patient-centred Intervention Network (SPIN) Cohort Study

Sophia J Sommer,^{1,2} Daphna Harel,^{1,2} Linda Kwakkenbos,³ Marie-Eve Carrier,⁴ Shadi Gholizadeh,^{5,6} Karen Gottesman,⁷ Catarina Leite,⁸ Vanessa L Malcarne,^{5,9} Brett D Thombs ,^{4,6,10,11,12,13} SPIN Investigators

To cite: Sommer SJ, Harel D, Kwakkenbos L, *et al.* Assessing differential item functioning for the Social Appearance Anxiety Scale: a Scleroderma Patient-centred Intervention Network (SPIN) Cohort Study. *BMJ Open* 2020;**10**:e037639. doi:10.1136/bmjopen-2020-037639

► Prepublication history for this paper is available online. To view these files, please visit the journal online (<http://dx.doi.org/10.1136/bmjopen-2020-037639>).

Received 12 February 2020
Revised 15 July 2020
Accepted 26 August 2020



© Author(s) (or their employer(s)) 2020. Re-use permitted under CC BY-NC. No commercial re-use. See rights and permissions. Published by BMJ.

For numbered affiliations see end of article.

Correspondence to

Dr Brett D Thombs;
brett.thombs@mcgill.ca

ABSTRACT

Objectives The Social Appearance Anxiety Scale (SAAS) is a 16-item questionnaire developed to evaluate fear of appearance-based evaluation by others. The primary objective of this research was to investigate the existence of differential item functioning (DIF) for the 16 SAAS items, comparing patients who completed the SAAS in English and French, either to confirm that scores are comparable or provide guidance on calculating comparable scores. A secondary research objective was to investigate the existence of DIF based on sex and disease status. A tertiary research objective was to assess DIF related to language, sex, and disease status on the recently developed SAAS-5.

Design This was a cross-sectional analysis using baseline data from patients enrolled in the Scleroderma Patient-centred Intervention Network (SPIN).

Setting SPIN patients included in the present study were enrolled at 43 centres in Canada, USA, UK, France and Australia, with questionnaires completed in April 2014 to July 2019.

Participants 1640 SPIN patients completed the SAAS in French (n=600) or English (n=1040).

Primary and secondary measures The SAAS was collected along with demographic and disease characteristics.

Results Six items were identified with statistically significant language-based DIF, four with sex-based DIF and one with disease type-based DIF. However, factor scores before and after accounting for DIF were similar (Pearson correlation >0.99), and individual score differences were small. This was true for both the full and shortened versions of the SAAS.

Conclusion SAAS and SAAS-5 scores are comparable across language, sex, and disease-type, despite small differences in how patients respond to some items.

INTRODUCTION

A desire to improve the patient-centred focus of healthcare research has led to the development and increased use of patient-reported

Strengths and limitations of this study

- This study uses a large cohort of patients which provides robust results, allowing for insights into the larger population of adults with scleroderma.
- Patients in the sample were required to have internet access in order to complete study questionnaires, which may limit generalisability of these findings due to selection bias.
- These findings are only generalisable to adults with scleroderma and should be confirmed for other populations.

outcome (PRO) measures aimed at a wide range of human experiences, including patient-perceived health, well-being and psychological status.¹ This is particularly important in chronic diseases that lead to symptoms that are not directly measurable.² Many PRO measures have been translated into multiple languages, which is relevant in treatment centres where more than one language is common, as well as in rare disease research, which often involves collaboration and communication across sites in multiple countries.³ In these situations, outcomes measured in more than one language are commonly combined in analyses.

In order to compare PROs across language and cultural groups, it is important to ensure that all patients interpret and respond to the questionnaire items in equivalent ways, and not based on idiosyncratic differences due to differing cultural norms, systematic differences in interpretation or indirect translations of some items.⁴ If this is not the case, then items or questions are said to have differential item functioning (DIF). When DIF is present, patients with equal underlying



levels of the construct, or latent trait, measured by that scale will respond differently to the same item.⁵

Systemic sclerosis (SSc) is a rare, multisystem autoimmune disorder with heterogeneous symptomatology characterised by microvascular damage and fibrosis in multiple organs.^{6,7} Changes in appearance are common and can include telangiectasias, hypopigmentation and hyperpigmentation, loss of skin folds, loss of flexibility of the lips, digital ulcers and hand contractures.^{6,8} These changes in appearance are often in socially relevant areas of the body, such as the face and hands, and can lead to problems with social interactions and increased social appearance anxiety.⁹

The Social Appearance Anxiety Scale (SAAS) is a 16-item scale, which aims to measure patients' fear of appearance-based evaluation.¹⁰ Among people with SSc, the SAAS may be used for both individual-level treatment plans and larger scale research, evaluating the impact of potential interventions. The Scleroderma Patient-centred Intervention Network (SPIN) Cohort is a web-based, international cohort designed to collect PROs at regular intervals and as a framework to conduct trials of psychosocial and rehabilitation interventions for patients with SSc.¹¹ Depending on their native language, participants enrolled in SPIN may complete the SAAS in French, English or Spanish; however, no research has yet confirmed that SAAS scores are comparable across these language groups.

A recent study developed a shortened version of the SAAS consisting of five items (SAAS-5) for use in patients with SSc.¹² The use of shortened versions, such as the SAAS 5, has the potential to decrease patient burden and increase data quality.¹³ However, it is of interest to determine whether the shortened version exhibits DIF.

Therefore, the primary purpose of this analysis is to investigate the comparability of responses to versions of the SAAS administered in different languages. As a secondary research objective, comparability of SAAS scores with respect to disease type and sex were also assessed. A tertiary research objective was to assess the comparability of SAAS scores on the 5-item shortened version.

MATERIALS AND METHODS

Patients and procedures

The sample consisted of patients enrolled in the SPIN Cohort with complete data study questionnaires from initial enrollment sessions between April 2014 and July 2019. Participants in the SPIN Cohort were enrolled at 43 centres in Canada, USA, UK, France and Australia. To be eligible for the SPIN Cohort, participants must be classified as having SSc according to the 2013 ACR/EULAR classification criteria,¹⁴ confirmed by a SPIN physician, be at least 18 years of age, have the ability to give informed consent, and be fluent in English, French, or Spanish. However, the present study only included patients who completed study questionnaires in English

or French, as the sample size of Spanish patients was too small to be included at the time of the analyses. Exclusion criteria for participation in the SPIN Cohort include not having access to the internet or otherwise not being able to respond to questionnaires via the internet. The SPIN sample is a convenience sample. Eligible participants are invited by the attending physician or a supervised nurse coordinator to participate in the SPIN Cohort, and written informed consent is obtained. The local SPIN physician or supervised nurse coordinator then completes a medical data form that is submitted online to initiate participants registration in the SPIN Cohort. After completion of online registration, an automated welcoming email is sent to participants with instructions on how to activate their SPIN online account and how to complete the SPIN Cohort patient measures online. SPIN Cohort participants complete outcome measures via the internet on enrollment and subsequently every 3 months.

Measures

Demographics and disease characteristics

Demographic and disease variables included age, sex, race/ethnicity, marital status, education level, time since diagnosis, and SSc subtype. Disease subtypes were classified as limited or diffuse. Limited disease was defined as skin involvement distal to the knees and elbows only, whereas diffuse disease included more extensive skin involvement.¹⁵ The country of patient recruitment and language of assessment were also recorded.

Social Appearance Anxiety Scale

The SAAS consists of 16 items assessing patients' self-reported anxiety about appearance-based evaluation. The SAAS was initially validated among three samples of undergraduate students ($n=512$, 853, and 541, respectively).¹⁰ In this population, the SAAS was shown to have unifactorial structure, high internal consistency, high test-retest reliability, and was positively correlated with other social anxiety measures.¹⁰ A recent study of 938 participants enrolled in the SPIN Cohort demonstrated that the SAAS is a unidimensional, reliable, and valid measurement of social appearance anxiety among people with SSc.¹⁶ The SAAS was initially written in English. The French version used in this study was translated by SPIN investigators using the forward-backward method.¹⁷ For both versions, item responses are recorded on a five point scale (1=not at all, 5=extremely). Item 1 ('I feel comfortable with the way I appear to others') is reverse coded before summing across items to produce a total score ranging from 16 to 80. Higher scores indicate higher levels of appearance anxiety.

The SAAS-5, consisting of items 6, 7, 12, 13 and 14 from the original SAAS, was recently developed and validated for use in patients with SSc.¹² Scores on the SAAS-5 range from 5 to 25, with higher scores indicating higher levels of appearance anxiety.

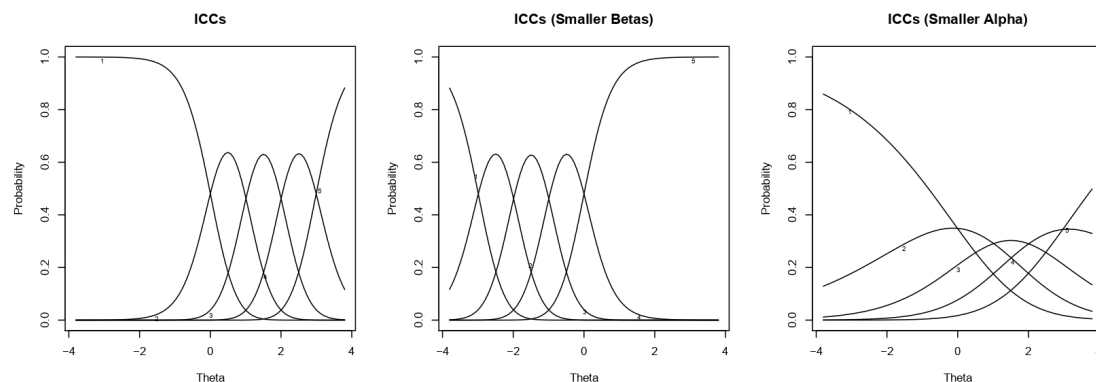


Figure 1 Three possible ICC for a five-category item. The left and middle panels show ICCs for items with the same approximate discrimination parameters (alphas) but different item-level thresholds (betas). The left and right panels show ICCs for items with the same approximate item-level thresholds (betas) but different discrimination parameters (alphas). ICC, item-characteristic curve.

Statistical analysis

The English-speaking and French-speaking samples were compared based on demographic and disease characteristics to identify possible differences between the two language groups.

A generalised partial credit model (GPCM) was then used to model the latent factor (social anxiety with appearance) underlying the SAAS. For each item, a GPCM was used to estimate two types of item-level parameters: (1) thresholds (betas) for the level of the latent factor (theta) at which respondents are more likely to endorse a given response category instead of the category below and (2) a discrimination parameter (alpha) that measures the strength of the relationship between that item and the underlying latent factor.¹⁸

Item-characteristic curves (ICCs) are often used to visualise these parameters, and figure 1 shows three examples of ICCs for a hypothetical 5-category item. Each curve in an ICC plot corresponds to a possible categorical response. Along the latent spectrum, the height of each curve indicates the estimated probability that a respondent with a particular level of the latent factor will respond with the corresponding category. Item-level thresholds are visualised as the intersections between consecutive curves; discrimination parameters are visualised as the peaked-ness of the curves. When item-level thresholds vary across observed groups, items are said to display uniform DIF. Uniform DIF could be visualised as a horizontal shift of ICC for one demographic group compared with the other. Meanwhile, when the discrimination parameter varies across observed groups, items are said to display non-uniform DIF. Non-uniform DIF could be visualised as a change in the peaked-ness of the curves for one demographic group compared with the other.

The *lordif* package in R^{19 20} was used to identify items with language-based DIF through an iterative procedure. The algorithm implemented by *lordif* iteratively fits three ordinal logistic models for each item and uses these models to flag items with potential DIF. The first model predicts the probability of each response category using estimated latent factor scores alone, while

the second and third models test for uniform and non-uniform DIF, respectively. Once a set of items is flagged, the algorithm then re-estimates latent factor scores using another GPCM that accounts for DIF on those items. DIF

Table 1 Demographic and disease characteristics by assessment language

Variable	All patients (n=1640)	English-speaking patients (n=1040)	French-speaking patients (n=600)
Mean age, years (SD)*	55.1 (12.5)	55.7 (11.7)	54.0 (13.8)
Female (%)	87.2	87.6	86.5
Mean SAAS summed score (SD)	29.1 (13.7)	28.3 (13.2)	30.5 (14.5)
Diffuse disease type (%)	39.0	42.4	33.2
Mean time since diagnosis, years (SD)	9.2 (7.9)	9.7 (8.0)	8.5 (7.6)
Married or common law (%)	71.2	73.3	67.5
At least 12 years of education (%)	85.7	94.2	70.8
Race†			
White (%)	83.6	83.9	83.0
Black (%)	7.1	6.1	8.8
Other (%)	9.3	10.0	8.2
Country of patient recruitment			
Canada (%)	24.9	28.7	18.5
USA (%)	35.5	55.9	0.2
UK (%)	9.7	15.3	0.0
France (%)	29.8	0.1	81.3
Australia (%)	0.1	0.1	0.0

Due to missing values.

*n=1036 for the English cohort.

†n=1038 for the English cohort.

SAAS, Social Appearance Anxiety Scale.

**Table 2** SAAS items

Item*	Item text
1	I feel comfortable with the way I appear to others.
2	I feel nervous when having my picture taken.
3	I get tense when it is obvious people are looking at me.
4	I am concerned people won't like me because of the way I look.
5	I worry that others talk about flaws in my appearance when I am not around.
6	I am concerned people will find me unappealing because of my appearance.
7	I am afraid people find me unattractive.
8	I worry that my appearance will make life more difficult for me.
9	I am concerned that I have missed out on opportunities because of my appearance.
10	I get nervous when talking to people because of the way I look.
11	I feel anxious when other people say something about my appearance.
12	I am frequently afraid that I won't meet others' standards of how I should look.
13	I worry people will judge the way I look negatively.
14	I am uncomfortable when I think others are noticing flaws in my appearance.
15	I worry that a romantic partner will/would leave me because of my appearance.
16	I am concerned that people think I am not good looking.

SAAS, Social Appearance Anxiety Scale.

is accounted for by allowing item level parameters to vary across groups. The process stops once the same items are repeatedly flagged for DIF.²⁰

During the iterative search for items with DIF, items were flagged using a χ^2 test comparing the first and third models (alpha=0.01 significance level). Flagged items were then re-examined to distinguish between uniform and non-uniform DIF. This was done by separately comparing the first and second models (to ascertain uniform DIF) and second and third models (to ascertain non-uniform DIF), again using a χ^2 test (alpha=0.01 significance level).

Items with DIF were further investigated by comparing item-level parameters from a GPCM for patients who completed the SAAS in English and French. To visualise and understand differences among the two groups on each item, item true score functions for English-speakers and French-speakers were compared, which show expected responses for items with DIF as a function of estimated latent social appearance anxiety accounting for DIF.

The questionnaire-level impact of DIF on estimated latent factor scores was assessed by plotting test characteristic curves, which show expected summed scores on the SAAS as a function of patients' GPCM scores accounting for DIF. As per previous guidelines, impact was numerically assessed by comparing initial scores (not accounting

for DIF) to final scores (accounting for DIF), using the Pearson correlation of the two scores and by comparing individual score differences to the SEs of initial scores.^{21 22}

To assess whether the correlation significantly differed from 1, a randomisation null distribution and p values were obtained by randomly permuting group labels 1000 times and re-estimating scores and statistics holding the measurement model fixed across permutations, but re-estimating the item-level parameters based on the permuted dataset.

Lastly, the median and range of score differences (of the difference between scores accounting for and not accounting for DIF) were also calculated, and score differences were plotted against initial scores to find areas of the latent spectrum with highest DIF impact. Before comparison, scores were placed on the same scale using a transformation by Stocking and Lord.²³ This was also done using the *lordif* package, which equates final scores accounting for DIF to initial GPCM estimates using the non-DIF items as anchors.

The same process was repeated to identify and investigate DIF related to sex and disease status, respectively, and additionally for the SAAS-5.

Patient involvement

SPIN was conceived by a collaboration of investigators and patients. SPIN's Patient Advisory Board advises the SPIN Steering Committee on priorities for investigation. Patients were included in the SPIN Publication Committee, which reviewed the proposal for the present study and its methods. Two patients were coauthors of the present report.

RESULTS

The English and French samples included 1040 and 600 patients, respectively. **Table 1** presents descriptive statistics for the full sample, as well as the English and French samples separately.

DIF analysis

Six of the 16 SAAS items (**table 2**) were identified as having statistically significant ($p < 0.01$) language-based DIF: items 2, 5, 8, 11, 12 and 13. Only item 11 was identified as having non-uniform DIF. Item true score functions for these six items are shown in **figure 2**. For most items with uniform DIF, French speakers' expected item level responses were slightly higher than their English-speaking counterparts with equal levels of appearance anxiety. This pattern is reversed for item 2.

Test characteristic curves for the English and French cohorts are plotted in **figure 3**, while **figure 4** shows score differences based on GPCMs that do and do not account for DIF. At the questionnaire level, French speakers are expected to have slightly larger summed scores on the SAAS as compared with English speakers with the same level of appearance anxiety. The correlation between the two sets of GPCM scores was 0.99977 (95% CI 0.99975

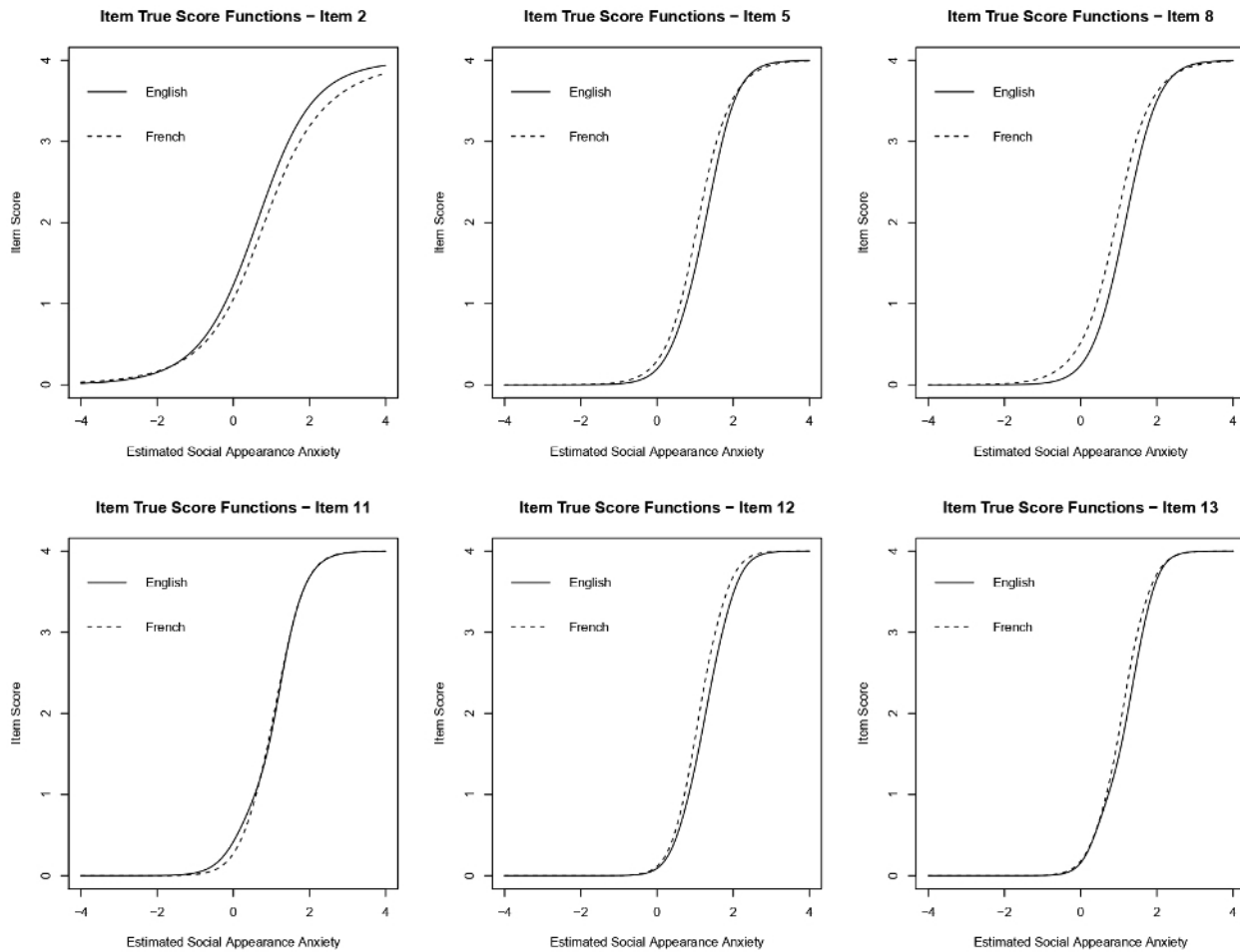


Figure 2 Item true score functions for the six items identified as having language-based DIF. For items 5, 8, 12 and 13, these plots demonstrate that French speakers are expected to give larger categorical responses than English speakers with equal levels of appearance anxiety. This trend is reversed for item 2, while item 11 demonstrates non-uniform DIF (ie, the true score functions for English and French speakers cross each other). DIF, differential item functioning.

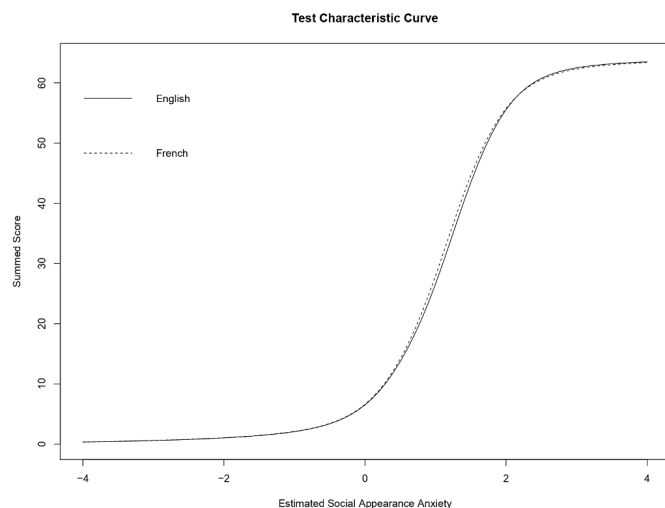


Figure 3 Test characteristic curve showing expected summed scores on the SAAS as a function of estimated social appearance anxiety accounting for DIF. Thus, among French and English speakers with the same estimated level of social appearance anxiety, French speakers are expected to have slightly larger summed scores. DIF, differential item functioning; SAAS, Social Appearance Anxiety Scale.

to 0.99979, $p < 0.001$). At the individual level, the median GPCM score difference (scores accounting for DIF minus scores that do not account for DIF) was 0.0049, and differences in factor scores ranged from -0.078 to 0.065 . No individual score differences exceeded the standard errors of initial estimates. Patients with the largest score differences had initial GPCM scores around -0.5 and 1.0 , whereas individuals whose initial estimated anxiety level was extreme (low or high) or average had smaller DIF impact.

Four items were identified as having sex-based DIF (all uniform): items 2, 4, 9 and 14. Only item two exhibited both language and sex based DIF. Item true score functions suggest that females tend to give slightly higher categorical responses than equally anxious males on items 2 and 14 and lower responses on items 4 and 9. Meanwhile, the test characteristic curves for males and females were practically indistinguishable, suggesting that equally anxious males and females have almost identical expected summed scores. The correlation between the two sets of GPCM scores was 0.99985 (95% CI 0.99983 to 0.99986, $p = 0.003$). At the individual level, the median score difference based on a GPCM was 0.0020, and differences in

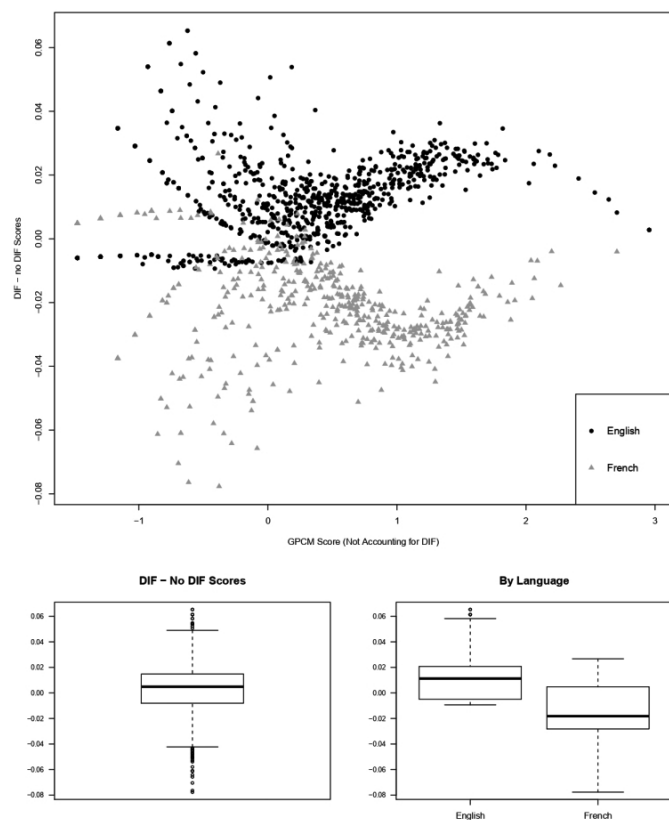


Figure 4 The top plot shows GPCM score differences at the questionnaire level (accounting for DIF— not accounting for DIF) compared with factor scores accounting for DIF. The largest score differences occur at estimated appearance anxiety levels .5 SD below average and 1 SD above average. The figure on the bottom left shows a box plot of these score differences among all respondents. The figure on the bottom right shows these differences by language. Overall differences are small and are mostly negative for English speakers and positive for French speakers, suggesting that pooled scores from a GPCM will tend to overestimate appearance anxiety for French speakers and underestimate it for English speakers. DIF, differential item functioning; GPCM, generalised partial credit model.

factor scores ranged from -0.047 to 0.135 . No individual score differences exceeded the SEs of initial estimates. The largest score differences were observed for individuals whose initial GPCM score was low (around -1.0 in this dataset); individuals with average or high estimated anxiety levels had comparatively low DIF impact.

Only one item (item 9) was identified as having DIF related to disease-type (non-uniform). This item was also identified as having sex-based DIF, but not language-based DIF. Among patients with low appearance anxiety, those with limited disease are expected to give smaller categorical responses to item nine than patients with diffuse disease and equal levels of appearance anxiety; this pattern is reversed at the higher end of the latent spectrum. At the questionnaire level, expected summed scores were nearly identical across disease-type groups. The correlation between the two sets of GPCM scores was 0.99996 (95% CI 0.99996 to 0.99997 , $p < 0.001$). At the

individual level, the median GPCM score difference was 0.001 and these differences in factor scores ranged from -0.101 to 0.080 . No individual score differences exceeded the SEs of initial estimates. The largest score differences were observed for individuals whose initial GPCM estimate was around 0, or slightly below.

For the SAAS-5, only item 12 was flagged for language based DIF, while item 14 was still flagged for gender-based DIF. In both cases, the correlation between factor scores was still high: 0.99995 for language-based DIF (95% CI 0.99995 to 0.99996 , $p = 0.017$) and 0.99971 for gender-based DIF (95% CI 0.99969 to 0.99974 , $p = 0.018$).

DISCUSSION

This study investigated whether the SAAS displays DIF across language, sex, and disease subtype groups among people with SSc. Nine items were flagged for language-based DIF (eight uniform, one non-uniform), four were flagged with sex-based DIF (all uniform), and only one was flagged with disease-type based DIF (non-uniform). In reviewing translations of the items flagged with language-based DIF, we did not observe any clear differences. Similarly negligible levels of DIF were found for the SAAS-5.

For all three analyses on the full-length SAAS, the high (>0.99) Pearson correlations between the two GPCM estimates imply that accounting for DIF does not provide much additional information about respondents' comparative levels of social appearance anxiety. The near-zero (<0.01) associated p values nonetheless suggest that observed correlations are lower than what would be expected by random chance in a no-DIF null condition under identical measurement models. While previous analyses have used Pearson correlations^{21 22} to compare GPCM scores that do and do not account for DIF, other research has cautioned against this.²⁴ Our findings imply that very large correlations between initial and final GPCM estimates may still be smaller than simulated values under a no-DIF condition. Thus, we caution that correlations alone may not be particularly interpretable as a measure of DIF impact.

The relatively small ranges of GPCM score differences in all three analyses nonetheless support the conclusion that accounting for DIF has limited impact on individual estimated scores. In all cases, no individual differences exceeded initial SEs. Thus, estimated scores accounting for DIF were all within the range of inherent uncertainty in naïve GPCM estimates. The median score difference was largest for language-based DIF and smallest for disease-type-based DIF; however, the range of score differences was smallest in the language-based analysis, due to the existence of a few outliers in the other two cases.

Scatter plots of GPCM score differences as a function of naïve GPCM estimates (see figure 4 for language-based DIF) show that language-type, sex-type and disease-type-based DIF impact is not constant across the latent spectrum. Naïve GPCM estimates near values where GPCM score differences are larger (ie, near -0.5 and 1

for language-based DIF, -1 for sex-based DIF, and 0 for disease-type-based DIF) may therefore be slightly less certain.

While DIF impact was found to be small for both simple summed scores and naïve GPCM estimates, it is important to note that the choice between these two scoring methods is also relevant.^{25 26} This paper explored three different methods for estimating social appearance anxiety levels based on responses to the SAAS: simple summed scores, naïve GPCM factor scores, and GPCM factor scores accounting for DIF. Our analysis aimed to assess comparability of scores across demographic groups, and therefore, mainly focused comparison between the two sets of GPCM factor scores; however, much more confidence in individual scores is gained in the jump from simple summed scores to a GPCM factor score, than in the jump from a naïve GPCM factor score to a GPCM factor score accounting for DIF. For example, in this dataset, individuals with the same summed score had naïve GPCM estimates of social appearance anxiety differing by up to 0.92 standardised units. Thus, regardless of whether DIF is accounted for in score calculations, a GPCM-based score or weighted summed score would be preferable over a simple summed score.

This study has several limitations. First, DIF was only investigated in the population of adults with scleroderma and results may not generalise to the general population. Second, in order to complete study questionnaires, patients were required to have access to the internet, which may bias the sample. Specifically, those with most severe disease may not be able to type due to the inability to use their fingers or hands. As well, it is possible that the oldest patients would be unable to participate. However, although the SPIN Cohort constitutes a convenience sample of SSc patients receiving treatment at a SPIN recruiting centre, and patients at these centres may differ from those in other settings, a comparison between SPIN Cohort participants and the European Scleroderma Trials and Research and Canadian Scleroderma Research Group cohorts showed broad comparability.²⁷ This increases confidence that insights gained from the SPIN Cohort should be generalisable.

CONCLUSION

In conclusion, this study used an iterative algorithm implemented via the *lordif* package in R to flag items on the SAAS for DIF related to language of test administration, sex and disease type. After flagging items with DIF, impact was assessed primarily by looking at GPCM score correlations and differences before and after accounting for DIF. While at least one item was flagged for DIF in each analysis, DIF impact was assessed to be small, supporting the conclusion that GPCM scores are comparable across groups produced by these three demographic variables.

Author affiliations

- ¹Applied Statistics, Social Science and Humanities, New York University, New York, New York, USA
- ²PRIISM Applied Statistics Center, New York University, New York, New York, USA
- ³Behavioral Science Institute, Radboud University, Nijmegen, The Netherlands
- ⁴Lady Davis Institute for Medical Research, Jewish General Hospital, Montreal, Quebec, Canada
- ⁵San Diego Joint Doctoral Program in Clinical Psychology, San Diego State University/University of California, San Diego, California, USA
- ⁶Psychiatry, McGill University, Montreal, Quebec, Canada
- ⁷Scleroderma Foundation, Los Angeles, California, USA
- ⁸School of Psychology, University of Minho, Braga, Portugal
- ⁹Psychology, San Diego State University, San Diego, California, USA
- ¹⁰Medicine, McGill University, Montreal, Quebec, Canada
- ¹¹Epidemiology, Biostatistics, and Occupational Health, McGill University, Montreal, Quebec, Canada
- ¹²Educational and Counselling Psychology, McGill University, Montreal, Quebec, Canada
- ¹³Psychology, McGill University, Montreal, Quebec, Canada

Collaborators

SPIN Investigators: Murray Baron (McGill University, Montreal, Quebec, Canada); Susan J Bartlett (McGill University, Montreal, Quebec, Canada); Daniel E Furst (Division of Rheumatology, Geffen School of Medicine, University of California, Los Angeles, California, USA); Maureen D Mayes (University of Texas McGovern School of Medicine, Houston, Texas, USA); Luc Mouthon (Université Paris Descartes, Paris, France); Warren R Nielson (St. Joseph's Health Care, London, Ontario, Canada); Robert Riggs (Scleroderma Foundation, Danvers, Massachusetts, USA); Maureen Sauve (Scleroderma Society of Ontario, Hamilton, Ontario); Fredrick Wigley (Johns Hopkins University School of Medicine, Baltimore, Maryland, USA); Shervin Assassi (University of Texas McGovern School of Medicine, Houston, Texas, USA); Andrea Benedetti (McGill University, Montreal, Quebec, Canada); Isabelle Boutron (Université Paris Descartes, and Assistance Publique - Hôpitaux de Paris, Paris, France); Angela Costa Maia (University of Minho, Braga, Portugal); Lindsay Cronin (Scleroderma Foundation, Western Pennsylvania Chapter, Pittsburgh, Pennsylvania, USA); Ghassan El-Baalbaki (Université du Québec à Montréal, Montreal, Quebec, Canada); Carolyn Ells (McGill University, Montreal, Quebec, Canada); Stephen Elrod (Scleroderma Foundation, Southern California Chapter, Los Angeles, California, USA); Cornelia van den Ende (Sint Maartenskliniek, Nijmegen, The Netherlands); Kim Fligelstone (Scleroderma & Raynaud's UK, London, UK); Catherine Fortune (Scleroderma Society of Ontario, Hamilton, Ontario, Canada); Tracy Frech (University of Utah, Salt Lake City, Utah, USA); Amy Gietzen (Scleroderma Foundation, Tri-State Chapter, Binghamton, New York, USA); Dominique Godard (Association des Sclérodermiques de France, Sorel-Moussel, France); Geneviève Guillot (Sclérodermie Québec, Montreal, Quebec, Canada); Shirley Haslam (Scleroderma Society of Ontario, Hamilton, Ontario, Canada); Monique Hinchcliff (Yale School of Medicine, New Haven, Connecticut, USA); Marie Hudson (McGill University, Montreal, Quebec, Canada); Ann Impens (Midwestern University, Downers Grove, Illinois, USA); Yeona Jang (McGill University, Montreal, Quebec, Canada); Sindhu R Johnson (Toronto Scleroderma Program, Mount Sinai Hospital, Toronto Western Hospital, and University of Toronto, Toronto, Ontario, Canada); Ann Tyrell Kennedy (Federation of European Scleroderma Associations, Dublin, Ireland); Annett Körner (McGill University, Montreal, Quebec, Canada); Maggie Larche (McMaster University, Hamilton, Ontario, Canada); Carlo Marra (Memorial University, St. John's, Newfoundland, Canada); Christelle Nguyen (Université Paris Descartes, and Assistance Publique - Hôpitaux de Paris, Paris, France); Karen Nielsen (Scleroderma Society of Ontario, Hamilton, Ontario, Canada); Janet Pope (University of Western Ontario, London, Ontario, Canada); Alexandra Portales (Asociación Española de Esclerodermia, Madrid, Spain); François Rannou (Université Paris Descartes, and Assistance Publique - Hôpitaux de Paris, Paris, France); Michelle Richard (Scleroderma Society of Nova Scotia, Halifax, Nova Scotia, Canada); Tatiana Sofia Rodriguez Reyna (Instituto Nacional de Ciencias Médicas y Nutrición Salvador Zubirán, Mexico City, Mexico); Ken Rozee (Scleroderma Society of Nova Scotia, Halifax, Nova Scotia, Canada); Anne A. Schouffoer (Leiden University Medical Center, Leiden, The Netherlands); Russell J Steele (Jewish General Hospital and McGill University, Montreal, Quebec, Canada); Nancy Stephens (Scleroderma Foundation, Michigan Chapter, Southfield, Michigan, USA); Maria E Suarez-Almazor (University of Texas MD Anderson Cancer Center, Houston, Texas, USA); Joep Welling (NVLE Dutch patient organization for systemic autoimmune diseases, Utrecht, The Netherlands, and Federation of European Scleroderma Associations, Brussel, Belgium); Durhane Wong-Rieger (Canadian Organization for Rare Disorders, Toronto, Ontario, Canada); Christian Agard (Centre Hospitalier Universitaire - Hôtel-Dieu de

Nantes, Nantes, France); Alexandra Albert (Université Laval, Quebec, Quebec, Canada); Marc André (Centre Hospitalier Universitaire Gabriel-Montpied, Clermont-Ferrand, France); Guylaine Arsenault (Université de Sherbrooke, Sherbrooke, Quebec, Canada); Ilham Benzidia (Assistance Publique - Hôpitaux de Paris, Hôpital St-Louis, Paris, France); Elana J Bernstein (Columbia University, New York, New York, USA); Sabine Berthier (Centre Hospitalier Universitaire Dijon Bourgogne, Dijon, France); Lyne Bissonnette (Université de Sherbrooke, Sherbrooke, Quebec, Canada); Gilles Boire (Université de Sherbrooke, Sherbrooke, Quebec, Canada); Alessandra Bruns (Université de Sherbrooke, Sherbrooke, Quebec, Canada); Patricia Carreira (Servicio de Reumatología del Hospital 12 de Octubre, Madrid, Spain); Marion Casadevall (Assistance Publique - Hôpitaux de Paris, Hôpital Cochin, Paris, France); Benjamin Chaigne (Assistance Publique - Hôpitaux de Paris, Hôpital Cochin, Paris, France); Lorinda Chung (Stanford University, Stanford, California, USA); Pascal Cohen (Assistance Publique - Hôpitaux de Paris, Hôpital Cochin, Paris, France); Chase Correia (Northwestern University, Chicago, Illinois, USA); Pierre Dagenais (Université de Sherbrooke, Sherbrooke, Quebec, Canada); Christopher Denton (Royal Free London Hospital, London, UK); Robyn Domsic (University of Pittsburgh, Pittsburgh, Pennsylvania, USA); Sandrine Dubois (Centre Hospitalier Régional Universitaire de Lille, Hôpital Claude Huriez, Lille, France); James V Dunne (St. Paul's Hospital and University of British Columbia, Vancouver, British Columbia, Canada); Bertrand Dunogue (Assistance Publique - Hôpitaux de Paris, Hôpital Cochin, Paris, France); Alexia Esquinca (Instituto Nacional de Ciencias Médicas y Nutrición Salvador Zubirán, Mexico City, Mexico); Regina Fare (Servicio de Reumatología del Hospital 12 de Octubre, Madrid, Spain); Dominique Farge-Bancel (Assistance Publique - Hôpitaux de Paris, Hôpital St-Louis, Paris, France); Paul R Fortin (CHU de Québec - Université Laval, Quebec, Quebec, Canada); Anna Gill (Royal Free London Hospital, London, UK); Jessica Gordon (Hospital for Special Surgery, New York City, New York, USA); Brigitte Granel-Rey (Aix Marseille Université, and Assistance Publique - Hôpitaux de Marseille, Hôpital Nord, Marseille, France); Claire Grange (Centre Hospitalier Lyon Sud, Lyon, France); Genevieve Gyger (Jewish General Hospital and McGill University, Montreal, Quebec, Canada); Eric Hachulla (Centre Hospitalier Régional Universitaire de Lille, Hôpital Claude Huriez, Lille, France); Pierre-Yves Hatron (Centre Hospitalier Régional Universitaire de Lille, Hôpital Claude Huriez, Lille, France); Sabrina Hoa (CHUM - Centre hospitalier de l'Université de Montréal, Montreal, Quebec, Canada); Ariane L Herrick (University of Manchester, Salford Royal NHS Foundation Trust, Manchester, UK); Adrian Hij (Assistance Publique - Hôpitaux de Paris, Hôpital St-Louis, Paris, France); Alena Ilic (Université Laval, Quebec, Quebec, Canada); Niall Jones (University of Alberta, Edmonton, Alberta, Canada); Artur Jose de B. Fernandes (Université de Sherbrooke, Sherbrooke, Quebec, Canada); Suzanne Kafaja (University of California, Los Angeles, California, USA); Nader Khalidi (McMaster University, Hamilton, Ontario, Canada); Marc Lambert (Centre Hospitalier Régional Universitaire de Lille, Hôpital Claude Huriez, Lille, France); David Launay (Centre Hospitalier Régional Universitaire de Lille, Hôpital Claude Huriez, Lille, France); Patrick Liang (Université de Sherbrooke, Sherbrooke, Quebec, Canada); Héliène Maillard (Centre Hospitalier Régional Universitaire de Lille, Hôpital Claude Huriez, Lille, France); Nancy Maltez (University of Ottawa, Ottawa, Ontario, Canada); Joanne Manning (Salford Royal NHS Foundation Trust, Salford, UK); Isabelle Marie (CHU Rouen, Hôpital de Bois-Guillaume, Rouen, France); Maria Martin (Servicio de Reumatología del Hospital 12 de Octubre, Madrid, Spain); Thierry Martin (Les Hôpitaux Universitaires de Strasbourg, Nouvel Hôpital Civil, Strasbourg, France); Ariel Masetto (Université de Sherbrooke, Sherbrooke, Quebec, Canada); François Maurier (Hôpitaux Privés de Metz, Hôpital Belle-Isle, Metz, France); Arsene Mekinian (Assistance Publique - Hôpitaux de Paris, Hôpital St-Antoine, Paris, France); Sheila Melchor (Servicio de Reumatología del Hospital 12 de Octubre, Madrid, Spain); Mandana Nikpour (St Vincent's Hospital and University of Melbourne, Melbourne, Victoria, Australia); Vincent Poindron (Les Hôpitaux Universitaires de Strasbourg, Nouvel Hôpital Civil, Strasbourg, France); Susanna Proudman (Royal Adelaide Hospital and University of Adelaide, Adelaide, South Australia, Australia); Alexis Régent (Assistance Publique - Hôpitaux de Paris, Hôpital Cochin, Paris, France); Sébastien Rivière (Assistance Publique - Hôpitaux de Paris, Hôpital St-Antoine, Paris, France); David Robinson (University of Manitoba, Winnipeg, Manitoba, Canada); Esther Rodriguez (Servicio de Reumatología del Hospital 12 de Octubre, Madrid, Spain); Sophie Roux (Université de Sherbrooke, Sherbrooke, Quebec, Canada); Perrine Smets (Centre Hospitalier Universitaire Gabriel-Montpied, Clermont-Ferrand, France); Doug Smith (University of Ottawa, Ottawa, Ontario, Canada); Vincent Sobanski (Centre Hospitalier Régional Universitaire de Lille, Hôpital Claude Huriez, Lille, France); Robert Spiera (Hospital for Special Surgery, New York City, New York, USA); Virginia Steen (Georgetown University, Washington, DC, USA); Wendy Stevens (St Vincent's Hospital and University of Melbourne, Melbourne, Victoria, Australia); Evelyn Sutton (Dalhousie University, Halifax, Nova Scotia, Canada); Benjamin Terrier (Assistance Publique - Hôpitaux de Paris, Hôpital Cochin, Paris, France); Carter Thorne (Southlake

Regional Health Centre, Newmarket, Ontario, Canada); John Varga (University of Michigan, Ann Arbor, Michigan, USA); Pearce Wilcox (St. Paul's Hospital and University of British Columbia, Vancouver, British Columbia, Canada); Angelica Bourgeault (Jewish General Hospital, Montreal, Quebec, Canada); Andrea Carboni Jiménez (Jewish General Hospital, Montreal, Quebec, Canada); Sami Harb (Jewish General Hospital, Montreal, Quebec, Canada); Lydia Tao (Jewish General Hospital, Montreal, Quebec, Canada).

Contributors DH and BDT were responsible for the study conception. LK, M-EC, VLM, BDT and the SPIN Investigators contributed to data collection. SJS, DH, LK, M-EC, SG, KG, CL, VLM, BDT contributed to data analysis and interpretation. SJS and DH drafted the manuscript. All authors provided a critical revision of the manuscript and approved the final version of the manuscript. DH is the guarantor.

Funding SPIN has been funded by grants from the Canadian Institutes of Health Research (TR3-119192, PJT-148504, PJT-149073, SCT-162963), the Canadian Initiative for Outcomes in Rheumatology Care, and the Arthritis Society. In addition, SPIN has received institutional contributions from the Female Davis Institute for Medical Research of the Jewish General Hospital, Montreal, Canada and from McGill University, Montreal, Canada. SPIN has also received support from Scleroderma Canada, the Scleroderma Society of Ontario, Sclérodermie Québec, the Scleroderma Society of Nova Scotia, the Scleroderma Association of British Columbia, Scleroderma Manitoba, and the Scleroderma Society of Saskatchewan.

Competing interests None declared.

Patient and public involvement Patients and/or the public were involved in the design, or conduct, or reporting, or dissemination plans of this research. Refer to the Methods section for further details.

Patient consent for publication Not required.

Ethics approval The SPIN Cohort study was approved by the Research Ethics Committee of the Jewish General Hospital, Montreal, Canada (MP-05-2013-150, 12–123) and by the Institutional Review Boards of each participating centre.

Provenance and peer review Not commissioned; externally peer reviewed.

Data availability statement Data are available on reasonable request. Data from the SPIN Cohort can be requested from the corresponding author.

Open access This is an open access article distributed in accordance with the Creative Commons Attribution Non Commercial (CC BY-NC 4.0) license, which permits others to distribute, remix, adapt, build upon this work non-commercially, and license their derivative works on different terms, provided the original work is properly cited, appropriate credit is given, any changes made indicated, and the use is non-commercial. See: <http://creativecommons.org/licenses/by-nc/4.0/>.

ORCID iD

Brett D Thombs <http://orcid.org/0000-0002-5644-8432>

REFERENCES

- 1 Weldring T, Smith SMS, Outcomes P-R. Pros and patient-reported outcome measures (PROMs). *Health Serv Insights* 2013;6:61–8.
- 2 Thombs BD, van Lankveld W, Bassel M, *et al*. Psychological health and well-being in systemic sclerosis: state of the science and consensus research agenda. *Arthritis Care Res* 2010;62:1181–9.
- 3 Kwakkenbos L, Willems LM, Baron M, *et al*. The comparability of English, French and Dutch scores on the functional assessment of chronic illness Therapy-Fatigue (FACIT-F): an assessment of differential item functioning in patients with systemic sclerosis. *PLoS One* 2014;9:e91979.
- 4 Teresi JA. Overview of quantitative measurement methods. equivalence, invariance, and differential item functioning in health applications. *Med Care* 2006;44:S39–49.
- 5 Holland PW, Wainer H. *Differential item functioning*. Taylor & Francis, 2012.
- 6 Allanore Y, Simms R, Distler O, *et al*. Systemic sclerosis. *Nat Rev Dis Primers* 2015;1:15002.
- 7 Clements PJ, Furst DE. *Systemic sclerosis*. 2 ed. Baltimore: Lippincott Williams & Wilkins, 2003.
- 8 Gholizadeh S, Fox RS, Mills SD, *et al*. *Coping with the disfigurement of scleroderma: facial, skin, and hand changes*. Scleroderma. . Springer international publishing, 2017: 2. 713–21.
- 9 Kwakkenbos L, Delisle VC, Fox RS, *et al*. Psychosocial aspects of scleroderma. *Rheum Dis Clin North Am* 2015;41:519–28.
- 10 Hart TA, Flora DB, Palyo SA, *et al*. Development and examination of the social appearance anxiety scale. *Assessment* 2008;15:48–59.

- 11 Kwakkenbos L, Jewett LR, Baron M, *et al.* The scleroderma patient-centered intervention network (spin) cohort: protocol for a cohort multiple randomised controlled trial (cmRCT) design to support trials of psychosocial and rehabilitation interventions in a rare disease context. *BMJ Open* 2013;3:e003563.
- 12 Harel D, Mills SD, Kwakkenbos L, *et al.* Shortening patient-reported outcome measures through optimal test assembly: application to the social appearance anxiety scale in the scleroderma patient-centered intervention network cohort. *BMJ Open* 2019;9:e024010.
- 13 Harel D, Baron M. Methods for shortening patient-reported outcome measures. *Stat Methods Med Res* 2019;28:2992–3011.
- 14 van den Hoogen F, Khanna D, Fransen J, *et al.* 2013 classification criteria for systemic sclerosis: an American College of rheumatology/ European League against rheumatism collaborative initiative. *Ann Rheum Dis* 2013;72:1747–55.
- 15 LeRoy EC, Black C, Fleischmajer R, *et al.* Scleroderma (systemic sclerosis): classification, subsets and pathogenesis. *J Rheumatol* 1988;15:202–5.
- 16 Mills SD, Kwakkenbos L, Carrier M-E, *et al.* Validation of the social appearance anxiety scale in patients with systemic sclerosis: a scleroderma patient-centered intervention network cohort study. *Arthritis Care Res* 2018;70:1557–62.
- 17 Wild D, Grove A, Martin M, *et al.* Principles of good practice for the translation and cultural adaptation process for patient-reported outcomes (pro) measures: report of the ISPOR Task force for translation and cultural adaptation. *Value Health* 2005;8:94–104.
- 18 Muraki E. A generalized partial credit model: application of an em algorithm. *Appl Psychol Meas* 1992;16:159–76.
- 19 R Core Team. *R: a language and environment for statistical computing*. Vienna, Austria: R Foundation for Statistical Computing, 2019.
- 20 Choi SW, Gibbons LE, Crane PK. lordif: an R package for detecting differential item functioning using iterative hybrid ordinal logistic Regression/Item response theory and Monte Carlo simulations. *J Stat Softw* 2011;39:1–30.
- 21 Chung H, Kim J, Cook KF, *et al.* Testing measurement invariance of the patient-reported outcomes measurement information system pain behaviors score between the US general population sample and a sample of individuals with chronic pain. *Qual Life Res* 2014;23:239–44.
- 22 Cook KF, Bombardier CH, Bamer AM, *et al.* Do somatic and cognitive symptoms of traumatic brain injury confound depression screening? *Arch Phys Med Rehabil* 2011;92:818–23.
- 23 Stocking ML, Lord FM. Developing a common metric in item response theory. *Appl Psychol Meas* 1983;7:201–10.
- 24 Rupp AA, Zumbo BD. A note on how to quantify and report whether irt parameter invariance holds: when Pearson correlations are not enough. *Educ Psychol Meas* 2004;64:588–99.
- 25 Harel D. *The effect of model misspecification for polytomous logistic adjacent category item response theory models*. McGill University Libraries, 2014.
- 26 van der Ark LA. Stochastic ordering of the latent trait by the sum score under various Polytomous irt models. *Psychometrika* 2005;70:283–304.
- 27 Dougherty DH, Kwakkenbos L, Carrier M-E, *et al.* The scleroderma patient-centered intervention network cohort: baseline clinical features and comparison with other large scleroderma cohorts. *Rheumatology* 2018;57:1623–31.