# Differential Item Functioning in the Unified Dyskinesia Rating Scale (UDysRS)

**Sheng Luo, PhD**[1], **Yuanyuan Liu, MS**[1], **Jeanne A. Teresi, Ed.D., Ph.D.**[2,3], **Glenn T. Stebbins, PhD**[4], and **Christopher G. Goetz, MD**[4]

[1]Department of Biostatistics, University of Texas Health Science Center, School of Public Health, Houston, TX, USA

[2]Columbia University Stroud Center at New York State Psychiatric Institute, New York, NY, USA

[3]Research Division, Hebrew Home at Riverdale, Bronx, NY, USA

[4]Department of Neurological Sciences, Rush University Medical Center, Chicago, IL USA

## Abstract

**Objective**—Test if differential item functioning due to gender, age, race/ethnicity or education impacts Unified Dyskinesia Rating Scale scores.

**Background**—Testing rating scales for differential item functioning is a core validation step. If differential item functioning exists, interpretation of item scores must consider secondary influences on dyskinesia ratings.

**Methods**—Using Unified Dyskinesia Rating Scale translation databases (N=3,132), we tested uniform and non-uniform differential item functioning. We required confirmation by two independent methods and considered differential item functioning pertinent if McFadden pseudo $R^2$ magnitude statistics exceeded negligible ratings.

**Results**—No age, race/ethnicity or education non-uniform differential item functioning was identified. Gender non-uniform differential item functioning occurred for two items, both with negligible magnitude. Gender, race, and education uniform differential item functioning was observed for multiple items, all with negligible magnitude.

**Conclusions**—The Unified Dyskinesia Rating Scale items effectively capture dyskinesia severity without pertinent gender, age, race/ ethnicity or education influence.

## Keywords

Parkinson's disease; dyskinesia; Rating Scales; Clinimetrics; Differential Item Functioning

## Introduction

The Unified Dyskinesia Rating Scale (UDysRS) was developed as a comprehensive rating tool of dyskinesia in Parkinson's Disease (PD) [1]. The scale was developed in English with a

Corresponding author: Sheng Luo, PhD, The University of Texas Health Science Center at Houston, 1200 Herman Pressler Dr, Rm E815, Houston, TX 77030, USA, Telephone: 713-500-9554; FAX: 713-500-9525; sheng.t.luo@uth.tmc.edu.

clinimetric program to provide validated non-English translations[2, 3]. Testing a rating scale for Differential Item Functioning (DIF) [4] is a core step to determine if covariates (e.g., age, gender) substantially bias any item score. Among people with similar severity levels of dyskinesia and the same probability of responding, DIF occurs for the UDysRS if the probability of an item-score differs according to selected covaraites. For example, gender-based DIF exits for item 4.1 (Communication disability related to dyskinesia) if men and women with the same severity level of dyskinesia responded differently on this item. Two kinds of DIF can occur. In non-uniform DIF (NU-DIF), covariate influences on item-scores vary across levels of the dyskinesia trait, while in uniform DIF (U-DIF), influences on item-scores by the covariate are constant across all trait levels (Figure S1 of Supplementary Material) [5].

We conducted both U-DIF and NU-DIF assessments on UDysRS items on the gender, age, race/ethnicity, and education level [6]. The absence of clinically relevant NU-DIF or U-DIF allows it to be used confidently as a true measure of dyskinesia.

## Methods

### The UDysRS dataset

We accessed the cross-sectional combined translation dataset of fully completed UDysRS scores from 13 languages (Chinese [N=250], English [N=70], French [N=250], German [N=284], Greek [N=260], Hungarian [N=256], Italian [N=252], Japanese [N=250], Korean [N=250], Portuguese [N=256], Slovak [N=251], Spanish [N=253], Turkish [N=250])[3].

### Assessing unidimensionality of the UDysRS

DIF analyses are anchored in the unidimensionality assumption, i.e., the items measure a single pertinent trait. To test unidimensionality of UDysRS, we conducted confirmatory factor analysis, requiring that the Confirmatory Fit Index (CFI) was 0.90 with Root Means Square Error of Approximation (RMSEA) $< 0.10$[7].

### Sample sizes for each analysis

DIF analyses require that for each item, all possible rating values must have some representation. Because there were no patients scoring in the most severe rating option (4) in many UDysRS items, we combined scores of 3 and 4 as a collapsed designation, termed 3/4. Further, we required at least 5 subject samples in each of the 0, 1, 2, and 3/4 categories for each UDysRS item.

### DIF determinations

We conducted DIF analysis using two independent latent variable models, the iterative hybrid ordinal logistic regression/item response theory (graded response model) [8] approach as realized in the *R* package *lordif* [9] and the MIMIC model [10, 11]. For an item to qualify for DIF designation, we required that both methods independently identify DIF at a significance level corrected for multiple comparisons using a Bonferroni correction [12].

All items were studied first for NU-DIF and those without NU-DIF were then analyzed for U-DIF[12]. For items identified with DIF, to determine clinical pertinence (DIF magnitude) , we used the McFadden pseudo $R^2$ magnitude estimate from the *R* package *lordif* and applied the recommended cut-offs of <0.035=negligible; 0.035–0.07=moderate; >0.07=large [13]. We considered an item with DIF to be clinically relevant if it exceeded negligible rating. Finally, we examined the combined impact (Scale Level Impact) of multiple identified items with DIF on the UDysRS using the Differential Test Function (DTF) index that compared the Test Characteristic Curves with and without DIF items [14]. The magnitude of the DTF [14] was assessed by a conservative threshold based on Monte Carlo simulations [15, 16] (cutoff DTF value = 1.404).

### Comparisons

For gender, the analyses compared males and females. For the age-based DIF analyses, we chose three age groups (ages 27–51, ages 52–75, and ages 76–93) to result in at least 280 cases in each age group. We chose this age divisions to reflect our age ranges (27–93), and they are similar to other reports examining age divisions in PD[17, 18]. Based on years of education, we divided the sample into three groups (<7, 7–12, >12), which resulted in 680 cases in each education group.[19] We chose race/ethnicity categories according to published divisions adopted by the US Office of Management and Budget [6]. Possible categories were: White (non-Hispanic), Hispanic, African descent, Asian, Pacific Islander, Native or Endogenous, and Other. Whereas the *lordif* model can accommodate multinomal options, MIMIC is restricted to binary comparisons. Therefore, we first conducted comparisons using *lordif,* and*,* if overall DIF was identified with this strategy*,* follow-up pairwise comparisons were conducted in *lordif* and MIMIC independently.

## Results

### Sample Sizes

The full dataset included UDysRS scores for 3,132 subjects, but missing data on isolated items or demographic information reduced the samples. In all assessments, however, the sample exceeded 2,500 UDysRS complete scores (Table S1 of Supplementary Material).

### Unidimensionality

The confirmatory factor analysis of the full dataset confirmed unidimensionality of the UDysRS. The scale met the criteria of a CFI    0.90 and a RMSEA < 0.10, allowing conduct of the DIF analyses[7] (CFI = 0.97, RMSEA = 0.08).

### Gender-based DIF (Upper part of Table 1)

NU-DIF for gender was identified in one Historical Disability item (Speech) and one Objective Disability item (Communication). U-DIF for gender was identified in one Historical Disability item (Time with Dyskinesia). In all cases, the magnitude of the DIF was "negligible." In assessing the combined effects of multiple "negligible" impacts, we did not detect an overall Scale Level Impact on UDysRS from gender-based DIF using the DTF index score (DTF=0.0214). (Supplementary Material provides all results for identified DIF).

### Education-based DIF (lower part of Table 1)

None of the items exhibited NU-DIF for education. Education-based U-DIF was found for Historical Disability ratings for Time Spent with On-Dyskinesia, Chewing/Swallowing, Eating Tasks, Dressing and Hygiene, though in all cases the magnitude of the DIF was "negligible." We did not detect an overall Scale Level Impact on UDysRS from education-based DIF using the DTF index score (<7 vs. all others=0.4285; 7–12 vs. all others=0.0144; >12 vs. all others=1.1297). (Supplementary Material provides all results for identified DIF).

### Age-based DIF

For age-based DIF, none of the Items was identified as having NU-DIF or U-DIF.

### Race/ethnicity-based DIF (Table 2)

The racial/ethnic groups under consideration were White non-Hispanic, Hispanic, and Asian. We did not have a sufficiently large score representation from other groups. For race/ethnicity-based DIF, none of the items was identified as having NU-DIF. Fourteen items exhibited race/ethnicity-based U-DIF for White vs other (Historical Disability ratings for Exciting or Emotional Settings, Effects of Pain from Off-Dystonia and Dystonia Pain; Objective Impairment ratings for Face and Right Leg/Hip; and, Objective Disability ratings for Drinking and Ambulation), Asian vs other (Historical Disability ratings for Exciting and Emotional Settings, Time Off Dystonia, Effects of Off-Dystonia Separate from Pain and Effects of Pain from Off-Dystonia and Dystonia Pain; Objective Impairment ratings for Face, Right Arm/Shoulder, Left Arm/Shoulder, Right Leg/Hip, and Left Leg/Hip; and Objective Disability ratings of Drinking), and Hispanic vs other (Historical Disability ratings for Eating Tasks and Public/Social Settings; and Objective Disability ratings for Ambulation). In all cases, the impact of U-DIF was negligible. We did not detect an overall Scale Level Impact on UDysRS using the DTF index score when comparing White vs. non-White (DTF = 0.2036) and Hispanic vs. non-Hispanic (DTF = 1.0501). The DTF simulation-based threshold was exceeded for Asian vs. non-Asian (DTF=5.0038). (Supplementary Material provides all results for identified DIF).

## Discussion

DIF, often termed "measurement bias",[12, 14–16, 20] is essential to test for a full validation of a rating scale and the confident conclusion that the scale is truly measuring the conceptual trait, in this case, dyskinesia severity. The fact that we did not detect DIF of moderate or large magnitude for any item relative to any of the studied demographic elements strongly argues that the UDysRS is effectively capturing dyskinesia severity and is not strongly influenced by gender, age, race/ethnicity or education. The conclusion is reinforced by our inability to detect a significant combined Scale Level Impact when multiple "negligible" DIF items occur in the scale. The DTF value above threshold observed for the Asian subsample indicated a small level of impact as evidenced in the graphs of the test characteristic curves for Asians and non-Asians. Although the level of aggregate impact was not sufficient to warrant concern, it is recommended that this finding be investigated further with other datasets.

There are two major differences between this study and our prior MDS-UPDRS DIF analysis[21]. First, because of the unidimensionality of UDysRS, we could justify performing DIF using the total UDysRS score as the index of dyskinesia severity. In the MDS-UPDRS, because the scale is unidimensional for each Part, but not as a total score, our approach necessitated DIF analysis for each Part. Second, although MDS-UPDRS items were not assessed for education-based DIF due to the lack of education information, we can add to our conclusions that the UDysRS scale item performance is not influenced by education level. We acknowledge that educational systems differ by culture, so the interpretation of DIF absence based on education is limited to conclusions regarding number of years of formal education and not knowledge base.

Although the sample sizes were very large, we were limited by the paucity of item-scores in the severe impairment and disability category (4), because all assessments were acquired in outpatient settings where the most severe patients are rarely seen. Hence, we collapsed 3 and 4 categories into a single designation, which may not achieve DIF analysis of the UDysRS as constructed. Moreover, DIF may exist from other covariates such as source of information for Parts 1 and 2 (patient, caregiver, or combined patient/caregiver) and rater- or site-based DIF. Our current dataset precluded such additional DIF analysis.

The strengths of our study include the very large dataset with worldwide representation across cultures using one validated scale. We have been rigorous in our clinimetric approach, requiring that designated items with DIF be identified by two independent statistical methods with correction for multiple comparisons. Using the McFadden's $R^2$ allows us to interpret the magnitude of identified DIF. The results suggest that the items composing the full UDysRS are highly specific to dyskinesia severity. With the negligible contributions from age, gender, race/ethnicity, and education level, the scale can be viewed as widely applicable and not impacted by these demographic indices.

## Supplementary Material

Refer to Web version on PubMed Central for supplementary material.

## Acknowledgments

## Author roles

Sheng Luo:

Research project - conception, organization, execution

Statistical analysis – design, supervision, conduct and review

Manuscript preparation – writing of the first draft, review and critique

Yuanyuan Liu:

Statistical analysis – conduct and review

Manuscript preparation – writing of the first draft, review and critique

Jeanne Teresi:

Research project-conceptual and methodological approach

Statistical analysis –conception and review

Manuscript preparation-review and critique

Glenn T. Stebbins:

Research project - conception, organization, execution

Statistical analysis - design, review and critique

Manuscript preparation – writing of first draft, review and critique

Christopher G. Goetz:

Research project - conception, organization, execution

Statistical analysis - design, review and critique

Manuscript preparation - writing of the first draft, critique and review

## Financial Disclosers for the past 12 months

**Sheng Luo, PhD**

**Consulting**: none

**Grants/Research**: NIH grants (R01NS091307, 5U01NS043127), grants from CHDI Foundation, Parkinson Disease Foundation and International Parkinson and Movement Disorder Society

**Honoraria:** none

**Intellectual Property Rights**: none

**Ownership interests**: none

**Royalties**: none

**Salary:** University of Texas School of Public Health: (MDS support for project included)

**Yuanyuan Liu, MS**

**Consulting:** none

**Honoraria:** none

**Intellectual Property Rights**: none

**Ownership interests**: none

**Royalties**: None

**Salary:** University of Texas School of Public Health (MDS support for project included); UTHealth Innovation for Cancer Prevention Research Training Program Predoctoral Fellowship (Cancer Prevention and Research Institute of Texas grant # RP160015)

### Jeanne A. Teresi, Ed.D., Ph.D.

**Consulting or Advisory Board Membership with honoraria:** none

**Grants/Research:** National Institute on Aging, National Institute of Minority Health Disparities, Patient Reported Outcomes Research Institute, Health Services Organization and Delivery

**Honoraria:** none

**Intellectual Property Rights:** none

**Ownership interests:** none

**Royalties:** none

**Salary:** New York State Office of Mental Health; Hebrew Home at Riverdale

### Glenn T. Stebbins, PhD

**Consulting and Advisory Board Membership with honoraria**: Acadia, Pharmaceuticals, Adamas Pharmaceuticals, Inc., Ceregene, Inc., CHDI Management, Inc., Ingenix Pharmaceutical Services (i3 Research), Neurocrine Biosciences, Inc., Pfizer, Inc., Ultragenyx.

**Grants and Research**: National Institutes of Health, Michael J. Fox Foundation for Parkinson's Research, Dystonia Coalition, CHDI, International Parkinson and Movement Disorder Society, CBD Solutions.

**Honoraria**: International Parkinson and Movement Disorder Society, American Academy of Neurology, Michael J. Fox Foundation for Parkinson's Research, Food and Drug Administration, National Institutes of Health.

**Intellectual Property Rights**: none

**Ownership interests**: none

**Royalties**: none

**Expert Testimony**: none

**Salary**: Rush University Medical Center

**Christopher G. Goetz, MD**

**Consulting or Advisory Board Membership with honoraria:** Acadia, Addex, Avanir, Boston Scientific, Neurocrine, Oxford Biomedica, WebMD.

**Grants/Research:** Funding to Rush University Medical Center from NIH, Michael J. Fox Foundation for research conducted by Dr. Goetz. Dr. Goetz directs the Rush Parkinson's Disease Research Center that receives support from the Parkinson's Disease Foundation and some of these funds support Dr. Goetz's salary as well as his research efforts. He directed the translation program for the MDS-UPDRS and UDysRS and received funds directed to Rush University Medical Center from the International Parkinson and Movement Disorder Society (IPMDS) for this effort.

**Honoraria: Oregon Health and Science University**

**Intellectual Property Rights: none**

**Ownership interests: none**

**Royalties: Elsevier Publishers, Oxford University Press, Wolters Kluwer,**

**Salary: Rush University Medical Center**

# References

1. Goetz CG, Nutt JG, Stebbins GT. The Unified Dyskinesia Rating Scale: presentation and clinimetric profile. Movement disorders : official journal of the Movement Disorder Society. 2008; 23(16): 2398–2403. [PubMed: 19025759]

2. Colosimo C, Martínez-Martín P, Fabbrini G, et al. Task force report on scales to assess dyskinesia in Parkinson's disease: critique and recommendations. Movement Disorders. 2010; 25(9):1131–1142. [PubMed: 20310033]

3. Goetz CG, Stebbins GT, Wang L, LaPelle NR, Luo S, Tilley BC. IPMDS-Sponsored Scale Translation Program: Process, Format, and Clinimetric Testing Plan for the MDS-UPDRS and UDysRS. Movement disorders clinical practice. 2014; 1(2):97–101. [PubMed: 27747259]

4. Hambleton RK. Good practices for identifying differential item functioning. Medical care. 2006; 44(11):S182–S188. [PubMed: 17060826]

5. Mellenbergh GJ. Contingency table models for assessing item bias. Journal of educational statistics. 1982; 7(2):105–118.

6. Office of Management and Budget. DIRECTIVE NO. 15 Race and Ethnic Standards for Federal Statistics and Administrative Reporting. 1977.

7. Brown, T. Confirmatory factor analysis for applied research. Kenny, DA., editor. New York, NY: The Guilford Press; 2006.

8. Samejima F. Estimation of latent ability using a response pattern of graded scores. Psychometrika monograph supplement. 1969

9. 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. Journal of statistical software. 2011; 39(8):1.

10. Muthén B. A general structural equation model with dichotomous, ordered categorical, and continuous latent variable indicators. Psychometrika. 1984; 49(1):115–132.

11. Jöreskog KG, Goldberger AS. Estimation of a model with multiple indicators and multiple causes of a single latent variable. journal of the American Statistical Association. 1975; 70(351a):631–639.

12. Teresi JA. Different approaches to differential item functioning in health applications: Advantages, disadvantages and some neglected topics. Medical care. 2006; 44(11):S152–S170. [PubMed: 17060822]

13. Jodoin MG, Gierl MJ. Evaluating type I error and power rates using an effect size measure with the logistic regression procedure for DIF detection. Applied Measurement in Education. 2001; 14(4): 329–349.

14. Roju NS, Van der Linden WJ, Fleer PF. IRT-based internal measures of differential functioning of items and tests. Applied Psychological Measurement. 1995; 19(4):353–368.

15. Flowers CP, Oshima T, Raju NS. A description and demonstration of the polytomous-DFIT framework. Applied Psychological Measurement. 1999; 23(4):309–326.

16. Kleinman M, Teresi JA. Differential item functioning magnitude and impact measures from item response theory models. Psychological Test and Assessment Modeling. 2016; 58(1):79–98. [PubMed: 28706769]

17. van Rooden S, Verbaan D, Stijnen T, Marinus J, Van Hilten J. The influence of age and approaching death on the course of nondopaminergic symptoms in Parkinson's disease. Parkinsonism & related disorders. 2016; 24:113–118. [PubMed: 26774535]

18. Keezer MR, Wolfson C, Postuma RB. Age, gender, comorbidity, and the MDS-UPDRS: results from a population-based study. Neuroepidemiology. 2016; 46(3):222–227. [PubMed: 26967747]

19. Statistics UIf. International Standard Classification of Education: ISCED 2011: UIS. Montreal, Quebec: 2012.

20. Embretson, SE., Reise, SP. Item response theory. Psychology Press; 2013.

21. Goetz CG, Liu Y, Stebbins GT, et al. Gender-, age-, and race/ethnicity-based differential item functioning analysis of the movement disorder society–sponsored revision of the Unified Parkinson's disease rating scale. Movement Disorders. 2016; 31(12):1865–1873. [PubMed: 27943473]

Author Manuscript

Author Manuscript

Author Manuscript

Author Manuscript

**Table 1**

Gender- and Education-Based Statistically Significant DIF

| Item | MIMIC p-values | LORDIF p-values | R² | Magnitude |
|------|---------------|-----------------|-----|-----------|
| **Gender-Based Non-Uniform DIF** | | | | |
| Historical Disability Speech | <0.0005 | 0.0001 | 0.0019 | Negligible |
| Objective Disability Communication | <0.0005 | <0.00005 | 0.0025 | Negligible |
| **Gender-Based Uniform DIF** | | | | |
| Historical Disability Time Spent with On-Dyskinesia | <0.0005 | 0.0005 | 0.0015 | Negligible |
| **Education-Based Non-Uniform DIF** | | | | |
| None | NA | NA | NA | NA |
| **Education-Based Uniform DIF** | | | | |
| Historical Disability Time Spent with On-Dyskinesia | | | | |
|    <7 vs. all others | <0.0005 | <0.00005 | 0.003 | Negligible |
| Historical Disability Chewing/Swallowing | | | | |
|    >12 vs. all others | <0.0005 | <0.00005 | 0.0035 | Negligible |
| Historical Disability Eating tasks | | | | |
|    >12 vs. all others | <0.0005 | 0.00001 | 0.0021 | Negligible |
| Historical Disability Dressing | | | | |
|    >12 vs. all others | <0.0005 | <0.00005 | 0.0031 | Negligible |
| Historical Disability Hygiene | | | | |
|    <7 vs. all others | <0.0005 | <0.00005 | 0.0057 | Negligible |
|    >12 vs. all others | <0.0005 | <0.00005 | 0.0079 | Negligible |

Legend: Most of the UDysRS items did not meet the minimal statistical criteria for DIF (see text). The Table lists items with DIF identified by both *lordif* and *MIMIC* as independent approaches (p values shown) with McFadden's $R^2$ ($R^2$) indicating the impact of the DIF.

**Table 2**

Race/Ethnicity-Based Statistically Significant DIF

| Race/Ethnicity-Based Non-Uniform DIF | | | | |
|---|---|---|---|---|
| Item | MIMIC p-values | LORDIF p-values | $R^2$ | Magnitude |
| None | NA | NA | NA | NA |
| **Race/Ethnicity-Based Uniform DIF** | | | | |
| Item | MIMIC p-values | LORDIF p-values | $R^2$ | Magnitude |
| Historical Disability - Eating Tasks | | | | |
|     Hispanic vs. all others | <0.0005 | <0.00005 | 0.0023 | Negligible |
| Historical Disability Public/Social Settings | | | | |
|     Hispanic vs. all others | <0.0005 | 0.0003 | 0.0015 | Negligible |
| Historical Disability Exciting or Emotional Settings | | | | |
|     White vs. all others | <0.0005 | <0.00005 | 0.0047 | Negligible |
|     Asian vs. all others | <0.0005 | <0.00005 | 0.0072 | Negligible |
| Historical Disability Time with Off-Dystonia | | | | |
|     Asian vs. all others | <0.0005 | <0.00005 | 0.0023 | Negligible |
| Historical Disability Effects of Off-Dystonia Separate from Pain | | | | |
|     Asian vs. all others | <0.0005 | <0.00005 | 0.0028 | Negligible |
| Historical Disability Effects of Pain from Off-Dystonia | | | | |
|     White vs. all others | <0.0005 | <0.00005 | 0.0018 | Negligible |
|     Asian vs. all others | <0.0005 | <0.00005 | 0.0099 | Negligible |
| Historical Disability Dystonia Pain | | | | |
|     White vs. all others | <0.0005 | <0.00005 | 0.0027 | Negligible |
|     Asian vs. all others | <0.0005 | <0.00005 | 0.0149 | Negligible |
| Objective Impairment Face | | | | |
|     White vs. all others | <0.0005 | <0.00005 | 0.0078 | Negligible |
|     Asian vs. all others | <0.0005 | <0.00005 | 0.0085 | Negligible |
| Objective Impairment Right Arm/Shoulder | | | | |
|     Asian vs. all others | <0.0005 | <0.00005 | 0.0027 | Negligible |
| Objective Impairment Left Arm/Shoulder | | | | |
|     Asian vs. all others | <0.0005 | <0.00005 | 0.0052 | Negligible |
| Objective Impairment Right Leg/Hip | | | | |
|     White vs. all others | <0.0005 | <0.00005 | 0.0031 | Negligible |
|     Asian vs. all others | <0.0005 | <0.00005 | 0.0119 | Negligible |

| Race/Ethnicity-Based Non-Uniform DIF | | | | |
|---|---|---|---|---|
| Item | MIMIC p-values | LORDIF p-values | $R^2$ | Magnitude |
| **Objective Impairment Left Leg/Hip** | | | | |
| Asian vs. all others | <0.0005 | <0.00005 | 0.0091 | Negligible |
| **Objective Disability Drinking** | | | | |
| White vs. all others | <0.0005 | <0.00005 | 0.0044 | Negligible |
| Asian vs. all others | <0.0005 | <0.00005 | 0.0046 | Negligible |
| **Objective Disability Ambulation** | | | | |
| White vs. all others | <0.0005 | 0.0001 | 0.002 | Negligible |
| Hispanic vs. all others | <0.0005 | <0.00005 | 0.0037 | Negligible |

Legend: Most of the UDysRS items did not meet the minimal statistical criteria for DIF (see text). The Table lists items with DIF identified by both *lordif* and *MIMIC* as independent approaches (p values shown) with McFadden's $R^2$ ($R^2$) indicating the impact of the DIF.