1    Comparison of reliability, construct validity and responsiveness of the IPAQ-SF and PASE in adults

2    with osteoarthritis.

3    **Abstract**

4    *Introduction*

5    This study assessed the measurement properties of two commonly used self-report physical activity

6    (PA) measures: the International Physical Activity Questionnaire- Short Form (IPAQ-SF) and the

7    Physical Activity Scale for the elderly (PASE) in adults with osteoarthritis.

8    *Methods*

9    Secondary analysis of the MOSAICS cluster randomised controlled trial baseline and 3-month follow

10    up questionnaires; total scores and subdomains of the IPAQ-SF and PASE were compared. Intra-class

11    correlations (ICC) were used to assess test-retest reliability, measurement error was assessed using

12    standard error of measurement (SEM), smallest detectable change (SDC) and 95% limits of

13    agreement (LoA). Responsiveness was assessed using effect size (ES), standard responsive

14    measurement (SRM) and response ratio (RR).

15    *Results*

16    There was moderate correlation (r=0.56) between the total IPAQ-SF scores (score ranges 0-16398)

17    and the total PASE scores (score ranges 0-400). Subdomain correlations were also moderate (ranges

18    0.39-0.57). The PASE showed greater reliability compared to the IPAQ-SF (ICC=0.68; 0.61-0.74 95%CI

19    & ICC=0.64; 0.55-0.72, respectively). Measurement error in both measures were large: PASE

20    SEM=46.7, SDC=129.6 and 95% LoA ranges=-117 to 136, the IPAQ-SF SEM=3532.2 METS$^{-1mins-1week}$,

21    SDC= 9790.8 and 95% LoA ranges=-5222 to 5597. Responsiveness was poor: ES -0.14 and -0.16, SRM

22    -0.21 and -0.21 and RR 0.12 and 0.09 for the IPAQ-SF and PASE respectively.

23    *Discussion*

1    The IPAQ-SF and PASE appear limited in reliability, measurement error and responsiveness.

2    Researchers and clinicians should be aware of these limitations, particularly when comparing

3    different levels of PA and monitoring PA levels changes over time in those with osteoarthritis.

4    **Keywords**

5    Measurement properties, osteoarthritis, physical activity, IPAQ, PASE

6    **Introduction**

7    Osteoarthritis (OA) is the most common cause of peripheral joint pain in adults aged 45 years and

8    over (Felson, 2009) . It is a clinical syndrome of joint pain causing varying degrees of limitation in

9    physical function and reduced quality of life (National Institute for Health and Care Excellence (NICE)

10   2014). Physical activity (PA) interventions are effective in reducing pain and improving physical

11   functioning among those with lower limb OA (Fransen et al., 2015; Fransen, McConnell, Hernandez-

12   Molina, & Reichenbach, 2014; Holden et al., 2015; Uthman et al., 2013). In addition, a physically

13   active lifestyle has wider health benefits reducing the risk of premature mortality, disability, chronic

14   diseases and mental health conditions (Warburton, Nicol, & Bredin, 2006). Due to its benefits, PA is

15   recommended as a core treatment for all adults with OA (NICE, 2014). However, current levels of

16   physical activity in people with OA are low (Herbolsheimer et al., 2016).

17   Despite PA recommendations, it is not clear what frequency, duration or intensity of PA is required

18   for those with OA to gain clinically important benefits in pain and physical functioning (Quicke,

19   Foster, Croft, Ogollah, & Holden, 2018). In order to draw inferences regarding these parameters of

20   PA, self-report measures also need to accurately capture these elements of PA. Valid and reliable

21   self-report measures of PA are required to establish current PA levels in people with OA, whilst

22   responsive measures are necessary to detect changes over time or following interventions.

23   There are two main approaches available for measuring levels of PA: objective measures (doubly

24   labelled water, indirect calorimetry or activity monitors such as pedometers and accelerometers)

1    and subjective measures (self-report questionnaires and activity diaries). Self-report measures of PA

2    are an attractive approach as they are inexpensive to administer, have the potential to measure all

3    forms of PA, can be self-completed and so used in large population level studies (Prince et al., 2008).

4    The limitations of self-report measures are that they can both over and underestimate levels of PA

5    due to reporting bias, recall bias and social desirability bias which may affect the reliability and

6    validity of their measurement of PA. Poor validity has been most notably identified at higher

7    intensities of PA where greater differences were found between self-report and objective measures

8    (Prince et al., 2008; Silsbury, Goldsmith, & Rushton, 2015).

9    Self-report PA measures are commonly used in OA research. A recent systematic review showed that

10   self-report measures of PA have been used in 91 studies since 2018 (Smith et al., 2019). The

11   International Physical Activity Questionnaire - Short Form (IPAQ-SF) and the Physical Activity Scales

12   for the Elderly (PASE) have been most commonly used. However, recent reviews identified a lack of

13   evidence for the validity and reliability of these measures in adults with OA (Healey et al., 2020;

14   Smith et al., 2019). Importantly, to date, neither the IPAQ-SF nor PASE have been assessed for

15   responsiveness (ability to detect changes in levels of PA over time). The aim of this study was to

16   evaluate and compare the test-retest reliability, construct validity and responsiveness of the IPAQ-SF

17   and PASE in adults with OA in the hand, hip, knee or foot.

18   **Methods**

19   *Design*

20   This study was conducted using secondary analysis of data from a cluster-randomised controlled trial

21   (RCT): Managing OSteoArthritis In ConsultationS (MOSAICS) study (ISRCTN number:

22   ISRCTN06984617). The protocol and results of the MOSAICS study are reported elsewhere (Dziedzic

23   et al., 2018; Dziedzic et al., 2014). Briefly, the MOSAICS study was a two-arm cluster Randomised

24   Controlled Trial (RCT) conducted in eight UK general practices. Adults aged 45 years and over, who

25   were registered with participating general practices were mailed a health survey. Participants that

1 reported peripheral joint pain who consented to follow-up and consulted their General Practitioner

2 (GP) for joint pain were invited to take part in the cluster RCT. Participants in the intervention

3 practices received an enhanced GP consultation, an OA self-management guidebook and were

4 offered follow-up practice nurse consultations, where core-recommended treatments were

5 delivered (Grime & Dudley, 2014). Participants in the control arm practices received usual care.

6 *Participants*

7 All 525 participants from the MOSAICS study cluster trial (288 from intervention practices and 237

8 from control practices), were included in this study. Inclusion criteria for the MOSAICS study were

9 adults aged 45 and over, consulting with joint pain in the hand, hip, knee or foot, at a participating

10 practice. Exclusion criteria comprised those who were: screened as ineligible by GP screening of

11 practice list; unable to give consent; a resident in a nursing home; had a history of serious disease

12 (malignancy, terminal illness), unable to consult to their GP or flagged on the practice list as

13 excluded from research (Dziedzic et al., 2014). Data from the baseline and 3-month follow-up of the

14 MOSAICS questionnaires were used for this study.

15 *Physical activity measures*

16 *The PASE*

17 The PASE is designed specifically to measure levels of PA in adults aged 65 and over. The PASE gives

18 an output score ranging from 0-400. The scoring of the PASE does not represent a quantifiable

19 amount of activity time which could indicate if a participant is participating in low, moderate or high

20 amounts of physical activity or meeting guideline-recommended activity levels, instead higher scores

21 represent higher levels of PA.  The PASE is a short self-report measure with items on PA in leisure,

22 occupational and household settings (Washburn, Smith, Jette, & Janney, 1993). The measure was

23 scored according to the instrument scoring guideline. Outliers of the PASE were checked if

24 individuals' total score exceeded 400 as recommended within the scoring guideline.

1    *The IPAQ-SF*

2    The IPAQ-SF was developed as an outcome to measure levels of PA for comparisons across

3    international populations. The IPAQ-SF measures energy expenditure per week (METS$^{-1mins-1week}$) and

4    can give a continuous or categorical score rating of an individual's weekly PA level; low (<600 METS$^{-}$

5    $^{1mins-1week}$), moderate (≥600-2999 METS$^{-1mins-1week}$) and high PA levels (≥3000 METS$^{-1mins-1week}$). The

6    IPAQ-SF contains four items which assess sedentary activities, walking activities, moderate intensity

7    activities and vigorous intensity activities (Booth et al., 2003). The measure was scored according to

8    the instrument scoring guideline and was conducted for both the categorical output and continuous

9    output. While there is not an upper limit on the scoring range of the IPAQ-SF; truncation of the

10   IPAQ-SF data was conducted in participants that scored higher than a total of 3 hours of either

11   walking, moderate or vigorous activities per day, to allow a maximum of 3 hours per day in each

12   activity. Outliers that where the sum of walking, moderate and vigorous activity totalling at greater

13   than 21 hours were excluded as recommended in the IPAQ-SF scoring manual.

14   *Muscle strengthening exercises or general fitness exercise*

15   Uptake of NICE core exercise recommendations was measured using a previously validated

16   questionnaire (Jinks, Jordan, Ong, & Croft, 2003; NICE, 2014). Participants were asked to report if

17   they had tried muscle strengthening exercises for their joint pain or general aerobic fitness exercise

18   for their joint pain in the last three months.

19   *Participant characteristics*

20   Self-reported participant characteristics and longitudinal descriptive statistics (baseline and three

21   months) were collected for this study. Participant baseline characteristics included age and gender.

22   Longitudinal descriptive statistics included: body mass index (BMI, kg/m$^2$), pain intensity (calculated

23   using a 0-10 numerical rating scale for each peripheral joint site, for those with multiple sites of pain

24   we took the score from the joint with the highest rated pain (Finney, Dziedzic, Lewis, & Healey,

1    2017)), health status was measured using the SF-12 (Ware Jr, Kosinski, & Keller, 1996), both the

2    physical component scale (PCS) and mental component scale (MCS) of the SF-12 were used, with a

3    score range of 0-100 with lower score indicating lower levels of health, and quality of life (QoL)

4    (measured using the EQ-5D with a score range from 0-1 with lower score indicating lower QoL (3-

5    level response) (EuroQol Group, 1990)). At 3 months follow-up participants recorded a global

6    assessment of change from baseline with ranges from 1= completely recovered to 6= much worse (K.

7    S. Dziedzic et al., 2014).

8    ***Procedure***

9    *Reliability & measurement error sub-sample*

10   Reliability and measurement error of the IPAQ-SF and PASE were assessed between baseline and 3-

11   month follow-up in a sub-group of participants who appeared to have remained stable in terms of

12   their PA level during the study period. The reliability sub-group of participants included: those who

13   completed either the IPAQ-SF or PASE at baseline and 3-month follow-up and self-reported no

14   change in their self-reported physical activity behaviour questions from baseline to follow-up. For

15   example, reported not trying muscle strengthening exercises or general fitness exercises at both

16   baseline and 3-month follow-up, or reported trying muscle strengthening exercises or general fitness

17   exercises at both baseline and three months follow-up. A 3-month follow-up period was selected as

18   an appropriate second measurement time-period for evaluating the reliability and measurement

19   error of the IPAQ-SF /PASE as most RCT investigating PA interventions for people with lower limb OA

20   are of 2-3 months in duration (Juhl et al., 2014).  It is of clinical importance to understand the test-

21   retest reliability of the PASE and IPAQ SF over a similar time-period. This knowledge would help

22   researchers and clinicians to understand whether changes in scores are due to real change or

23   measurement error.

24   *Responsiveness sub-sample*

1    Responsiveness of the IPAQ-SF and PASE was assessed between baseline and 3-month follow-up in a

2    sub-group of participants, who reported they had increased their levels of PA between the baseline

3    data collection and the 3-month follow-up on the self-reported physical activity behaviour

4    questions. The responsiveness subsample included participants who completed either the IPAQ-SF

5    or PASE at both baseline and 3-month follow-up and self-reported a positive change in their self-

6    reported physical activity behaviour questions from baseline to follow-up.  For example, reported

7    not trying muscle strengthening exercises or general fitness exercises at baseline, but reported

8    trying muscle strengthening exercises or general fitness exercises at 3-month follow-up). To allow for

9    the analysis of responsiveness, a 3-months follow-up was selected as the most suitable second

10   measurement time-period to evaluate whether the IPAQ-SF or PASE could detect changes in PA

11   behaviours that occur over a sufficient intervention duration with a component that targets physical

12   activity.

13   *Statistical analysis*

14   The IPAQ-SF total METS$^{-1mins-1week}$ scores were positively skewed and so a logarithmic transformation

15   was used to allow for a parametric statistical model when evaluating measurement properties of the

16   IPAQ-SF. Baseline descriptive statistics were reported in all participants who completed a baseline

17   IPAQ-SF or PASE and also in the reliability and responsiveness sub-samples using frequencies with

18   proportions or mean values with standard deviation (SD). Median and interquartile ranges (IQR)

19   were reported for skewed data. Changes in longitudinal descriptive statistics were reported as mean

20   change from baseline to 3-month follow-up in the reliability and responsiveness sub-samples.

21   Changes in scores from baseline to 3-month follow-up in the reliability and responsiveness sub-

22   samples were tested using paired sample t-tests with an α level of 5%. All analyses were conducted

23   using Stata Statistical Software: Release 13 (StataCorp. 2013. College Station, TX: StataCorp LP).

24   *Comparison of IPAQ-SF and PASE*

1   Baseline data of all participants were used to compare levels of PA determined by the IPAQ-SF and

2   PASE. Total scores of both measures, using IPAQ-SF logarithmic transformation, were compared

3   using Pearson's correlations. Sub-domains of PA (sitting, walking, moderate intensity and vigorous

4   intensity activities) were compared using a Spearman's rank coefficient. A priori hypothesis were

5   made, as recommended, on the comparisons of IPAQ-SF and PASE scores (Terwee et al., 2007). We

6   hypothesised that the IPAQ-SF and PASE would correlate in terms of total PA score and sub-domains

7   of PA with a strong association (0.6-0.9) based on recommendations that correlations between self-

8   report measures of PA have previously been shown to be strong (Svege, Kolle, & Risberg, 2012;

9   Terwee et al., 2010).

10  *Reliability and measurement error*

11  To assess reliability of both measures, a two-way random effects intraclass correlation for absolute

12  agreement was used (ICC$_{agreement}$)(Shrout & Fleiss, 1979). An a priori cut-off of 0.7 was selected to

13  represent adequate reliability for both the IPAQ-SF and PASE (Terwee et al., 2007). Reliability of the

14  IPAQ-SF categorical scoring was assessed using a quadratic weighted Kappa (Cohen, 1968).

15  Weightings were assigned as follows; 0 for same category, 1 for adjacent categories, and 4 for 2

16  categories apart. Reliability was assessed in sub-domains of the IPAQ-SF and PASE. For the IPAQ-SF,

17  vigorous, moderate, light, walking and sitting activities were tested. For the PASE, strenuous,

18  moderate, light, walking, sitting and strengthening exercises were tested. For dichotomous items of

19  the PASE's household and work-related activities, reliability was tested using Kappa tests and 95%

20  confidence intervals (95%CI). Interpretation of Kappa values followed recommended reference

21  values (Landis & Koch, 1977).

22  Measurement error was assessed using standard error measurement (SEM), smallest detectable

23  change (SDC) and Bland and Altman plots. The SEM was calculated for total scores of both measures

24  using SEM absolute agreement, to account for systematic difference between time points and

25  residual variance. The SDC was calculated using the SEM absolute agreement. Bland and Altman

1    plots were used to display measurement error using mean scores at baseline and 3-month follow-up

2    and difference in scores between baseline and 3-month follow-up. For the Bland and Altman plot,

3    95%CI of mean values were calculated to represent 95% limits of agreement (Bland & Altman, 1986).

4    *Responsiveness*

5    To assess responsiveness, change scores in the PASE and IPAQ-SF were compared using effect sizes

6    (ES), standardised responsiveness ratio (SRM) and response ratio (RR). ES were calculated by the

7    difference in mean change divided by the baseline SD. SRM was calculated as mean change divided

8    by the SD of change. RR was calculated as the mean change divided by the SD of the baseline of the

9    reliability subsample. ES and SRM were interpreted using cut-off values; small (0.2-0.5), moderate

10   (0.5-0.8) and large (0.8 or greater) (Guyatt, Walter, & Norman, 1987; Kazis, Anderson, & Meenan,

11   1989; Liang, Fossel, & Larson, 1990). A RR of above 1 indicated good responsiveness.

12   **Ethical approval**

13   The MOSAICS study was approved by the North West 1 Research Ethics Committee, Cheshire (REC

14   reference: 10/H1017/76).

15   **Results**

16   Of the 525 participants who returned the MOSAICS baseline questionnaire; 489 (93%) completed

17   either the IPAQ-SF (n=371, 70%) or the PASE (n=432, 82%) at baseline and were included in this

18   secondary data analysis.  314 (60%) completed both the IPAQ-SF and the PASE questionnaire (60% of

19   the total sample). Of the 470 participants who returned the MOSAICS 3-months follow-up

20   questionnaire, there were 401 participants classified into the reliability sub-sample, who were

21   deemed to have been stable during the study period, and had completed either the PASE or IPAQ-SF.

22   Of these, 360 (90%) completed the PASE and 312 (78%) completed the IPAQ-SF at baseline and 3-

23   month follow-up. There were 90 participants classified into the responsiveness sub-sample, 66 of

24   those completed the IPAQ-SF (73%) and 83 PASE (86%) at both baseline and 3 months follow-up

1    (**figure 1**). Baseline descriptive statistics for the whole sample and each subsample are displayed in

2    **table 1**. We compared baseline descriptive statistics of those who completed either the IPAQ-SF or

3    PASE at baseline to those that did not complete either PA measure to assess any differences. Those

4    who did not complete either the IPAQ-SF or PASE were older in age, had poorer mental health (as

5    indicated by a lower MCS score in the SF-12), but were not significantly different in terms of gender

6    distributions, pain intensity, PCS and QoL.

7    *Comparison of IPAQ-SF and PASE*

8    Total scores of the IPAQ-SF and PASE were moderately associated with each other ($r=0.56$,

9    $p=<0.001$). Comparisons of the sub-domains showed moderate strength associations in sitting

10    activities ($r_s=0.46$, $p=<0.001$), walking activities ($r_s=0.57$, $p=<0.001$), moderate intensity activities

11    ($r_s=0.34$, $p=<0.001$) and vigorous/strenuous activities ($r_s=0.39$, $p=<0.001$). While all associations were

12    statistically significant, neither the total scores nor subdomains of the IPAQ-SF and PASE

13    demonstrated a correlation coefficient above 0.6.

14    *Reliability and measurement error*

15    Changes in longitudinal descriptive statistics for the reliability sub-sample are displayed in **table 2**.

16    Within the reliability sub-sample, there were statistically significant changes in joint pain intensity,

17    PCS and QoL. Although changes in these outcomes were not of a magnitude commonly considered

18    to be of a minimal clinically important change (Jenkinson & Layte, 1997; Salaffi, Stancati, Silvestri,

19    Ciapetti, & Grassi, 2004; Walters & Brazier, 2005). The mean PASE scores also changed significantly

20    at 3-month follow-up compared to baseline in the reliability sub-sample. The intraclass correlation

21    between baseline and 3-month follow-up for the total score of the PASE was below the 0.7 cut-off

22    value; $ICC_{agreement}=0.68$ (0.61-0.73 95%CI, $p=<0.001$), SEM was 46.7 and SDC was 129.6. **Figure 2**

23    displays the Bland and Altman plot with the lower 95% limit of agreement -117 and upper 95% limit

24    of agreement 136, representing large measurement error and limits of agreement when considering

25    the score range of the PASE (0-400). The intraclass correlation between baseline and 3-month

1    follow-up for the total score of the IPAQ-SF was below the 0.7 cut-off value; $ICC_{agreement}$=0.62 (0.55-

2    0.71 95%CI, p=<0.001), SEM was 3532.2 METS$^{-1mins-1week}$ and SDC was 9790.8 METS$^{-1mins-1week}$. **Figure 3**

3    displays the Bland and Altman plot with the lower 95% limit of agreement -5222 METS$^{-1mins-1week}$ and

4    upper limits of agreement 5597METS$^{-1mins-1week}$, representing large measurement error and limits of

5    agreement, considering 3000METS$^{-1mins-1week}$ equates to 6-8 hours of running in a week. A quadratic

6    weighted Kappa showed agreement between baseline and 3-month follow-up was below 0.7 cut-off:

7    K=0.56 (0.43-0.67 95%CI). **Table 3** displays the reliability of the sub-domains within the IPAQ-SF and

8    PASE, Spearman's rank coefficients and Kappa values ranged from 0.33-0.74 across domains in the

9    IPAQ-SF and PASE.

10   *Responsiveness*

11   Changes in longitudinal descriptive statistics for the responsiveness sub-sample are displayed in

12   **table 2**. There were no statistically significant changes in the descriptive statistics between baseline

13   and 3-month follow-up in the responsiveness subsample. Mean changes in PASE and median change

14   in IPAQ-SF suggested the responsiveness sub-group reduced their levels of PA between baseline and

15   3 months follow-up. Low scores in indicators of responsiveness were observed; the ES was -0.14 and

16   -0.16, SRM was -0.21 and -0.21 and RR was 0.12 and 0.09 for the IPAQ-SF and PASE respectively.

17   **Discussion**

18   *Main findings*

19   This study has investigated the measurement properties of the IPAQ-SF and PASE in those aged 45

20   and over, consulting primary care with OA of the hand, hip, knee or foot. We assessed reliability,

21   measurement error, responsiveness and compared scoring for the IPAQ-SF and PASE, assessing the

22   ability of both measures in detecting changes in physical activity levels. When comparing total scores

23   of the IPAQ-SF and PASE we found that both instruments correlated moderately with each other,

24   suggesting the IPAQ-SF and PASE are moderately similar in terms of measuring total PA levels in our

1  sample, however, the correlation strength was below the a priori cut-off of 0.6. When exploring

2  subdomains of PA, walking activities also had a moderate correlation to each other, but sitting,

3  moderate and vigorous activities all had weaker correlations in comparison to those found for the

4  total scores. The PASE contains several sub-domains that the IPAQ-SF does not, differences in

5  categorisation of activities between the IPAQ-SF and PASE may explain the higher magnitude of

6  association in total scores but not in matched sub-domains.

7  While there is no complete consensus on appropriate cut-off values for ICC to demonstrate good

8  reliability of a measurement instrument, an ICC=0.7 has been generally recommended (Terwee et

9  al., 2007). Neither the IPAQ-SF nor PASE achieved this in 3-month test-retest assessment, although

10  the PASE was closest. For measurement error, the IPAQ-SF statistics represent a relatively large SEM,

11  SDC and limits of agreement, meaning an extremely large change in weekly PA levels would be

12  required to be detected by the IPAQ-SF outside of the measurement error. The SEM, SDC and limits

13  of agreement findings were relatively large in relation to the PASE total scale range (0-400)

14  suggesting large measurement error. This suggests large measurement error in proportion to

15  possible range of total scores. In IPAQ-SF and PASE responsiveness findings suggest low

16  responsiveness. However, it is difficult to ascertain whether this was due to a true lack of change in

17  the responsiveness subsample or due to the instruments' inability to detect change.

18  Comparison with other studies

19  Our findings comparing the PASE and IPAQ-SF total scores showing moderate correlations were

20  lower than previously shown strong correlations (Svege et al., 2012), although previous studies

21  showed low correlations to activity monitors (Casartelli et al., 2015; Svege et al., 2012). The findings

22  in measurement error of the IPAQ-SF and PASE are comparable with three previous investigations

23  (Svege et al., 2012) and PASE (Bolszak, Casartelli, Impellizzeri, & Maffiuletti, 2014; Casartelli, Bolszak,

24  Impellizzeri, & Maffiuletti, 2015) in adults with OA, which identified that neither the IPAQ-SF nor

25  PASE had adequate test-retest reliability (greater than or equal to 0.7) (Terwee et al., 2007). A study

1    using the Dutch version of the IPAQ-SF in a sample after joint replacement demonstrated test-retest

2    reliability below 0.7 (Blikman, Stevens, Bulstra, van den Akker-Scheek, & Reininga, 2013). However,

3    there is inconsistency in the literature with some previous studies reporting adequate test-retest

4    reliability in both measures which could potentially be explained by a shorter time gap between

5    testing (7-10 days) or a more active approach to keeping participants stable in their levels of PA

6    (Bolszak et al., 2014; Naal, Impellizzeri, & Leunig, 2009; Svege et al., 2012). The findings in this study

7    on measurement error were similar to other studies on the IPAQ-SF (Blikman et al., 2013; Bolszak et

8    al., 2014) and PASE (Casartelli et al., 2015), suggesting that comparisons of levels of PA between

9    individuals or time-points are at risk of large measurement error when using the IPAQ-SF and PASE.

10   The general findings on the IPAQ-SF and PASE in terms of reliability and measurement error in our

11   study are in line with a previous systematic review on all self-report instruments measuring PA in

12   adults with OA, showing some evidence for acceptable reliability, but large measurement error

13   (Smith et al., 2019).

14   *Strengths and weakness of study*

15   A strength of this study was the large sample size and investigation of the two most commonly used

16   self-report PA instruments in OA research (Smith et al., 2019). Sample sizes above n=100 are

17   generally considered adequate for statistical precision in evaluating test-retest reliability and

18   measurement error and comparing scores of the IPAQ-SF and PASE (Terwee et al., 2007). The

19   characteristics of primary care practices included within the MOSAICS cluster RCT are similar to

20   those within the wider UK, allowing for our findings to be generalisable to other UK general practices

21   (K. Dziedzic et al., 2018). In those that did not complete either the PASE or IPAQ-SF, we found only

22   small differences in descriptive characteristics to those that completed either instrument, suggesting

23   that the findings in our sample are generalisable to those aged 45 and over with OA in the wider

24   primary care MOSAICS sample.

Despite the strengths of this study, the evaluation of test-retest reliability and responsiveness does have limitations. For example, our criteria for stability in the reliability subsample could not guarantee that the whole sub-sample remained stable with regards to levels of PA during the study period. This may have led to an underestimation on the instruments' reliability and measurement, if participants who were not stable in their level of PA were included into the subsample. Conversely, our criteria for changes in PA behaviours in our responsiveness subsample may not have been sensitive or specific enough to identify those that experienced true changes in their PA levels during the study period. There is a risk of underestimation of the instruments' responsiveness if those with true changes to their PA levels were not detected. An anchor measurement, such as an objective measure in levels of PA would have been desirable to accurately estimate stability or changes in true levels of PA in the study sample. However objective measures such as activity monitors are costly for population level studies and so were not viable within the MOSAICS study. The 'Gold Standard' measure of PA, doubly-labelled water (Schuit, Schouten, Westerterp, & Saris, 1997), is costly and requires specialised laboratory expertise and equipment, which were not viable as part of the MOSAICS study. Because of this, we were unable to establish the degree to which the IPAQ-SF or PASE represent true levels of PA.

*Implications for research and clinical practice*

Currently, there are no international recommendations to guide the selection of instruments to measure PA in adults with OA. Due to the large measurement error of the PASE and IPAQ-SF, relative to their scoring range or quantity of PA measured in METS respectively, both instruments may perform poorly when comparing individual or population PA levels or evaluating change in PA levels over time. For example, biased associations or inferences may occur in research studies that investigate the associations between PA level and other clinical outcomes or studies that investigate change in PA over time.

1 More evaluations of measurement properties of self-report physical activity instruments are

2 warranted, particularly those that assess criterion validity in adults with OA. Further evaluations on

3 test-retest reliability and responsiveness using objective measures as an anchor for detecting

4 stability or changes in PA would also provide important information on their measurement

5 properties. Such evaluations are necessary before recommendations can be made on the selection

6 of instrument for measuring levels of PA in adults with OA.

7 *Conclusion*

8 Despite their low cost and ease to administer in population level research, the IPAQ-SF and PASE

9 appear to be limited in their reliability and measurement error for measuring levels of PA among

10 adults with OA. Their ability to accurately compare PA levels between populations and to detect

11 changes in PA over time is questionable. Researchers and clinicians should be aware of the

12 limitations of the IPAQ-SF and PASE's measurement properties for assessing PA levels among adults

13 aged 45 and over with OA.

14 **References**

15 Bland, J. M., & Altman, D. (1986). Statistical methods for assessing agreement between two methods
16     of clinical measurement. *The Lancet, 327*(8476), 307-310.
17 Blikman, T., Stevens, M., Bulstra, S. K., van den Akker-Scheek, I., & Reininga, I. H. (2013). Reliability
18     and validity of the Dutch version of the International Physical Activity Questionnaire in
19     patients after total hip arthroplasty or total knee arthroplasty. *journal of orthopaedic &*
20     *sports physical therapy, 43*(9), 650-659.
21 Bolszak, S., Casartelli, N. C., Impellizzeri, F. M., & Maffiuletti, N. A. (2014). Validity and reproducibility
22     of the Physical Activity Scale for the Elderly (PASE) questionnaire for the measurement of the
23     physical activity level in patients after total knee arthroplasty. *BMC musculoskeletal*
24     *disorders, 15*(1), 1.
25 Booth, M. L., Ainsworth, B. E., Pratt, M., Ekelund, U., Yngve, A., Sallis, J. F., & Oja, P. (2003).
26     International physical activity questionnaire: 12-country reliability and validity. *Med Sci*
27     *Sports Exerc, 195*(9131/03), 3508-1381.
28 Casartelli, N. C., Bolszak, S., Impellizzeri, F. M., & Maffiuletti, N. A. (2015). Reproducibility and
29     Validity of the Physical Activity Scale for the Elderly (PASE) Questionnaire in Patients After
30     Total Hip Arthroplasty. *Physical therapy, 95*(1), 86-94. doi:10.2522/ptj.20130557
31 Cohen, J. (1968). Weighted kappa: Nominal scale agreement provision for scaled disagreement or
32     partial credit. *Psychological bulletin, 70*(4), 213.
33 Dziedzic, K., Healey, E., Porcheret, M., Afolabi, E., Lewis, M., Morden, A., . . . Finney, A. (2018).
34     Implementing core NICE guidelines for osteoarthritis in primary care with a model

consultation (MOSAICS): a cluster randomised controlled trial. *Osteoarthritis and Cartilage, 26*(1), 43-53.

Dziedzic, K. S., Healey, E. L., Porcheret, M., Ong, B. N., Main, C. J., Jordan, K. P., . . . Morden, A. (2014). Implementing the NICE osteoarthritis guidelines: a mixed methods study and cluster randomised trial of a model osteoarthritis consultation in primary care-the Management of OsteoArthritis In Consultations (MOSAICS) study protocol. *Implementation Science, 9*(1), 95.

EuroQol Group. (1990). EuroQol--a new facility for the measurement of health-related quality of life. *Health policy (Amsterdam, Netherlands), 16*(3), 199.

Felson, D. T. (2009). Developments in the clinical understanding of osteoarthritis. *Arthritis research & therapy, 11*(1), 203.

Finney, A., Dziedzic, K. S., Lewis, M., & Healey, E. (2017). Multisite peripheral joint pain: a cross-sectional study of prevalence and impact on general health, quality of life, pain intensity and consultation behaviour. *BMC musculoskeletal disorders, 18*(1), 535.

Fransen, M., McConnell, S., Harmer, A. R., Van der Esch, M., Simic, M., & Bennell, K. L. (2015). Exercise for osteoarthritis of the knee: a Cochrane systematic review. *Br J Sports Med, 49*(24), 1554-1557.

Fransen, M., McConnell, S., Hernandez-Molina, G., & Reichenbach, S. (2014). Exercise for osteoarthritis of the hip. *Cochrane Database of Systematic Reviews*(4).

Grime, J., & Dudley, B. (2014). Developing written information on osteoarthritis for patients: facilitating user involvement by exposure to qualitative research. Health Expectations, 17(2), 164-173.

Guyatt, G., Walter, S., & Norman, G. (1987). Measuring change over time: assessing the usefulness of evaluative instruments. *Journal of chronic diseases, 40*(2), 171-178.

Healey, E. L., Allen, K. D., Bennell, K., Bowden, J. L., Quicke, J. G., & Smith, R. (2020). Self-Report Measures of Physical Activity. *Arthritis Care & Research, 72*, 717-730.

Herbolsheimer, F., Schaap, L. A., Edwards, M. H., Maggi, S., Otero, Á., Timmermans, E. J., . . . Cooper, C. (2016). Physical Activity Patterns Among Older Adults With and Without Knee Osteoarthritis in Six European Countries. *Arthritis Care & Research, 68*(2), 228-236.

Holden, M. A., Nicholls, E. E., Young, J., Hay, E. M., & Foster, N. E. (2015). Exercise and physical activity in older adults with knee pain: a mixed methods study. Rheumatology, 54(3), 413-423.

Jenkinson, C., & Layte, R. (1997). Development and testing of the UK SF-12. *Journal of health services research & policy, 2*(1), 14-18.

Jinks, C., Jordan, K., Ong, B., & Croft, P. (2003). A brief screening tool for knee pain in primary care (KNEST). 2. Results from a survey in the general population aged 50 and over. *Rheumatology, 43*(1), 55-61.

Juhl, C., Christensen, R., Roos, E. M., Zhang, W., & Lund, H. (2014). Impact of exercise type and dose on pain and disability in knee osteoarthritis: a systematic review and meta-regression analysis of randomized controlled trials. *Arthritis & Rheumatology, 66*(3), 622-636.

Kazis, L. E., Anderson, J. J., & Meenan, R. F. (1989). Effect sizes for interpreting changes in health status. *Medical care*, S178-S189.

Landis, J. R., & Koch, G. G. (1977). The measurement of observer agreement for categorical data. *Biometrics*, 159-174.

Liang, M. H., Fossel, A. H., & Larson, M. G. (1990). Comparisons of five health status instruments for orthopedic evaluation. *Medical care*, 632-642.

Naal, F. D., Impellizzeri, F. M., & Leunig, M. (2009). Which is the best activity rating scale for patients undergoing total joint arthroplasty? *Clinical orthopaedics and related research, 467*(4), 958-965.

National Institute for Health and Care Excellence (2014). Osteoarthritis: care and management.

1   Prince, S. A., Adamo, K. B., Hamel, M. E., Hardt, J., Gorber, S. C., & Tremblay, M. (2008). A
2       comparison of direct versus self-report measures for assessing physical activity in adults: a
3       systematic review. *International Journal of Behavioral Nutrition and Physical Activity, 5*(1), 1.
4   Quicke, J. G., Foster, N. E., Croft, P. R., Ogollah, R. O., & Holden, M. A. (2018). Change in physical
5       activity level and clinical outcomes in older adults with knee pain: a secondary analysis from
6       a randomised controlled trial. *BMC musculoskeletal disorders, 19*(1), 59.
7   Salaffi, F., Stancati, A., Silvestri, C. A., Ciapetti, A., & Grassi, W. (2004). Minimal clinically important
8       changes in chronic musculoskeletal pain intensity measured on a numerical rating scale.
9       *European Journal of Pain, 8*(4), 283-291.
10  Schuit, A. J., Schouten, E. G., Westerterp, K. R., & Saris, W. H. (1997). Validity of the Physical Activity
11      Scale for the Elderly (PASE): according to energy expenditure assessed by the doubly labeled
12      water method. *Journal of clinical epidemiology, 50*(5), 541-546.
13  Shrout, P. E., & Fleiss, J. L. (1979). Intraclass correlations: uses in assessing rater reliability.
14      *Psychological bulletin, 86*(2), 420.
15  Silsbury, Z., Goldsmith, R., & Rushton, A. (2015). Systematic review of the measurement properties
16      of self-report physical activity questionnaires in healthy adult populations. *BMJ open, 5*(9),
17      e008430.
18  Smith, R. D., Dziedzic, K. S., Quicke, J. G., Holden, M. A., McHugh, G. A., & Healey, E. L. (2019).
19      Identification and Evaluation of Self-Report Physical Activity Instruments in Adults With
20      Osteoarthritis: A Systematic Review. *Arthritis Care & Research, 71*(2), 237-251.
21  Svege, I., Kolle, E., & Risberg, M. A. (2012). Reliability and validity of the Physical Activity Scale for the
22      Elderly (PASE) in patients with hip osteoarthritis. *BMC musculoskeletal disorders, 13*(1), 1.
23  Terwee, C. B., Bot, S. D., de Boer, M. R., van der Windt, D. A., Knol, D. L., Dekker, J., . . . de Vet, H. C.
24      (2007). Quality criteria were proposed for measurement properties of health status
25      questionnaires. *Journal of clinical epidemiology, 60*(1), 34-42.
26  Terwee, C. B., Mokkink, L. B., van Poppel, M. N., Chinapaw, M. J., van Mechelen, W., & de Vet, H. C.
27      (2010). Qualitative attributes and measurement properties of physical activity
28      questionnaires. *Sports Medicine, 40*(7), 525-537.
29  Uthman, O. A., van der Windt, D. A., Jordan, J. L., Dziedzic, K. S., Healey, E. L., Peat, G. M., & Foster,
30      N. E. (2013). Exercise for lower limb osteoarthritis: systematic review incorporating trial
31      sequential analysis and network meta-analysis. *Bmj, 347*, f5555.
32  Walters, S. J., & Brazier, J. E. (2005). Comparison of the minimally important difference for two
33      health state utility measures: EQ-5D and SF-6D. *Quality of Life Research, 14*(6), 1523-1532.
34  Warburton, D. E., Nicol, C. W., & Bredin, S. S. (2006). Health benefits of physical activity: the
35      evidence. CMAJ, 174(6), 801-809.
36  Ware Jr, J. E., Kosinski, M., & Keller, S. D. (1996). A 12-Item Short-Form Health Survey: construction
37      of scales and preliminary tests of reliability and validity. *Medical care, 34*(3), 220-233.
38  Washburn, R. A., Smith, K. W., Jette, A. M., & Janney, C. A. (1993). The Physical Activity Scale for the
39      Elderly (PASE): development and evaluation. *Journal of clinical epidemiology, 46*(2), 153-162.
40

Table 1 Baseline descriptive statistics of whole sample, those that completed the IPAQ-SF and PASE at baseline, reliability subsample and responsiveness subsample.

| | Baseline IPAQ-SF and PASE (n=314) | Reliability subsample (n=401) | Responsiveness subsample (n=90) |
|---|---|---|---|
| **Gender, n (%)** | | | |
| Females | 172 (55) | 530 (57) | 50 (56) |
| **Site of peripheral joint pain, n (%)** | | | |
| Knee | 259 (82) | 335 (84) | 71 (79) |
| Hip | 181 (57) | 232 (58) | 46 (51) |
| Feet | 150 (48) | 194 (48) | 42 (47) |
| Hand | 167 (53) | 220 (55) | 54 (60) |
| **Mean age, years (SD)** | 67.2 (10.5) | 68.3 (10.6) | 66.8 (10.0) |
| **Mean BMI, kg/m$^2$ (SD)** | 27.0 (8.1) | 28.4 (4.9) | 29.7 (6.2) |
| **Mean pain intensity (SD)** | 7.4 (2.1) | 7.4 (2.0) | 6.7 (2.3) |
| **Mean health Status, SF-12 (SD)** | | | |
| PCS | 36.7 (11.3) | 36.8 (11.5) | 38.1 (9.2) |
| MCS | 51.5 (11.3) | 51.2 (10.9) | 51.8 (11.3) |
| **Mean QoL, EQ-5D (SD)** | 0.6 (0.3) | 0.6 (0.3) | 0.6 (0.3) |
| **Mean PASE score (SD)** | 142.3 (78.8) | 140.3 (76.3) | 153.4 (89.3) |
| **Median IPAQ-SF, total METS$^{-1mins-1week}$ (interquartile range)** | 1527 (462-3732) | 1386 (198-3452) | 2574 (305-5153) |
| **IPAQ-SF categories, n (%)** | | | |
| Low | 113 (36) | 112 (36) | 17 (26) |
| Moderate | 103 (33) | 104 (33) | 21 (32) |
| High | 98 (31) | 96 (31) | 28 (42) |

Note: percentages may not equal 100% due to rounding up. Abbreviations: SD, standard deviation, IPAQ-SF, International physical activity questionnaire - short form, PASE, physical activity scale for the elderly, BMI, body mass index, SF-12, 12-Item Short Form Survey, PCS, physical component score, MCS, mental component score, QoL, quality of life, METS$^{-1mins-1week}$, metabolic equivalent per minute per week.

Table 2 Changes in longitudinal descriptive statistics from baseline to 3 months follow-up for reliability and responsiveness subsample.

| | Change in reliability subsample (n=401) | p-value | Change in responsiveness subsample (n=90) | p-value |
|---|---|---|---|---|
| **Change in mean BMI, kg/m$^2$** | -0.1 | 0.19 | -0.5 | 0.18 |
| **Change in mean pain intensity** | -1.3 | <0.001* | -0.5 | 0.06 |
| **Change in mean health Status, SF-12** | | | | |
| PCS | 1.0 | 0.03* | -0.4 | 0.7 |
| MCS | 0.4 | 0.81 | -0.6 | 0.6 |
| **Change in mean QoL, EQ-5D** | 0.1 | <0.001* | -0.1 | 0.3 |
| **Change in mean PASE score** | -8.8 | 0.04* | -14.3 | 0.2 |
| **Change in mean IPAQ-SF, total METS$^{-1mins-1week}$** | 323 | 0.52 | -347 | 0.1 |
| **Global assessment of change, n (%)** | | | | |
| Missing | 9 (2) | | 0 (0) | |
| Completely recovered | 5 (1) | | 1 (1) | |
| Much better | 40 (10) | | 10 (11) | |
| Better | 80 (20) | | 20 (22) | |
| No change | 141 (35) | | 39 (43) | |
| Worse | 108 (27) | | 16 (18) | |
| Much worse | 18 (4) | | 4 (4) | |

Note: percentages may not equal 100% due to rounding up. Abbreviations: IPAQ-SF, International physical activity questionnaire - short form, PASE, physical activity scale for the elderly, BMI, body mass index, SF-12, 12-Item Short Form Survey, PCS, physical component score, MCS, mental component score, QoL, quality of life, METS$^{-1mins-1week}$, metabolic equivalent per minute per week.

Table 3 Reliability of PASE and IPAQ-SF subdomains in the reliability subsample.

| Instrument subdomain | Coefficient |
|---|---|
| **PASE – Leisure time activities** | **Spearman's Rank, *p* value** |
| Sitting activities | r=0.54, *p*=<0.001 |
| Walking activities | r=0.57, *p*=<0.001 |
| Light intensity activities | r=0.33, *p*=<0.001 |
| Moderate intensity activities | r=0.47, *p*=<0.001 |
| Strenuous intensity activities | r=0.64, *p*=<0.001 |
| Muscle endurance or strengthening activities | r=0.56, *p*=<0.001 |
| **PASE - Household and work-related activities** | **Weighted Kappa (95%CI)** |
| Light housework | K=0.49 (0.08-0.89) |
| Heavy housework | K=0.60 (0.48-0.72) |
| Home repairs | K=0.49 (0.32-0.65) |
| Garden care | K=0.46 (0.34-0.58) |
| Outdoor gardening | K=0.37 (0.25-0.50) |
| Caring others | K=0.52 (0.39-0.64) |
| Work or volunteer | K=0.65 (0.53-0.77) |
| Working physical activity | K=0.50 (0.29-0.71) |
| **IPAQ-SF – Subdomains** | **Spearman's Rank, *p* value** |
| Sitting activities | r=0.74, *p*=<0.001 |
| Walking activities | r=0.59, *p*=<0.001 |
| Moderate intensity activities | r=0.45, *p*=<0.001 |
| Vigorous intensity activities | r=0.46, *p*=<0.001 |

Abbreviations: IPAQ-SF, International physical activity questionnaire - short form, PASE, physical activity scale for the elderly.

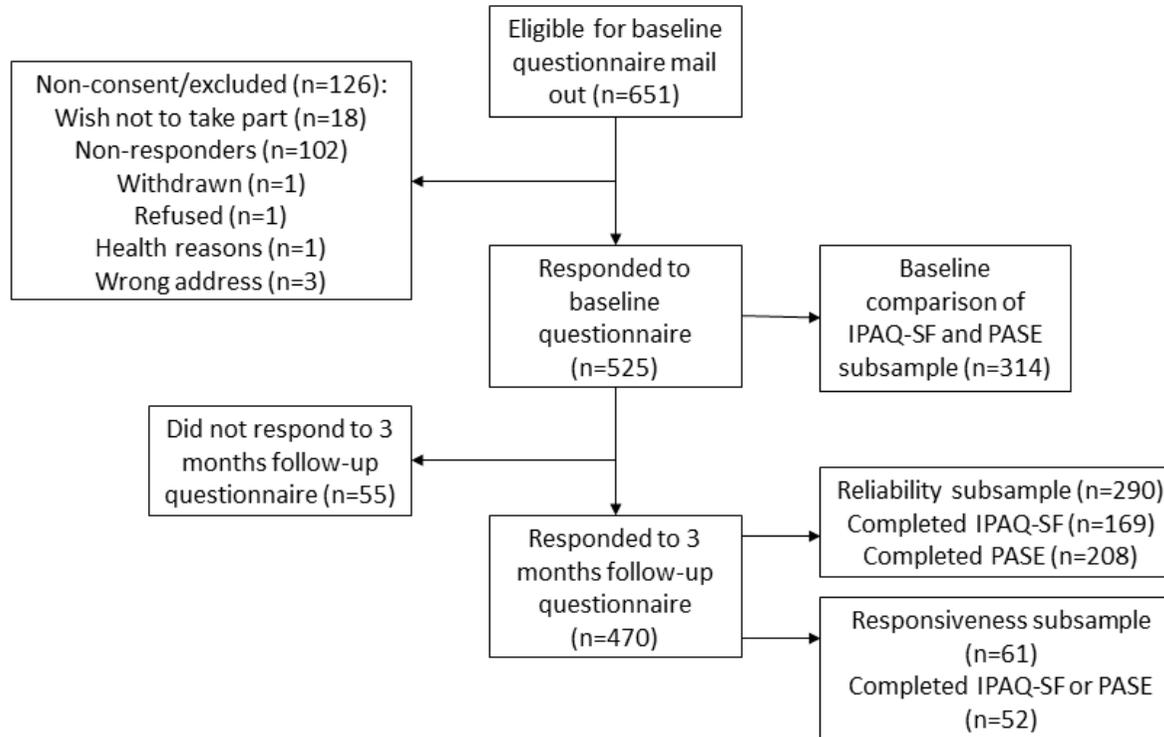Figure 1 Flowchart of study participants for analysis taken from the MOSAICS study.

Figure 2 Bland and Altman plot of the PASE in the reliability subsample from baseline to 3-months follow-up.
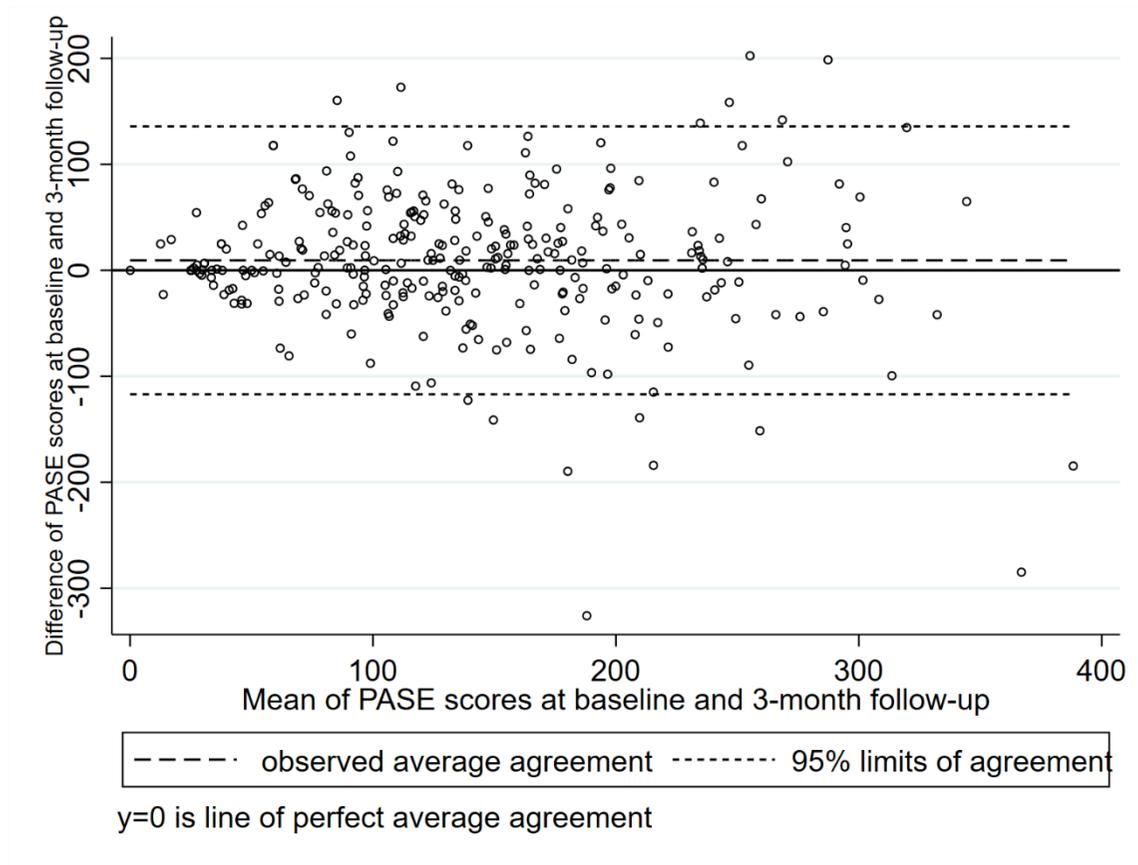
Figure 3 Bland and Altman plot of the IPAQ-SF in the reliability subsample from baseline to 3-months follow-up.