

Cover Page



Universiteit Leiden



The handle <http://hdl.handle.net/1887/19033> holds various files of this Leiden University dissertation.

Author: Lukas, Cédric

Title: Methodological aspects of outcome assessment in inflammatory rheumatic diseases

Date: 2012-05-31



Summary and general discussion

1 The development of outcome measures in the field of rheumatology has shown
2 a rapid and constructive advance in the last decades. It has been achieved to
3 obtain more uniformity in the evaluation of inflammatory rheumatic disorders,
4 especially concerning long-term outcomes and disease activity assessments. In
5 the field of RA, clinical experts and researchers have done a lot of work under
6 the umbrella of OMERACT (Outcome Measures in Rheumatology Clinical
7 Trials) [1]. In the field of AS, major advances in the outcome and therapeutic
8 advances in the last years have been possible because of projects piloted within
9 ASAS (the Assessment of Spondyloarthritis international Society).

10 The analyses and studies presented in this thesis were part of this process,
11 with most of the work performed in international working groups under
12 auspices of either OMERACT or ASAS or both.

13 14 **ASSESSING RADIOGRAPHIC PROGRESSION IN RA IN CLINICAL** 15 **PRACTICE** 16

17
18 In **chapter 1**, we have assessed the metrological properties of a scoring method
19 for radiographic damage in RA, the SENS method, which aimed at making
20 the objective and numerical evaluation of structural damage in a given patient
21 feasible for clinical rheumatologists without specific training [2]. Only the
22 recognition of an erosive lesion or of any joint space narrowing in the joints
23 of interest is required, as no further assessment of the degree of destruction is
24 necessary to complete the scoring: While comprehensive scoring methods, in
25 particular the Sharp and Larsen methods and their modifications, attribute a
26 numerical value to each of the scored joint according to the level of damage as
27 appreciated by the reader, SENS only acknowledges the presence (1 point) or
28 absence (zero point) of erosions and joint space narrowing separately [3]. Base-
29 line- and 52 weeks radiographies of 680 RA patients included in the TEMPO
30 trial (Trial of Etanercept and Methotrexate with Radiographic Patient Out-
31 come) were scored according to the Sharp-van der Heijde (SHS) method by 2
32 independent readers, and SENS was derived from their results [4, 5]. Reliability
33 of status scores was found high for both methods, with intraclass correlation
34 coefficients (ICC) for SHS of 0.81 and 0.77 at baseline and 52 weeks and 0.91
35 and 0.89 for SENS respectively. ICCs of change scores were somewhat lower.
36 Both scores performed equally well with regard to sensitivity to change, and
37 the optimal cut-off levels for progression vs. no progression based on smallest
38 detectable change were, as determined by ROC-curve analysis, at 2 units for
39 SHS and 1 unit for SENS. With regard to their discriminatory ability, both

1 methods also performed very well. Additionally, it was suggested by further
2 analysis that the SENS does not suffer a lot from ceiling effect. In conclu-
3 sion it was proven that SENS has a good reliability, sensitivity to change and
4 discriminatory ability, and can be considered for use in clinical practice also in
5 patients with established disease.
6

7 8 **AUTOMATED MEASUREMENT OF JOINT SPACE WIDTH IN** 9 **RHEUMATOID ARTHRITIS**

10
11 In **chapter 2**, five (semi-)automated methods to measure joint space width
12 and its change over time in joints of patients with RA have been compared.
13 One set of radiographs from 107 patients included in the COBRA (*Combina-*
14 *tietherapie Bij Reumatoïde Artritis*) trial was used, and different aspects of the
15 5 methods, such as feasibility, efficiency, reliability and discriminatory ability
16 were tested [6]. Most methods showed efficiency problems: With a 50% suc-
17 cess rate (i.e. at least half of the evaluated joints successfully measured both at
18 baseline- and follow-up visits) as a benchmark, 4 of the 5 methods passed the
19 efficiency test. But when higher requirements were asked, efficiency of most
20 methods fell short: Only 3 methods passed the requirement of at least 75% of
21 joints successfully measured. Repeatability, on the other hand, as assessed by
22 ICC, was very acceptable ($ICC > 0.80$). The interpretation of discriminatory
23 ability of the different methods was more challenging, and direct comparison
24 across methods is even unreasonable, because different joints are assessed by
25 the various systems, missing (i.e. unsuccessfully assessed) joints differed across
26 methods. The same films scored by the Sharp van der Heijde method served
27 as external standard. The manual scoring of the joint space narrowing score
28 did not accurately distinguish the treatment arms in the COBRA trial, while
29 the difference was highly significant for the erosion score. However, some of
30 the (semi)automated measurements were able to pick up a difference between
31 the treatment arms, even with the limited number of joints included and the
32 limited success rate. It was concluded that (semi-)automated measurement of
33 joint space width in RA may help in better discriminating between treatment
34 arms in the context of a clinical trial, provided that their efficiency is improved
35 and that they are combined with the assessment of radiographic erosions.
36
37
38
39

1 **RADIOGRAPHIC PROGRESSION IN RHEUMATOID ARTHRITIS,** 2 **PLAUSIBILITY OF REPAIR**

3
4 The research described in **Chapters 3 and 4** was aiming at getting a better and
5 more detailed insight in radiographic progression in clinical trials of RA. In
6 particular, the occurrence of so-called negative change scores, which suggest
7 improvement, was investigated. An important finding was that negative scores
8 were not only due to reading error but partly to true improvement (repair).
9 Somewhat to our surprise we found that multiple and independent scoring of
10 the same radiographs under blinded time sequence led to seemingly different
11 results on a per joint basis: From the 7255 individual joints with 4 change
12 scores available, pertaining to 178 patients in TEMPO trial, only a minority
13 (1.3 to 5.8% across reads) showed change over 1 year. Most surprisingly, the
14 absolute agreement across readers about the direction of change (either posi-
15 tive, i.e. worsening of erosive damage, or negative, i.e. apparent improvement
16 of lesions) was very low: In only 12 joints a consistently positive or negative
17 change pattern was found in all 4 readings. Discrepant change patterns, on the
18 other hand, were also very rare. We explained the paradox of reliable patient
19 total change scores made up by unreliable single joint scores by introducing
20 the concept of conservative vs. sensitive readers: In this concept, the direction
21 of the score is always determined by the sensitive reader. The conservative
22 reader scores zero change in case of doubt and does change the magnitude of
23 the signal but not the direction. These subtleties underscore the importance of
24 perfect blinding of time order in readings of RCTs in RA.

25 Further exploration of radiographic results at the individual joint level in
26 **chapter 4** showed that the negative change scores can be regarded as a sur-
27 rogate of repair, namely as a decrease in the size of erosive lesions. The finding
28 that a negative change score was statistically more likely to occur in a joint of a
29 patient treated with a TNF blocker, if that joint also showed improvement or
30 resolution of clinical swelling, adds to the belief that repair is a real event and
31 not just an assessment artifact.

32 33 34 **THERAPEUTIC APPROACH OF EARLY INFLAMMATORY ARTHRITIS:** 35 **BEHAVIORS AND CONSEQUENCES**

36
37 **Chapter 5** reported and analyzed the therapeutic approach of French rheuma-
38 tologists facing a patient with early inflammatory arthritis. It was demonstrated
39 that both the time-to-treatment-start and the choice of the drug were hetero-

1 geneous across geographical regions. The tendency of initiating a DMARD in
2 a patient with high disease activity, abnormal acute phase reactants, polyar-
3 thritis (i.e. more than 3 joint groups) and the presence of anti-CCP antibodies,
4 however, was paramount. Interestingly, a significant interaction between the
5 study center and the result of the anti-CCP test on one hand, and swollen
6 joint count on the other hand, pointed to differences in the interpretation of
7 clinical data in light of the therapeutic decision: Some rheumatologists find
8 the presence of bad prognostic markers less important than others.

9
10 In **chapter 6**, the therapeutic behavior of rheumatologists in early inflam-
11 matory arthritis and its consequences with regard to radiographic progres-
12 sion were further investigated. Because the ESPOIR cohort is observational
13 and does not dictate a treatment protocol, time to initiate a DMARD and
14 choice of the drug were left to the discretion of the treating rheumatologist.
15 Consequently, as was previously described in chapter 5, the time elapsed from
16 the arthritis onset to the start of a DMARD was variable across patients, and
17 could thus be modeled as a factor of its efficacy [7, 8]. The propensity analysis
18 approach was chosen in these circumstances, because in observational studies,
19 the “treatment”-groups (here early treatment versus delayed treatment) by de-
20 fault show imbalances on prognostic factors [9]. These prognostic imbalances
21 may confound the relation between the chosen DMARD treatments and the
22 outcome (radiographic progression): It is difficult to attribute differences in
23 responses to the treatment itself because the patient- and disease characteristics
24 are also believed to influence the response. The propensity score method aims
25 to reduce the confounding effects of such covariates, and allows differences of
26 responses to be attributed to differences in therapeutic strategies (early versus
27 delayed). In the case of the ESPOIR study, we were able to demonstrate that
28 in daily clinical practice the very early initiation of a DMARD, which means
29 within the first 3 months after the first occurrence of joint swelling, had a sig-
30 nificant impact on the radiographic progression after 12 months of follow up.
31 This demonstration adds to the credibility of the current recommendations
32 regarding the management of patients with early arthritis, as well as to the ap-
33 preciation that both a fast referral to a rheumatologist and an early treatment
34 start are required to improve the long term prognosis of the disease [10-13].

35
36
37
38
39

1 MAGNETIC RESONANCE IMAGING IN AS: OUTCOME MEASURES

2
 3 **Chapter 7** was a multi-reader experiment to compare feasibility, discrimina-
 4 tory ability, responsiveness and reliability of 3 different methods developed to
 5 measure inflammation on MRI in the spine of AS patients. The Ankylosing
 6 Spondylitis spine MRI score for activity (ASspiMRI-a), the Berlin method (a
 7 modification of the ASspiMRI-a), and the Spondyloarthritis Research Consor-
 8 tium of Canada Magnetic Resonance Imaging Index for Assessment of Spinal
 9 Inflammation in AS (SPARCC) were compared using a set of MRIs from 30
 10 patients participating in a clinical trial that compared a TNF-blocking drug
 11 with placebo [14-16]. The MRI sets were presented in randomized order to 9
 12 readers who scored them by all three methods. Repeatability and inter-reader
 13 reliability for every reader pair were highly heterogeneous across methods,
 14 although consistently higher ICCs (relative agreement) were found for the
 15 SPARCC method. Absolute agreement, however, was worst for the SPARCC
 16 method. This apparent contradiction can be explained by the fact that the
 17 scoring methods are based on different possible score ranges [17]. With regard
 18 to sensitivity to change and discrimination between active drug and placebo,
 19 all methods were appropriate. In conclusion, a definite preference for one of
 20 the 3 methods could not be made. While all methods would probably show
 21 similar performances in the context of a clinical trial, where sensitivity to
 22 change and discriminatory ability are of major interest, assessment of MRIs
 23 of AS patients included and followed up in a longitudinal cohort would make
 24 the SPARCC method less useful, because only 6 discovertebral units are taken
 25 into account (the most severely affected).
 26

27 ASSESSING DISEASE ACTIVITY IN AS: THE ASDAS

28
 29 **Chapter 8** describes the results of an analytical process to derive a composite
 30 score with optimal metrological performance in assessing disease activity in
 31 AS patients, the ASDAS (Ankylosing Spondylitis Disease Activity Score). This
 32 work resulted in 4 draft indices that combine patient-reported assessments
 33 of disease activity and acute phase reactants. All these 4 indexes showed high
 34 discriminatory ability, both in the cohort they were derived from (ISASS:
 35 *International Study on Starting tumour necrosis factor-blocking agents in Ankylos-*
 36 *ing Spondylitis*) and in an independent cohort that was chosen: the OASIS
 37 (*Outcome in Ankylosing Spondylitis International Study*) cohort. They showed a
 38 consistently good performance in distinguishing patients with high versus low
 39

1 disease activity, regardless of the external reference that was used. Most impor-
2 tantly, this discriminatory ability was systematically higher in comparison with
3 the current reference measure, the BASDAI (Bath Spondylitis Disease Activity
4 Index), which is often criticized because it is a strictly patient-reported outcome
5 measure [18]. Follow-up research, conducted by others, has since advanced
6 the concept: The ASDAS has further been validated in relevant subgroups,
7 confirming a stable performance in patients with AS with peripheral arthritis
8 or in patients with normal levels of acute phase reactants [19]. Definition of
9 disease activity levels have also been derived, as well as cutoff values to deter-
10 mine whether a change in ASDAS observed after therapeutic adjustment can
11 be regarded as clinically relevant: To separate inactive disease, moderate, high
12 and very high disease activity, cutoff values were derived: 1.3, 2.1 and 3.5 units.
13 A change ≥ 1.1 unit was considered clinically important improvement and a
14 change ≥ 2.0 units a major improvement [20].

15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39