



Data Harmonisation across Databases for Machine Learning Model Development to Predict Knee Osteoarthritis Progression

AUTHORS: S. Danso¹, P. Widera¹, P. M. Welsing², S. Peelen³, A. Tielmann⁴, F. Berenbaum⁵, F. J. Blanco⁶, I. K. Haugen⁷, L. Hussaarts³, M. Kloppenburg⁸, F. Lafeber², J. Larkin⁹, M. C. Levesque¹⁰, A. Mobasher¹¹, L. Paolozzi¹², F. Petit-Dop¹², J. Sellam⁵, W. E. van Spil², H. Weinans², C. Ladel⁴, J. Loughlin¹, J. Bacardit (Jaume.Bacardit@Newcastle.ac.uk)¹, for the APPROACH consortium

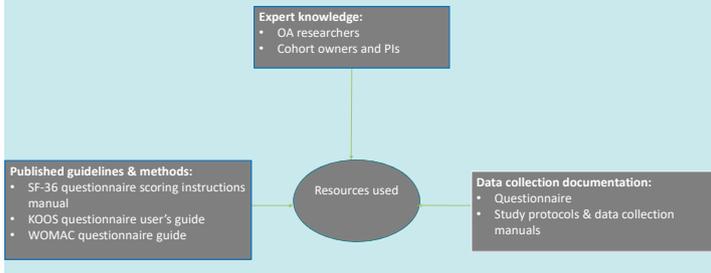
¹Newcastle Univ., Newcastle upon Tyne, United Kingdom, ²Univ. ir Medisch Centrum Utrecht, Utrecht, Netherlands, ³Lygature, Utrecht, Netherlands, ⁴Merck, Darmstadt, Germany, ⁵Hosp. Saint-Antoine, Paris, France, ⁶Servizo Galego de Saude, A Coruña, Spain, ⁷Diakonhjemmet Hosp. AS, Oslo, Norway, ⁸Leids Univ. ir Medisch Centrum, Leiden, Netherlands, ⁹GlaxoSmithKline, Philadelphia, PA, ¹⁰AbbVie, Chicago, IL, ¹¹Univ. of Surrey, Guildford, United Kingdom, ¹²Inst. de Recherches Servier, Suresnes, France

Background

- APPROACH is an Innovative Medicines Initiative/European Union funded consortium of scientists and clinicians from academia and industry
- APPROACH aims to use biomarkers to assess which patients are at risk of knee osteoarthritis (OA) and most likely to progress in a 2 year timeframe and prioritise them for inclusion to OA clinical trials
- The consortium will recruit subjects to APPROACH from existing OA cohorts using machine learning algorithms, hence harmonisation of the data is needed so we can use common models on different cohorts

Inherent structural differences

- Cross-population data harmonisation is needed to align data from the different OA cohorts prior to application of algorithms to recruit patients
- This process involves transforming the data from multiple sources
- This can be challenging due to variability in data collection protocols used resulting in inherent structural differences (ISD) in datasets

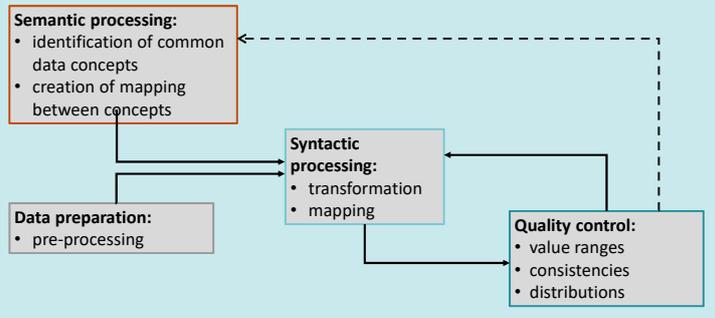


Methods

- We obtained data from five European OA cohorts: CHECK, DIGICOD, PROCOAC, MUST and HOSTAS
- The identified ISDs were categorised into semantics, syntactics and transformation differences.
- Developed specific strategies to address each type of difference identified based on questionnaire guidelines and expert knowledge.
- Semantics: conceptually similar data elements that could be mapped to common higher-level data elements, e.g. 'pain present on knee flexion' and 'pain present on knee extension' were mapped to 'pain on movement of knee'
- Syntactics – mapping : differences in coding practices between datasets, e.g. {1=positive, 2=negative} vs. {1=no, 2=yes}
- Syntactics - transformations: manipulations of representation schemes, where functions were employed to map from one type to another e.g. 'bodily pain score' = round (mean((6 - SF7) * 20, (5 - SF8) * 25)), where SF7 = 'bodily pain in the past 4 weeks' with values ranging from 1 = 'no pain' to 6 = 'very severe' and SF8 = 'pain interference in your work in the past 4 weeks' with values ranging from 1 = 'not at all' to 5='extremely'.

Harmonisation Pipeline

- The mapping rules were further executed as an automated software pipeline
- Quality control methods were employed to identify errors in patterns and outliers in the harmonised datasets



Results

- Our harmonisation process resulted in eight different datasets. For each cohort, we have one harmonised version of CHECK to train machine learning models and the actual cohort data to select patients
- Over 5520 records were harmonised across the five cohorts
- Overall, harmonisation of the demographic variables was easiest, and due to variations in data collection protocols, harmonisation of lifestyle and radiographic variables was most difficult

Dataset	No. of variables	No. of variables mapped
MUST	884	80
HOSTAS	124	50
DIGICOD	772	57
PROCOAC	301	53

Conclusions

- Data harmonisation was used to facilitate recruitment for APPROACH
- Transforming the data from existing cohorts to a common set of variables enabled the use of machine learning technology to identify fast OA progressors
- We believe this strategy will enhance the power and efficiency of OA clinical trials

