

VALIDITY, AGREEMENT, CONSENSUALITY AND ANNOTATED DATA QUALITY

Anaëlle Baledent, Yann Mathet, Antoine Widlöcher, Christophe Couronne, Jean-Luc Manguin

Normandie Univ, UNICAEN, ENSICAEN, CNRS, GREYC, Caen, FRANCE

{anaelle.baledent, yann.mathet, antoine.widlocher, christophe.couronne, jean-luc.manguin}@unicaen.fr



1. Context

In NLP, reference annotated datasets are required for various common tasks. But their production is known as a difficult problem, from both a theoretical and practical point of view. The multi-annotation is most of the time a necessity to build these reference corpora. We focus here on the complex relations between agreement and reference (of which agreement among annotators is supposed to be an indicator), and the emergence of consensus.

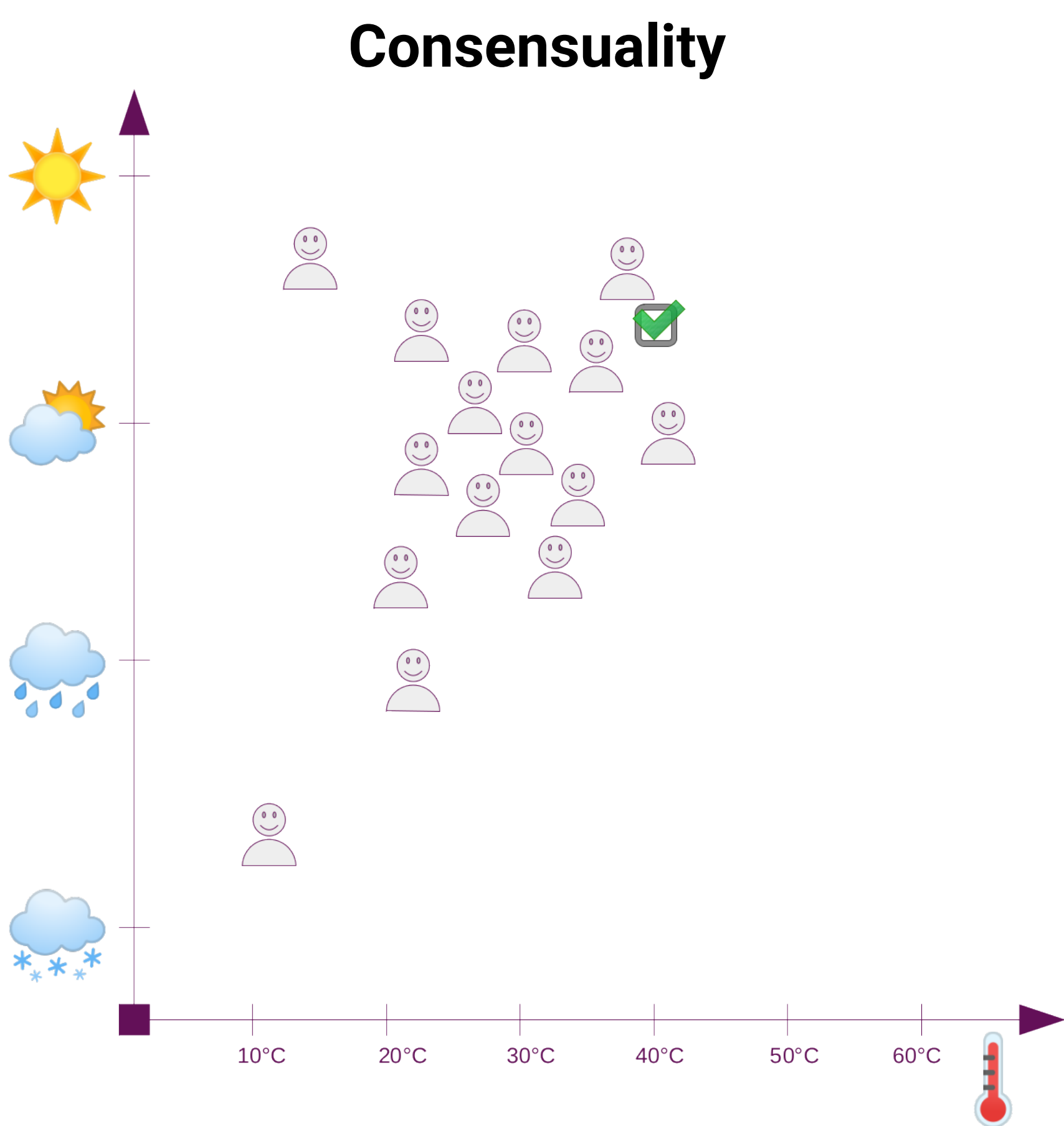
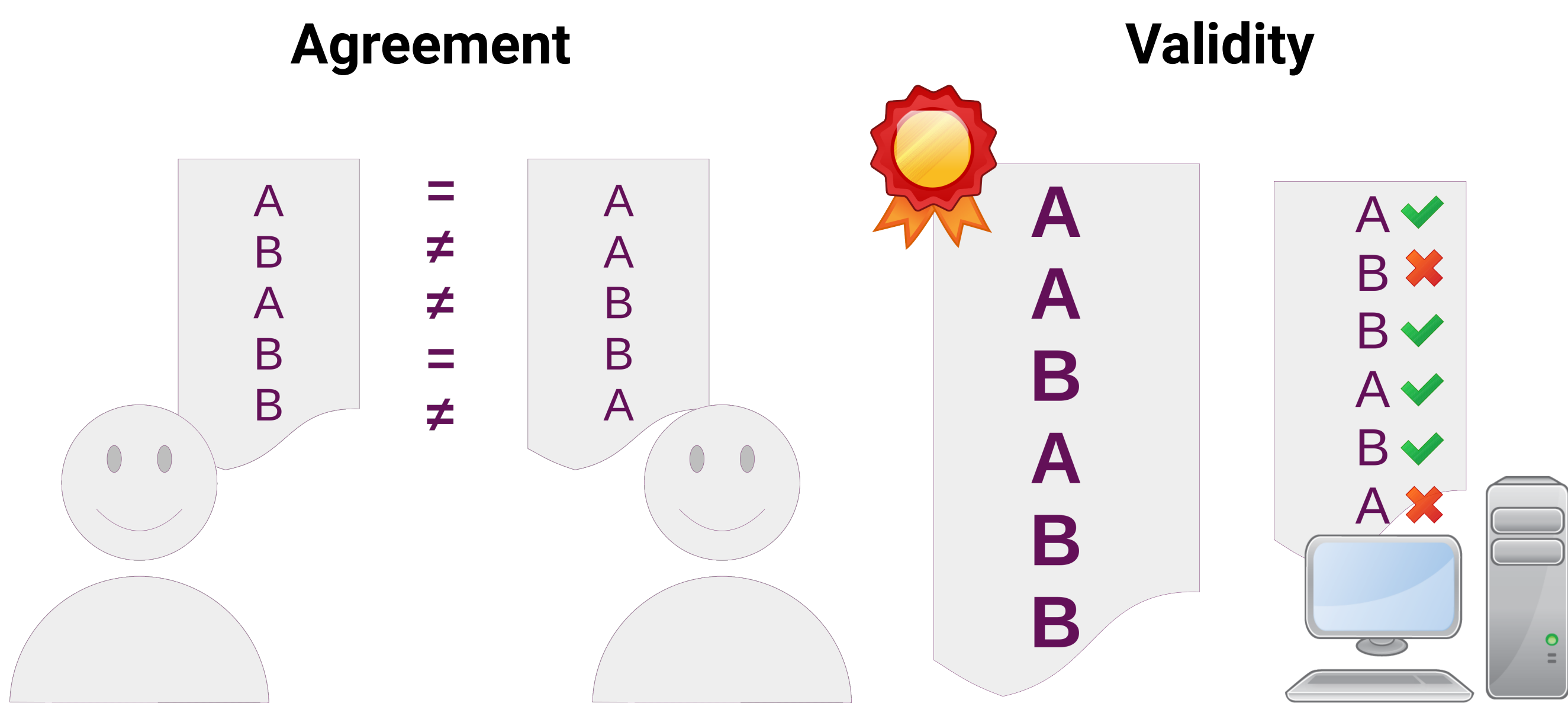
2. How old are these people?

Corpus : 100 photographs (collected from WIKIMEDIA COMMONS) of persons whose ages are within the range 3 months – 97 years.



Estimated age	17.84 (\pm 5.99)	41.98 (\pm 9.69)	53.57 (\pm 6.3)
Real age	11	64	54

3. Agreement, validity and consensuality



Computing the consensuality

The **annotator's consensuality degree regarding a group** is given by the difference between the disagreement of this group deprived of this annotator, and the disagreement of this group.

4. Consensuality ranking versus performance ranking

There is no strong correlation between consensuality and annotators performance.

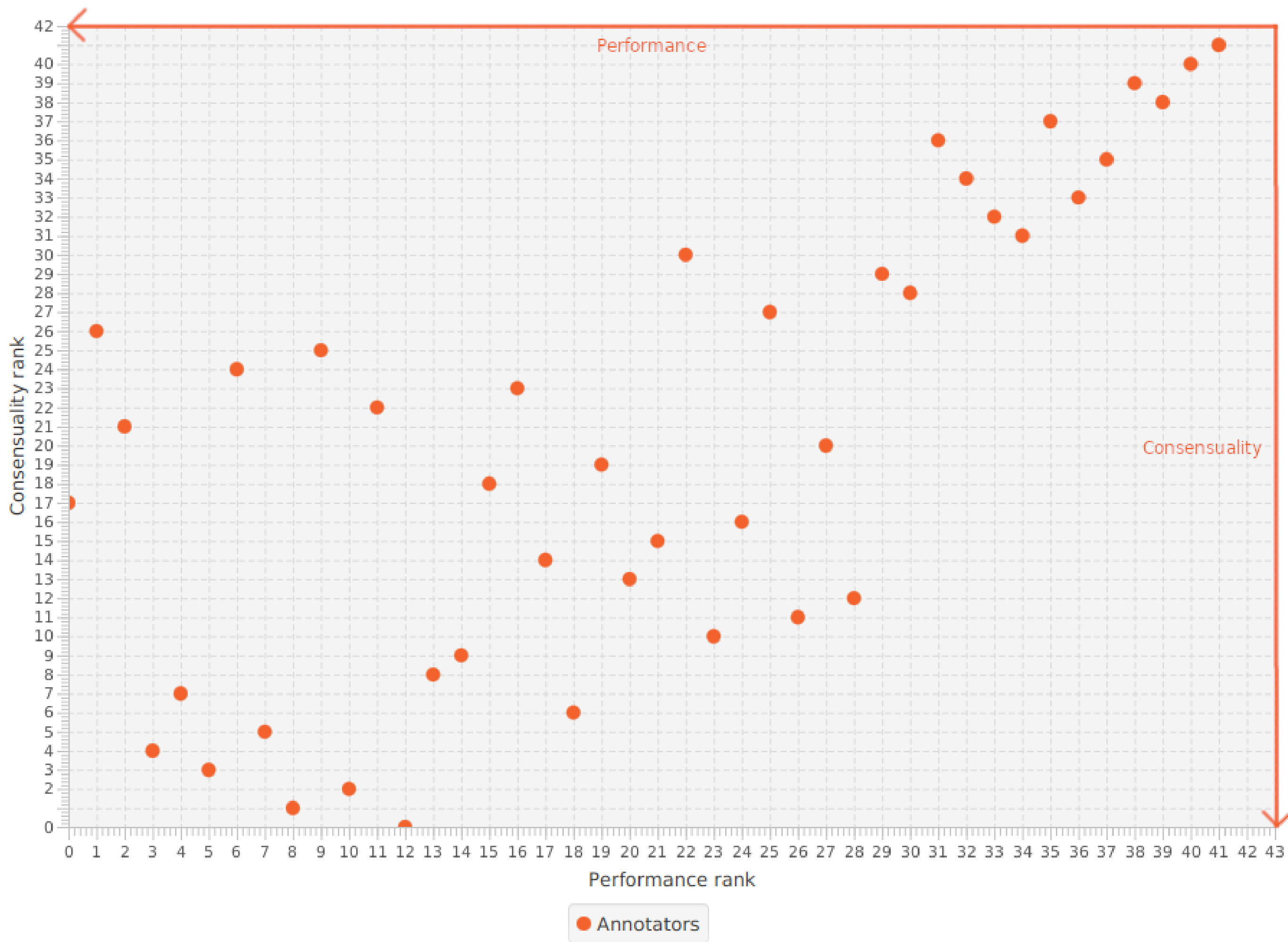


FIGURE 1 : Annotators' ranks according to their performance and their (progressive) consensuality.

5. Distinguishing between initial and progressive consensuality

Two types of consensuality

Initial : Ranking the annotators from the initial group.
Progressive : Recomputing the consensuality ranks after removing the least consensual annotator, and repeat.

Progressive consensuality is better than initial one. Removing the less consensual annotators improves the group performance.

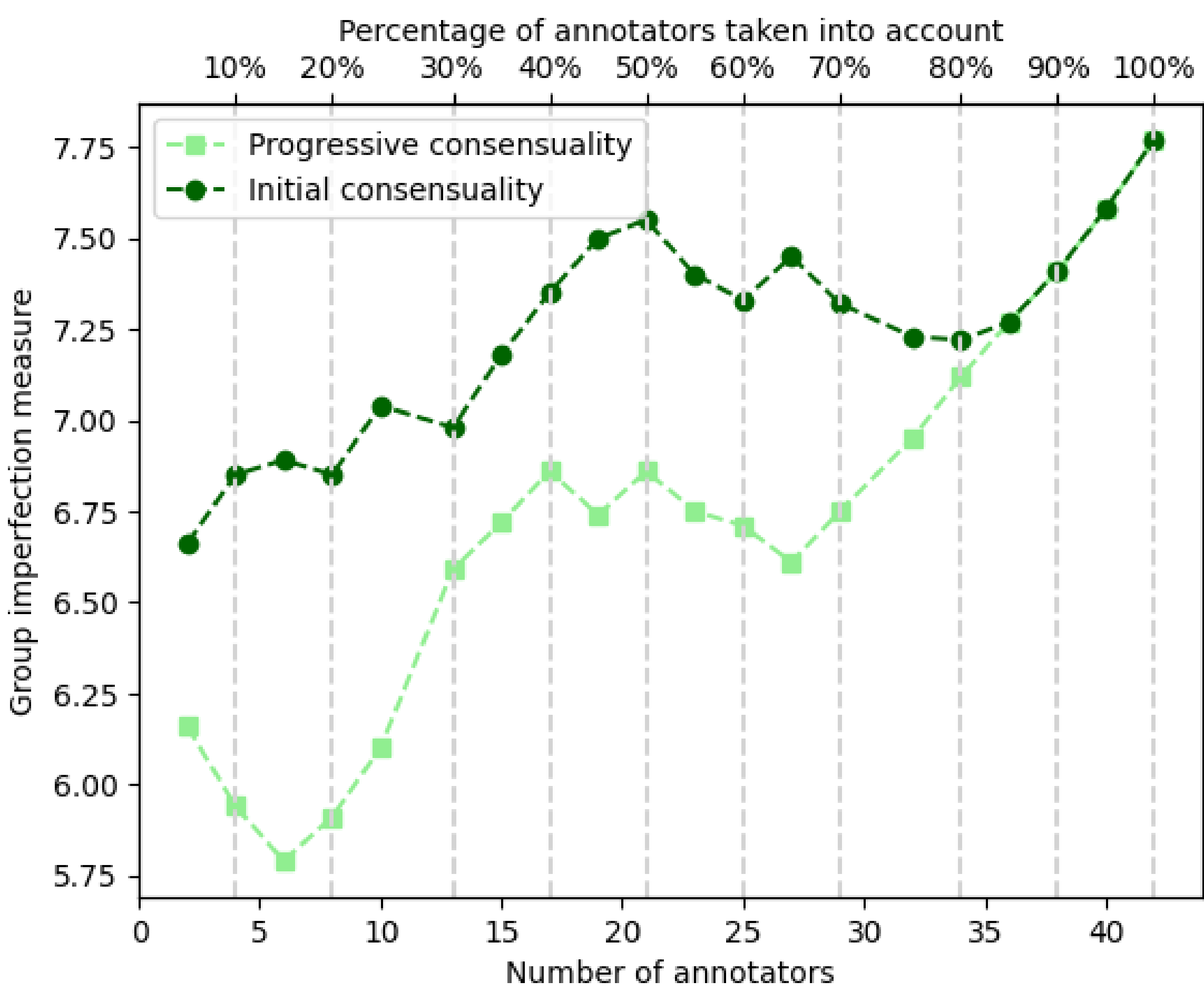


FIGURE 2 : Imperfection of most consensual annotators