# Statistical Methods for Annotation Analysis

# Synthesis Lectures on Human Language Technologies

## Editor
**Graeme Hirst,** *University of Toronto*

Synthesis Lectures on Human Language Technologies is edited by Graeme Hirst of the University of Toronto. The series consists of 50- to 150-page monographs on topics relating to natural language processing, computational linguistics, information retrieval, and spoken language understanding. Emphasis is on important new techniques, on new applications, and on topics that combine two or more HLT subfields.

## Statistical Methods for Annotation Analysis
Silviu Paun, Ron Artstein, and Massimo Poesio
2021

## Validity, Reliability, and Significance: Empirical Methods for NLP and Data Science
Stefan Riezler and Michael Hagmann
2021

## Pretrained Transformers for Text Ranking: BERT and Beyond
Jimmy Lin, Rodrigo Nogueira, and Andrew Yates
2021

## Automated Essay Scoring
Beata Beigman Klebanov and Nitin Madnani
2021

## Explainable Natural Language Processing
Anders Søgaard
2021

## Finite-State Text Processing
Kyle Gorman and Richard Sproat
2021

## Semantic Relations Between Nominals, Second Edition
Vivi Nastase, Stan Szpakowicz, Preslav Nakov, and Diarmuid Ó Séagdha
2021

Statistical Language Models for Information Retrieval
ChengXiang Zhai
2008

Statistical Methods for Annotation Analysis

Silviu Paun, Ron Artstein, and Massimo Poesio

# Statistical Methods for Annotation Analysis

Silviu Paun
Queen Mary University of London

Ron Artstein
University of Southern California

Massimo Poesio
Queen Mary University of London and Turing Institue

# ABSTRACT

Labelling data is one of the most fundamental activities in science, and has underpinned practice, particularly in medicine, for decades, as well as research in corpus linguistics since at least the development of the Brown corpus. With the shift towards Machine Learning in Artificial Intelligence (AI), the creation of datasets to be used for training and evaluating AI systems, also known in AI as corpora, has become a central activity in the field as well.

Early AI datasets were created on an *ad-hoc* basis to tackle specific problems. As larger and more reusable datasets were created, requiring greater investment, the need for a more systematic approach to dataset creation arose to ensure increased quality. A range of statistical methods were adopted, often but not exclusively from the medical sciences, to ensure that the labels used were not subjective, or to choose among different labels provided by the coders. A wide variety of such methods is now in regular use. This book is meant to provide a survey of the most widely used among these statistical methods supporting annotation practice.

As far as the authors know, this is the first book attempting to cover the two families of methods in wider use. The first family of methods is concerned with the development of labelling schemes and, in particular, ensuring that such schemes are such that sufficient agreement can be observed among the coders. The second family includes methods developed to analyze the output of coders once the scheme has been agreed upon, particularly although not exclusively to identify the most likely label for an item among those provided by the coders.

The focus of this book is primarily on Natural Language Processing, the area of AI devoted to the development of models of language interpretation and production, but many if not most of the methods discussed here are also applicable to other areas of AI, or indeed, to other areas of Data Science.

# KEYWORDS

*In memory of Janyce Wiebe*

# Contents

# Preface

When, almost 15 years ago, two of us (Massimo and Ron) completed what was to become our (Artstein and Poesio, 2008) paper, we felt elated at reaching a milestone after three years of hard work. But we were also aware that the stage we had reached, on the one hand, left a number of open questions, particularly about the interpretability of coefficients of agreement. On the other hand, it failed to cover important areas of research within the field of statistical methods for annotation analysis, such as latent models of agreement or probabilistic annotation models, which we felt could make important contributions to data creation and use within Computational Linguistics even though, at the time, it had only begun being experimented with in such pioneering work as Beigman Klebanov and Beigman (2009), Bruce and Wiebe (1999), Carpenter (2008), and Reidsma and Carletta (2008). 2008 was also the year of the seminal paper by Snow et al. (2008) that started the crowdsourcing revolution in Natural Language Processing (NLP), in which such methods were to play an essential role. We were therefore delighted when, many years later, Graeme Hirst and Mike Morgan offered us the opportunity to make further progress along the path begun in that paper. This book is the result of that progress.

As the 2008 paper covered the material in its scope—coefficients of agreement—in a, we thought, reasonably thorough way; and, as even more thorough works covering that material have appeared since, such as Gwet (2014), we didn't attempt to expand our coverage of that topic much in this book. We merely aimed to incorporate some material that had been left out in the original paper and to update the presentation to take into account more recent proposals, e.g., on unitizing. Our effort focused instead primarily on covering methods that had not been covered at all in the 2008 paper, also because meanwhile those techniques have become much more widespread in NLP. We hope this book offers as accessible an introduction to latent models of agreement, probabilistic models of aggregation, and learning directly from multiple coders, as the 2008 paper did for coefficients of agreement.

Silviu Paun, Ron Artstein, and Massimo Poesio
September 2021

# Acknowledgements

We are extremely grateful to Graeme Hirst and Mike Morgan for offering us the opportunity to continue our work in this area, and to the many colleagues who collaborated with us on this effort at various stages and/or offered their comments and support. We would especially like to thank Bob Carpenter, Jon Chamberlain, Janosch Haber, Dirk Hovy, Tommaso Fornaciari, Udo Kruschwitz, Chris Madge, Becky Passonneau, Barbara Plank, Edwin Simpson, Alexandra Uma, and Juntao Yu, for much inspiration, feedback, and help over the years. We would also like to thank Jacopo Amidei, Lora Aroyo, Valerio Basile, Beata Beigman-Klebanov, Raffaella Bernardi, Alex Braylan, Chris Callison-Burch, Jean Carletta, Chris Cieri, Barbara di Eugenio, Anca Dumitrache, James Fiumara, Iryna Gurevych, Nancy Ide, Klaus Krippendorff, Walter Lasecki, Matt Lease, Mark Liberman, Yann Mathet, Sabine Schulte im Walde, Yannick Versley, Renata Vieira, Aline Villavicencio, and Janyce Wiebe.

Silviu Paun, Ron Artstein, and Massimo Poesio
September 2021