

Project Essay: History of R in Linguistics

Introduction:

Quantitative methods are frequently relevant to contemporary research in the field of linguistics. Quantitative analysis allows for advancements in the various sub-fields of linguistics, such as in phonetics, sociolinguistics, corpus linguistics, and syntax studies, as well as language acquisition studies. Use of programming languages facilitates the advancement of empirical and data-driven lines of research; accordingly, the need for robust statistical analysis tools has grown. The R programming language has come to be considered one of the best currently available tools to facilitate the statistical analysis of linguistics research. This is in large part due to R offering a comprehensive statistical ecosystem that is capable of handling the complexities, variabilities, and theoretical nuances in linguistic data analysis. This essay aims to examine the relevance of R to statistical analysis in linguistics, particularly focusing on the history and applications of the R language, including its status as part of free open source software (FOSS). The essay is divided into three sections: the first covers a history of the R language, including its roots in the S language; the second moves to examine its use and applications in linguistics; the third compares its use in linguistic data analysis against other possible alternative statistical analysis tools, such as Python.

History of R:

The history of the R programming language can be traced back to the S programming language, which was created between the 1970s to the 1980s at Bell Labs by John Chambers et al (source). S was initially developed as an interactive statistical computing environment that would allow researchers to manipulate data in a way that was both flexible and iterative. It was different from compiled languages at the time since, unlike contemporary compiled languages, S placed emphasis on syntax that was readable by humans and also on the exploratory nature of statistical work. The focus on flexibility, interactivity, and exploration anticipated the methodological needs of various scientific fields that would later deal with increasingly complex empirical datasets. Among these fields, linguistic data analysis required a programming language that would be particularly sensitive to the nuances, complexities, and variabilities

intrinsic to the field, e.g., high inter-speaker and intra-speaker variation, and context-sensitive linguistic behaviours.

R was developed as a free open source software re-implementation of the existing S language (unsurprisingly, ‘S’ stood for ‘Statistics’), with the aim of retaining the statistical strengths of the S language whilst expanding user access through the GNU General Public License. It was developed in the early 1990s, dating it approximately thirty years, by Ross Ihaka and Robert Gentleman in the University of Auckland. Open source models facilitate the growth of a community of contributors, which naturally leads to continual improvement of the given language as part of a vast ecosystem of usable packages. The collaborative and continually evolving statistical environment fostered by the R language is thus intrinsic to its use in the ever-evolving and inherently collaborative field of linguistic analysis.

The usefulness of the R language to linguistic data analysis merely begins at this stage. The establishment of the Comprehensive R Archive Network (CRAN) enabled systematic distribution of contributed packages, which allowed the language to evolve into a modular and extensible environment capable of accommodating specialised analytic techniques.

This evolution—driven by a global community of statisticians, computer scientists, and domain experts—positioned R as an ideal platform for linguistics, a field whose analytic needs are both interdisciplinary and rapidly developing.