SciDataCon 2025



Contribution ID: 124

Type: Presentation

## Data-Driven Risk Identification in Supervision Reports of the Ministry of Health

Monday 13 October 2025 15:47 (11 minutes)

In response to inefficiencies in governmental regulation, such as excessive regulation, an overabundance of laws and procedures, lack of flexibility, and disregard for the costs, countries around the world began efforts to optimize regulation, in part by privatization. The trend toward privatization of social services necessitated substantial development of governmental supervision practices. Risk management in regulation emerged as an efficient approach to supervising public sector services, assisting regulators in deciding the extent of intervention necessary to prevent harm to the public interest and to ensure that service recipients are protected and safe. Today, risk management in supervision is a critical component of decision-making processes under conditions of uncertainty and is recognized as one of eleven principles of best practices in supervision and enforcement by the OECD.

This study explores the potential of artificial intelligence (AI) in identifying and categorizing risks from unstructured open text, using advanced natural language processing (NLP) architectures such as Dicta and HeBERT. The research aimed to develop a methodology for analyzing supervision reports from the healthcare sector, enabling risk detection and classification into predefined categories.

The study's results indicate high performance of the Dicta model in identifying and classifying risks from unstructured text, achieving an accuracy of 93.3%, a recall of 85.9%, and an F1 score of 92.3%. In comparison, the HeBERT model yielded lower results across all metrics. In the multi-class classification task, Dicta also outperformed HeBERT, with an accuracy of 74.4% versus 65.1%, respectively. These differences were statistically significant (p < 0.05), underscoring the advantages of using Hebrew-adapted models, particularly those tailored to the healthcare domain.

The study highlights the critical role of semantic features and keywords in risk identification. It also addresses challenges associated with ambiguous sentences and overlapping categories, emphasizing the need for future research to develop multi-category classification algorithms. While Dicta showed superior performance in identifying key categories such as "Infrastructure, Equipment, and Logistics" and "Medical Services and Quality of Care," HeBERT exhibited limitations in distinguishing mid-range categories, resulting in higher error rates.

The findings suggest practical applications for regulatory bodies, such as optimizing resource allocation, enhancing decision-making through data-driven insights, and improving transparency and service quality. Despite its promising results, the study acknowledges limitations, including the reliance on a single corpus of healthcare supervision reports and the constrained sample size. Future research should expand the corpus and explore AI techniques for less structured texts.

This research provides a foundational framework for applying AI to risk detection in healthcare and other domains, offering valuable insights for improving supervision, monitoring, and service delivery.

Primary author: Dr ZADOK, Avital (University of Haifa)

Co-author: RABAN, Daphne (University of Haifa, CODATA Israel NC)

Presenter: RABAN, Daphne (University of Haifa, CODATA Israel NC)

Session Classification: Presentations Session 2: Data and Research & Data Science and Data Analysis Track Classification: SciDataCon Persistent Themes: Data Science and Data Analysis