# Exploring quality dimensions in trustworthy Machine Learning in the context of official statistics: model explainability, accuracy and uncertainty

Saeid Molladavoudi, Wesley Yung (Statistics Canada)

_saeid.molladavoudi@statcan.gc.ca_

*Abstract*

Despite the fact that National Statistical Organizations (NSOs) continue to embrace and adopt Machine Learning (ML) methods and tools in a variety of areas of their operations, including data collection, integration, and processing, it is still not clear how these complex and prediction-oriented approaches can be incorporated into the quality standards and frameworks within NSOs. This presentation focuses on two of the quality dimensions: model explainability and accuracy (including uncertainty) and it builds upon and extends the previous Quality Framework for Statistical Algorithms (QF4SA). The implications of the current methods for explainable ML and uncertainty quantification will be examined in further detail, as well as their possible uses in statistical production, such as continuous model monitoring in intermediate ML classifications and auto-coding phases. This strategy will ensure that human subject-matter experts, who are an essential component of every statistical programme, are effectively integrated into the life cycle of ML projects. It will also guarantee to maintain the quality of ML models in production, adhere to the current quality frameworks within NSOs, and ultimately boost confidence and trust in these emerging technologies.