

Credit scoring models are commonly trained and evaluated on data of accepted applicants whose repayment behavior has already been observed. However, a representative sample should be drawn from a population that applies for credit including both accepted and rejected cases. This issue degrades the model’s real-world performance due to the difference between the training (sample-based) and the real (population-based) probability density functions. Reject inference refers to techniques that assign labels to the rejected samples.

The contributions of this paper are two-fold. First, we introduce a self-learning framework with distinct training regimes for its iterative labeling and training stages. This idea results in having a highly calibrated model for labeling, which minimizes the noise introduced on newly labeled rejected cases, and a highly discriminative model for training, which maximizes the performance of a scoring model.

Second, we propose a novel evaluation strategy to guide holdout-based model selection in the presence of sample bias without labeling rejected cases in the test set, which is subject to noise. We show that the standard practice of doing model selection based on the model’s performance on the accepted cases may lead to a suboptimal bias-variance trade-off.

Experiments on a large-scale, high-dimensional real-world credit scoring data set demonstrate the effectiveness of the proposed approach when compared to reject inference techniques conventionally used in this domain. The data set used in this paper includes a labeled unbiased sample containing both accepts and rejects, which gives us a unique opportunity to evaluate the real-world performance gains from reject inference.