

Normalizing Pandemic Data for Credit Scoring

Abstract

The COVID-19 pandemic created abnormal credit risk conditions that did not align well with pre-2020 credit scores. Since the pandemic, most organizations have either excluded the period 2020-2021 from their modeling or included it without adjustment, leaving it as noise in the data. Model validators and examiners have been divided about requiring one of these approaches or defaulting to model developer judgment. None of this is ideal from a model development perspective.

We have found that a technical solution is available. Our analysis uses lifecycle and environment outputs from an Age-Period-Cohort analysis as fixed offsets to the credit score development. Panel data is used, so the credit score is developed with a discrete time survival model approach. We tested both logistic regression and stochastic gradient boosted regression trees as estimators with the panel data and APC inputs.

For this research, we used Fannie Mae data. The APC model was estimated on the full available history, from 2005 through 2024. The origination scores were estimated on two-year periods from 2016 through 2024 and tested on all other periods, including a score that was developed on the full period. All models were also tested on comparably prepared data from Freddie Mac for cross-validation.

The origination scoring variables available from both Fannie Mae and Freddie Mac were:

BORROWER_CREDIT_SCORE_AT_ORIGINATION
LOAN_PURPOSE
NUMBER_OF_BORROWERS
ORIGINAL_COMBINED_LOAN_TO_VALUE_RATIO_CLTV
ORIGINAL_DEBT_TO_INCOME_RATIO
ORIGINAL_INTEREST_RATE
PROPERTY_TYPE

Doing this, we find that pandemic data can be used in the credit score development process without degrading the rank ordering performance of the score. In fact, the pandemic period jumps from being the least predictable in rank ordering to the most predictable. This holds true for both logistic regression estimation and stochastic gradient boosted regression trees, and is significantly better than the same models created without APC inputs. Notably, the models did not degrade significantly when tested out-of-time. This is different from commonly reported experiences with machine learning models but consistent with prior experiences integrating survival modeling with machine learning.

Authors & Affiliations

Joseph Breeden¹

¹Deep Future Analytics LLC, Santa Fe, USA