

Executive Summary: Regression Analysis and Modeling

TikTok claims classification project

OVERVIEW

The TikTok data team aims to create a machine learning model to aid in classifying claims for user submissions. Notably, the team observed a higher likelihood of users with verified accounts posting opinions. Given the ultimate objective of classifying claims and opinions, it becomes crucial to develop a model that predicts the behavior of verified accounts, which tend to lean towards posting more opinions. In this phase of the project, the data team constructed a logistic regression model specifically designed to predict `verified_status`.

PROJECT STATUS

The choice of the variable `verified_status` for this regression model was driven by the observed relationship between the verified account type and the nature of video content. Opting for a logistic regression model was influenced by the data type and distribution characteristics.

A LOOK AT THE MODEL RESULTS

The logistic regression model demonstrated a precision of 69% and a recall of 66% (weighted averages). The model also achieved an f1 accuracy of 66%. These results provide valuable insights into video features, as further discussed in the "Key Insights" section."

NEXT STEPS

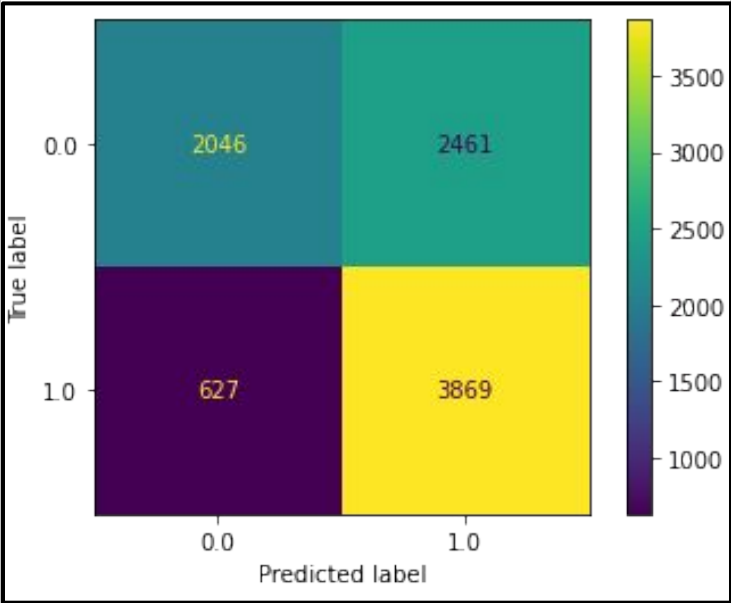
The subsequent phase involves building a classification model to predict the status of claims made by users. This constitutes the final project and aligns with the original expectations from the TikTok team. With the accumulated information, there is now sufficient context to thoroughly analyze the results of this model, incorporating insights into user behavior.

KEY INSIGHTS

According to the estimated model coefficients from the logistic regression, there is a tendency for longer videos to be associated with higher odds of the user being verified.

However, other video features display small estimated coefficients in the model, suggesting that their association with verified status is minimal. Consequently, besides video length, other video features do not appear to be significantly linked to verified status.

Confusion matrix for logistic regression model



Upper-left: the number of videos posted by unverified accounts. Upper-right: the number of videos posted by unverified accounts. Lower-left: the number of videos posted by verified accounts. Lower-right: the number of videos posted by verified accounts.