

TikTok Claims Classification Project

Machine Learning Model

Project Overview

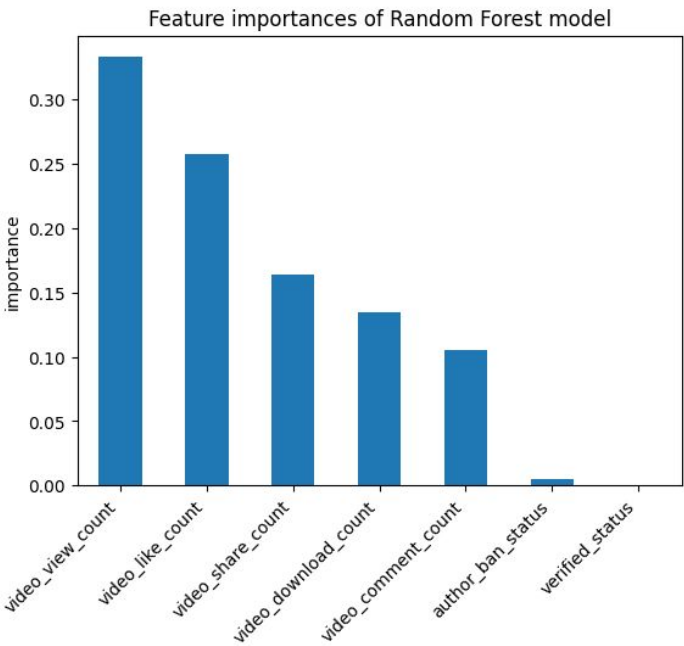
The data team is developing a machine learning algorithm to precisely detect whether a video contains claims or opinions. After Hypothesis Testing and Feature Engineering it is time to build the machine learning model.

Key Insights

1. To find the best model, both Random Forest and XGBoost were run in a GridSearch with a predefined set of hyperparameters. The decisive metric is recall because it measures the number of actual claims that are falsely predicted as opinions.
2. The Random Forest model beats the XGBoost with a slightly better recall of 98.9%.
3. The Random Forest model performs a recall of 99.0% and an accuracy of 99.5% on the test dataset.
4. Following features are listed in descending order of importance for the ML model prediction: view count (33%), like count (26%), share count (16%), download count (13%), comment count (10%), author ban status (0.5%), author verified status (0.2%)

Details

model	recall	accuracy
Random Forest	98.9%	99.4%
XGBoost	98.8%	99.3%



Next Steps:

The Data Team recommends the use of the Random Forest model for incoming reports. If the backlog is still overloaded then an automated claim check algorithm that checks all TikTok videos at given interval (e.g. every 24h) until the video has been checked a certain number of times, should be developed.