

# Video Classification Project | Machine Learning Results

Prepared For: TikTok Data Team

## ISSUE / PROBLEM

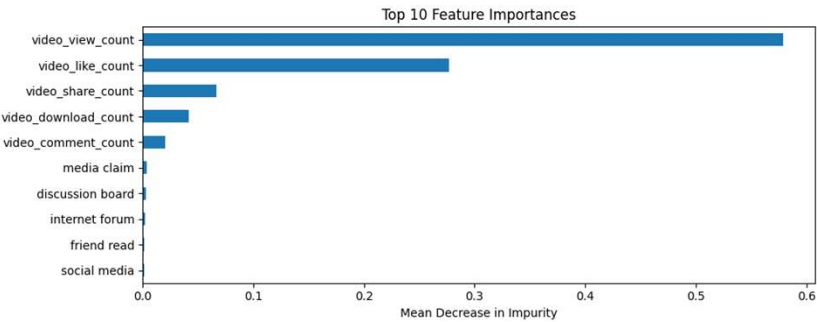
- The TikTok data team is developing a data analytics project to identify whether a TikTok video expresses a claim or an opinion. Prior analysis indicated that video engagement metrics were highly related to claim status. A model capable of this distinction will significantly reduce moderators' workload.

## IMPACT

- Historically, videos that make claims tend to receive higher views, likes, and shares, often featuring harmful content. Consequently, it is ethically crucial to develop a model that can identify these videos early, preventing their widespread dissemination across the platform.
- In order to achieve this, the data team created two classification models based on tree algorithms: Random Forest and XGBoost. These models were utilized to predict outcomes on a validation dataset, and the one demonstrating the highest recall score was chosen as the final model. Subsequently, this final model was evaluated against the held-out (test) dataset to assess its expected performance on new, unseen data.

## RESPONSE

- Both models showed strong performance; however, the random forest model achieved a superior recall score of 99.5% and was consequently chosen as the top-performing model.



## KEY INSIGHTS

- The model showed excellent performance on the hold-out dataset, misclassifying only 9 out of 3,817 samples.
- Analysis revealed that the main predictive features were video engagement metrics, such as view count, likes, shares, comments, and downloads, which greatly enhanced the dataset's predictive power.
- We can confidently state that the level of video engagement correlates with whether the video presents a claim or an opinion.
- Given the model's strong results, it's recommended to test it on more user data subsets to ensure its effectiveness remains consistent across different variations before deployment.

