

# TikTok Claims Classification Project

Logistic regression model - Executive Summary

## ISSUE / PROBLEM

The TikTok data team seeks to develop a machine learning model to assist in the classification of claims for user submissions. In this part of the project, We decided to analyze the verified status. We've noticed that not-verified users tend to upload more claim videos.

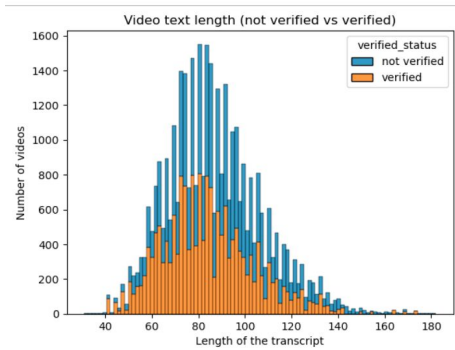
## RESPONSE

The Tik Tok data team decided to conduct a logistic regression model. With the purpose of figuring out which variables are more related to the verified status of the users. In order to accomplish this, the team will choose the verified status as the outcome and the different count continuous variables, and the claim and verified statuses as the features or independent variables of the model.

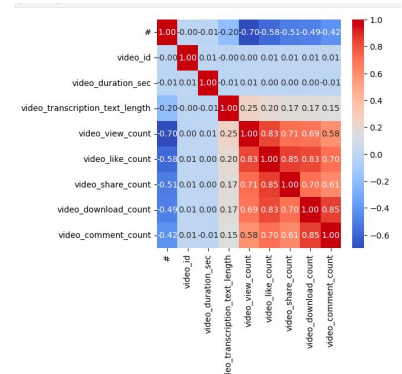
## IMPACT

During past stages, I discovered there are a majority of not-verified users, so I should balance the data in order to build the model. On the other hand, the model had a success rate of 65% in predicting true values. Besides, we discovered that video duration is the variable more closely related to whether or not a user is verified.

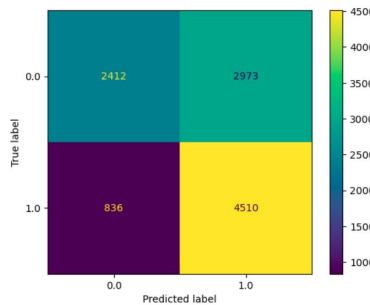
A key component of this project's logistic regression model involves the creation of a correlation matrix, a comparison between how long the transcriptions are made for not-verified users vs. the ones made by verified users, and a confusion matrix about the predictions made by the model. These are illustrated in the next visuals.



As we can see in this visual. Not-verified users tend to write longer transcripts than verified users.



All the counts used as features for the train and the test sample (likes, comments, shares, downloads, and views) are highly correlated between them. This makes sense because the more popular a video is, the higher the numbers in it counts.



After checking the confusion matrix, we can say it was capable of predicting 0.65, or 65%, of the outcomes. The majority of these outcomes were true/false positives, which confirms that the majority of the users are unverified users.

## KEY INSIGHTS

As part of the model, we did a classification report of it. Here are the results:

### Verified

- The model had a precision of 75% in its prediction.
- 47% of the positive values were predicted with success.
- Even when the recall percentage is low and the f1-score is 58%. Under the goals of the project, it is more important to be able to predict the not-verified users. Which we had more success with.

### Not verified

- The model had a precision of 61% in its prediction.
- 84% of the positive values were predicted with success.
- The model had an f1-score of 71%.