

# Executive Summary: Regression Analysis

## Milestone 5

### Overview

The objective of this milestone was to build and evaluate a logistic regression model to predict whether a TikTok user is verified based on video characteristics and account attributes. This analysis aims to uncover which factors are most associated with verified status, support TikTok’s broader claims classification project, and provide insights that will guide the development of a final predictive model distinguishing claims from opinions.

### Data Preparation

Initial data exploration showed that the dataset was highly imbalanced, with 93.7% of posts from unverified accounts and 6.3% from verified accounts.

To address the class imbalance, we applied resampling techniques, resulting in a balanced dataset with equal representation (50% verified, 50% unverified).

Outliers were identified in **video\_like\_count** and **video\_comment\_count**. Using the IQR method, we capped extreme values to minimize their influence on the model.

We developed a **logistic regression model** to predict **verified\_status** using video engagement metrics. This model reveals which video characteristics most influence the likelihood of verification and serves as a diagnostic tool in the larger claims detection framework.

### Details

**Features Used:** video\_duration\_sec, video\_view\_count, transcription\_length, claim\_status, author\_ban\_status

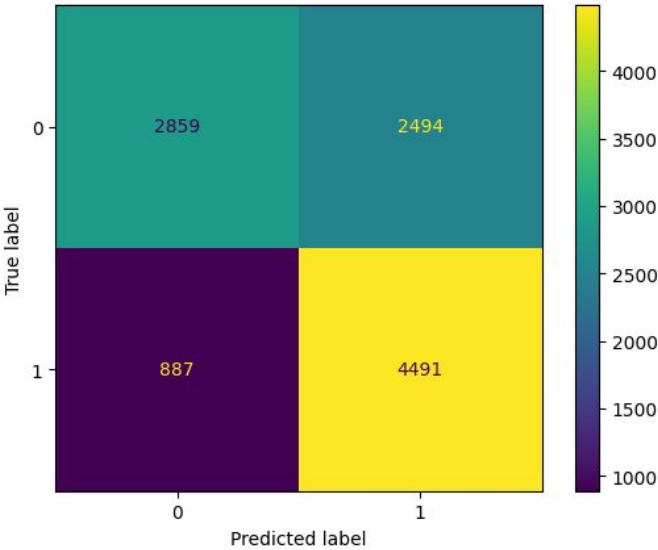
**Key Correlations:** Strong correlation observed between engagement metrics (e.g., likes, shares, downloads)

**Top Predictors (Log-Odds Coefficients):**

- claim\_status\_opinion: +1.47506
- author\_ban\_status\_banned: -0.43319
- author\_ban\_status\_under\_review: -0.37031
- transcription\_length: -0.00192
- video\_duration\_sec: -0.00174

**Model Performance:**

- Accuracy: 68.5%
- Precision: 64.3%
- Recall: 83.5%
- F1 Score: 72.7%



### Next Steps

The next step is to build the final classification model to predict claim status. With insights from the regression analysis, we are now better equipped to interpret user behavior and improve model performance.