# TikTok project: detect claims and opinions
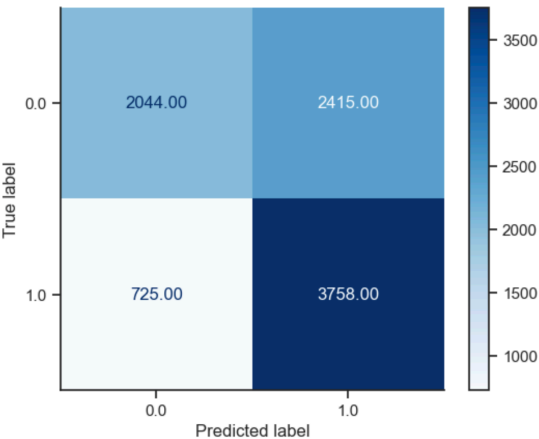
## Milestone #4: Logistic regression

## Overview

The objetive of this project is to implement a machine learning model to classify the videos uploaded to TikTok into claims (unsourced statements) and opinions (personal thoughts) in order to reduce moderating times and reduce economical costs.

## Objective

The objective of this milestone was to evaluate **how video features relate to the verification status** of the users. For this purpose, we implemented a logistic regression model and we evaluated the coefficients associated to each independent variable. We followed a simple train-test split approach. We evaluated the multicollinearity of the variables to comply with the logistic regression model assumptions.

## Results

- From the coefficients of the logistic regression model: the video duration and the comment count are the two variables which most influence the outcome variable.

- Every second of the video duration increase the log odds of the probability of an user being verified by 0.009 (considering the rest of the independent variables are held constant).

- Therefore, **longer videos tend to be associated with higher odds of the user being verified.** In the case of the comment count, each additional comment reduces the former log odds by 0.0004.

- The logistic regression models predicts the outcome variable *verified_status* with **higher accuracy than random choice** (0.65 vs 0.5)



```
                  precision    recall  f1-score   support

      verified       0.74      0.46      0.57      4459
  not verified       0.61      0.84      0.71      4483

      accuracy                           0.65      8942
     macro avg       0.67      0.65      0.64      8942
  weighted avg       0.67      0.65      0.64      8942
```

## Next Steps

The next step is to **construct a classification model that will predict the status of claims made by users.** That is the final project and original expectation from the TikTok team. Now, there is enough information to analyze the results of that model with helpful context around user behavior.