

TikTok Claims Classification Project

Hypothesis Testing

Project Overview

The data team is developing a machine learning algorithm to precisely detect whether a video contains claims or opinions. After conducting Exploratory Data Analysis we performed hypothesis testing on all variables to determine which variables are statistically significant for future predicting model.

Details

Key Insights

- 1. Following variables/features are statistically significant and will be used for further steps: views, likes, shares, downloads, comments, author ban status and author verified status
- 2. The video duration and the video transcription text are not used for further steps because they have no statistically significant difference between claims and opinions
- 3. For further steps the original data will be used and not the log-transformed data, because it will not be necessary anymore.

Next Steps:

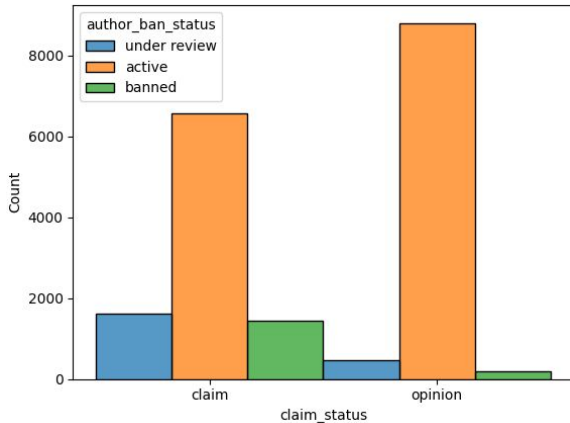
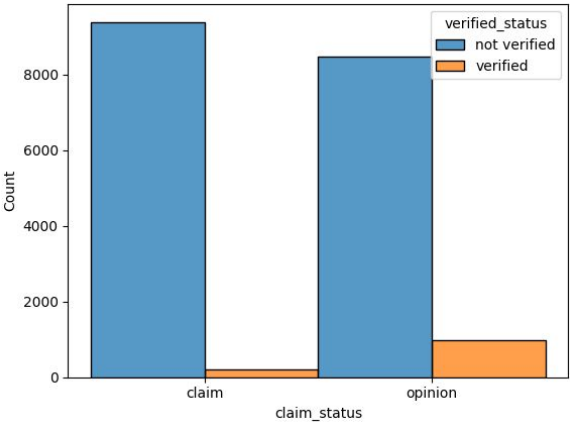
The Data Team will build and compare XGBoost and Random Forest models, evaluating which model predicts the claim status best.

Rates have p-value close to 0

```
1 for rate in rates_list:
2     p_list = []
3     for i in range(1000):
4         p_value = stats.ttest_ind(claim_data[rate].sample(30), opinion_data[rate].sample(30), equal_var=False).pvalue
5         p_list.append(p_value)
6
7     p_value = np.array(p_list).mean()
8     decision = ('definitely statistically significant' if p_value * 1000 < significance_level
9               else 'statistically significant' if p_value < significance_level else 'by chance')
10    print(f"The difference in {rate} between claims and opinions is {decision} with a p-value of {p_value}.")
```

✓ 5.7s

The difference in video_duration_sec between claims and opinions is by chance with a p-value of 0.5035878414226542.
The difference in video_view_count_log between claims and opinions is definitely statistically significant with a p-value of 8.619851150615286e-18.
The difference in video_like_count_log between claims and opinions is definitely statistically significant with a p-value of 5.3139185797998696e-15.
The difference in video_share_count_log between claims and opinions is definitely statistically significant with a p-value of 8.62872035285769e-11.
The difference in video_download_count_log between claims and opinions is definitely statistically significant with a p-value of 7.53311123124776e-11.
The difference in video_comment_count_log between claims and opinions is definitely statistically significant with a p-value of 3.92053287654959e-11.



The difference in author verified/ban status between claims and opinions is statistically significant.