

Phase 2: Exploratory Data Analysis & Engagement Insights for TikTok Claims Classification

Email from Rosie Mae Bradshaw, TikTok's Data Science Manager

Subject: Help with coding notebook?

From: "Bradshaw, Rosie Mae" —rosiemaebadshaw@tiktok

Cc: "Rainier, Orion" —orionrainier@tiktok

Good morning,

I have a couple of updates on our latest project. The leadership team has approved the project proposal that we completed previously. Thanks for all of your great work so far. Additionally, I just received an email from our Project Management Officer, Mary Joanna Rodgers that the data team is clear to proceed.

Before we begin the process of Exploratory Data Analysis (EDA), we could really use your help with coding and prepping the data. During your interview you mentioned that you worked with Python specifically in the Google certificate program you completed. That experience sounds applicable here.

Orion Rainier (Cc'd above) started a Jupyter notebook with the relevant dataset (attached). Orion is very involved in the final stages of another project. I'm sure your assistance in completing the coding and setting up the notebook for the project would be greatly appreciated.

Orion, do you mind sharing the details?

Humblest regards,

Rosie Mae Bradshaw

Data Science Manager

TikTok

Email from Orion Rainier, Data Scientist

Subject: RE: Help with coding notebook?

From: "Rainier, Orion"—orionrainier@tiktok

Cc: "Bradshaw, Rosie Mae"—rosiemaebradshaw@tiktok

Nice to meet you (virtually)!

Hope you have enjoyed your first few weeks!

With the project proposal approved, we are ready to begin the process of preparing the claim classification data. The goal of this project is to ultimately build a machine learning model that can streamline the claims process by identifying whether statements made in videos are claims or opinions.

A claim refers to information that is either unsourced or from an unverified source. For example, "The news reported that someone revealed that around 50% of the mined gold on Earth comes from one source."

Opinions refer to the personal beliefs or thoughts of a group or an individual. Here's an example, "In my opinion the most productive work day of the week is Tuesday."

There are a number of data team members committed to adjusting the machine learning developed for the last project, so your help is greatly appreciated!

Until we finish the prior project, there is no need to do a full EDA on this data. We will get to that soon. Do you mind importing the data (attached) and reviewing it for the team? It would be fantastic if you could include a summary of the column Data types, data value nonnull counts, relevant and irrelevant columns, along with anything else code related you think is worth sharing/showing in the notebook? You'll need to select a couple of variables to focus on. Include their minimum and maximum values. I haven't looked closely at the data yet, but it would be really helpful if you can create meaningful variables by combining or modifying the structures given.

Thanks,

Orion Rainier

Data Scientist

TikTok

—

"Big data isn't about bits, it's about talent." — Douglas Merrill