# 

### **GOALS AND OBJECTIVES**

- Expand and operationalize Twitter data sets
- Expand and explore new platforms for data collection
- Develop pipeline for automated topic modeling and sentiment analysis of Twitter data
- Collect and analyze YouTube comments on Minecraft updates
- Investigate additional social media platforms for data collection • Establish automated weekly analysis of topics and sentiment for efficient processing

### DATA SCIENCE

DRTH SCIENCE

During the first semester, we chose to focus on Twitter because there was already an existing Twitter scraper from last year. This means that we could build on previous work and expand on the previous pipeline.

During the second semester, we chose to focus on YouTube because it is one of the most popular platforms for Minecraft. The previous scraper collected data from specific channels, and we wanted create something more general to Minecraft as a whole for the platform.

# DATA PREPROCESSING

Data preprocessing work was done to clean up and make the data easier to interpret and work with the ML models.

The first preprocessing pass was done to translate emojis to text that could be using in NLP.

The second preprocessing pass was done to lemmatize and remove stop words from the data. The purpose was the further optimize the data for interpretation by removing low level information and normalizing text.

Lastly, work was done to translate text from other languages into English so that data from other countries could be used and interpreted.

# **TOPIC MODELING - LATENT** DIRICHLET ALLOCATION (LDA)

For topic modeling, we used Latent Dirichlet Allocation (LDA). LDA is an unsupervised clustering technique that is commonly used for text analysis. It's a type of topic modeling in which words are represented as topics, and documents are represented as a collection of these word topics.

Why LDA?

- Can find latent topic inside documents
- Supervised learning requires a true label, which may not be available
- LDA is easy to train
- LDA give interpretable topics
- > face\_with\_tears\_of\_joy

0.0.00

-> smiling\_face\_with\_smiling\_eyes -> pouting\_face

Example of emoji translation

### Leaves -> Leaf

Ponies -> Pony

Lemmatization to reduce words to their base forms

"An apple a day keeps the doctor away:" -> "apple day keeps doctor away"

"This is the best" -> "best" Removing stop words that do not add meaning

# 



- Existing listener from the previous team built with python library Tweepy
- Listener collected Minecraft related Tweets everyday
- Previously, listener lacked a
- pipeline to perform topic
- modeling and sentiment analysis



### YOUTUBE

- Existing scraper from the previous team
- Old scraper relied on YouTube API
- Old scraper limited by YouTube API quota
- Old scraper limited to only a select few channels and creators of interest
- New scraper is more "generalized" and collects data about what is being most talked about

## TWITTER TOPIC MODELING RESULTS - OCT 19TH, 2022 – OCT 26TH, 2022



between circles is how closely related each topic is. The size of each circle is proportional to the number of words in each topic.

The bars represent the frequency of terms in a topic. This is showing terms for topic 1.

### MINECRAFT YOUTUBE CHANNEL TOPIC MODELING RESULTS - MAR 21ST, 2023

With the data from the YouTube scraper, we used the same LDA model that we used for the Twitter data. The word clouds below are of the most brought up topics from the comments of the 28 most recent YouTube videos on the official Minecraft channel using our LDA model. While we started by using comments from general Minecraft videos, we found that too many creator-specific topics showed up. Instead, we decided to pivot to using videos from the official Minecraft channel as we felt that the videos stayed more on topic and we would get less creator-specific comments. However, we are still able to performing topic modeling using comments from general Minecraft videos.







Rohit Kannan, Siddharth Singh, Hridhay Monangi