Data Scraping, Data Analysis, predictions model Load the Data: Read the CSV file into a DataFrame using a library like pandas. Preprocess the Data: Clean and preprocess the data as necessary (e.g., handling missing values, encoding categorical variables, etc.). Feature Engineering: Create features that will be useful for the model. Split the Data: Divide the data into training and testing sets. Select a Model: Choose a machine learning model appropriate for your prediction task (e.g., regression, classification). Train the Model: Fit the model to the training data. Make Predictions: Use the trained model to make predictions on the test set or new data. Evaluate the Model: Assess the model's performance using appropriate metrics. 1. Data Scraping Data Collection Scrape or extract data from social media platforms. Platforms: Twitter (using Twitter API or scrapers like snscrape) Reddit (using Python's PRAW library or Pushshift API) Telegram (using libraries like Telethon or scraping tools) Key Data Points to Extract: Text content (posts, tweets, comments) Timestamps User engagement (likes, shares, comments) Stock symbols or mentions (e.g., \$AAPL, GME, etc.) Tools and Libraries: tweepy, snscrape for Twitter PRAW, psaw for Reddit Telethon for Telegram Preprocessing Data: Remove irrelevant posts (spam, off-topic). Normalize text (e.g., removing URLs, hashtags, special characters) import praw import pandas as pd import datetime # Reddit API setup reddit = praw.Reddit( client\_id="LZT12Q37jVYJc-aOMx8xpw", # Replace with your Reddit app's client ID client\_secret="tVdnlBF1pR6rxd9U5mpPmQ4dNG6P1Q", # Replace with your Reddit app's client secret user\_agent="StockScraper/1.0 by RiteshKumarYadav" # Replace with your Reddit app's user agent # Function to scrape stock market-related subreddits def scrape\_subreddit(subreddit\_name, limit=100, keywords=None): subreddit = reddit.subreddit(subreddit\_name) data = [] for post in subreddit.new(limit=limit): if keywords: # Only scrape posts that match the keywords if any(keyword.lower() in post.title.lower() for keyword in keywords): 'title': post.title, 'text': post.selftext, 'score': post.score, 'created': datetime.datetime.fromtimestamp(post.created), 'url': post.url, 'comments': post.num\_comments }) else: data.append({ 'title': post.title, 'text': post.selftext, 'score': post.score, 'created': datetime.datetime.fromtimestamp(post.created), 'url': post.url, 'comments': post.num\_comments return data # Example: Scrape from subreddits related to stock market discussions subreddits = ['stocks', 'WallStreetBets', 'investing'] keywords = ['stock', 'investment', 'market', 'trade', 'prediction'] # Example keywords for filtering  $all_data = []$ for sub in subreddits: print(f"Scraping {sub} subreddit...") scraped\_data = scrape\_subreddit(sub, limit=100, keywords=keywords) all\_data.extend(scraped\_data) # Convert to DataFrame df = pd.DataFrame(all\_data) # Save to CSV for further analysis df.to\_csv('stock\_market\_discussions.csv', index=False) # Preview the data print(df.head()) Scraping stocks subreddit... Scraping WallStreetBets subreddit... Scraping investing subreddit... 0 /r/Stocks Weekend Discussion Saturday - Dec 07... 1 Should I talk my friend out of a \$50k investme... Pros and cons of this investment strategy? These are the stocks on my watchlist (12/6) 4 r/Stocks Daily Discussion & Fundamentals Frida... text score \ O This is the weekend edition of our stickied di... 1 First, me and my friend talk every few week ab... 2 I have Roth IRA with 150k in it and like every... 3 Hi! I am an ex-prop shop equity trader.\n\nThi... 4 This is the daily discussion, so anything stoc... created 0 2024-12-07 16:00:18 https://www.reddit.com/r/stocks/comments/1h8pr... 1 2024-12-06 23:22:55 https://www.reddit.com/r/stocks/comments/1h879... 2 2024-12-06 22:28:08 https://www.reddit.com/r/stocks/comments/1h85y... 3 2024-12-06 19:46:34 https://www.reddit.com/r/stocks/comments/1h82b... 4 2024-12-06 16:00:10 https://www.reddit.com/r/stocks/comments/1h7yg... comments 0 53 2 13 20 3 311 4 In [2]: #Data Cleaning and Preprocessing import re import nltk from nltk.corpus import stopwords from nltk.tokenize import word\_tokenize from textblob import TextBlob # Download NLTK data for stopwords nltk.download('punkt') nltk.download('stopwords') # Function to clean and preprocess text def clean\_text(text): # Remove URLs text =  $re.sub(r'http\S+', '', text)$ # Remove non-alphanumeric characters text = re.sub(r'[ $^A-Za-z0-9$ ]+', '', text) # Convert to lowercase text = text.lower() return text # Function to remove stopwords def remove\_stopwords(text): stop\_words = set(stopwords.words('english')) words = word\_tokenize(text) return " ".join([word for word in words if word not in stop\_words]) # Function for sentiment analysis using TextBlob def get\_sentiment(text): analysis = TextBlob(text) # Classify polarity: Positive (> 0), Negative (< 0), Neutral (0)</pre> if analysis.sentiment.polarity > 0: return 'positive' elif analysis.sentiment.polarity < 0:</pre> return 'negative' else: return 'neutral' # Apply text cleaning to the 'title' and 'text' columns df['cleaned\_title'] = df['title'].apply(clean\_text) df['cleaned\_text'] = df['text'].apply(clean\_text) # Apply stopword removal to 'cleaned\_title' and 'cleaned\_text' df['cleaned\_title'] = df['cleaned\_title'].apply(remove\_stopwords) df['cleaned\_text'] = df['cleaned\_text'].apply(remove\_stopwords) # Perform sentiment analysis on the cleaned text df['sentiment'] = df['cleaned\_title'].apply(get\_sentiment) # Handle missing data (remove rows with missing titles or texts) df = df.dropna(subset=['cleaned\_title', 'cleaned\_text']) # Preview the cleaned data print(df.head()) [nltk\_data] Downloading package punkt to C:\Users\rites\AppData\Roaming\nltk\_data... [nltk\_data] [nltk\_data] Package punkt is already up-to-date! [nltk\_data] Downloading package stopwords to [nltk\_data] C:\Users\rites\AppData\Roaming\nltk\_data... [nltk\_data] Package stopwords is already up-to-date! 0 /r/Stocks Weekend Discussion Saturday - Dec 07... 1 Should I talk my friend out of a \$50k investme... Pros and cons of this investment strategy? These are the stocks on my watchlist (12/6) 4 r/Stocks Daily Discussion & Fundamentals Frida... text score \ O This is the weekend edition of our stickied di... 1 First, me and my friend talk every few week ab... 2 I have Roth IRA with 150k in it and like every... 3 Hi! I am an ex-prop shop equity trader. $\n\$ 4 This is the daily discussion, so anything stoc... created 0 2024-12-07 16:00:18 https://www.reddit.com/r/stocks/comments/1h8pr... 1 2024-12-06 23:22:55 https://www.reddit.com/r/stocks/comments/1h879... 2 2024-12-06 22:28:08 https://www.reddit.com/r/stocks/comments/1h85y... 3 2024-12-06 19:46:34 https://www.reddit.com/r/stocks/comments/1h82b... 4 2024-12-06 16:00:10 https://www.reddit.com/r/stocks/comments/1h7yg... cleaned\_title \ comments rstocks weekend discussion saturday dec 07 2024 53 talk friend 50k investment xrt pros cons investment strategy 13 stocks watchlist 126 311 rstocks daily discussion fundamentals friday d... cleaned\_text sentiment 0 weekend edition stickied discussion thread dis... neutral 1 first friend talk every week trades hundred fi... neutral 2 roth ira 150k like everyone want grow know mem... neutral 3 hi exprop shop equity traderthis daily watchli... neutral 4 daily discussion anything stocks related fine ... neutral In [3]: # Save the cleaned data to a new CSV file df.to\_csv('cleaned\_stock\_market\_discussions.csv', index=False) # Preview the cleaned data print(df.head()) 0 /r/Stocks Weekend Discussion Saturday - Dec 07... 1 Should I talk my friend out of a \$50k investme... Pros and cons of this investment strategy? These are the stocks on my watchlist (12/6)4 r/Stocks Daily Discussion & Fundamentals Frida... text score \ O This is the weekend edition of our stickied di... 1 First, me and my friend talk every few week ab... 2 I have Roth IRA with 150k in it and like every... 3 Hi! I am an ex-prop shop equity trader.\n\nThi... 4 This is the daily discussion, so anything stoc... created 0 2024-12-07 16:00:18 https://www.reddit.com/r/stocks/comments/1h8pr... 1 2024-12-06 23:22:55 https://www.reddit.com/r/stocks/comments/1h879... 2 2024-12-06 22:28:08 https://www.reddit.com/r/stocks/comments/1h85y... 3 2024-12-06 19:46:34 https://www.reddit.com/r/stocks/comments/1h82b... 4 2024-12-06 16:00:10 https://www.reddit.com/r/stocks/comments/1h7yg... comments cleaned\_title \ rstocks weekend discussion saturday dec 07 2024 talk friend 50k investment xrt 13 pros cons investment strategy stocks watchlist 126 311 rstocks daily discussion fundamentals friday d... cleaned\_text sentiment 0 weekend edition stickied discussion thread dis... neutral 1 first friend talk every week trades hundred fi... neutral 2 roth ira 150k like everyone want grow know mem... neutral 3 hi exprop shop equity traderthis daily watchli... neutral 4 daily discussion anything stocks related fine ... In [4]: # using matplotlib and seaborn import matplotlib.pyplot as plt import seaborn as sns # Plot the distribution of sentiment sns.countplot(x='sentiment', data=df) plt.title('Sentiment Distribution of Stock Market Posts') Sentiment Distribution of Stock Market Posts 40 35 30 25 20 15 10 5 0 positive neutral negative sentiment 2. Data Analysis Sentiment Analysis Analyze the sentiment of user-generated content to gauge market sentiment. Approaches: Pre-trained models: Hugging Face's BERT (e.g., FinBERT specialized for financial sentiment). Sentiment analysis APIs (e.g., Google NLP, AWS Comprehend). Custom model: Train a sentiment classifier using labeled financial sentiment datasets (e.g., FinancialPhraseBank). Sentiment Categories: Positive, Neutral, Negative Optionally, assign scores (e.g., [-1, 1]). Tools and Libraries: spaCy, NLTK, Transformers (Hugging Face) from textblob import TextBlob # Function to get sentiment polarity (positive, negative, neutral) def get\_sentiment(text): analysis = TextBlob(text) # Sentiment polarity: >0 is positive, <0 is negative, ==0 is neutral if analysis.sentiment.polarity > 0: return 'positive' elif analysis.sentiment.polarity < 0:</pre> return 'negative' return 'neutral' # Apply sentiment analysis on the 'cleaned\_text' column df['sentiment'] = df['cleaned\_text'].apply(get\_sentiment) # Optionally, calculate the polarity score directly df['polarity'] = df['cleaned\_text'].apply(lambda text: TextBlob(text).sentiment.polarity) # Preview the sentiment results print(df[['title', 'sentiment', 'polarity']].head()) # Save the sentiment data to CSV df.to\_csv('sentiment\_analyzed\_stock\_discussions.csv', index=False) title sentiment polarity 0 /r/Stocks Weekend Discussion Saturday - Dec 07... negative -0.031250 1 Should I talk my friend out of a \$50k investme... positive 0.200962 Pros and cons of this investment strategy? negative -0.0285713 These are the stocks on my watchlist (12/6) positive 0.1450524 r/Stocks Daily Discussion & Fundamentals Frida... positive 0.018981 Topic Modeling Using LDA In [6]: import nltk from nltk.corpus import stopwords from nltk.tokenize import word\_tokenize from sklearn.feature\_extraction.text import CountVectorizer from sklearn.decomposition import LatentDirichletAllocation # Download NLTK data (tokenizer, stopwords) nltk.download('punkt') nltk.download('stopwords') # Function to preprocess the text (tokenization and stopword removal) def preprocess\_text(text): stop\_words = set(stopwords.words('english')) words = word\_tokenize(text.lower()) return " ".join([word for word in words if word.isalnum() and word not in stop\_words]) # Apply preprocessing to the 'cleaned\_text' column df['processed\_text'] = df['cleaned\_text'].apply(preprocess\_text) # Vectorize the processed text using CountVectorizer vectorizer = CountVectorizer(max\_features=1000) # You can adjust the number of features X = vectorizer.fit\_transform(df['processed\_text']) # Perform LDA (Topic Modeling) lda = LatentDirichletAllocation(n\_components=5, random\_state=42) # 5 topics, you can adjust lda.fit(X) # Print the top words for each topic  $num\_words = 10$ for index, topic in enumerate(lda.components\_): print(f"Topic {index+1}:") print([vectorizer.get\_feature\_names\_out()[i] for i in topic.argsort()[-num\_words:]]) print("\n") [nltk\_data] Downloading package punkt to [nltk\_data] C:\Users\rites\AppData\Roaming\nltk\_data... [nltk\_data] Package punkt is already up-to-date! [nltk\_data] Downloading package stopwords to [nltk\_data] C:\Users\rites\AppData\Roaming\nltk\_data... [nltk\_data] Package stopwords is already up-to-date! ['dont', 'meme', 'shares', 'know', 'ive', 'fund', 'investment', 'im', 'price', 'would'] Topic 2: ['money', 'think', 'ago', 'good', 'much', 'like', 'stocks', 'years', 'stock', 'market'] Topic 3: ['roi', 'list', 'companies', 'price', 'company', 'value', 'intrinsic', 'like', 'market', 'stock'] Topic 4: ['im', 'trading', 'question', 'daily', 'fundamentals', 'earnings', 'stock', 'market', 'news', 'stocks'] Topic 5: ['financial', 'production', 'market', 'optimus', 'competitors', 'level', 'like', 'autonomy', 'tesla', 'teslas'] Frequency of Mentions (Keyword Frequency) In [7]: from collections import Counter # Define stock-related keywords keywords = ['stock', 'investment', 'market', 'trade', 'prediction', 'bull', 'bear', 'portfolio'] # Function to count keyword frequency in each post def count\_keywords(text, keywords): words = text.lower().split() keyword\_counts = {keyword: words.count(keyword) for keyword in keywords} return keyword\_counts # Apply the keyword count function to each post df['keyword\_counts'] = df['cleaned\_text'].apply(count\_keywords, keywords=keywords) # Convert keyword counts into separate columns for analysis keyword\_df = pd.DataFrame(df['keyword\_counts'].tolist()) # Merge the keyword counts with the original DataFrame df = pd.concat([df, keyword\_df], axis=1) # Preview the DataFrame with keyword frequencies print(df[['title', 'keyword\_counts']].head()) # You can also calculate the total frequency of mentions for each keyword across all posts total\_keyword\_counts = keyword\_df.sum() print("Total keyword frequencies across all posts:") print(total\_keyword\_counts) title \ 0 /r/Stocks Weekend Discussion Saturday - Dec 07... 1 Should I talk my friend out of a \$50k investme... Pros and cons of this investment strategy? These are the stocks on my watchlist (12/6) 4 r/Stocks Daily Discussion & Fundamentals Frida... keyword\_counts 0 {'stock': 0, 'investment': 0, 'market': 3, 'tr... 1 {'stock': 0, 'investment': 0, 'market': 0, 'tr... 2 {'stock': 1, 'investment': 4, 'market': 0, 'tr... 3 {'stock': 2, 'investment': 0, 'market': 1, 'tr... 4 {'stock': 1, 'investment': 0, 'market': 3, 'tr... Total keyword frequencies across all posts: stock 54 20 investment 75 market trade 11 prediction 0 bull 0 0 portfolio dtype: int64 Visualization In [8]: #Sentiment Distribution import seaborn as sns import matplotlib.pyplot as plt # Plot the sentiment distribution sns.countplot(x='sentiment', data=df) plt.title('Sentiment Distribution of Stock Market Posts') plt.show() Sentiment Distribution of Stock Market Posts 40 30 10 negative positive neutral sentiment In [11]: #Keyword Frequency Distribution # Plot the total frequency of keywords total\_keyword\_counts.sort\_values(ascending=False).plot(kind='bar', figsize=(10, 6)) plt.title('Frequency of Stock Market Keywords') plt.ylabel('Frequency') plt.show() Frequency of Stock Market Keywords 70 60 50 Frequency 30 20 10 stock IInq investment portfolio prediction Saving the Final Processed Data df.to\_csv(r'C:\Users\rites\OneDrive\Desktop\final\_processed\_stock\_discussions.csv', index=False) 3. Prediction Model Prepare Data for the Model Load and Inspect Data: Load the CSV and inspect the available columns to understand the data. Feature Engineering: Create relevant features from the data (e.g., sentiment analysis, keyword counts, etc.). Model Training: Train a machine learning model to predict stock movements (e.g., up or down). Model Evaluation: Evaluate the model's performance using appropriate metrics. Stock Movement Prediction i am using RandomForestClassifier, which is a good choice for many classification tasks. However, you might want to experiment with other models or tune hyperparameters for better performance. import pandas as pd from sklearn.model\_selection import train\_test\_split from sklearn.ensemble import RandomForestClassifier from sklearn.metrics import accuracy\_score, classification\_report # Load the Data file\_path = r"C:\Users\rites\OneDrive\Desktop\final\_processed\_stock\_discussions.csv" data = pd.read\_csv(file\_path) # Check the column names print("Column names:", data.columns.tolist()) # Check the first few rows of the DataFrame print(data.head()) # Step 2: Preprocess the Data # Handle missing values data.fillna(method='ffill', inplace=True) # Encode categorical variables if necessary data = pd.get\_dummies(data) Column names: ['title', 'text', 'score', 'created', 'url', 'comments', 'cleaned\_title', 'cleaned\_text', 'sentiment', 'polarity', 'processed\_text', 'keyword\_counts', 'stock', 'inves tment', 'market', 'trade', 'prediction', 'bull', 'bear', 'portfolio'] title 0 /r/Stocks Weekend Discussion Saturday - Dec 07... Should I talk my friend out of a \$50k investme... Pros and cons of this investment strategy? These are the stocks on my watchlist (12/6) 4 r/Stocks Daily Discussion & Fundamentals Frida... text score \ O This is the weekend edition of our stickied di... 1 First, me and my friend talk every few week ab... 2 I have Roth IRA with 150k in it and like every... 3 Hi! I am an ex-prop shop equity trader.\n\nThi... 4 This is the daily discussion, so anything stoc... created 0 2024-12-07 16:00:18 https://www.reddit.com/r/stocks/comments/1h8pr... 1 2024-12-06 23:22:55 https://www.reddit.com/r/stocks/comments/1h879... 2 2024-12-06 22:28:08 https://www.reddit.com/r/stocks/comments/1h85y... 3 2024-12-06 19:46:34 https://www.reddit.com/r/stocks/comments/1h82b... 4 2024-12-06 16:00:10 https://www.reddit.com/r/stocks/comments/1h7yg... cleaned\_title \ comments rstocks weekend discussion saturday dec 07 2024 53 talk friend 50k investment xrt 2 13 pros cons investment strategy 3 stocks watchlist 126 4 311 rstocks daily discussion fundamentals friday d... cleaned\_text sentiment polarity \ 0 weekend edition stickied discussion thread dis... negative -0.031250first friend talk every week trades hundred fi... positive 0.200962 roth ira 150k like everyone want grow know mem... negative -0.028571 3 hi exprop shop equity traderthis daily watchli... positive 0.145052 4 daily discussion anything stocks related fine ... positive 0.018981 processed\_text \ 0 weekend edition stickied discussion thread dis... 1 first friend talk every week trades hundred fi... 2 roth ira 150k like everyone want grow know mem... 3 hi exprop shop equity traderthis daily watchli... 4 daily discussion anything stocks related fine ... keyword\_counts stock investment \ 0 {'stock': 0, 'investment': 0, 'market': 3, 'tr... 0 0 1 {'stock': 0, 'investment': 0, 'market': 0, 'tr... 0 {'stock': 1, 'investment': 4, 'market': 0, 'tr... 4 3 {'stock': 2, 'investment': 0, 'market': 1, 'tr... 0 4 {'stock': 1, 'investment': 0, 'market': 3, 'tr... 0 market trade prediction bull bear portfolio 0 2 0 0 0 3 1 0 0 0 Feature Importance: After training the model, you might want to check which features are most important for predictions. This can help you understand your model better and refine your feature set. Hyperparameter Tuning: Consider using techniques like Grid Search or Random Search to find the best hyperparameters for your model. Cross-Validation: Instead of a single train-test split, you might want to use cross-validation to get a better estimate of your model's performance.