<pre>om sklearn.feature_extraction.text import TfidfVectorizer port folium om folium.plugins import HeatMap port plotly.express as px om geopy.geocoders import Nominatim port matplotlib.pyplot as plt port seaborn as sns port warnings rnings.filterwarnings('ignore') Download necessary NLTK data tk.download('stopwords') tk.download('punkt')</pre>			
tk_data] Downloading package stopwords to tk_data] C:\Users\sriva\AppData\Roaming\nltk_data tk_data] Package stopwords is already up-to-date! tk_data] Downloading package punkt to tk_data] C:\Users\sriva\AppData\Roaming\nltk_data tk_data] Package punkt is already up-to-date! Task 1: Social Media Data Extraction & Preprocessing this section, we'll: Define crisis-related keywords for filtering relevant posts Configure Reddit API access Extract mental health-related posts from multiple subreddits Clean and preprocess the text data for further analysis			
<pre>isis_keywords = ["depressed", "depression", "anxiety", "suicidal", "suicide", "addiction", "overwhelmed", "hopeless", "self-harm", "cutting", "lonely", "mental health", "help me", "can't cope", "end my life" Reddit API Configuration f setup_reddit_api(): reddit = praw.Reddit(client_id="rUrpqf9QSk8Pt110AcUaNg", client_secret="p3he-liblvaBqrFGM3swNnEsI7yp8A", user_agent="mental_health_data_collection_script_v1.0") return reddit Extract posts from Reddit based on keywords f extract_reddit_posts(reddit, keywords, limit=1000):</pre>			
<pre>posts_data = [] # Search across multiple subreddits focused on mental health subreddits = ["depression", "selfharm", "anxiety", "mentalhealth", "MMFB"] for subreddit_name in subreddits: subreddit = reddit.subreddit(subreddit_name) for submission in subreddit.hot(limit=limit//len(subreddits)): # Check if any keyword is in the post if any(keyword.lower() in submission.title.lower() or</pre>			
<pre>'content': submission.selftext,</pre>			
<pre># Remove special characters and digits text = re.sub(r'[^\w\s]', '', text) text = re.sub(r'\d+', '', text) # Convert to lowercase text = text.lower() # Remove stop words stop_words = set(stopwords.words('english')) word_tokens = word_tokenize(text) filtered_text = [w for w in word_tokens if not w in stop_words] # Remove emojis text = re.sub(r':[a-zA-Z_]+:', '', ' '.join(filtered_text)) return text</pre>			
<pre>ddit = setup_reddit_api() sts_df = extract_reddit_posts(reddit, crisis_keywords) Apply text cleaning sts_df['cleaned_content'] = posts_df['content'].apply(clean_text) Display the preprocessed data int(f"Total posts collected: {len(posts_df)}") sts_df.head() Save to CSV sts_df.to_csv('crisis_posts.csv', index=False) al posts collected: 480 sts_df.head()</pre>			
1frqlk0 2024-09-29 04:35:43 Regular check-in post, with information about Welcome to 1j8hemq 2025-03-11 08:44:33 I died and was revived. I'm angry I'm alive Every day I 1j88rqx 2025-03-11 02:09:14 Some people don't understand animals and their Of course the	pes not apply to everyone but depression 88 17 Beautifu	author cleaned_content SQLwitch understand people reply immediately op invitat SQLwitch welcome rdepressions checkin post place take m every day see sun get angry hate things living course apply everyone male cat quite literally Xbonespooky im years old year literally nothing ongoing ad	
<pre># Initialize TF-IDF vectorizer tfidf_vectorizer = TfidfVectorizer(max_features=1000, stop_words='english', ngram_range=(1, 2)) # Fit and transform the corpus tfidf_matrix = tfidf_vectorizer.fit_transform(corpus) # Get feature names (terms) feature_names = tfidf_vectorizer.get_feature_names_out() # Create a dictionary to store the average TF-IDF score for each term term_scores = {} # Calculate the average TF-IDF score for each term</pre>			
<pre># Calculate the average TF-IDF score for each term for i in range(len(feature_names)): term_scores[feature_names[i]] = tfidf_matrix[:, i].mean() # Sort terms by their average TF-IDF score sorted_terms = sorted(term_scores.items(), key=lambda x: x[1], reverse=True # Return the top 50 terms return [term for term, score in sorted_terms[:50]] Risk level classification using both keyword approach and TF-IDF insights f classify_risk(text, high_risk_terms=None, moderate_risk_terms=None, tfidf_ti if high_risk_terms is None: # Base high-risk keywords high_risk_terms = ["kill myself", "end my life", "suicide", "suicidal", "don't want to live", "want to die", "better off dead"]</pre>	=None):		
<pre>if moderate_risk_terms is None: # Base moderate concern keywords moderate_risk_terms = ["depressed", "depression", "anxiety", "help me", "can't cope", "hopeless", "lost", "overwhelmed"] # Use TF-IDF terms to enhance keyword lists if available if tfidf_terms is not None: # Add high-risk TF-IDF terms high_risk_indicators = ["suicide", "kill", "die", "end", "pain"] for term in tfidf_terms: if any(indicator in term for indicator in high_risk_indicators): high_risk_terms.append(term) # Add moderate-risk TF-IDF terms moderate_risk_indicators = ["depress", "anxi", "help", "cope", "overwhe</pre>			
<pre>moderate_risk_indicators = ["depress", "anxi", "help", "cope", "overwhe for term in tfidf_terms: if any(indicator in term for indicator in moderate_risk_indicators) moderate_risk_terms.append(term) # Remove duplicates high_risk_terms = list(set(high_risk_terms)) moderate_risk_terms = list(set(moderate_risk_terms)) if any(term in text.lower() for term in high_risk_terms): return "High Risk" elif any(term in text.lower() for term in moderate_risk_terms): return "Moderate Concern" else: return "Low Concern" Extract crisis terms using TF-IDF int("Extracting high-risk terms using TF-IDF analysis")</pre>			
<pre>idf_crisis_terms = extract_crisis_terms(posts_df['cleaned_content'].fillna('' int(f"Top crisis terms identified via TF-IDF: {', '.join(tfidf_crisis_terms[: Apply sentiment analysis sts_df['sentiment'] = posts_df['content'].apply(analyze_sentiment) Apply enhanced risk classification with TF-IDF insights sts_df['risk_level'] = posts_df['content'].apply(lambda x: classify_risk(x, t Display risk term statistics int("\nRisk Classification Results:") int(posts_df['risk_level'].value_counts()) Visualize sentiment distribution t.figure(figsize=(10, 6)) ntiment_counts = posts_df['sentiment'].value_counts()</pre>			
<pre>s.barplot(x=sentiment_counts.index, y=sentiment_counts.values) t.title('Sentiment Distribution of Crisis-Related Posts') t.xlabel('Sentiment') t.ylabel('Count') t.show() Visualize risk level distribution t.figure(figsize=(10, 6)) sk_counts = posts_df['risk_level'].value_counts() s.barplot(x=risk_counts.index, y=risk_counts.values, palette='YlOrRd') t.title('Risk_Level Distribution of Crisis-Related Posts') t.xlabel('Risk_Level') t.ylabel('Count') t.show() Cross-tabulation of sentiment and risk level oss_tab = pd.crosstab(posts_df['sentiment'], posts_df['risk_level'])</pre>			
Risk Level Distribution of Crisis-Related Post Risk Level Distribution of Crisis-Related Post 200	169 - 140 - 120 - 100 - 100 - 80 - 60 - 40 - 87 - 20 - 20 - Moderate Concern Moderate Concern		
<pre>Geocoding function f geocode_location(location_name, cache={}): if not location_name: return None # Check if we've already geocoded this location if location_name in cache: return cache[location_name] try: geolocator = Nominatim(user_agent="crisis_mapping") location = geolocator.geocode(location_name, timeout=10) if location: # Store in cache for future use cache[location_name] = (location.latitude, location.longitude)</pre>			
<pre>return cache[location_name] # Try adding "city" to improve geocoding results location = geolocator.geocode(f"{location_name} city", timeout=10) if location: cache[location_name] = (location.latitude, location.longitude) return cache[location_name] return None except Exception as e: print(f"Geocoding error for {location_name}: {e}") return None Extract and geocode locations from posts sts_df['location_mentioned'] = posts_df['content'].apply(extract_location) sts_df['coordinates'] = posts_df['location_mentioned'].apply(lambda x: geocode_location(x) if x else None</pre>			
Create a clean dataframe for mapping p_df = posts_df.dropna(subset=['coordinates']).copy() Count occurrences of each location cation_counts = map_df['location_mentioned'].value_counts().reset_index() cation_counts.columns = ['location', 'count'] p_locations = location_counts.head(5) int("Top 5 locations with highest crisis discussions:") int(top_locations) Create a Folium map centered on the US = folium.Map(location=[39.8283, -98.5795], zoom_start=4) Add points to the map r idx, row in map_df.iterrows(): lat, lon = row['coordinates']			
<pre>lat, lon = row['coordinates'] risk_color = { 'High Risk': 'red', 'Moderate Concern': 'orange', 'Low Concern': 'green' }.get(row['risk_level'], 'blue') popup_text = f""" Risk Level: {row['risk_level']} < br > Sentiment: {row['sentiment']} < br > Location: {row['location_mentioned']} < br > """ folium.CircleMarker(</pre>			
<pre>location=[lat, lon], radius=5, color=risk_color, fill=True, fill_opacity=0.7,</pre>			

1. Data Collection: Extracting relevant mental health posts from Reddit using targeted keywords

2. Text Processing: Cleaning and preparing text data for analysis

3. Sentiment Analysis: Using VADER to determine emotional tone of posts4. Risk Classification: Categorizing posts into risk levels for prioritization

5. Geographic Analysis: Extracting location information and visualizing crisis hotspots

Crisis Detection and Mental Health Risk Analysis System

intervention and response planning.

Table of Contents

1. Setup and Dependencies

1. Setup and Dependencies

4. Task 3: Crisis Geolocation & Mapping

Task 1: Social Media Data Extraction & Preprocessing
 Task 2: Sentiment Analysis & Crisis Risk Classification

This notebook implements a comprehensive system to analyze mental health crisis signals from social media data. The solution covers data collection from Reddit API, text processing, sentiment analysis, risk classification, and geospatial visualization - making it suitable for early