```
1 mkdir -p raw_data
```

```
1 cd raw_data
```

```
1 !wget -q https://clinicaltrials.gov/AllPublicXML.zip
```

```
1 !unzip -q AllPublicXML.zip
```

```python
1  import csv
2  import xml.etree.ElementTree as ET
3
4
5  def pull_tag_content(root_node, tag_path):
6      content_list = []
7      stack = [(root_node, '')]
8
9      while stack:
10         node, prefix = stack.pop()
11         path = prefix + '/' + node.tag
12
13         if path == tag_path and node.text:
14             content_list.append(node.text)
15
16         for child in node:
17             stack.append((child, path))
18
19     return content_list
20
21 def extract_criteria(criteria_textblock):
22     # Find the positions of "Inclusion Criteria:" and "Exclusion Criteria:"
23     c1 = criteria_textblock.find("Inclusion Criteria:",0)
24     c2 = criteria_textblock.find("Exclusion Criteria:",0)
25
26     # Extract the Inclusion Criteria and handle missing criteria
27     if c1 >= 0:
28         if c2 >= 0:
29             inclusion_criteria = criteria_textblock[c1 + len("Inclusion Criteria:"):c2].strip()
30         else:
31             inclusion_criteria = criteria_textblock[c1 + len("Inclusion Criteria:"):].strip()
32     else:
33         inclusion_criteria = ""
34
35     # Extract the Exclusion Criteria and handle missing criteria
36     if c2 >= 0:
37         if c1 >= 0:
38             exclusion_criteria = criteria_textblock[c2 + len("Exclusion Criteria:"):].strip()
39         else:
40             exclusion_criteria = criteria_textblock[c2 + len("Exclusion Criteria:"):].strip()
41     else:
42         exclusion_criteria = ""
43     return inclusion_criteria, exclusion_criteria
44
45 def read_xml_file(file_path):
46     try:
47         # Parse the XML file
48         tree = ET.parse(file_path)
49         root = tree.getroot()
50
51         clinical_data = {
52             'nct_id': root.findtext('id_info/nct_id',''),
```

```python
            nct_id : root.findtext('id_info/nct_id',''),
            'brief_title': root.findtext('brief_title',''),
            'official_title': root.findtext('official_title',''),
            'agency': root.findtext('sponsors/lead_sponsor/agency',''),
            'agency_class': root.findtext('sponsors/lead_sponsor/agency_class',''),
            'collaborator_agency': root.findtext('sponsors/collaborator/agency',''),
            'brief_summary': root.findtext('brief_summary/textblock',''),
            'detailed_description': root.findtext('detailed_description/textblock',''),
            # 'conditions': root.findtext('condition',''),
            'overall_status': root.findtext('overall_status',''),
            'phase': root.findtext('phase',''),
            'study_type': root.findtext('study_type',''),
            'has_expanded_access': root.findtext('has_expanded_access',''),
            'intervention': root.findtext('intervention',''),
            'intervention_type': root.findtext('intervention/intervention_type',''),
            'intervention_name': root.findtext('intervention/intervention_name',''),
            'lead_sponsor_agency': root.find('sponsors/lead_sponsor/agency',''),
            'primary_completion_date': root.findtext('primary_completion_date',''),
            'start_date': root.findtext('start_date',''),
            'completion_date': root.findtext('completion_date',''),
            'gender': root.findtext('eligibility/gender',''),
            'minimum_age': root.findtext('eligibility/minimum_age',''),
            'maximum_age': root.findtext('eligibility/maximum_age',''),
            'healthy_volunteers': root.findtext('eligibility/healthy_volunteers',''),
            'why_stopped': root.findtext('why_stopped',''),
        }
        # check for multipal marks
        conditions = pull_tag_content(root, '/clinical_study/condition')
        keywords = pull_tag_content(root, '/clinical_study/keyword')
        clinical_data['conditions'] = conditions
        clinical_data['keywords'] = keywords
        # Extract Inclusion and Exclusion Criteria
        criteria_textblock = root.findtext('eligibility/criteria/textblock','')
        inclusion_criteria, exclusion_criteria = extract_criteria(criteria_textblock)
        clinical_data['inclusion_criteria'] = inclusion_criteria
        clinical_data['exclusion_criteria'] = exclusion_criteria
        return clinical_data

    except FileNotFoundError:
        print(f"Error: File '{file_path}' not found.")
        return None
    except ET.ParseError:
        print(f"Error: Invalid XML format in '{file_path}'.")
        return None


def save_to_csv(data_list, csv_file):
    with open(csv_file, 'w', newline='', encoding='utf-8') as f:
        fieldnames = data_list[0].keys()
        writer = csv.DictWriter(f, fieldnames=fieldnames)
        writer.writeheader()
        writer.writerows(data_list)
```

```python
raw_data_dir = '/content/raw_data'
import os
# # Create an empty list to store all the file paths
all_file_paths = []

# # Use os.walk to traverse through all subdirectories in 'raw_data'
for dirpath, _, filenames in os.walk(raw_data_dir):
    # Concatenate the directory path with the filenames to get the full file paths
    file_paths = [os.path.join(dirpath, filename) for filename in filenames]
    # Extend the all_file_paths list with the file_paths list for each subdirectory
```

```
11    all_file_paths.extend(file_paths)
12
```

```
1 data_list = []
2 for file_path in all_file_paths:
3     if file_path.endswith(".xml"):
4         clinical_data = read_xml_file(file_path)
5         if clinical_data is not None:
6             data_list.append(clinical_data)
7
8
```

```
1 mkdir -p data_output
```

```
1
2 csv_file= '/content/data_output/combine_output.csv'
3 if data_list:
4     save_to_csv(data_list, csv_file)
5     print(f"Data from {len(data_list)} XML files saved to '{csv_file}' successfully.")
6 else:
7     print("No valid data found in XML files.")
```

```
1 from google.colab import drive
2 drive.mount('/content/drive')
```
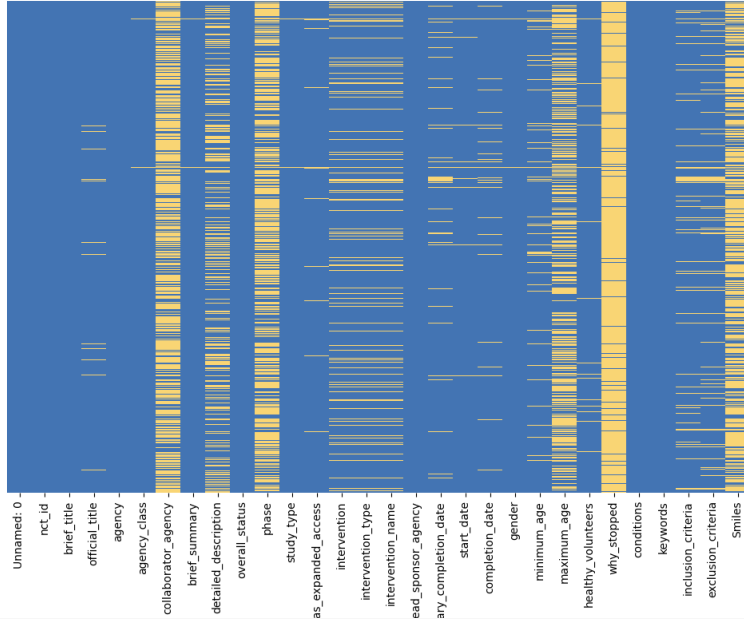
    Mounted at /content/drive

```
1 import pandas as pd
2 df=pd.read_csv('/content/drive/MyDrive/Mtech /Dissertation/data_output/combine_output_with_smiles.csv')
```

```
1 import pandas as pd
2 import seaborn as sns
3 import matplotlib.pyplot as plt
```

```
1 # Assuming you have loaded your data into a pandas DataFrame called 'df'
2 # If not, load your data into 'df' using pandas read_csv or any other method.
3
4 # Create the null value DataFrame (True for null, False for non-null)
5 null_df = df.isnull()
6
7 # Set a custom color palette for the heatmap (diverging color map with blue and yellow)
8 colors = ['#4374B3', '#F9D574']
9
10 # Create the heatmap using seaborn
11 plt.figure(figsize=(12, 8))
12 sns.heatmap(null_df, cmap=sns.color_palette(colors), cbar=False, yticklabels=False)
13
14 # Add a title and labels to the heatmap
15 plt.title('Null Values Heatmap', fontsize=20)
16 plt.xlabel('Columns', fontsize=14)
17
18 # Show the plot
19 plt.show()
20
```
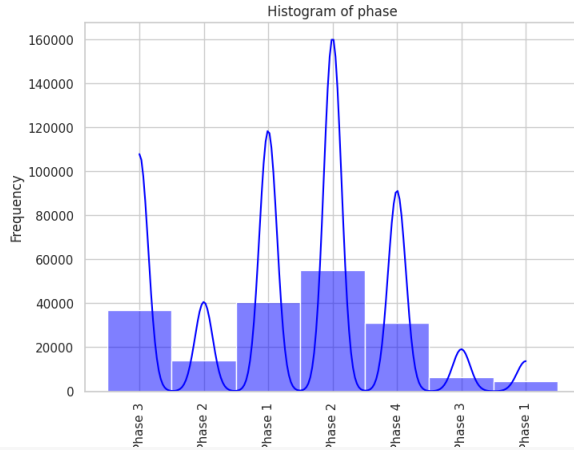
Null Values Heatmap

Column labels (left to right): Unnamed: 0, nct_id, brief_title, official_title, agency, agency_class, collaborator_agency, brief_summary, detailed_description, overall_status, phase, study_type, as_expanded_access, intervention, intervention_type, intervention_name, ead_sponsor_agency, ary_completion_date, start_date, completion_date, gender, minimum_age, maximum_age, healthy_volunteers, why_stopped, conditions, keywords, inclusion_criteria, exclusion_criteria, Smiles

```
1  sns.set(style="whitegrid")
2
3  # Data Visualization
4  # Histograms for Numerical Columns
5  numerical_cols = ['phase']#, 'minimum_age', 'maximum_age']
6  for col in numerical_cols:
7      plt.figure(figsize=(8, 6))
8      sns.histplot(df[col].dropna(), kde=True, color='blue', bins=30)
9      plt.title(f'Histogram of {col}')
10     plt.xlabel(col)
11     plt.ylabel('Frequency')
12     plt.xticks(rotation=90)
13     plt.show()
```
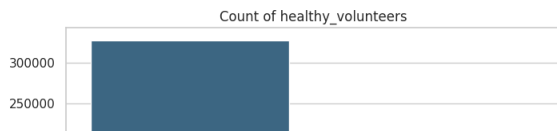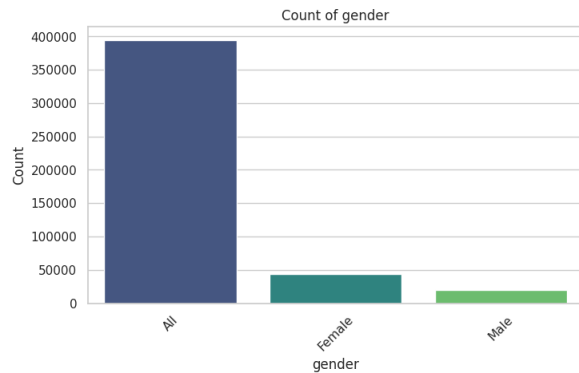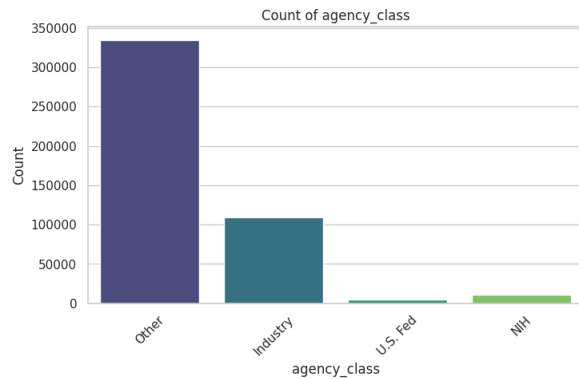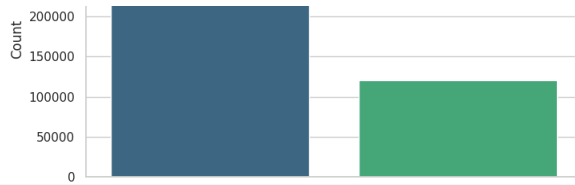
Histogram of phase

```
1 categorical_cols = ['agency_class', 'gender', 'healthy_volunteers']
2
3 num_plots = len(categorical_cols)
4 num_rows = num_plots
5 num_cols = 1
6
7 fig, axes = plt.subplots(num_rows, num_cols, figsize=(8, 6*num_rows))
8
9 # Plot settings
10 sns.set_palette("viridis")
11
12 for i, col in enumerate(categorical_cols):
13     sns.countplot(data=df, x=col, palette='viridis', ax=axes[i])
14     axes[i].set_title(f'Count of {col}')
15     axes[i].set_xlabel(col)
16     axes[i].set_ylabel('Count')
17     axes[i].tick_params(axis='x', rotation=45)
18
19 plt.subplots_adjust(hspace=0.5)  # Adjust the vertical spacing between subplots
20
21 plt.show()
```

Count of agency_class

Count of gender

Count of healthy_volunteers

```
1  categorical_cols = ['overall_status', 'study_type']
2
3  num_plots = len(categorical_cols)
4  num_rows = 1
5  num_cols = num_plots
6
7  fig, axes = plt.subplots(num_rows, num_cols, figsize=(12,6))
8
9  # Plot settings
10 sns.set_palette("viridis")
11
12 for i, col in enumerate(categorical_cols):
13     sns.countplot(data=df, x=col, palette='viridis', ax=axes[i])
14     axes[i].set_title(f'Count of {col}')
15     axes[i].set_xlabel(col)
16     axes[i].set_ylabel('Count')
17     axes[i].tick_params(axis='x', rotation=90)
18
19 plt.tight_layout()  # Automatically adjust spacing between subplots
20
21 plt.show()
```
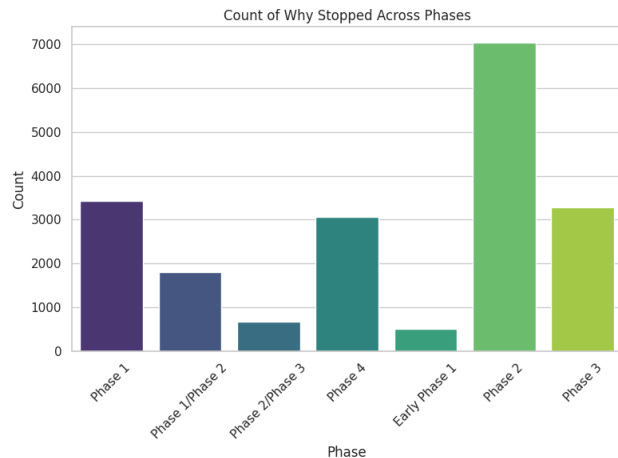
```
1 plt.figure(figsize=(8, 6))
2 sns.countplot(data=df[df['why_stopped'].notnull()], x='phase', palette='viridis')
3 plt.title('Count of Why Stopped Across Phases')
4 plt.xlabel('Phase')
5 plt.ylabel('Count')
6 plt.xticks(rotation=45)
7 plt.tight_layout()
8 plt.show()
```



```
1 # Text Data Analysis (if relevant, for text columns)
2 # Example: Word Cloud for 'brief_summary'
3 # from wordcloud import WordCloud
4
5 # text_column = 'why_stopped'
6 # text_data = " ".join(text for text in df[text_column].dropna())
7 # wordcloud = WordCloud(width=800, height=400, background_color='white').generate(text_data)
8 # plt.figure(figsize=(10, 6))
9 # plt.imshow(wordcloud, interpolation='bilinear')
10 # plt.title(f'Word Cloud for {text_column}')
11 # plt.axis('off')
12 # plt.show()
13
```

```
1
```

✓ 1s    completed at 4:45 PM