

Ans to question 1

Task1:

```
Q1_Task_1.py > ...
1  import pandas as pd
2
3  # Define the file paths
4  csv_files = [
5      r'CSV1.csv',
6      r'CSV2.csv',
7      r'CSV3.csv',
8      r'CSV4.csv'
9  ]
10
11  columns_to_check = ['SHORT-TEXT', 'TEXT']
12  output_txt_file = 'Q1_extracted_texts.txt'
13
14  def extract_and_save_text(df, column_name, output_file, mode='a'):
15      """Extracts text from the specified column and saves it to a text file."""
16      if column_name in df.columns:
17          texts = df[column_name].dropna().tolist() # Drop NaN values
18          with open(output_file, mode, encoding='utf-8') as outfile:
19              for text in texts:
20                  outfile.write(text + '\n')
21          return True
22      return False
23
24  # Process each CSV file and extract text from the appropriate columns
25  with open(output_txt_file, 'w', encoding='utf-8') as outfile:
26      outfile.write("") # Clear the file at the beginning
27
28  for i, csv_file in enumerate(csv_files):
29      try:
30          # Read the CSV file into a DataFrame
31          df = pd.read_csv(csv_file)
32
33          # Check for both possible column names and save text if the column exists
34          saved = False
35          for column in columns_to_check:
36              if extract_and_save_text(df, column, output_txt_file, 'a'):
37                  print(f"Extracted text from '{column}' in {csv_file}.")
38                  saved = True
39                  break # Stop checking other columns if one is found
40
41          if not saved:
42              print(f"None of the specified columns found in {csv_file}.")
43
44      except FileNotFoundError:
45          print(f"Error: The file '{csv_file}' does not exist.")
46      except pd.errors.EmptyDataError:
47          print(f"Error: The file '{csv_file}' is empty.")
48      except Exception as e:
49          print(f"An unexpected error occurred while processing '{csv_file}': {str(e)}")
```

Task2: Research

```
Q1_Task_2.py > ...
1  #imports
2  import spacy
3  from transformers import AutoTokenizer, AutoModelForTokenClassification, pipeline
4
5  #Installation Command:
6  ##Install the libraries(SpaCy - 'en_core_sci_sm').
7  ###pip install spacy==2.3.5
8  ###python -m spacy download en_core_sci_sm
9  ##Transformer Install Command
10  ###pip install transformers
11
12  #Sample Usage in my code
13
14  ## spacy
15  def load_spacy_model(model_path):
16      try:
17          return spacy.load(model_path)
18      except OSError as e:
19          print(f"Error: Failed to load the SpaCy model '{model_path}'. {str(e)}")
20          return None
21
22      nlp_sci_sm = load_spacy_model("en_core_sci_sm")
23
24
25  ## transformers
26  def load_biobert_pipeline(model_name):
27      try:
28          tokenizer = AutoTokenizer.from_pretrained(model_name)
29          model = AutoModelForTokenClassification.from_pretrained(model_name)
30          return pipeline("ner", model=model, tokenizer=tokenizer)
31      except Exception as e:
32          print(f"Error: Failed to load BioBERT model '{model_name}'. {str(e)}")
33          return None
34
35  biobert_pipeline = load_biobert_pipeline("dmis-lab/biobert-base-cased-v1.1")
36  tokenizer = AutoTokenizer.from_pretrained("bert-base-uncased", use_fast=True)
37
```

Task3: Programming and Research

```
Q1_Task_3.py > ...
1  import pandas as pd
2  from collections import Counter
3  from transformers import AutoTokenizer
4  import warnings
5
6  # Suppress warnings
7  warnings.filterwarnings('ignore')
8
9  output_txt_file = 'Q1_extracted_texts.txt'
10 top_words_csv = 'Q1_top_30_words.csv'
11
12 # Task 3.1: Count the Top 30 Most Common Words and Store in a CSV File
13 def count_top_words(text_file, output_csv):
14     try:
15         with open(text_file, 'r', encoding='utf-8') as file:
16             text = file.read()
17
18             if not text.strip(): # Check if the file is empty
19                 print(f"Error: The file '{text_file}' is empty.")
20                 return
21
22             # Split the text into words and count occurrences
23             words = text.split()
24             word_counts = Counter(words)
25
26             # Get the 30 most common words
27             top_30_words = word_counts.most_common(30)
28
29             # Store the result in a CSV file
30             df_top_words = pd.DataFrame(top_30_words, columns=['Word', 'Count'])
31             df_top_words.to_csv(output_csv, index=False)
32
33             print(f"\nTask 3.1: Top 30 words saved to {output_csv}\n")
34     except FileNotFoundError:
35         print(f"Error: The file '{text_file}' was not found.")
36     except Exception as e:
37         print(f"An error occurred while counting words: {str(e)}")
38
39 # Task 3.2: Count Unique Tokens Using Transformers AutoTokenizer (with Chunking and Truncation)
40 def count_unique_tokens(text_file, chunk_size=5000, max_length=500):
41     try:
42         tokenizer = AutoTokenizer.from_pretrained("bert-base-uncased", use_fast=True)
43
44         token_counts = Counter()
45
46         with open(text_file, 'r', encoding='utf-8') as file:
47             while True:
```

```
48     # Read the file in chunks
49     chunk = file.read(chunk_size)
50     if not chunk:
51         break # End of file
52
53     # Tokenize the chunk with truncation to handle long sequences
54     tokens = tokenizer.tokenize(chunk, truncation=True, max_length=max_length)
55
56     # Count token occurrences
57     token_counts.update(tokens)
58
59     # Get the 30 most common tokens
60     top_30_tokens = token_counts.most_common(30)
61
62     print(f"\nTask 3.2: Top 30 tokens:\n {top_30_tokens}")
63     return top_30_tokens
64
65 except FileNotFoundError:
66     print(f"Error: The file '{text_file}' was not found.")
67 except Exception as e:
68     print(f"An error occurred while counting tokens: {str(e)}")
69
70 # Run the function for Task 3.1
71 count_top_words(output_txt_file, top_words_csv)
72
73 # Run the function for Task 3.2
74 count_unique_tokens(output_txt_file)
```

Task 4: Named-Entity Recognition (NER)

```
Q1_Task_4.py > ...
1  import pandas as pd
2  import spacy
3  from collections import Counter
4  from transformers import AutoTokenizer, AutoModelForTokenClassification, pipeline
5  import warnings
6
7  # Suppress specific warnings related to truncation
8  warnings.filterwarnings("ignore")
9
10 output_txt_file = 'Q1_extracted_texts.txt'
11
12 # Define maximum token length for BioBERT (500 tokens, plus 2 for [CLS] and [SEP])
13 MAX_TOKEN_LENGTH = 500
14
15 # Try to load the SpaCy models with error handling
16 def load_spacy_model(model_path):
17     try:
18         return spacy.load(model_path)
19     except OSError as e:
20         print(f"Error: Failed to load the SpaCy model '{model_path}'. {str(e)}")
21         return None
22
23 # Try to load BioBERT model and tokenizer with error handling
24 > def load_biobert_pipeline(model_name): ...
25
26 # Function to extract entities using SpaCy models
27 def extract_entities_spacy(text_file, model, label_filter):
28     try:
29         with open(text_file, 'r', encoding='utf-8') as file:
30             text = file.read()
31
32         if model is None:
33             print("Error: No model available for entity extraction.")
34             return []
35
36         doc = model(text)
37         entities = [ent.text for ent in doc.ents if ent.label_ in label_filter]
38         return entities
39     except FileNotFoundError:
40         print(f"Error: The file '{text_file}' was not found.")
41         return []
42     except Exception as e:
43         print(f"An error occurred during entity extraction: {str(e)}")
44         return []
45
46
47
48
49
50
51
52
```

```

53 # Function to extract entities using BioBERT with correct truncation and chunking
54 def extract_entities_biobert(text_file, biobert_pipeline, max_length=MAX_TOKEN_LENGTH):
55     try:
56         with open(text_file, 'r', encoding='utf-8') as file:
57             text = file.read()
58
59             if biobert_pipeline is None:
60                 print("Error: No BioBERT pipeline available for entity extraction.")
61                 return [], []
62
63             # Tokenize the text and apply truncation and padding
64             tokenizer = biobert_pipeline.tokenizer
65             tokens = tokenizer(text, truncation=True, padding='max_length', max_length=max_length, return_tensors='pt')
66
67             # Ensure the total number of tokens is within the 512 limit, including [CLS] and [SEP]
68             input_ids = tokens['input_ids'].tolist()
69
70             diseases, drugs = [], []
71             for chunk in input_ids:
72                 entities = biobert_pipeline(chunk) # Pass text chunks to the BioBERT pipeline
73                 diseases += [entity['word'] for entity in entities if entity['entity'] == 'B-Disease']
74                 drugs += [entity['word'] for entity in entities if entity['entity'] == 'B-Drug']
75
76             return diseases, drugs
77
78     except FileNotFoundError:
79         print(f"Error: The file '{text_file}' was not found.")
80         return [], []
81     except Exception as e:
82         print(f"An error occurred during BioBERT entity extraction: {str(e)}")
83         return [], []
84
85 # Load models with error handling
86 nlp_sci_sm = load_spacy_model("en_core_sci_sm")
87 biobert_pipeline = load_biobert_pipeline("dmis-lab/biobert-base-cased-v1.1")
88
89 # Extract diseases and drugs using SpaCy models with error handling
90 diseases_sci_sm = extract_entities_spacy(output_txt_file, nlp_sci_sm, ['DISEASE'])
91
92 # Extract diseases and drugs using BioBERT with correct chunking
93 diseases_biobert, drugs_biobert = extract_entities_biobert(output_txt_file, biobert_pipeline)
94
95 # Compare the results of SpaCy models and BioBERT
96 print(f"SpaCy (en_core_sci_sm) detected {len(diseases_sci_sm)} diseases")
97 print(f"BioBERT detected {len(diseases_biobert)} diseases and {len(drugs_biobert)} drugs")
98
99 # Output most common entities if found
100 if diseases_biobert:
101     print("Most common diseases detected by BioBERT:", Counter(diseases_biobert).most_common(5))
102 else:
103     print("No diseases detected by BioBERT.")
104
105 if drugs_biobert:
106     print("Most common drugs detected by BioBERT:", Counter(drugs_biobert).most_common(5))
107 else:
108     print("No drugs detected by BioBERT.")

```