

TASK 3 - Development History

Regarding to divide our workload, each group member will incharge each task in this assignment, and if there any challenge happens, one will support another.

Here is the information about student main incharge:

- Task 1: Anh Huy Phung
- Task 2: Wei Yu Su

Task 2

Date: 05/04/2023

Today, we went through the assignment specification, and started exploring task 2 data.

Contribution:

Group member 1:

- Process concatenating all the data in the Excel sheets, and remove duplicated records and null columns. Extract the final DataFrame with 2 columns 'ChannelID' and 'textOriginal'.

Group member 2:

- Remove all the special characters and filter all English comments only, create a summary for each ChannelID that has at least 15 English comments.

proof

```
def remove_special_characters(text):
    """Remove text using regular expressions.

    Parameters:
    text (str): The input text.

    Returns:
    List: List of tokens.

    """
    tokenizer = RegexpTokenizer(re.sub(r'[^\w\s]', '', text))
    tokens = tokenizer.tokenize(text)
    return tokens

def load_data(channel_id, df):
    """Load data for a specific channel ID.

    Parameters:
    channel_id (int): Channel ID.
    df (DataFrame): Dataframe containing text data.

    Returns:
    DataFrame: Dataframe with word lists aggregated by ChannelID.

    """
    df['tokens'] = df['tokens'].apply(lambda x: remove_special_characters(x))
    df['tokens'] = df['tokens'].apply(lambda x: [token.lower() for token in x])
    df['tokens'] = df['tokens'].apply(lambda x: [token for token in x if token.isalpha()])
    df['tokens'] = df['tokens'].apply(lambda x: [token for token in x if token.isalpha()])

def channel_summary(df):
    """Generates a summary DataFrame containing aggregated information about channels.

    Args:
    df (DataFrame): Dataframe containing channel data.

    Returns:
    DataFrame: Dataframe with channel summary information including total comments and English comments count.

    """
    # Create a copy of the DataFrame to avoid modifying the original
    df_copy = df.copy()

    # Replace null values in 'textOriginal' with 1
    df_copy['textOriginal'].fillna(1, inplace=True)

    # Perform aggregation
    summary_df = df_copy.groupby('ChannelID').agg(
        total_comments=('textOriginal', 'count'), # Count non-null values
        eng_comments=('textOriginal', lambda x: (x == True).sum())
    ).reset_index()

    # Rename columns
    summary_df.columns = ['channel_id', 'all_comment_count', 'eng_comment_count']

    return summary_df

def english_df(summary_df, df):
    """Filters DataFrame to include only English comments from channels with at least 15 English comments.

    Args:
    summary_df: Summary DataFrame containing channel information.
    df: Original DataFrame containing comments.

    Returns:
    DataFrame: Filtered DataFrame containing English comments from selected channels.

    """
    # Filter channels with at least 15 English comments
    filter_id = summary_df['channel_id'][summary_df['eng_comment_count'] >= 15]

    # Filter DataFrame to include only selected channels and English comments
    filtered_channels = df[df['ChannelID'].isin(filter_id)]
    filtered_channels = filtered_channels[filtered_channels['is_english'] == True]

    return filtered_channels[['ChannelID', 'textOriginal']]
```

Step 2: Text extraction and cleaning

```
def json_txt_catch(dataframe, idx):
    """Extract text from JSON data.

    Parameters:
    dataframe (DataFrame): DataFrame containing JSON data.
    idx (int): Index of the JSON data.

    Returns:
    str: Extracted text.

    """
    try:
        snippet_txt = total_df.iloc[idx]['snippet']
        final_dictionary = eval(snippet_txt)
        textOriginal = final_dictionary['channelId'], final_dictionary['topLevelComment']
    except:
        log.append(str(x))
    finally:
        return textOriginal

def load_emoji_txt(filename):
    """Load emoji text from a file.

    Parameters:
    filename (str): File name.

    Returns:
    List: List of emoji words.

    """
    try:
        with open(filename, encoding='utf-8') as f:
            emojiwords = f.read().splitlines()
    except:
        print('The file is not found')
    finally:
        return emojiwords

def remove_emoji(string, emojiwords):
    """Remove emoji from a string.

    Parameters:
    string (str): Input string.
    emojiwords (List): List of emoji words.

    Returns:
    str: String with emoji removed.

    """
    emoji_pattern = re.compile(''.join(re.escape(emoji) for emoji in emojiwords))
    tem_txt = emoji_pattern.sub('', string)
    print('Something went wrong when removing emoji')
    finally:
        return tem_txt
```

Date 13/04/2024

Today, we started deciding to allocate steps for processing tokens and came with this flow:

- Tokenization
- Bigram
- Remove stopwords
- Remove rare tokens
- Remove less than 3 words tokens
- Stem

Contribution:

Group member 1:

- Develop a function to tokenize all the English comments by comment level as well as the bigram function to extract bigram functions from these..

Group member 2:

- Develop other functions (remove stopwords, rare tokens and less than 3 words tokens).

After this, we also plan to do other functions for extracting vocabulary and sparse presentation.

Proof

```
def load_stopwords_txt(filename):
    """Load stopwords from a file.

    Parameters:
    filename (str): File name.

    Returns:
    List: List of stopwords.

    """
    try:
        with open(filename) as f:
            stopwords = f.read().splitlines()
    except:
        print('The file is not found')
    finally:
        return stopwords

def remove_indep_stopwords(string, stopwords):
    """Remove independent stopwords from a string.

    Parameters:
    string (str): Input string.
    stopwords (List): List of stopwords.

    Returns:
    str: String with independent stopwords removed.

    """
    tem_txt = ''
    try:
        sw_pattern = re.compile(''.join(re.escape(stp) for stp in stopwords))
        pattern = re.compile(''.join(re.escape(stp) for stp in stopwords))
        regex = re.compile(pattern, re.IGNORECASE)
        tem_txt = regex.sub('', string)
        # tem_txt = sw_pattern.sub('', string)
    except:
        print('Something went wrong when removing stopwords')
    finally:
        return tem_txt

def remove_indep_stopwords(tokens, stopwords, is_uni):
    """Remove independent stopwords from tokens.

    Parameters:
    tokens (List): List of tokens.
    stopwords (List): List of stopwords.
    is_uni (bool): Whether to use unigram or bigram.

    Returns:
    List: List of tokens with independent stopwords removed.

    """
    stopwords_set = set(stopwords)
    if is_uni:
        stopped_tokens = [w for w in tokens if w not in stopwords_set]
    else:
        stopped_tokens = [w for w in tokens if w.split('_')[0] not in stopwords_set and w.split('_')[1] not in stopwords_set]
    return stopped_tokens

def remove_dep_stopwords(df, threshold):
    """Remove dependent stopwords from a DataFrame.

    Parameters:
    df (DataFrame): DataFrame containing token data.
    threshold (int): Threshold for removing dependent stopwords.

    Returns:
    DataFrame: DataFrame with dependent stopwords removed.

    """
    # Step 1: Calculate the number of unique ChannelIDs each word appears in
    word_id_counts = df.explode('word_list').groupby('word_list')['channel_id'].nunique()

    # Step 2: Calculate the threshold count
    ch_id_set = set(df['channel_id'])
    threshold_count = threshold * len(ch_id_set)

    # Step 3: Set the threshold
    context_dependent_stopwords = word_id_counts[word_id_counts >= threshold_count].index.tolist()

    return context_dependent_stopwords

def remove_dep_stopwords_list(df, list, threshold):
    """Remove dependent stopwords from a DataFrame.

    Parameters:
    df (DataFrame): DataFrame containing token data.
    list (List): List of stopwords.
    threshold (int): Threshold for removing dependent stopwords.

    Returns:
    DataFrame: DataFrame with dependent stopwords removed.

    """
    # Step 1: Calculate the number of unique ChannelIDs each word appears in
    word_id_counts = df.explode('word_list').groupby('word_list')['channel_id'].nunique()

    # Step 2: Calculate the threshold count
    ch_id_set = set(df['channel_id'])
    threshold_count = threshold * len(ch_id_set)

    # Step 3: Set the threshold
    temp_list = word_id_counts[word_id_counts >= threshold_count].index.tolist()
    context_dependent_stopwords = [w for w in list if w not in temp_list]

    return context_dependent_stopwords

# Removing the Less Frequent Words
def remove_less_freq_words(df, threshold):
    """Remove less frequent words from a DataFrame.

    Parameters:
    df (DataFrame): DataFrame containing token data.
    threshold (int): Threshold for removing less frequent words.

    Returns:
    DataFrame: DataFrame with less frequent words removed.

    """
    # Step 1: Calculate the number of unique ChannelIDs each word appears in
    word_id_counts = df.explode('word_list').groupby('word_list')['channel_id'].nunique()

    # Step 2: Calculate the threshold count
    ch_id_set = set(df['channel_id'])
    threshold_count = threshold * len(ch_id_set)

    # Step 3: Set the threshold
    rare_tokens = word_id_counts[word_id_counts >= threshold_count].index.tolist()

    return rare_tokens

def remove_less_len(word_list, length):
    """Remove less than length words from a list.

    Parameters:
    word_list (List): List of words.
    length (int): Minimum length.

    Returns:
    List: List of words with length greater than or equal to length.

    """
    tem_list = [token for token in word_list if len(token) >= length]
    return tem_list
```

Date 15/04/2024

Today, we reviewed process, step 4 processing flow offline.

Contribution:

Group member 1:

- Program testing.

Group member 2:

- Check that the program results are consistent with the excel data.

Proof

```
[ ] import matplotlib.pyplot as plt
def element_check(input1, input2):
    return [i for i in input1 if i not in input2]

for element in test_list:
    combination, name, input1, input2 = element[0], element[1], element[2], element[3]
    check, len_check = element_check(input1, input2), len(element_check(input1, input2))
    combination = ' & '.join(a + '_' + b for a, b in combination)
    print(f'Combination function: {combination} \n',
          f'Bigram extracting place: {name}, \n',
          f'final bigram: {len(input2)}',
          f'remove bigram: {len(element_check(input1, input2))} \n',
          f'total vocab: {len(total_vocab_sorted)}')

import pickle

# Combination function: rare_tokens_set & less_three_tokens_remove & eliminate_or_choose_
# Bigram extracting place: tokens_in_channelID_df,
# final bigram: 194,
# remove bigram: 6
# total vocab: 2991
# Combination function: rare_tokens_set & less_three_tokens_remove & eliminate_or_choose_
# Bigram extracting place: eliminate_or_choose_words,
# final bigram: 195,

# Step 4-a
#token_df = tokenization_df.drop_new_df
#token_df = tokenization_df(only_eng_df)
token_df = tokenization_df(eng_fifteen_df)
unigram_df = token_df.copy()
bigram_df = generate_ngrams_df(token_df, 2)
bigram_txt_df = generate_ngrams_txtpair_df(token_df, 2)
#token_df
#type(token_df)
#DataFrame
#channel_id, word_list

unigram_agg_df = token_agg_by_chid(unigram_df)
bigram_agg_df = token_agg_by_chid(bigram_df)

# Step4-d
#Unigram
rare_unigram_df = unigram_agg_df.copy()
rare_unigram_list = remove_less_freq_words(rare_unigram_df, 0.01)
#Bigram
rare_bigram_df = bigram_agg_df.copy()
rare_bigram_list = remove_less_freq_words(rare_bigram_df, 0.01)
#List

n = 0
for i in range(len(rare_bigram_list)):
    if len(rare_bigram_list[i].split('_')) < 2:
        print(rare_bigram_list[i])
    else:
        n = n + 1

# Step4-e
less_uni_token_list = remove_less_len(rare_unigram_list, 3)
less_bi_token_list = remove_less_bi_len(rare_bigram_list, 3)

# Step4-f
#Unigram
less_uni_token_list_cp = less_uni_token_list.copy()
stopped_uni_list = remove_dep_stopwords_list(rare_unigram_df, less_uni_token_list_cp, 0.99)
#Bigram
less_bi_token_list_cp = less_bi_token_list.copy()
stopped_bi_list = remove_dep_stopwords_list(rare_bigram_df, less_bi_token_list_cp, 0.99)
#stopped_bi_df
#context-independent
#Unigram
stopped_uni_ind_list = remove_indep_stopwords(stopped_uni_list, stopwords, True)
#Bigram
stopped_bi_ind_list = remove_indep_stopwords(stopped_bi_list, stopwords, False)
```

Date 18/04/2024

Today, we started review all of the functions and generated file submission file.

Contribution:

Group member 1:

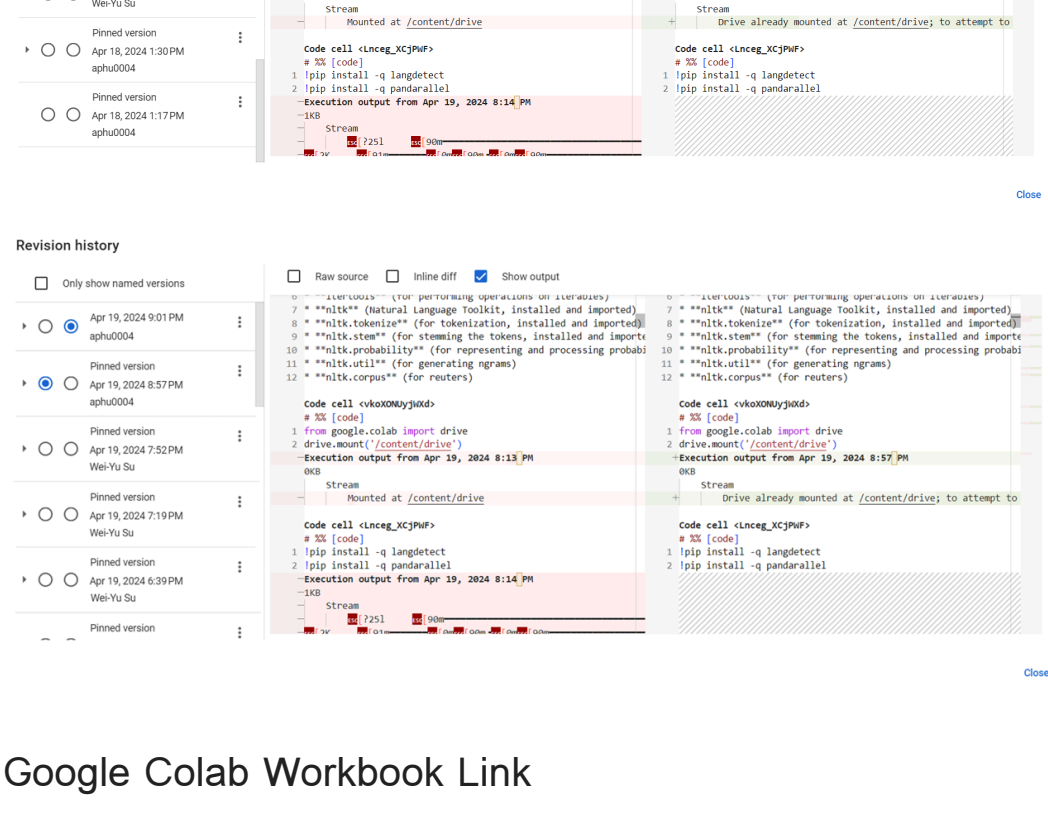
- Review all of the generated functions (handling tokens and deliver sparse representation), give example in each function for easier documentation.

Group member 2:

- Documentation.

Proof

These are pictures to show our explanation to some of funtions and examples.



Google Colab Workbook Link

Task 2