## HOMEWORK 6: TEXT CLASSIFICATION

In this homework, you will create models to classify texts from TRUE call-center. There are two classification tasks:

1. Action Classification: Identify which action the customer would like to take (e.g. enquire, report, cancle)
2. Object Classification: Identify which object the customer is referring to (e.g. payment, truemoney, internet, roaming)

We will focus only on the Object Classification task for this homework.

In this homework, you are asked compare different text classification models in terms of accuracy and inference time.

You will need to build 3 different models.

1. A model based on tf-idf
2. A model based on MUSE
3. A model based on wangchanBERTa

**You will be ask to submit 3 different files (.pdf from .ipynb) that does the 3 different models. Finally, answer the accuracy and runtime numbers in MCV.**

This homework is quite free form, and your answer may vary. We hope that the processing during the course of this assignment will make you think more about the design choices in text classification.

```
1 !wget --no-check-certificate https://www.dropbox.com/s/37u83g55p19kvrl/clean-phone-data-for-students.csv -O ./clean-phone-
```

```
--2025-02-15 13:57:38--  https://www.dropbox.com/s/37u83g55p19kvrl/clean-phone-data-for-students.csv
Resolving www.dropbox.com (www.dropbox.com)... 2620:100:6035:18::a27d:5512, 162.125.85.18
Connecting to www.dropbox.com (www.dropbox.com)|2620:100:6035:18::a27d:5512|:443... connected.
HTTP request sent, awaiting response... 302 Found
Location: https://www.dropbox.com/scl/fi/8h8hvsw9uj6o0524lfe4i/clean-phone-data-for-students.csv?rlkey=lwv5xbf16jerehnv3
--2025-02-15 13:57:38--  https://www.dropbox.com/scl/fi/8h8hvsw9uj6o0524lfe4i/clean-phone-data-for-students.csv?rlkey=lw
Reusing existing connection to [www.dropbox.com]:443.
HTTP request sent, awaiting response... 302 Found
Location: https://uc7df6d5a23f54bd3592a4a478d5.dl.dropboxusercontent.com/cd/0/inline/CkIpCV84EQav89R1pHJhwTj_FwQ5tzBwdvR
--2025-02-15 13:57:39--  https://uc7df6d5a23f54bd3592a4a478d5.dl.dropboxusercontent.com/cd/0/inline/CkIpCV84EQav89R1pHJh
Resolving uc7df6d5a23f54bd3592a4a478d5.dl.dropboxusercontent.com (uc7df6d5a23f54bd3592a4a478d5.dl.dropboxusercontent.com
Connecting to uc7df6d5a23f54bd3592a4a478d5.dl.dropboxusercontent.com (uc7df6d5a23f54bd3592a4a478d5.dl.dropboxusercontent
HTTP request sent, awaiting response... 200 OK
Length: 2518977 (2.4M) [text/plain]
Saving to: './clean-phone-data-for-students.csv'

./clean-phone-data- 100%[===================>]   2.40M  4.04MB/s    in 0.6s

2025-02-15 13:57:41 (4.04 MB/s) - './clean-phone-data-for-students.csv' saved [2518977/2518977]
```

```
1 !pip install -q pythainlp
```

## Import Libs

```
1 %matplotlib inline
2 import pandas
3 import sklearn
4 import numpy as np
5 import time
6 import matplotlib.pyplot as plt
7 import pandas as pd
8 from pprint import pprint
9
10 from torch.utils.data import Dataset
11 from IPython.display import display
12 from collections import defaultdict
13 from sklearn.metrics import accuracy_score
```

## Loading data

First, we load the data from disk into a Dataframe.

A Dataframe is essentially a table, or 2D-array/Matrix with a name for each column.

```
1 data_df = pd.read_csv('clean-phone-data-for-students.csv')
```

Let's preview the data.

```
1 # Show the top 5 rows
2 display(data_df.head())
3 # Summarize the data
4 data_df.describe()
```

| | Sentence Utterance | Action | Object |
|---|---|---|---|
| 0 | <PHONE_NUMBER_REMOVED> ผมไปจ่ายเงินที่ Counte... | enquire | payment |
| 1 | internet ยังความเร็วอยู่เท่าไหร ครับ | enquire | package |
| 2 | ตะกี้ไปชำระค่าบริการไปแล้ว แต่ยังใช้งานไม่ได้... | report | suspend |
| 3 | พี่ค่ะยังใช้ internet ไม่ได้เลยค่ะ เป็นเครื่อ... | enquire | internet |
| 4 | ฮาโหล คะ พอดีว่าเมื่อวานเปิดซิมทรูมูฟ แต่มันโ... | report | phone_issues |

| | Sentence Utterance | Action | Object |
|---|---|---|---|
| count | 16175 | 16175 | 16175 |
| unique | 13389 | 10 | 33 |
| top | บริการอื่นๆ | enquire | service |
| freq | 97 | 10377 | 2525 |

## Data cleaning

We call the DataFrame.describe() again. Notice that there are 33 unique labels/classes for object and 10 unique labels for action that the model will try to predict. But there are unwanted duplications e.g. Idd,idd,lotalty_card,Lotalty_card

Also note that, there are 13389 unqiue sentence utterances from 16175 utterances. You have to clean that too!

### #TODO 0.1:

You will have to remove unwanted label duplications as well as duplications in text inputs. Also, you will have to trim out unwanted whitespaces from the text inputs. This shouldn't be too hard, as you have already seen it in the demo.

```
1 display(data_df.describe())
2 display(data_df.Object.unique())
3 display(data_df.Action.unique())
```

| | Sentence Utterance | Action | Object |
|---|---|---|---|
| count | 16175 | 16175 | 16175 |
| unique | 13389 | 10 | 33 |
| top | บริการอื่นๆ | enquire | service |
| freq | 97 | 10377 | 2525 |

```
array(['payment', 'package', 'suspend', 'internet', 'phone_issues',
       'service', 'nonTrueMove', 'balance', 'detail', 'bill', 'credit',
       'promotion', 'mobile_setting', 'iservice', 'roaming', 'truemoney',
       'information', 'lost_stolen', 'balance_minutes', 'idd',
       'TrueMoney', 'garbage', 'Payment', 'IDD', 'ringtone', 'Idd',
       'rate', 'loyalty_card', 'contact', 'officer', 'Balance', 'Service',
       'Loyalty_card'], dtype=object)
array(['enquire', 'report', 'cancel', 'Enquire', 'buy', 'activate',
       'request', 'Report', 'garbage', 'change'], dtype=object)
```

```
 1 clean_data_time = time.time()
 2
 3 # Group the duplicate label
 4 data_df.dropna(subset=['Object'], inplace=True)
 5 data_df['Object'] = data_df['Object'].apply(lambda x: x.lower())
 6
 7 # Clean the data
 8 data_df['Sentence Utterance'] = data_df['Sentence Utterance'].apply(lambda x: str(x).strip())
 9 data_df['Sentence Utterance'] = data_df['Sentence Utterance'].apply(lambda x: x.lower())
10 data_df.drop_duplicates(subset=['Sentence Utterance'], inplace=True)
11
12 # Drop the unused columns
13 data_df.drop(columns=['Action'], inplace=True)
14
15 clean_data_time = time.time() - clean_data_time
```

Split data into train, valdation, and test sets (normally the ratio will be 80:10:10 , respectively). We recommend to use train_test_spilt from scikit-learn to split the data into train, validation, test set.

In addition, it should split the data that distribution of the labels in train, validation, test set are similar. There is **stratify** option to handle this issue.

Make sure the same data splitting is used for all models.

```
1 from sklearn.model_selection import train_test_split
2 from collections import Counter
3
4 split_data_time = time.time()
5
6 # For the object column, we will only keep the object that has more than 2% of the total data
7 object_counter = Counter(data_df['Object'])
8 stratify_col = data_df['Object'].apply(lambda x: 'other' if object_counter[x] < 0.02*len(data_df) else x)
9 train_df, test_df = train_test_split(data_df, test_size=0.2, random_state=4242, stratify=stratify_col)
10
11 object_counter = Counter(test_df['Object'])
12 stratify_col = test_df['Object'].apply(lambda x: 'other' if object_counter[x] < 0.02*len(test_df) else x)
13 test_df, val_df = train_test_split(test_df, test_size=0.5, random_state=4242, stratify=stratify_col)
14
15 train_df = train_df.reset_index(drop=True)
16 test_df = test_df.reset_index(drop=True)
17 val_df = val_df.reset_index(drop=True)
18
19 print(f"Train size: {len(train_df)}")
20 print(f"Test size: {len(test_df)}")
21 print(f"Val size: {len(val_df)}")
22
23 split_data_time = time.time() - split_data_time
```

```
   Train size: 10689
   Test size: 1336
   Val size: 1337
```

```
1 # Save the data
2 # train_df.to_csv('train.csv', index=False)
3 # test_df.to_csv('test.csv', index=False)
4 # val_df.to_csv('val.csv', index=False)
```

## ⌄ Model 1 TF-IDF

Build a model to train a tf-idf text classifier. Use a simple logistic regression model for the classifier.

For this part, you may find this tutorial helpful.

Below are some design choices you need to consider to accomplish this task. Be sure to answer them when you submit your model.

What tokenizer will you use? Why?

**Ans:** Use `newmm` tokenizer from PyThaiNLP. It is a dictionary-based tokenizer that is suitable for Thai text. Additionally, I use `ngram_range=(1, 2)` to include bigram in the feature set. This is because bigram can capture the relationship between words that are close to each other.

Will you ignore some stop words (a, an, the, to, etc. for English) in your tf-idf? Is it important? PythaiNLP provides a list of stopwords if you want to use (https://pythainlp.org/docs/2.0/api/corpus.html#pythainlp.corpus.common.thai_stopwords)

**Ans:** I use the list of stopwords provided by PyThaiNLP. It is important to remove stopwords because they are common words that do not carry much meaning. Removing them can help the model focus on more important words.

The dictionary of TF-IDF is usually based on the training data. How many words in the test set are OOVs?

**Ans:** 324 words

```
1 import pandas as pd
2 train_df, test_df, val_df = pd.read_csv('train.csv'), pd.read_csv('test.csv'), pd.read_csv('val.csv')
```

```
1 from sklearn.feature_extraction.text import TfidfVectorizer
2 from sklearn.linear_model import LogisticRegression
3 from sklearn.metrics import classification_report, accuracy_score
4 import pythainlp
```

```
1 tokenize_and_vectorize_time = time.time()
2
3 # Create the tf-idf vectorizer
4 tf_idf = TfidfVectorizer(tokenizer=pythainlp.word_tokenize,
5                          ngram_range=(1, 2),
6                          stop_words=list(pythainlp.corpus.thai_stopwords()))
7 tf_idf.fit(train_df['Sentence Utterance'])
```

```
 8
 9 # Transform the data
10 train_x = tf_idf.transform(train_df['Sentence Utterance'])
11 test_x = tf_idf.transform(test_df['Sentence Utterance'])
12 val_x = tf_idf.transform(val_df['Sentence Utterance'])
13
14 tokenize_and_vectorize_time = time.time() - tokenize_and_vectorize_time
15
16 # Get y label
17 train_y = train_df['Object']
18 test_y = test_df['Object']
19 val_y = val_df['Object']
```

```
/Users/jirayuwat/anaconda3/envs/nlp/lib/python3.11/site-packages/sklearn/feature_extraction/text.py:517: UserWarning: Th
  warnings.warn(
  /Users/jirayuwat/anaconda3/envs/nlp/lib/python3.11/site-packages/sklearn/feature_extraction/text.py:402: UserWarning: Yo
  warnings.warn(
```

```
 1 train_time = time.time()
 2
 3 # Create the model
 4 model = LogisticRegression(max_iter=1000, penalty='elasticnet', solver='saga', l1_ratio=0.5, C=1.0, n_jobs=-1)
 5 model.fit(train_x, train_y)
 6
 7 train_time = time.time() - train_time
```

```
 1 all_inference_time = time.time()
 2
 3 # Predict the data
 4 train_pred = model.predict(train_x)
 5 test_pred = model.predict(test_x)
 6 val_pred = model.predict(val_x)
 7
 8 all_inference_time = time.time() - all_inference_time
 9
10 # Calculate the accuracy
11 train_acc = accuracy_score(train_y, train_pred)
12 test_acc = accuracy_score(test_y, test_pred)
13 val_acc = accuracy_score(val_y, val_pred)
14
15 print(f"Train accuracy: {train_acc}")
16 print(f"Test accuracy: {test_acc}")
17 print(f"Val accuracy: {val_acc}")
```

```
Train accuracy: 0.7176536626438395
Test accuracy: 0.6803892215568862
Val accuracy: 0.6798803290949887
```

```
 1 # Print the classification report
 2 print("Val classification report")
 3 print(classification_report(val_y, val_pred))
 4
 5 import pickle
 6
 7 # Save the classificaion report
 8 classification_report_dict = classification_report(val_y, val_pred, output_dict=True)
 9 with open('classification_report_tfidf.pkl', 'wb') as f:
10     pickle.dump(classification_report_dict, f)
```

```
Val classification report
                 precision    recall  f1-score   support

        balance       0.75      0.78      0.76       148
balance_minutes       0.80      0.80      0.80         5
           bill       0.65      0.48      0.55        54
         credit       1.00      0.76      0.87        17
         detail       0.89      0.24      0.38        33
        garbage       0.00      0.00      0.00         6
            idd       0.82      0.64      0.72        14
    information       1.00      0.31      0.47        29
       internet       0.66      0.82      0.73       179
       iservice       0.00      0.00      0.00         3
    lost_stolen       0.93      0.90      0.92        30
   loyalty_card       1.00      0.62      0.77         8
 mobile_setting       0.85      0.39      0.54        28
     nontruemove       0.89      0.38      0.53        21
        officer       0.00      0.00      0.00         1
        package       0.62      0.72      0.67       179
        payment       0.58      0.66      0.61        64
   phone_issues       0.63      0.57      0.60        58
      promotion       0.56      0.55      0.56       115
           rate       0.00      0.00      0.00         6
       ringtone       1.00      0.91      0.95        11
```

```
        roaming          0.85      0.65      0.73         17
        service          0.66      0.82      0.73        211
        suspend          0.70      0.67      0.69         73
      truemoney          0.83      0.74      0.78         27

       accuracy                              0.68       1337
      macro avg          0.67      0.54      0.57       1337
   weighted avg          0.69      0.68      0.67       1337
```

```
/Users/jirayuwat/anaconda3/envs/nlp/lib/python3.11/site-packages/sklearn/metrics/_classification.py:1565: UndefinedMetri
  _warn_prf(average, modifier, f"{metric.capitalize()} is", len(result))
/Users/jirayuwat/anaconda3/envs/nlp/lib/python3.11/site-packages/sklearn/metrics/_classification.py:1565: UndefinedMetri
  _warn_prf(average, modifier, f"{metric.capitalize()} is", len(result))
/Users/jirayuwat/anaconda3/envs/nlp/lib/python3.11/site-packages/sklearn/metrics/_classification.py:1565: UndefinedMetri
  _warn_prf(average, modifier, f"{metric.capitalize()} is", len(result))
/Users/jirayuwat/anaconda3/envs/nlp/lib/python3.11/site-packages/sklearn/metrics/_classification.py:1565: UndefinedMetri
  _warn_prf(average, modifier, f"{metric.capitalize()} is", len(result))
/Users/jirayuwat/anaconda3/envs/nlp/lib/python3.11/site-packages/sklearn/metrics/_classification.py:1565: UndefinedMetri
  _warn_prf(average, modifier, f"{metric.capitalize()} is", len(result))
/Users/jirayuwat/anaconda3/envs/nlp/lib/python3.11/site-packages/sklearn/metrics/_classification.py:1565: UndefinedMetri
  _warn_prf(average, modifier, f"{metric.capitalize()} is", len(result))
```

```python
1 print(f'''
2 All preprocessing time: {clean_data_time + split_data_time + tokenize_and_vectorize_time:.2f} seconds
3 - Clean data time: {clean_data_time:.2f} seconds
4 - Split data time: {split_data_time:.2f} seconds
5 - Tokenize and vectorize time: {tokenize_and_vectorize_time:.2f} seconds
6 Train time: {train_time:.2f} seconds
7 Inference time: {all_inference_time:.2f} seconds
8 '''.strip())
```

```
All preprocessing time: 1.20 seconds
 - Clean data time: 0.02 seconds
 - Split data time: 0.03 seconds
 - Tokenize and vectorize time: 1.15 seconds
Train time: 34.22 seconds
Inference time: 0.01 seconds
```

```python
 1 # Count OOV words
 2 tokenizer = pythainlp.word_tokenize
 3
 4 train_words = set()
 5 for sentence in train_df['Sentence Utterance']:
 6     train_words.update(tokenizer(sentence))
 7
 8 test_words = set()
 9 for sentence in test_df['Sentence Utterance']:
10     test_words.update(tokenizer(sentence))
11
12 val_words = set()
13 for sentence in val_df['Sentence Utterance']:
14     val_words.update(tokenizer(sentence))
15
16 oov_words = test_words.union(val_words) - train_words
17 print(f"OOV words: {len(oov_words)}")
18 print("OOV words:", oov_words)
```

```
OOV words: 324
OOV words: {'เหล่านี้', 'ไทร', 'เหตุการณ์', 'เชียงคาน', 'สเวิด', 'ประวัติการ', 'การผ่อน', 'lift', 'คล่องตัว', '4541843', 'ละคร', 'เลค
```