# Artificial Intelligence In Medical Applications and Services
# Homework3 Report

## 1. Score Result

```
1   f1score

0.4258144301165879
```

## 2. Process

As the TA mentioned, we can seperate the whole project into 2 parts which are **Preprocessing** and **Model_Training.**

### 1.  Preprocessing

First we use **loadInputFile** function to load the data given by TA (SampleData_deid.txt).

```
 1   def loadInputFile(path):
 2       trainingset = list()  # store trainingset [content,content,...]
 3       position = list()  # store position [article_id, start_pos, end
 4       mentions = dict()  # store mentions[mention] = Type
 5       with open(file_path, 'r', encoding='utf8') as f:
 6           file_text=f.read().encode('utf-8').decode('utf-8-sig')
 7       datas=file_text.split('\n\n--------------------\n\n')[:-1]
 8       for data in datas:
 9           data=data.split('\n')
10           content=data[0]
11           trainingset.append(content)
12           annotations=data[1:]
13           for annot in annotations[1:]:
14               annot=annot.split('\t') #annot= article_id, start_pos,
15               position.extend(annot)
16               mentions[annot[3]]=annot[4]
17
18       return trainingset, position, mentions
```

Then use CRFFormatData to change the input data into a corresponding format (clear the space, adding 'B' 'O' 'I' label)

```
1  def CRFFormatData(trainingset, position, path):
2      if (os.path.isfile(path)):
3          os.remove(path)
4      outputfile = open(path, 'a', encoding= 'utf-8')
5
6      # output file lines
7      count = 0 # annotation counts in each content
8      tagged = list()
9      for article_id in range(len(trainingset)):
10         trainingset_split = list(trainingset[article_id])
11         while '' or ' ' in trainingset_split:
12             if '' in trainingset_split:
13                 trainingset_split.remove('')
14             else:
15                 trainingset_split.remove(' ')
16         start_tmp = 0
17         for position_idx in range(0,len(position),5):
18             if int(position[position_idx]) == article_id:
19                 count += 1
20                 if count == 1:
21                     start_pos = int(position[position_idx+1])
22                     end_pos = int(position[position_idx+2])
23                     entity_type=position[position_idx+4]
24                     if start_pos == 0:
25                         token = list(trainingset[article_id][start_pos:end_pos])
26                         whole_token = trainingset[article_id][start_pos:end_pos]
27                         for token_idx in range(len(token)):
28                             if len(token[token_idx].replace(' ','')) == 0:
29                                 continue
30                             # BIO states
31                             if token_idx == 0:
32                                 label = 'B-'+entity_type
33                             else:
34                                 label = 'I-'+entity_type
```

## 2. Model_Training

Create a model to fit and predict, and don't forget to flatten to get final f1score
and print the prediction result

```
1  def CRF(x_train, y_train, x_test, y_test):
2      crf = sklearn_crfsuite.CRF(
3          algorithm='lbfgs',
4          c1=0.1,
5          c2=0.1,
6          max_iterations=100,
7          all_possible_transitions=True
8      )
9      crf.fit(x_train, y_train)
10     # print(crf)
11     y_pred = crf.predict(x_test)
12     y_pred_mar = crf.predict_marginals(x_test)
13
14     # print(y_pred_mar)
15
16     labels = list(crf.classes_)
17     labels.remove('O')
18     f1score = metrics.flat_f1_score(y_test, y_pred, average='weighted', labels=labels)
19     sorted_labels = sorted(labels,key=lambda name: (name[1:], name[0])) # group B and I results
20     print(flat_classification_report(y_test, y_pred, labels=sorted_labels, digits=3))
21     return y_pred, y_pred_mar, f1score
```

This part is very important. We need to prepare a pretrained word vector file to let
it find corresponding key value of the token that get from input data.

The first two tokens are the parameters which are vocabulary_size and it's dimension. And keep the rest tokens into **vec** array in float type, and return word_vecs dictionary.

```
1   # load pretrained word vectors
2   # get a dict of tokens (key) and their pretrained word vectors (value)
3   # pretrained word2vec CBOW word vector: https://fgc.stpi.narl.org.tw/activity/videoDe
4   dim = 0
5   word_vecs= {}
6   # open pretrained word vector file
7   with open('/content/drive/MyDrive/Colab Notebooks/NER/glove.42B.300d.txt') as f:
8       for line in f:
9           tokens = line.strip().split()
10
11          # there 2 integers in the first line: vocabulary_size, word_vector_dim
12          if len(tokens) == 2:
13              dim = int(tokens[1])
14              continue
15
16          word = tokens[0]
17          vec = np.array([ float(t) for t in tokens[1:] ])
18          word_vecs[word] = vec
19
```

Prepare dataset and split them into training_set and testing_set with test_size = 0.33

```
7   def Dataset(data_path):
8       with open(data_path, 'r', encoding='utf-8') as f:
9           data=f.readlines()#.encode('utf-8').decode('utf-8-sig')
10      data_list, data_list_tmp = list(), list()
11      article_id_list=list()
12      idx=0
13      for row in data:
14          data_tuple = tuple()
15          if row == '\n':
16              article_id_list.append(idx)
17              idx+=1
18              data_list.append(data_list_tmp)
19              data_list_tmp = []
20          else:
21              row = row.strip('\n').split(' ')
22              data_tuple = (row[0], row[1])
23              data_list_tmp.append(data_tuple)
24      if len(data_list_tmp) != 0:
25          data_list.append(data_list_tmp)
26
27      # here we random split data into training dataset and testing dataset
28      # but you should take `development data` or `test data` as testing data
29      # At that time, you could just delete this line,
30      # and generate data_list of `train data` and data_list of `development/test data` by this function
31      traindata_list, testdata_list, traindata_article_id_list, testdata_article_id_list=train_test_split(data_list,
32                                                                                          article_id_list,
33                                                                                          test_size=0.33,
34                                                                                          random_state=42)
35
36      return data_list, traindata_list, testdata_list, traindata_article_id_list, testdata_article_id_list
```

Use the pretrained word vector that we just imported as embedding_dict to find the input data_list's corresponding value with three for loops and keep them as a list.

```
 4   def Word2Vector(data_list, embedding_dict):
 5       embedding_list = list()
 6
 7       # No Match Word (unknown word) Vector in Embedding
 8       unk_vector=np.random.rand(*(list(embedding_dict.values())[0].shape))
 9
10       for idx_list in range(len(data_list)):
11           embedding_list_tmp = list()
12           for idx_tuple in range(len(data_list[idx_list])):
13               key = data_list[idx_list][idx_tuple][0] # token
14               # print(str(idx_tuple) + key)
15               if key in embedding_dict:
16                   value = embedding_dict[key]
17               else:
18                   value = unk_vector
19               embedding_list_tmp.append(value)
20           embedding_list.append(embedding_list_tmp)
21       return embedding_list
```

Extract all the features of embedding_list and keep also(three loops go through 3 dimension of embed_list)

```
 1   # input features: pretrained word vectors of each token
 2   # return a list of feature dicts, each feature dict corresponding to each token
 3   def Feature(embed_list):
 4       feature_list = list()
 5       for idx_list in range(len(embed_list)):
 6           feature_list_tmp = list()
 7           for idx_tuple in range(len(embed_list[idx_list])):
 8               feature_dict = dict()
 9               for idx_vec in range(len(embed_list[idx_list][idx_tuple])):
10                   feature_dict['dim_' + str(idx_vec+1)] = embed_list[idx_list][idx_tuple][idx_vec]
11               feature_list_tmp.append(feature_dict)
12
13           feature_list.append(feature_list_tmp)
14       return feature_list
```

This is how embed_list[idx_list][idx_tuple][idx_vec] looks like:

```
0.6180830816584256
0.23310316349906202
0.7140401709909245
0.19546616359681368
0.8567940396548351
0.7400688220122569
```

Prepare a Label list by set the third dimension of data_list to **1**

```
 1   # get the labels of each tokens in train.data
 2   # return a list of lists of labels
 3   def Preprocess(data_list):
 4       label_list = list()
 5       for idx_list in range(len(data_list)):
 6           label_list_tmp = list()
 7           for idx_tuple in range(len(data_list[idx_list])):
 8               label_list_tmp.append(data_list[idx_list][idx_tuple][1])
 9           label_list.append(label_list_tmp)
10       return label_list
```

Call the functions

```
 1   # Load Word Embedding
 2   trainembed_list = Word2Vector(traindata_list, word_vecs)
 3   testembed_list = Word2Vector(testdata_list, word_vecs)
 4
 5   # CRF - Train Data (Augmentation Data)
 6   x_train = Feature(trainembed_list)
 7   y_train = Preprocess(traindata_list)
 8
 9   # CRF - Test Data (Golden Standard)
10   x_test = Feature(testembed_list)
11   y_test = Preprocess(testdata_list)
```

Fit and predict

```
 1   y_pred, y_pred_mar, f1score = CRF(x_train, y_train, x_test, y_test)
```

Result

|             | precision | recall | f1-score | support |
|-------------|-----------|--------|----------|---------|
| B-location  | 0.000     | 0.000  | 0.000    | 15      |
| I-location  | 0.000     | 0.000  | 0.000    | 41      |
| B-med_exam  | 0.200     | 0.030  | 0.053    | 33      |
| I-med_exam  | 1.000     | 0.075  | 0.140    | 80      |
| B-money     | 0.364     | 0.333  | 0.348    | 12      |
| I-money     | 0.353     | 0.171  | 0.231    | 35      |
| B-name      | 0.500     | 0.143  | 0.222    | 7       |
| I-name      | 0.333     | 0.100  | 0.154    | 10      |
| B-time      | 0.667     | 0.450  | 0.538    | 111     |
| I-time      | 0.825     | 0.532  | 0.647    | 265     |
|             |           |        |          |         |
| micro avg   | 0.714     | 0.345  | 0.465    | 609     |
| macro avg   | 0.424     | 0.184  | 0.233    | 609     |
| weighted avg| 0.661     | 0.345  | 0.426    | 609     |

F1Score

```
 1   f1score
```

0.4258144301165879

Change the output format and save into output.tsv file
Here we show the entity text's id, start_position, ending_position, content and type.

```
1  output="article_id\tstart_position\tend_position\tentity_text\tentity_type\n"
2  for test_id in range(len(y_pred)):
3      pos=0
4      start_pos=None
5      end_pos=None
6      entity_text=None
7      entity_type=None
8      for pred_id in range(len(y_pred[test_id])):
9          if y_pred[test_id][pred_id][0]=='B':
10             start_pos=pos
11             entity_type=y_pred[test_id][pred_id][2:]
12         elif start_pos is not None and y_pred[test_id][pred_id][0]=='I' and y_pred[test_id][pred_id+1][0]=='O':
13             end_pos=pos
14             entity_text=''.join([testdata_list[test_id][position][0] for position in range(start_pos,end_pos+1)])
15             line=str(testdata_article_id_list[test_id])+'\t\t'+str(start_pos)+'\t\t'+str(end_pos+1)+'\t\t'+entity_text+'\t\t'+entity_type
16             output+=line+'\n'
17         pos+=1
```

| article_id | start_position | end_position | entity_text | entity_type |
|---|---|---|---|---|
| 8 | 52 | 54 | 前天 | time |
| 8 | 68 | 70 | 昨天 | time |
| 8 | 189 | 193 | 二十分鐘 | time |
| 8 | 293 | 295 | 五年 | time |
| 8 | 540 | 544 | 兩個禮拜 | time |
| 8 | 726 | 728 | 前天 | time |
| 8 | 730 | 732 | 前天 | time |
| 8 | 858 | 860 | 前天 | time |
| 8 | 898 | 900 | 前天 | time |
| 8 | 1549 | 1551 | 五天 | time |
| 8 | 1622 | 1626 | 五天禮拜 | time |
| 8 | 2352 | 2354 | 去喬 | time |
| 8 | 2560 | 2563 | 兩個月 | time |
| 16 | 51 | 55 | 九、十點 | time |
| 16 | 60 | 64 | 九、十點 | time |
| 16 | 122 | 124 | 三年 | time |
| 16 | 130 | 132 | 三年 | time |
| 16 | 247 | 249 | 三年 | time |
| 0 | 1268 | 1271 | 8公分 | med_exam |
| 0 | 1358 | 1362 | 三多路上 | time |
| 0 | 2576 | 2578 | 五天 | time |

## 3. Features

　　這裡原本我是只使用助教給的抓取data_list的第一個([0])也就是token內容去比對embedding_dict的key value 作為features

```
for idx_tuple in range(len(data_list[idx_list])):
    key = data_list[idx_list][idx_tuple][0] # token
    # print(str(idx_tuple) + key)
    if key in embedding_dict:
        value = embedding_dict[key]
    else:
        value = unk_vector
```

```
for idx_list in range(len(embed_list)):
    feature_list_tmp = list()
    for idx_tuple in range(len(embed_list[idx_list])):
        feature_dict = dict()
        for idx_vec in range(len(embed_list[idx_list][idx_tuple])):
            feature_dict['dim_' + str(idx_vec+1)] = embed_list[idx_list][idx_tuple][idx_vec]
```

This is how embed_list[idx_list][idx_tuple][idx_vec] looks like:

```
0.6180830816584256
0.23310316349906202
0.7140401709909245
0.19546616359681368
0.8567940396548351
0.7400688220122569
```

1. 後來有試着先做了加入word_length的features，平均f1score有比較微微的上升(比原本的0.35~0.41)，固定維持在0.4以上
2. 再來是加入word_position(start_pos)，我這裡是選擇給start_pos，因為給end_pos的效果其實沒什麼影響，給了start_pos後只加了一點點
3. 再來是對於features的組織是要通過word2vector和features這兩個function後才能作成x_train, x_test的，我發現其實我的word2vector有成功把我想要的features加入到embedding_dict中並且return吐出，但是由於時間的關係我在features的dict中沒有成功整合，所以model其實沒有吃到我加入的features。不過就上述情況而言，我已經成功了解如果要訓練一個NER-CRF的模型的大概流程是怎麼樣的了。

```python
else:
    # print('Row:' + row)
    row = row.strip('\n').split(' ')
    data_tuple = (row[0], row[1], row[2], row[3])
    # print(data_tuple)
```

```python
for idx_list in range(len(data_list)):
    embedding_list_tmp = list()
    embed_dict = dict()
    for idx_tuple in range(len(data_list[idx_list])):
        key = data_list[idx_list][idx_tuple][0] # token
        # print(str(idx_tuple) + key)
        if key in embedding_dict:
            value = embedding_dict[key]
        else:
            value = unk_vector
        embed_dict[0] = value # token in embedding_dict's key value
        embed_dict[1] = data_list[idx_list][idx_tuple][2] # word_position(start_pos)
        embed_dict[2] = data_list[idx_list][idx_tuple][3] # word_length
        # embedding_list_tmp.append(value)
        embedding_list_tmp.append(embed_dict)
    embedding_list.append(embedding_list_tmp)
return embedding_list
```

```python
def Feature(embed_list):
    feature_list = list()
    for idx_list in range(len(embed_list)):
        feature_list_tmp = list()
        for idx_tuple in range(len(embed_list[idx_list])):
            feature_dict = dict()
            # print(embed_list[idx_list][idx_tuple][0])
            for idx_vec in range(len(embed_list[idx_list][idx_tuple])):
                feature_dict['dim_' + str(idx_vec+1)] = embed_list[idx_list][idx_tuple][idx_vec]
                # feature_dict[str(0+1)] = embed_list[idx_list][0]
                # feature_dict[str(1+1)] = embed_list[idx_list][2]
                # feature_dict[str(2+1)] = embed_list[idx_list][3]
            # print(feature_dict)
            feature_list_tmp.append(feature_dict)

        feature_list.append(feature_list_tmp)
    return feature_list
```

## 4. Experience

- From this homework, I had learned what are **NER** and **CRF model**,and how computer learn a high-level human language.
- As we know, **f1score** show our model performance, but the point is how to get higher score! This is definitely different to the typical ML model, because it is not only the input_data but also how we processing the input_data into tokens and label is decisively important !
- There is no so much parameters to let you adjust to get better model performance，so a great processing_data and completely pretrained word vector file will help your model to become better!
- Interesting topic for us to explore how great ML is in this generation.
- 雖然由於時間的關係我在features的dict中沒有成功整合不過就上述情況而言，我已經成功了解如果要訓練一個NER-CRF的模型的大概流程是怎麼樣的了。
- Thanks for Prof and TA! :)