

---

***BÁO CÁO BÀI TẬP***  
**MÔN HỌC: KHAI THÁC VĂN BẢN VÀ ỨNG DỤNG**

Giảng viên hướng dẫn: Nguyễn Trường Sơn

*Học viên: Nguyễn Huỳnh Trường Khang*

20C12006

## Mục lục

<b>1 MÔI TRƯỜNG CHẠY CHƯƠNG TRÌNH .....</b>	<b>3</b>
<b>2 BÀI 1.....</b>	<b>3</b>
2.1    Yêu cầu.....	3
2.2    Hướng dẫn.....	3
2.3    Kết quả .....	3
<b>3 BÀI 2.....</b>	<b>4</b>
3.1    Yêu cầu.....	4
3.2    Hướng dẫn.....	4
3.3    Kết quả .....	5
<b>4 BÀI 3.....</b>	<b>6</b>
4.1    Yêu cầu.....	6
4.2    Hướng dẫn.....	6
4.3    Kết quả .....	6

# 1 MÔI TRƯỜNG CHẠY CHƯƠNG TRÌNH

Các chương trình trong báo cáo được chạy với môi trường thiết lập như sau:

- Python 3.6
- Nltk 3.6.7
- Scikit-learn 0.23.2
- Sklearn-crfsuite 0.3.6

## 2 BÀI 1

### 2.1 YÊU CẦU

Xây dựng chương trình python cho phép rút trích thông tin email và số điện thoại từ văn bản.

### 2.2 HƯỚNG DẪN

Khi chạy chương trình chia làm hai trường hợp:

- Trường hợp 1: không truyền tham số vào chương trình

Chương trình yêu cầu nhập đoạn văn bản cần rút trích (có thể nhập nhiều đoạn) và xuất ra màn hình danh sách các email và số điện thoại có trong văn bản.

- Trường hợp 2: truyền tham số vào chương trình, bao gồm input là file văn bản và output là tên file chứa danh sách các email và số điện thoại có trong văn bản

Chương trình đọc văn bản cần rút trích từ input và ghi kết quả ra file output (Nếu file output chưa tồn tại thì sẽ tạo mới).

### 2.3 KẾT QUẢ

#### Trường hợp 1

*Input*

```
Enter/Paste your content. Ctrl-Z ( windows ) to save it.  
Địa chỉ: Số 28 Trần Hưng Đạo - Quận Hoàn Kiếm - Hà Nội. Điện thoại Cổng thông tin điện tử  
: (+84)24 2220 2828. Email: support@mof.gov.vn  
Công ty CP truyền thông Luật Việt Nam. Điện thoại: 093 836 1919. Email: cskh@luatvietnam.  
vn  
^Z
```

*Output*

```
Phone number list= ['(+84)24 2220 2828', '093 836 1919']  
Email list= ['support@mof.gov.vn', 'cskh@luatvietnam.vn']
```

#### Trường hợp 2

*Tập tin input\_text.txt*

```

1 <p><strong>Kuala Lumpur</strong><strong>:</strong> +60 (0)3 2723 7900</p>
2 <p><strong>Mutiar Damansara:</strong> +60 (0)3 2723 7900</p>
3 <p><strong>Penang:</strong> + 60 (0)4 255 9000</p>
4 <h2>Where we are </h2>
5 <strong>&nbsp;Call us on:</strong>&nbsp;+6 (03) 8924 8686
6 </p></div><div class="sys_two">
7 <h3 class="parentSchool">General enquiries</h3><p style="FONT-SIZE: 11px">
8 <strong>&nbsp;Call us on:</strong>&nbsp;+6 (03) 8924 8000
9 + 60 (7) 268-6200 <br />
10 Fax:<br />
11 +60 (7) 228-6202<br />
12 The standard American telephone number is ten digits, such as (555) 555-1234 .
13 0974007496 0911905428
14 Phone:</strong><strong style="color: #f00">+601-4228-8055</strong>
15 email: </strong><strong style="color: #f00">support123@gmail.com</strong>

```

#### Input

```
bai1.py input_text.txt result.txt
```

#### Output

```
The result of phone number and email extraction is saved in result.txt
```

#### Tập tin result.txt

```

1 Phone number list:
2 (7) 268-6200
3 (7) 228-6202
4 (555) 555-1234
5 723 7900
6 723 7900
7 255 9000
8 924 8686
9 924 8000
10 601-4228
11 Email list:
12 support123@gmail.com

```

## 3 BÀI 2

### 3.1 YÊU CẦU

Xây dựng chương trình python cho phép phân loại câu hỏi tiếng Việt sử dụng thuật toán SVM

### 3.2 HƯỚNG DẪN

Dữ liệu cung cấp sẵn được xử lý theo dạng: <class>:<câu hỏi> và lưu ở file *dataset.txt*

Sử dụng file *bai2\_config.py* để thiết lập các cấu hình.

Sử dụng file *bai2\_train.py* để huấn luyện mô hình theo thuật toán SVM và lưu thành tập tin tại thư mục *models*.

Sử dụng file *bai2\_main.py* để load mô hình lên và phân loại câu hỏi do người dùng nhập vào.

### 3.3 KẾT QUẢ

File *dataset.txt*

```
1  DESC.def:Con gà là con gì
2  DESC.def:Định nghĩa định lí pytago là gì
3  DESC.def:Thế nào gọi là bò sát
4  DESC.def:Thế nào gọi là gia cầm
5  DESC.def:Thế nào gọi là động vật có vú
6  DESC.def:Thế nào gọi là hiệu ứng nhà kính
7  DESC.def:Định nghĩa của môi trường sống
8  DESC.desc:Quả cam có vị chua như thế nào
9  DESC.desc:Trời hôm nay nóng như thế nào
10 DESC.desc:Com mẹ nấu ngon như thế nào
11 DESC.desc:Tỉ phú Phạm Nhật Vượng giàu như thế nào
12 DESC.desc:Giá đất đắt như thế nào
13 DESC.desc:Xe Mercedes mắc như thế nào
14 DESC.reason:Tại sao đèn xanh lại phải đi
15 DESC.reason:Tại sao mật ong lại ngọt
16 DESC.reason:Tại sao tránh uống cafe vào buổi tối
17 DESC.reason:Tại sao không nên ăn nhiều vào ban đêm
```

Chạy file *bai2\_train.py* để huấn luyện mô hình. Quá trình huấn luyện (bao gồm xây dựng và tối ưu tham số cho mô hình) sẽ được ghi nhận lại tại *bai2\_train.log* (được cấu hình mặc định trong *bai2\_config.py*)

Do tập dữ liệu ít nên cần xử lý lại, chỉ lấy theo các phân lớp lớn:

```
Class: {'SELECT', 'HUM', 'ABBR', 'LOC', 'ENTY', 'NUM', 'DESC'}
```

Kết quả đạt được với mô hình ban đầu:

Accuracy: **0.8640350877192983**

Tiếp theo, ta tiến hành tối ưu các tham số để đạt được kết quả tốt hơn:

Best score= **0.8738521752220383**

with parameter:

clf\_\_C: **2.0**

tfidf\_\_use\_idf: False

vect\_\_ngram\_range: (**1**, **2**)

Mô hình được lưu lại tại đường dẫn

The SVM model is saved in C:\models\svm\_text\_cls

Cuối cùng, ta chạy *bai2\_main.py* để load mô hình vừa lưu và nhập vào câu hỏi để kiểm tra khả năng phân loại của mô hình:

```
D:\Demo>python bai2_main.py
Enter/Paste your content. Ctrl-Z ( windows ) to save it.
Tại sao trái đất hình tròn?
Hồ Chí Minh là ai?
^Z
```

Kết quả thu được

```
Tại sao trái đất hình tròn? ==> DESC
Hồ Chí Minh là ai? ==> HUM
```

## 4 BÀI 3

### 4.1 YÊU CẦU

Xây dựng chương trình python cho nhận diện thực thể định danh (Named Entity Recogtion) sử dụng thuật toán CRF.

### 4.2 HƯỚNG DẪN

Dữ liệu cung cấp sẵn trong thư mục CoNLL-2003 với các file *train.txt*, *valid.txt* và *test.txt*  
Sử dụng file *bai3\_config.py* để thiết lập các cấu hình và cung cấp các hàm rút trích các feature của các từ.

Sử dụng file *bai3\_train.py* để huấn luyện mô hình theo thuật toán CRF và lưu thành tập tin tại thư mục *models*.

Sử dụng file *bai3\_main.py* để load mô hình lên và chạy chương trình chính.

### 4.3 KẾT QUẢ

File *train.txt*

```
1  -DOCSTART- -X- -X- O
2
3  EU NNP B-NP B-ORG
4  rejects VBZ B-VP O
5  German JJ B-NP B-MISC
6  call NN I-NP O
7  to TO B-VP O
8  boycott VB I-VP O
9  British JJ B-NP B-MISC
10 lamb NN I-NP O
11 . . O O
12
13 Peter NNP B-NP B-PER
14 Blackburn NNP I-NP I-PER
15
16 BRUSSELS NNP B-NP B-LOC
17 1996-08-22 CD I-NP O
```

*File valid.txt*

```
1  -DOCSTART- -X- -X- O
2
3  CRICKET NNP B-NP O
4  - : O O
5  LEICESTERSHIRE NNP B-NP B-ORG
6  TAKE NNP I-NP O
7  OVER IN B-PP O
8  AT NNP B-NP O
9  TOP NNP I-NP O
10 AFTER NNP I-NP O
11 INNINGS NNP I-NP O
12 VICTORY NN I-NP O
13 . . O O
14
15 LONDON NNP B-NP B-LOC
16 1996-08-30 CD I-NP O
17
```

*File test.txt*

```
1  -DOCSTART- -X- -X- O
2
3  SOCCER NN B-NP O
4  - : O O
5  JAPAN NNP B-NP B-LOC
6  GET VB B-VP O
7  LUCKY NNP B-NP O
8  WIN NNP I-NP O
9  , , O O
10 CHINA NNP B-NP B-PER
11 IN IN B-PP O
12 SURPRISE DT B-NP O
13 DEFEAT NN I-NP O
14 . . O O
15
16 Nadim NNP B-NP B-PER
17 Ladki NNP I-NP I-PER
```

Chạy file *bai3\_train.py* để huấn luyện mô hình. Quá trình huấn luyện (bao gồm xây dựng và tối ưu tham số cho mô hình) sẽ được ghi nhận lại tại *bai3\_train.log* (được cấu hình mặc định trong *bai3\_config.py*)

Ban đầu ta được huấn luyện mô hình với thông số sau:

L-BFGS optimization

c1: 0.100000

c2: 0.100000

num\_memories: 6

max\_iterations: 100

Ta thu được F1 score như sau:

F1 score: 0.8811508100893499

Ta sử dụng RandomizedSearchCV để tối ưu các tham số nhằm thu được mô hình tốt hơn:

Best score in train dataset= 0.8495688951646253

Best score in test dataset= 0.8856710914573505

with best\_params= {'c2': 0.01, 'c1': 0.0001}

Tiếp tục, ta tạo mô hình với tham số vừa tối ưu được

Create model with best params:

Và lưu tại đường dẫn

The CRF model is saved in D:\models\ner\_conll\_crf

Cuối cùng, ta chạy *bai3\_main.py* để load mô hình vừa lưu và nhập vào các câu input để chương trình xử lý

```
Enter/Paste your content. Ctrl-Z ( windows ) to save it.  
He is a German who works at Facebook Inc.  
The French President's remarks drew swift and angry condemnation as tensions rise  
over new vaccine pass  
^Z
```

Kết quả đạt được

```
He ==> 0  
is ==> 0  
a ==> 0  
German ==> B-MISC  
who ==> 0  
works ==> 0  
at ==> 0  
Facebook ==> B-ORG  
Inc ==> I-ORG  
. ==> 0  
The ==> 0  
French ==> B-MISC  
President ==> 0  
's ==> 0  
remarks ==> 0  
drew ==> 0  
swift ==> 0  
and ==> 0  
angry ==> 0  
condemnation ==> 0  
as ==> 0  
tensions ==> 0  
rise ==> 0  
over ==> 0  
new ==> 0  
vaccine ==> 0  
pass ==> 0
```