

IST 664 Natural Language Processing Homework 2

Harper He

xhe128@syr.edu

Part 1	2
Email Patterns	2
Epattern 1	2
Epattern 2	2
Epattern 3	2
Epattern 4	2
Epattern 5	3
Phone Number Patterns	3
Ppattern 1	3
Ppattern 2	3
Ppattern 3	4
Ppattern 4	4
Ppattern 5	4
Ppattern 6	4
Result	5
Part 2	5
a.	5
b.	7

Part 1

Email Patterns

Epattern 1

Expression

`([A-Za-z.]+)\sat\s([A-Za-z.]+\.[A-Za-z]+)\.EDU`

Description

This regular expression matches strings in the format: (1 or more characters) + whitespace + word “at” + whitespace + (1 or more characters + punctuation dot “.” + 1 or more characters) + punctuation dot “.” + word “EDU”.

The first parentheses here is used to capture “someone” while the second is to capture “somewhere”

Examples

This regular expression matches “uma at cs.Stanford.EDU” in the file “cheriton”.

Epattern 2

Expression

`([A-Za-z]+\sWHERE\s([A-Za-z]+\sDOM\sedu`

Description

This regular expression matches strings in the format: (1 or more characters) + whitespace + word “WHERE” + whitespace + (1 or more characters) + whitespace + word “DOM” + whitespace + word “edu”.

The first parentheses here is used to capture “someone” while the second is to capture “somewhere”

Examples

This regular expression matches “engler WHERE stanford DOM edu” in the file “engler”.

Epattern 3

Expression

`([A-Za-z]+\.[A-Za-z]+\s\([A-Za-z]+\s[A-Za-z]+\s"@([A-Za-z]+\).edu`

Description

This regular expression matches strings in the format: (1 or more characters + punctuation dot “.” + 1 or more characters) + whitespace + left parenthesis + 1 or more characters + whitespace + 1 or more characters + whitespace + punctuation quote “”” + punctuation @ + (1 or more characters) + punctuation dot “.” + word “edu”.

The first parentheses here is used to capture “someone” while the second is to capture “somewhere”

Examples

This regular expression matches [“teresa.lynn \(followed by “@stanford.edu”\)”](#) in the file “ouster”.

Epattern 4

Expression

`([A-Za-z]+\sAT\s([A-Za-z]+\sDOT\sedu`

Description

This regular expression matches strings in the format: (1or more characters) + whitespace + word “AT” + whitespace + (1or more characters) + word “DOT” + whitespace + word “edu”. The first parentheses here is used to capture “someone” while the second is to capture “somewhere”

Examples

This regular expression matches “subh AT stanford DOT edu” in the file “subh”.

Epattern 5

Expression

```
([A-Za-z]+\s)\s([A-Za-z]+\.[A-Za-z]+\.)\s\.
```

Description

This regular expression matches strings in the format: (1or more characters) + whitespace + word “at” + whitespace + (1or more characters + punctuation dot “.”+ 1or more characters) + punctuation dot “.”+ word “edu”.

The first parentheses here is used to capture “someone” while the second is to capture “somewhere”

Examples

This regular expression matches “lam at cs.stanford.edu” in the file “lam”.

Phone Number Patterns

Ppattern 1

Expression

```
\+\d\s(\d{3})\s(\d{3})\s(\d{4})
```

Description

This regular expression matches strings in the format: punctuation plus “+” + one digit + whitespace + (3 digits) + whitespace + (3 digits) + whitespace + (4 digits).

The first parentheses here is used to capture “area code”, the second is to capture “exchange part” while the third one is to capture “number part”.

Examples

One example that this regular expression matches is the number “+1 650 723 5666” and “+1 650 725 9046” in the file “jurafrsky”.

Ppattern 2

Expression

```
\[(\d{3})\]\s(\d{3})-(\d{4})
```

Description

This regular expression matches strings in the format: punctuation left bracket “[” + (3 digits) + punctuation right bracket “]” +whitespace + (3 digits) + punctuation hyphen “-” + (4 digits).

The first parentheses here is used to capture “area code”, the second is to capture “exchange part” while the third one is to capture “number part”.

Examples

One example that this regular expression matches is the number “[650] 723-5499” and “[650] 725-2472” in the file “nass”.

Ppattern 3

Expression

`\+\d\s((\d{3})\)\s(\d{3})-(\d{4})`

Description

This regular expression matches strings in the format: punctuation plus “+” + one digit + whitespace + left parenthesis “(“ + (3 digits) + right parenthesis “)” + whitespace + (3 digits) + punctuation hyphen “-” + (4 digits).

The first parentheses here is used to capture “area code”, the second is to capture “exchange part” while the third one is to capture “number part”.

Examples

One example that this regular expression matches is the number “+1 (650) 723-7683” and “+1 (650) 725-1449” in the file “manning”.

Ppattern 4

Expression

`\((\d{3})\)(\d{3})-(\d{4})`

Description

This regular expression matches strings in the format: left parenthesis “(“ + (3 digits) + right parenthesis “)” + (3 digits) + punctuation hyphen “-” + (4 digits).

The first parentheses here is used to capture “area code”, the second is to capture “exchange part” while the third one is to capture “number part”.

Examples

One example that this regular expression matches is the number “(650)814-1478” and “(650)723-1614” in the file “ashishg”.

Ppattern 5

Expression

`\((\d{3})\)\s(\d{3})-(\d{4})'`

Description

This regular expression matches strings in the format: left parenthesis “(“ + (3 digits) + right parenthesis “)” + whitespace + (3 digits) + punctuation hyphen “-” + (4 digits).

The first parentheses here is used to capture “area code”, the second is to capture “exchange part” while the third one is to capture “number part”.

Examples

One example that this regular expression matches is the number “(650) 724-6354” and “(650) 723-4539” in the file “bgirod”.

Ppattern 6

Expression

`\+\d\s(\d{3})\s(\d{3})-(\d{4})`

Description

This regular expression matches strings in the format: punctuation plus “+” + one digit + whitespace + (3 digits) + whitespace + (3 digits) + punctuation hyphen “-” + (4 digits).

The first parentheses here is used to capture “area code”, the second is to capture “exchange part” while the third one is to capture “number part”.

Examples

One example that this regular expression matches is the number “+1 650 723-3432” and “+1 650 725-1449” in the file “shoham”.

Result

The results before and after I added the expressions are listed below. The final output is 99 TP, 2 FP and 18 FN.

Expression Added (Cumulated)	Result		
	TP	FP	FN
Example Expression: ([A-Za-z.]+)@([A-Za-z.]+)\.edu ([A-Za-z.]+)s@([A-Za-z.]+)\.edu (\d{3})-(\d{3})-(\d{4})	41	0	76
Epattern 1	42	0	75
Epattern 2	43	0	74
Epattern 3	44	0	73
Epattern 4	45	0	72
Epattern 5	46	2	71
Pattern 1	48	2	69
Pattern 2	50	2	67
Pattern 3	53	2	64
Pattern 4	61	2	56
Pattern 5	97	2	20
Pattern 6	99	2	18

Part 2

a.

As the result shown above, there are 20 examples that I cannot match using regular expressions and all of them are email address.

Some examples can be matched by the regular expression on regex online tester but failed to be extracted. I list all the 18 FN examples with some regular expressions I tried but failed as well as my assumption of the reasons. For those regular expressions that didn’t work, I have no idea of the reasons because there is no extracted information for my reference.

(Please click the link for the source file of this screenshot:

https://drive.google.com/file/d/1JUQ1MVWYRvabXOBiKBa_Ftr-sipEVAu/view?usp=sharing)

There are also 2 False Positives examples, they are “Server at infolab.stanford.edu” and “Server at cs.stanford.edu”, I used the regular expression “([A-Za-z]+\s([A-Za-z]+\.[A-Za-z]+\.)\.edu” but failed to get TP results. I had no idea why they cannot match because the regular expression matched them on regex online tester, my assumption is that italic may affect the matching process, just like the strikethrough examples.

b

Based on the FN examples, we can see some characteristics of these example, for example, using word/punctuation/formatting to obscure the email addresses. Below is the email address I designed based on my observation:

nlp dot e-o_u[r]s;e @ is;ch-oo!>l dot s-%y-r dot edu
(nlp.course@ischool.syr.edu)