

Summary Report

Name:- Goutham Sri Vishwesh Bikkumalla

I gained knowledge of regular expressions and their applications through this project, which included a real-time example. I utilized regular expressions, or Reg Exp, on the Stanford University database's ContactFinder list. I chose option 3 and attempted to solve as many false negatives as I could; I was able to solve virtually all of them, with the exception of nine files.

Let me enumerate every regular expression pattern I created.

Patterns Used for the phone number: -

1. `ppatterns.append('(\\d{3})-(\\d{3})-(\\d{4})'):` -

A regular expression pattern is added to the list `ppatterns` by this line. To match a phone number in the format XXX-XXX-XXXX, where each X indicates a digit, utilize the pattern `(\\d{3})-(\\d{3})-(\\d{4})`.

`(\\d{3})` corresponds to three digits.

A literal hyphen is matched by the hyphen `-`.

Three repetitions of this pattern correspond to the three phone number segments, which are XXX-XXX-XXXX.

2. `ppatterns.append('\\(\\d{3})\\)(\\d{3})-(\\d{4})')`

This line appends a regular expression pattern to `ppatterns`.

The pattern `'\\(\\d{3})\\)(\\d{3})-(\\d{4})'` is used to match a phone number in the format of (XXX)XXX-XXXX.

`'\\('` and `'\\)'` match literal parentheses.

`(\\d{3})` matches three digits enclosed within the parentheses.

`(\\d{3})-(\\d{4})` matches the remaining part of the phone number in the format XXX-XXXX.

3. `ppatterns.append('\\(\\d{3})\\)\\s*(\\d{3})-(\\d{4})')`

This line appends a regular expression pattern to `ppatterns`.

The pattern `'\\(\\d{3})\\)\\s*(\\d{3})-(\\d{4})'` is used to match a phone number in the format of (XXX) XXX-XXXX or (XXX)XXX-XXXX.

`\\s*` matches zero or more whitespace characters.

The rest of the pattern is similar to the previous one.

4. `ppatterns.append('\\(?:\\d{3})\\)?\\s*-?\\s*(\\d{3})-(\\d{4})')`

A regular expression pattern is added to `ppatterns` by this line.

To match a phone number in formats like XXX-XXX-XXXX, (XXX)XXX-XXXX, XXX-XXX-XXXX, and XXXXXXXXXX, use the pattern `'\\(?:\\d{3})\\)?\\s*-?\\s*(\\d{3})-(\\d{4})'`. `\\-?` denotes a single or zero hyphen.

The remainder of the pattern is reminiscent of the earlier ones.

5. `ppatterns.append('(\\d{3})\\s(\\d{3})-(\\d{4})')`

The pattern `(\\d{3})\\s(\\d{3})-(\\d{4})` is used to match a phone number in the format XXX XXX-XXXX.

6. `ppatterns.append('\[(\d{3})\]\s(\d{3})-(\d{4})')`
The pattern `\[(\d{3})\]\s(\d{3})-(\d{4})'` is used to match a phone number in the format `[XXX] XXX-XXXX`.
7. `ppatterns.append('(\d{3})\s(\d{3})\s(\d{4})')`
The pattern `(\d{3})\s(\d{3})\s(\d{4})'` is used to match a phone number in the format `XXX XXX XXXX`.

Patterns Used for the Email: -

- `epatterns.append('[A-Za-z.]+@[A-Za-z.]+\.edu')`
Matches email addresses ending with ".edu".
Example: john.doe@example.edu
- `epatterns.append('[A-Za-z.]+\s@\s[A-Za-z.]+\.edu')`
Matches email addresses separated by "@", ending with ".edu".
Example: jane_smith@example.edu
- `epatterns.append('[a-zA-z.]+\sat\s[a-zA-Z]+\.[a-zA-Z.]+\.EDU')`
Matches email addresses with "at" and ".edu" spelled out.
Example: aliceatexample.com.edu
- `epatterns.append('[A-Za-z]+\sWHERE\s[A-Za-z]+\sDOM\sedu')`
Matches email addresses with "WHERE" and "DOM" spelled out, ending with ".edu".
Example: johndoe WHERE example DOM edu
- `epatterns.append('[A-Za-z.]+@[A-Za-z.]+\.edu')`
Matches email addresses with "@" representing "@" and ending with ".edu".
Example: info@example.edu
- `epatterns.append('[A-Za-z]+\s*at\s*!--.*?-->\s*[A-Za-z]+\s*!--.*?-->\s*dot\s*!--.*?-->\s*edu')`
Matches email addresses with "at" and "dot" spelled out, embedded in HTML comments, and ending with ".edu".
Example: user1<!--at-->example<!--dot-->edu
- `epatterns.append('[A-Za-z.]+\s<at symbol>\s[A-Za-z.]+\.edu')`
Matches email addresses with "<at symbol>" representing "@" and ending with ".edu".
Example: support <at symbol> example.edu
- `epatterns.append('\([A-Za-z.]+\s*\s*[A-Za-z.]+\.edu\)')`
Matches email addresses enclosed in parentheses, separated by "@" and ending with ".edu".
Example: (info@example.edu)
- `epatterns.append('[A-Za-z.]+\s*(?:@|AT)\s*[A-Za-z.]+\s*(?:\.|DOT)\s*edu')`
Matches email addresses with variations of "@" and "dot" spelled out, ending with ".edu".
Example: contact AT example DOT edu
- `epatterns.append('email to ([A-Za-z.]+) at ([A-Za-z.]+) dt com')`
Matches email addresses formatted as "email to <username> at <domain> dt com".
Example: email to info at example dt com
- `epatterns.append('E-mail: ([A-Za-z.]+) at ([A-Za-z.]+\.edu')`
Matches email addresses prefixed with "E-mail:" and separated by "at", ending with ".edu".


```

res.append((name,'e',email))

# phone pattern list
for ppat in ppatterns:
    # each ppat has 3 sets of parentheses so each match will have 3
items in a list
    matches = re.findall(ppat,line)
    for m in matches:
        # phone number has form areacode-exchange-number
        #phone = '%s-%s-%s' % m
        phone = '{}-{}-{}'.format(m[0],m[1],m[2])
        res.append((name,'p',phone))

return res

```

Output Screenshot :-

```

Assuming ContactFinder.py called in directory with data folder
True Positives (189):
[('ashishg', 'e', 'ashishg@stanford.edu'),
('ashishg', 'e', 'roze@stanford.edu'),
('ashishg', 'p', '650-723-1614'),
('ashishg', 'p', '650-723-4173'),
('ashishg', 'p', '650-814-1478'),
('balaji', 'e', 'balaji@stanford.edu'),
('bgirod', 'p', '650-723-4539'),
('bgirod', 'p', '650-724-3648'),
('bgirod', 'p', '650-724-6354'),
('cheriton', 'e', 'cheriton@cs.stanford.edu'),
('cheriton', 'e', 'uma@cs.stanford.edu'),
('cheriton', 'p', '650-723-1131'),
('cheriton', 'p', '650-725-3726'),
('dabo', 'e', 'dabo@cs.stanford.edu'),
('dabo', 'p', '650-725-3897'),
('dabo', 'p', '650-725-4671'),
('engler', 'e', 'engler@lcs.mit.edu'),
('engler', 'e', 'engler@stanford.edu'),
('eroberts', 'e', 'eroberts@cs.stanford.edu'),
('eroberts', 'p', '650-723-3642'),
('eroberts', 'p', '650-723-6092'),
('fedkin', 'e', 'fedkin@cs.stanford.edu'),
('hager', 'e', 'hager@cs.jhu.edu'),
('hager', 'p', '410-516-5521'),
('hager', 'p', '410-516-5553'),
('hager', 'p', '410-516-8000'),
('hanrahan', 'e', 'hanrahan@cs.stanford.edu'),
('hanrahan', 'p', '650-723-0033'),
('hanrahan', 'p', '650-723-8530'),
('horowitz', 'p', '650-725-3707'),
('horowitz', 'p', '650-725-6949'),
('jurafsky', 'p', '650-723-6666'),
('kosecka', 'e', 'kosecka@cs.gmu.edu'),
('kosecka', 'p', '703-993-1710'),
('kosecka', 'p', '703-993-1876'),
('kunle', 'e', 'darlene@csl.stanford.edu'),
('kunle', 'e', 'kunle@ogun.stanford.edu'),
('kunle', 'p', '650-723-1430'),

```

```

('kunle', 'p', '650-723-1430'),
('kunle', 'p', '650-725-3713'),
('kunle', 'p', '650-725-6949'),
('lam', 'e', 'lam@cs.stanford.edu'),
('lam', 'p', '650-725-3714'),
('lam', 'p', '650-725-6949'),
('latombe', 'p', '650-721-6625'),
('latombe', 'p', '650-723-0350'),
('latombe', 'p', '650-723-4137'),
('latombe', 'p', '650-725-1449'),
('levoy', 'e', 'ada@graphics.stanford.edu'),
('levoy', 'e', 'melissa@graphics.stanford.edu'),
('levoy', 'p', '650-723-0033'),
('levoy', 'p', '650-724-6865'),
('levoy', 'p', '650-725-3724'),
('levoy', 'p', '650-725-4089'),
('manning', 'e', 'dbarros@cs.stanford.edu'),
('manning', 'e', 'manning@cs.stanford.edu'),
('manning', 'p', '650-723-7683'),
('manning', 'p', '650-725-1449'),
('manning', 'p', '650-725-3358'),
('nass', 'e', 'nass@stanford.edu'),
('nass', 'p', '650-723-5499'),
('nass', 'p', '650-725-2472'),
('nick', 'e', 'nick.parlante@cs.stanford.edu'),
('nick', 'p', '650-725-4727'),
('ok', 'p', '650-723-9753'),
('ok', 'p', '650-725-1449'),
('pal', 'e', 'pal@cs.stanford.edu'),
('pal', 'p', '650-725-9046'),
('psyoung', 'e', 'patrick.young@stanford.edu'),
('rajeev', 'p', '650-723-4377'),
('rajeev', 'p', '650-723-6045'),
('rajeev', 'p', '650-725-4671'),
('rinard', 'e', 'rinard@lcs.mit.edu'),
('rinard', 'p', '617-253-1221'),
('rinard', 'p', '617-258-6922'),
('serafim', 'e', 'serafim@cs.stanford.edu'),
('serafim', 'p', '650-723-3334'),
('serafim', 'p', '650-725-1449'),

```

```

('shoham', 'p', '650-723-3432'),
('shoham', 'p', '650-725-1449'),
('subh', 'e', 'subh@stanford.edu'),
('subh', 'e', 'uma@cs.stanford.edu'),
('subh', 'p', '650-724-1915'),
('subh', 'p', '650-725-3726'),
('subh', 'p', '650-725-6949'),
('thm', 'e', 'pkrokel@stanford.edu'),
('thm', 'p', '650-725-3383'),
('thm', 'p', '650-725-3636'),
('thm', 'p', '650-725-3938'),
('tim', 'p', '650-724-9147'),
('tim', 'p', '650-725-2348'),
('tim', 'p', '650-725-4671'),
('ullman', 'e', 'support@gradiance.com'),
('ullman', 'e', 'ullman@cs.stanford.edu'),
('ullman', 'p', '650-494-8816'),
('ullman', 'p', '650-725-2888'),
('ullman', 'p', '650-725-4882'),
('vladlen', 'e', 'vladlen@stanford.edu'),
('widom', 'e', 'siroker@cs.stanford.edu'),
('widom', 'e', 'widom@cs.stanford.edu'),
('widom', 'p', '650-723-0872'),
('widom', 'p', '650-723-7690'),
('widom', 'p', '650-725-2588'),
('zelenski', 'e', 'zelenski@cs.stanford.edu'),
('zelenski', 'p', '650-723-6092'),
('zelenski', 'p', '650-725-8596'),
('zm', 'e', 'manna@cs.stanford.edu'),
('zm', 'p', '650-723-4364'),
('zm', 'p', '650-725-4671'))
False Positives (2):
('nass', 'p', '375-742-1353')
gold: ('nass', 'e', 'nass@stanford.edu')
gold: ('nass', 'p', '650-723-5499')
gold: ('nass', 'p', '650-725-2472')
('nass', 'p', '114-910-7699')
gold: ('nass', 'e', 'nass@stanford.edu')
gold: ('nass', 'p', '650-723-5499')
gold: ('nass', 'p', '650-725-2472')

```

```

('widom', 'p', '650-723-0872'),
('widom', 'p', '650-723-7690'),
('widom', 'p', '650-725-2588'),
('zelenski', 'e', 'zelenski@cs.stanford.edu'),
('zelenski', 'p', '650-723-6092'),
('zelenski', 'p', '650-725-8596'),
('zm', 'e', 'manna@cs.stanford.edu'),
('zm', 'p', '650-723-4364'),
('zm', 'p', '650-725-4671'))}

```

False Positives (2):

```

('nass', 'p', '375-742-1353')
gold: ('nass', 'e', 'nass@stanford.edu')
gold: ('nass', 'p', '650-723-5499')
gold: ('nass', 'p', '650-725-2472')
('nass', 'p', '114-910-7699')
gold: ('nass', 'e', 'nass@stanford.edu')
gold: ('nass', 'p', '650-723-5499')
gold: ('nass', 'p', '650-725-2472')

```

False Negatives (8):

```

{('dlwh', 'e', 'dlwh@stanford.edu'),
('jks', 'e', 'jks@robotics.stanford.edu'),
('jurafsky', 'e', 'jurafsky@stanford.edu'),
('latombe', 'e', 'asandra@cs.stanford.edu'),
('latombe', 'e', 'latombe@cs.stanford.edu'),
('latombe', 'e', 'liliana@cs.stanford.edu'),
('ouster', 'e', 'ouster@cs.stanford.edu'),
('ouster', 'e', 'teresa.lynn@stanford.edu'))}

```

Summary: tp=109, fp=2, fn=8

C:\Users\JYOTHIKA\Desktop\ContactFinder (2) (1)\ContactFinder>

Observations and learnings-

Regular expressions are frequently employed in data extraction and parsing. We were able to retrieve emails and phone numbers from a significant amount of contested data, as demonstrated by our instances.

➤ Text searching and manipulation are possible with Regular Expressions (Regex). Reg exp allows us to easily manipulate texts and use them as needed.

➤ Data validation: Reg exp programs are employed to validate data. It is typically used to confirm whether or not the user-provided inputs are in the correct format. It is easy to check for the sequence or patterns we want.