

CISS445: Programming Languages Assignment a07

OBJECTIVES

- Design DFA and NFA
- Write regular expressions using Perl-style regular expression syntax

Mathematically given a set of symbols S , a regular expression is a string made up of characters from S together with \cup (union), $($ (left parenthesis), $)$ (right parenthesis), $*$ (Kleene star), ϵ (empty string), \emptyset (empty set).

Although these characters are enough to mathematically construct regular expressions (and hence any regular language), programmers have found this regular expression syntax cumbersome. Different programming languages and different third party software (classes, modules, packages, etc) have included additional symbols to help in writing regular expressions more compactly. One such syntax is the Perl regular expression syntax. The objective of this assignment is to learn to write regular expression using Perl's regular expression syntax. Note that this is not the only regular expression syntax. In your professional life (in academia or industry) you will come across different regular expression engines.

The following includes a quick introduction to Python. You will need to read the documents listed below for more information. Note that I am using Python 3.7.9

HOW TO SPECIFY DFA AND NFA IN THE QUESTIONS

Suppose for a question, you have designed an NFA, say you call it N . The states are A, B, C, D, E . Suppose your start state is A and C, D are accept states. Suppose further that the alphabet of your N is $\{0, 1\}$. In your N , you have a transition from A to B labeled 0, another transition from A to A labeled 1, and a transition from B to C labeled ϵ . Then your NFA, N , is described in `main.py` as follows:

```
# Name: John Doe
# File: main.py

N = NFA(alphabet=["0", "1"],
        states=["A", "B", "C", "D", "E"],
        start="A",
        accepts=["C", "D"],
        transitions=[("A", "0", "B"),
                     ("A", "1", "A"),
                     ("B", "", "C")])
```

Note in particular you should use "" as ϵ .

Suppose for another question, you have designed a DFA, say you call it M . For instance if your DFA has only two states "q0" and "q1" where "q0" is the start state and "q1" is the accept state, the alphabet is "0" and "1". As for transitions, "q0" transitions to "q0" on "0", "q0" transitions to "q1" on "1", "q1" transitions to "q0" on "0", and "q1" transitions to "q1" on "1", then your DFA is described follows in `main.py`:

```
# Name: John Doe
# File: main.py

M = DFA(alphabet=["0", "1"],
        states=["q0", "q1"],
        start="q0",
        accepts=["q1"],
        transitions=[("q0", "0", "q0"),
                     ("q0", "1", "q1"),
                     ("q1", "0", "q0"),
                     ("q1", "1", "q1")])
```

WARNING: This is a DFA. Therefore out of every state, there are exactly 2 transitions, one labeled 0 and one labeled 1.

Q1. [Design NFA]

Design an NFA N that accept strings made up of symbols 0 and 1 where each string contains an even number of 0s or odd number of 1s.

Complete the following skeleton:

```
# Name: John Doe
```

```
# File: main.py
```

```
N = NFA(alphabet=[],  
        states=[],  
        start=None,  
        accepts=[],  
        transitions=[])
```

Q2. [Design DFA]

Design a DFA M that accept strings made up of symbols 0 and 1 where each string contains an even number of 0s or odd number of 1s.

Complete the following skeleton:

```
# Name: John Doe
# File: main.py

M = DFA(alphabet=[],
        states=[],
        start=None,
        accepts=[],
        transitions=[])
```

PROGRAMMING TOOL

To test your regular expressions, we will use the Python programming language. You should use one of my fedora virtual machines.

```
$ Python 3.7.9 (default, Aug 19 2020, 17:05:11)
[GCC 9.3.1 20200408 (Red Hat 9.3.1-2)] on linux
Type "help", "copyright", "credits" or "license" for more information.
>>>
```

(You can also download other Python implementation from Python's web site: <http://www.python.org>. If you're using the Windows version, you can get to the same Python shell when you select the menu item "Python (command line)" from your Start menu under Python.)

If you have a Linux box, Python is probably already installed. You can also run the python interpreter by typing `python` at the shell prompt.

After spring 2022, make sure you use Python 3.7. You can find all the releases for Python at <http://www.python.org>. Choose the right one.

After running the Python interpreter, you can issue Python statements at the Python prompt in the Python shell:

```
>>> 1 + 1
2
>>>
```

Note that the text in bold is entered by the user. Note that Python responded by printing the resulting value. You can also issue a `print` command if you like but the two have the same effect:

```
>>> print(1 + 1)
2
>>>
```

You can write your Python statements in a file, say `test.py`:

```
print(1 + 1)
```

After saving it, you can run it at your shell prompt by typing:

```
python test.py
```

Your shell will run the program with python and respond with:

2

For more information the Python programming language refer to the tutorials and guide at <http://www.python.org>. Here's one: <http://docs.python.org/tutorial/index.html>.

PYTHON STRINGS

Let's talk about Python strings. Python strings have the "usual" behavior. For instance you can concatenate strings. Run this:

```
s = "abc"
t = "def"
u = s + t
print(u)
```

Of course you can cut out a piece of the string. Run this:

```
s = "abcdef"
t = s[2]
print(t)
```

What do you get? Python does not have the concept of characters in the sense of C/C++. A character in Python is actually just a string of length 1. Therefore the variable `t` above is a string.

Somewhat more surprising is this (if you have not seen this before):

```
s = "abcdef"
t = s[1:3]
print(t)
```

Instead of `1:3` try different values such as `0:2`, `1:4`, etc.

Python strings are immutable. In other words you cannot change the contents of the string. For instance if you want to change the first character of a string to `x`, then you basically have to build a new string like this. Try this:

```
s = "abcdef"
t = "x" + s[1:]
print(t)
```

You can of course compare strings. Try this:

```
s = "abcdef"
t = "abcdef"
u = "abcdez"
print(s == t)
print(s == u)
```

And of course we must have the string length function:

```
s = ""  
t = "a"  
u = "abcdef"  
print(len(s), len(t), len(u))
```

Since python does not have the concept of characters, the single quote is also used to delimit strings:

```
s = ''  
t = 'a'  
u = 'abcdef'  
print(len(s), len(t), len(u))
```

For more on Python strings, refer to the tutorial stated above.

PYTHON FUNCTIONS

Writing a Python function is easy. Here's an example. Run it:

```
def f():  
    return 42  
  
print(f())
```

And here's another. Run it:

```
def g(x):  
    y = "you gave me " + x  
    return y  
  
print(g("hello world"))
```

And another. Run it:

```
def h(x):  
    if x < 0:  
        a = 0  
        b = 1  
        c = 2  
    elif x == 0:  
        a = 1  
        b = 2  
        c = 3  
    else:  
        a = 2  
        b = 3  
        c = 4  
    return a + b + c  
  
print(h(-1))  
print(h(0))  
print(h(2))
```

Note that Python determines blocks by whitespaces. The above also show you how to write branching statements. As for the for-loop, try this:

```
for a in [1,2,3,4,5]:  
    print(a)  
  
for a in range(1, 10):  
    print(a)  
  
for a in range(1, 10, 2):
```

```
print(a)
```

And for the while-loop try this:

```
s = input('gimme something ... ')
while s != '':
    print(s)
    s = input('gimme something ... ')
```

PYTHON REGULAR EXPRESSIONS

For more information, you will need to study <https://docs.python.org/3/howto/regex.html>.

Python comes with a regular expression module. With the library you can build a regular expression object. With this object you can carry out “matching” and “searching”: with matching you will get a match object and with searching you will get a search object. We’ll talk about matching first.

EXAMPLES. The following example shows you how to use the regular expression module in Python. Run this:

```
>>> import re
>>> p = re.compile("abc")
>>> p.match("abc")
<re.Match object; span=(0, 3), match='abc'>
>>> p.match("xyz")
>>>
```

In the above example `p` is a regular expression object corresponding to the regular expression `"abc"`.

Now we test the regular expression against some strings.

`p.match("abc")` returns an object `<re.Match object; span=(0, 3), match='abc'>` telling you that the string `"abc"` matches the regular expression object `"abc"`.

On the other hand when you test the string `"xyz"` against your regular expression, you get nothing. This tells you that `"xyz"` does not match the regular expression of `p`.

However there is a difference between Python’s regular expression match functionality and the mathematical regular expression matching that you should be aware of. The regular expression object’s `match()` method actually matches leftmost substrings in the test string. For instance try this (continuing the above example):

```
>>> p.match("abcxyz")
<re.Match object; span=(0, 3), match='abc'>
```

As you can see, `"abcxyz"` matches `"abc"` in the sense that `"abc"` appears as a leftmost substring in the test string `"abcxyz"`.

So what if do not want to match leftmost substrings - you want exact match? You

can force matching of end-of-string with `\Z`, the end-of-string syntax. In other words you use the regular expression `"abc\Z"`. Continuing the above example, let's build another regular expression object:

```
>>> p1 = re.compile("abc\Z")
>>> p1.match("abc")
<re.Match object; span=(0, 3), match='abc'>
>>> p1.match("abcxyz")
```

Note that now `"abcxyz"` does not match the regular expression `"abc\Z"`. That's because `x` does not match `\Z` since `\Z` is the end-of-string marker.

(Note that you only need to do `"import re"` once in your Python shell. You don't have to do it again in this example. The only time you need to do that is when you start the Python shell.)

Now let's look closely at a match object. Do this again:

```
>>> p = re.compile("abc")
>>> match = p.match("abcdef")
```

This gives you a match object, `match`. You can also view the attributes (including the methods) of the `match` object:

```
>>> print(dir(match))
['__class__', '__copy__', '__deepcopy__', '__delattr__', '__dir__', '__doc__',
 '__eq__', '__format__', '__ge__', '__getattribute__', '__getitem__', '__gt__',
 '__hash__', '__init__', '__init_subclass__', '__le__', '__lt__', '__ne__',
 '__new__', '__reduce__', '__reduce_ex__', '__repr__', '__setattr__', '__sizeof__',
 '__str__', '__subclasshook__', 'end', 'endpos', 'expand', 'group', 'groupdict',
 'groups', 'lastgroup', 'lastindex', 'pos', 're', 'regs', 'span', 'start',
 'string']
```

You print the string that was matched:

```
>>> print(match.group())
abc
```

You can print the starting index and ending index of the match within the string `"abcdef"`:

```
>>> print(match.start())
0
>>> print(match.end())
```

3

Refer to Python documentation for more information.

Now for searching. While matching attempts to find a pattern starting at the beginning of a string, `search` tries to match any substring. Try this

```
>>> p2 = re.compile("abc")
>>> search = p2.search("12 34abc56789")
>>> print(search)
<re.Match object; span=(5, 8), match='abc'>
```

This tells you that the regular expression "abc" is found in the string "12 34abc56789", not necessarily at the beginning. We can find the substring within "12 34abc56789" that matches our regular expression "abc" and we can also find the starting index and ending index of this substring within "12 34abc56789":

```
>>> print(search.group())
abc
>>> print(search.start())
5
>>> print(search.end())
8
```

As you can see the "abc" is found within "12 34abc56789" at starting index position 5 and ending index position 8.

PERL REGULAR EXPRESSION SYNTAX

There are special characters to denote special meanings. For instance you will see that `[` is used to mean something. But what if you want to match the character `[`? You use the escape character, i.e. you use `\[`.

The only time when you do not need to use the escape character `\` is in `[...]`. The `[...]` notation is similar to the set notation. You can also specify a range of values in `[...]`. Here are some examples:

```
[abcd]    matches 'a' or 'b' or 'c' or 'd'.
[a-d]     matches 'a' to 'd' (the order is taken from the ASCII table).
[0-9a-z]  matches '0' to '9' or 'a'-'z'

[^abc]    matches everything except a,b,c
```

All these special characters used in building a regular expression in Python actually comes from another language: Perl. (There are many other regular expression syntax, but Perl is probably the most famous.)

Therefore the regular expression `"[a-d][0-9]\Z"` should match strings such as `"b0"`, `"d9"`:

```
>>> p2 = re.compile("[a-d][0-9]\Z")
>>> p1.match("ad")
>>> p1.match("a0")
```

The regular expression `"[a-d][^a-z]"` should match `"a0"` but not `"ab"`:

```
>>> p3 = re.compile("[a-d][^a-z]")
>>> p3.match("a0")
<re.Match object; span=(0, 2), match='a0'>
>>> p3.match("ab")
```

Getting the hang of it yet?

```
a*      matches any number of a (including none)
a+      matches any positive number of a. This is the same as aa*
a{4}    same as aaaa
a{2,4}  aa or aaa or aaaa
a{2,}   aa or aaa or aaaa or aaaaa or aaaaaa or ....
a?      same as empty string or a, i.e. optional a. This is the same as a{0, 1}
```

Try these:

```
>>> p4 = re.compile("[01]{1,2}\Z")
>>> p4.match("0")
<re.Match object; span=(0, 1), match='0'>
>>> p4.match("1")
<re.Match object; span=(0, 1), match='1'>
>>> p4.match("00")
<re.Match object; span=(0, 2), match='00'>
>>> p4.match("01")
<re.Match object; span=(0, 2), match='01'>
>>> p4.match("10")
<re.Match object; span=(0, 2), match='10'>
>>> p4.match("11")
<re.Match object; span=(0, 2), match='11'>
>>> p4.match("000")
>>> p4.match("101")
>>>
```

And this:

```
>>> p5 = re.compile("[01]{2,}\Z")
>>> p5.match("0")
>>> p5.match("00")
<re.Match object; span=(0, 2), match='00'>
>>> p5.match("000")
<re.Match object; span=(0, 3), match='000'>
>>> p5.match("0000")
<re.Match object; span=(0, 4), match='0000'>
>>>
```

You should try a few search examples, printing the `group()`, `start()`, and `end()`.

This is not too surprising by now:

```
a|b          match a or b
```

Try this:

```
>>> p6 = re.compile("(ab)|(de)\Z")
>>> p6.match("ab")
<re.Match object; span=(0, 2), match='ab'>
>>> p6.match("de")
<re.Match object; span=(0, 2), match='de'>
```

Note that matches are *greedy*. For instance the string "aaaaa" will not match the regular expression "a{3,5}aa\Z" since "aaaaa" will match "a{3,5}" leaving nothing for "aa\Z".

FAQ

Q: “What about epsilon?”

A: Use an empty substring. Study the following example:

```
>>> import re
>>> p = re.compile("abc(1|2|)def")
>>> p.search("a")
>>> print(p.search("a"))
None
>>> print(p.search("ab"))
None
>>> print(p.search("abc"))
None
>>> print(p.search("abc1"))
None
>>> print(p.search("abc1d"))
None
>>> print(p.search("abc1de"))
None
>>> print(p.search("abc1def"))
<re.Match object; span=(0, 7), match='abc1def'>
>>> print(p.search("abc2def"))
<re.Match object; span=(0, 7), match='abc2def'>
>>> print(p.search("abcdef"))
<re.Match object; span=(0, 6), match='abcdef'>
>>>
```


Q3. [Regex using PERL Syntax]

Write down the mathematical description of a regular expression for strings representing polynomials with integer coefficients. The following are examples of such strings:

- 0
- -2
- 2
- 1+x
- 1 + x
- 1 + x
- x + 1
- x+1
- -1+x³
- 1-x³
- 1+x+x+x
- 1+x²³+x³ - 42 x¹⁰⁰
- 0x¹⁰⁰ + -1x
- x⁰

Integer coefficients and powers must be integers. An integer is either 0 or a non-zero digit followed by any number of digits. The following are strings which should not be accepted:

- x[^]
- ^3
- 2^{^2}
- x^{^x}

Here's a skeleton (skeleton code is to help you – it is not meant to be correct or error-free):

```
# Name: John Doe
# File: main.py

import re

p = re.compile("x*Z") // replace string with correct regex
s = input()
print(p.match(s))
```

(The regular expression is for any number of x. Therefore it's not correct although it does match "x" correctly.)

Q3. [Regex using PERL Syntax]

Write a regular expression for strings of 0's and 1's where all maximal substrings of 0's have even length. For instance

- 0000
- 1001100111
- 0011

are accepted by the regular expression. But the following are not:

- 000
- 1001100011
- 100110010

Skeleton code:

```
# Name: John Doe
# File: main.py

import re

p = re.compile("x*Z") // replace string with correct regex
s = input()
print(p.match(s))
```

Q4. [Extracting data from XML using regex]

An XML string is a string in a special format. Here's an XML fragment:

```
<firstname>John</firstname>
Here's another:
<lastname>Doe</lastname>
```

Here's another

```
<person>
  <firstname>John</firstname>
  <lastname>Doe</lastname>
</person>
```

or maybe

```
<person id=1231235>
  <firstname>John</firstname>
  <lastname>Doe</lastname>
</person>
```

XML data provides “meaning” and “structure” to otherwise flat data without meaning.

For this question, you will need to visit <https://openweathermap.org/>. Go ahead and create an account at that website. (You can also go to wikipedia and read up on openweathermap.) After confirming the account, login and click on API (near the top) and then subscribe to the weather data. You should get an email from the website that looks like this:

```
Dear Customer!

Thank you for subscribing to Free OpenWeatherMap!

API key:
- Your API key is 014f77b0g63693crffcbhab9dfvb7903
- Within the next couple of hours, it will be activated and ready to use
- You can later create more API keys on your account page
- Please, always use your API key in each API call

[etc.]
```

You will need the above API key.

Wait for a couple of minutes. Make sure you are connected to the internet. Then

run the following python program:

```
import urllib.request

city_id = 4575352
api_key = "27496f4d8de865143282c78b41979cc0"
url = "http://api.openweathermap.org/data/2.5/weather?id=%s&mode=xml&units=imperial&APPID=%s"
url = url % (city_id, api_key)
xml = urllib.request.urlopen(url).read()
xml = xml.decode("utf-8")
print(xml)
```

You will see this:

```
<?xml version="1.0" encoding="UTF-8"?>
<current>
<city id="4575352" name="Columbia"><coord lon="-81.0348" lat="34.0007"></coord>
<country>US</country><timezone>-14400</timezone>
<sun rise="2022-03-22T11:25:04" set="2022-03-22T23:36:51"></sun></city>
<temperature value="53.53" min="51.04" max="55.35" unit="fahrenheit"></temperature>
<feels_like value="51.64" unit="fahrenheit"></feels_like>
<humidity value="65" unit="%"></humidity><pressure value="1023" unit="hPa"></pressure>
<wind><speed value="4.61" unit="mph" name="Light breeze"></speed><gusts></gusts>
<direction value="180" code="S" name="South"></direction></wind>
<clouds value="0" name="clear sky"></clouds><visibility value="10000"></visibility>
<precipitation mode="no"></precipitation><weather number="800" value="clear sky" icon="01n"></weather>
<lastupdate value="2022-03-22T05:44:58"></lastupdate>
</current>
```

or after formatting:

```
<?xml version="1.0" encoding="UTF-8"?>
<current>
  <city id="4575352" name="Columbia">
    <coord lon="-81.03" lat="34">
      </coord>
    <country>US</country>
    <timezone>-18000</timezone>
    <sun rise="2020-03-07T11:44:40" set="2020-03-07T23:25:37">
      </sun>
    </city>
    <temperature value="47.82" min="43" max="53.6" unit="fahrenheit">
      </temperature>
    <feels_like value="41.41" unit="fahrenheit">
      </feels_like>
    <humidity value="43" unit="%"></humidity>
    <pressure value="1030" unit="hPa"></pressure>
    <wind>
      <speed value="3.71" unit="mph" name="Light breeze"></speed>
      <gusts></gusts>
      <direction value="32" code="NNE" name="North-northeast">
        </direction>
      </wind>
    <clouds value="1" name="clear sky"></clouds>
    <visibility value="16093"></visibility>
    <precipitation mode="no"></precipitation>
    <weather number="800" value="clear sky" icon="01n"></weather>
    <lastupdate value="2020-03-08T00:17:16"></lastupdate>
  </current>
```

Write a Python program so that when you run it, it displays the above information (which should be extracted using regular expressions). For instance the above XML data will result in the following output:

```
City: Columbia
Country: US
Sun rise: 2020-03-07T11:44:40
Sun set: 2020-03-07T23:25:37
Temperature: 47.82, min 43, max 53.6 fahrenheit
Feels like: 41.41 fahrenheit
Humidity: 43%
Pressure: 1030 hPa
Wind speed: 3.71 mph, Light breeze
Wind direction: North-northeast
Clouds: clear sky
Visibility: 16093
Precipitation: no
Last update: 2020-03-08T00:17:16
```

Of course if you run the program at different times, you will get different results.

Here's a skeleton:

```
# Name: John Doe
# File: main.py

import urllib.request
import re

city_id = "4575352"
api_key = "014f77b0g63693crffcbhab9dfvb7903"
url = "http://api.openweathermap.org/data/2.5/weather?id=%s&mode=xml&units=imperial&APPID=%s"
url = url % (city_id, api_key)
xml = urllib.request.urlopen(url).read()
xml = xml.decode("utf-8")

city = ""
country = ""
sun_rise = ""

print("City:", city)
print("Country:", country)
print("Sun rise:", sun_rise)
```

Your goal is to, of course, extract the city from the string xml using a regular expression:

```
# Name: John Doe
# File: main.py

import urllib.request
import re

city_id = "4575352"
api_key = "014f77b0g63693crffcbhab9dfvb7903"
url = "http://api.openweathermap.org/data/2.5/weather?id=%s&mode=xml&units=imperial&APPID=%s"
url = url % (city_id, api_key)
```

```
xml = urllib.request.urlopen(url).read()
xml = xml.decode("utf-8")

...
city = ""
country = ""
sun_rise = ""

print("City:", city)
print("Country:", country)
print("Sun rise:", sun_rise)
```

Then you need to extract the country.

```
# Name: John Doe
# File: main.py

import urllib.request
import re

city_id = "4575352"
api_key = "014f77b0g63693crffcbhab9dfvb7903"
url = "http://api.openweathermap.org/data/2.5/weather?id=%s&mode=xml&units=imperial&APPID=%s"
url = url % (city_id, api_key)
xml = urllib.request.urlopen(url).read()
xml = xml.decode("utf-8")

...
city = ...
...
country = ...
sun_rise = ""

print("City:", city)
print("Country:", country)
print("Sun rise:", sun_rise)
```

Etc. (It's even better if you can write one single regex to extract all the data.)

Once your program is working, change it so that it works for any city id:

```
# Name: John Doe
# File: main.py

import urllib.request
import re

city_id = input("enter city id: ")
api_key = "014f77b0g63693crffcbhab9dfvb7903"
url = "http://api.openweathermap.org/data/2.5/weather?id=%s&mode=xml&units=imperial&APPID=%s"
url = url % (city_id, api_key)
xml = urllib.request.urlopen(url).read()
xml = xml.decode("utf-8")

...
```