# Regular Expressions

2020年7月8日　17:12

| | |
|---|---|
| **^** | Matches the beginning of a line |
| **$** | Matches the end of the line |
| **.** | Matches any character |
| **\s** | Matches whitespace |
| **\S** | Matches any non-whitespace character |
| **\*** | Repeats a character zero or more times |
| **\*?** | Repeats a character zero or more times (non-greedy) |
| **+** | Repeats a character one or more times |
| **+?** | Repeats a character one or more times (non-greedy) |
| **[aeiou]** | Matches a single character in the listed set |
| **[^XYZ]** | Matches a single character $not$ in the listed set |
| **[a-z0-9]** | The set of characters can include a range |
| **(** | Indicates where string extraction is to start |
| **)** | Indicates where string extraction is to end |

# re.search()

Return True or False

## Using re.search() Like find()

```python
hand = open('mbox-short.txt')
for line in hand:
    line = line.rstrip()
    if line.find('From:') >= 0:
        print(line)
```

```python
import re

hand = open('mbox-short.txt')
for line in hand:
    line = line.rstrip()
    if re.search('From:', line) :
        print(line)
```

## Using re.search() Like startswith()

```python
hand = open('mbox-short.txt')
for line in hand:
    line = line.rstrip()
    if line.startswith('From:'):
        print(line)
```

```python
import re

hand = open('mbox-short.txt')
for line in hand:
    line = line.rstrip()
    if re.search('^From:', line) :
        print(line)
```

# re.findall()

Extract data

```
>>> import re
>>> x = 'My 2 favorite numbers are 19 and 42'
>>> y = re.findall('[0-9]+',x)
>>> print(y)
['2', '19', '42']
```

# Fine-Tuning String Extraction

Parentheses are not part of the match - but they tell where to start and stop what string to extract

`From stephen.marquard@uct.ac.za Sat Jan  5 09:14:16 2008`

```
>>> y = re.findall('\S+@\S+',x)
>>> print(y)
['stephen.marquard@uct.ac.za']
>>> y = re.findall('^From (\S+@\S+)',x)
>>> print(y)
['stephen.marquard@uct.ac.za']
```

`^From  (\S+@\S+)`

# Spam Confidence

```python
import re
hand = open('mbox-short.txt')
numlist = list()
for line in hand:
    line = line.rstrip()
    stuff = re.findall('^X-DSPAM-Confidence: ([0-9.]+)', line)
    if len(stuff) != 1 :  continue
    num = float(stuff[0])
    numlist.append(num)
print('Maximum:', max(numlist))
```

`X-DSPAM-Confidence: 0.8475`

```
python ds.py
Maximum: 0.9907
```

# Escape Character

If you want a special regular expression character to just behave normally (most of the time) you prefix it with '\'

```
>>> import re
>>> x = 'We just received $10.00 for cookies.'
>>> y = re.findall('\$[0-9.]+',x)
>>> print(y)
['$10.00']
```

At least one or more

`\$[0-9.]+`

A real dollar sign     A digit or period

# Networks and Sockets
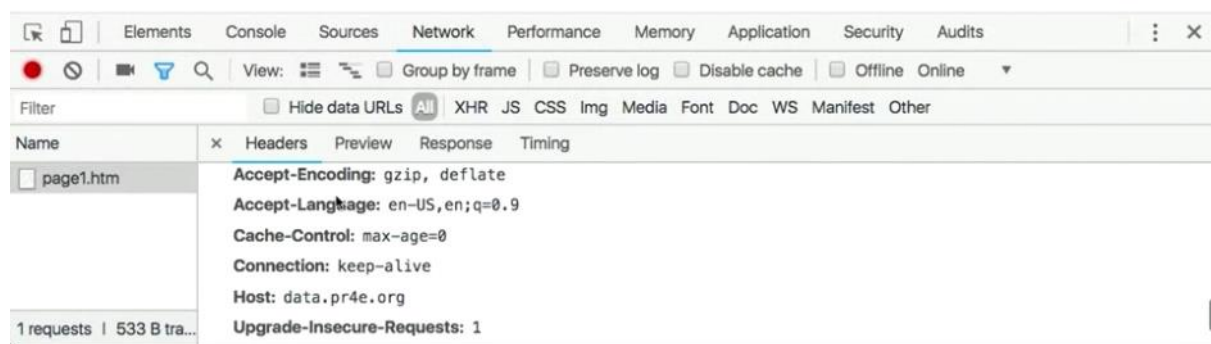
2020年7月13日　　23:43

## Python sockets

```python
1    import socket
2    mysock = socket.socket(socket.AF_INET, socket.SOCK_STREAM)
3    mysock.connect(('data.pr4e.org', 80))
4    cmd = 'GET http://data.pr4e.org/intro-short.txt HTTP/1.0\r\n\r\n'.encode()
5    mysock.send(cmd)
6
7    while True:
8        data=mysock.recv(512)
9        if len(data)<1:
10           break
11       print(data.decode(),end='')
12
13   mysock.close()
```
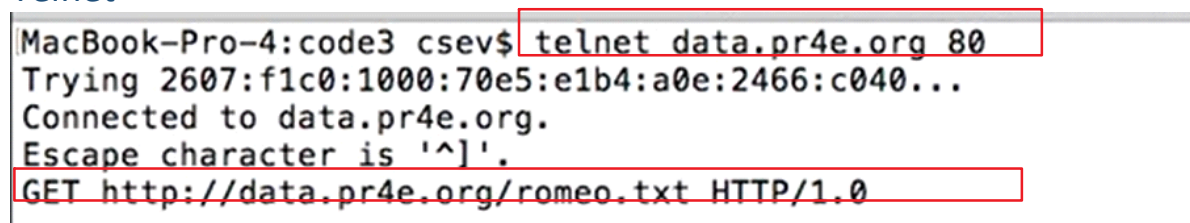
## Developer Console



## Telnet

# Unicode Characters and Strings

2020年7月15日 17:49

In python 3, all strings are Unicode.

# Retrieving Web Page

## Urllib

```python
import urllib.request, urllib.parse, urllib.error

fhand = urllib.request.urlopen('http://data.pr4e.org/romeo.txt')
for line in fhand:
    print(line.decode().strip())
```

# BeautifulSoup4

Susie Choe;
Lunhui Chen

2020年7月16日　　16:20

```python
# To run this, download the BeautifulSoup zip file
# http://www.py4e.com/code3/bs4.zip
# and unzip it in the same directory as this file
from urllib.request import urlopen
from bs4 import BeautifulSoup
import ssl
# Ignore SSL certificate errors
ctx = ssl.create_default_context()
ctx.check_hostname = False
ctx.verify_mode = ssl.CERT_NONE
url = input('Enter - ')
html = urlopen(url, context=ctx).read()
soup = BeautifulSoup(html, "html.parser")
# Retrieve all of the anchor tags
tags = soup('a')
for tag in tags:
    # Look at the parts of a tag
    print('TAG:', tag)
    print('URL:', tag.get('href', None))
    print('Contents:', tag.contents[0])
    print('Attrs:', tag.attrs)
```

```
C:\Users\15378\Desktop\AccessWebData\w4\bs4>python3 urllink2.py
Enter - http://www.dr-chuck.com/page1.htm
TAG: <a href="http://www.dr-chuck.com/page2.htm">
Second Page</a>
URL: http://www.dr-chuck.com/page2.htm
Contents:
Second Page
Attrs: {'href': 'http://www.dr-chuck.com/page2.htm'}
```

# Web Services and XML

2020年7月17日     11:42



Serialize the data to be XML or JSON

# XML

2020年7月17日　　12:29

White space/ indent is only for easy-understanding
It only matters in between a text area.

# SCHEMA

2020年7月17日    13:23

## XSD

XML Scheme Definition

```
<?xml version="1.0" encoding="ISO-8859-1" ?>
<xs:schema xmlns:xs="http://www.w3.org/2001/XMLSchema">
<xs:element name="shiporder">
 <xs:complexType>
  <xs:sequence>
   <xs:element name="orderperson" type="xs:string"/>
   <xs:element name="shipto">
    <xs:complexType>
     <xs:sequence>
      <xs:element name="name" type="xs:string"/>
      <xs:element name="address" type="xs:string"/>
      <xs:element name="city" type="xs:string"/>
      <xs:element name="country" type="xs:string"/>
     </xs:sequence>
    </xs:complexType>
   </xs:element>
   <xs:element name="item" maxOccurs="unbounded">
    <xs:complexType>
     <xs:sequence>
      <xs:element name="title" type="xs:string"/>
      <xs:element name="note" type="xs:string" minOccurs="0"/>
      <xs:element name="quantity" type="xs:positiveInteger"/>
      <xs:element name="price" type="xs:decimal"/>
     </xs:sequence>
    </xs:complexType>
   </xs:element>
  </xs:sequence>
  <xs:attribute name="orderid" type="xs:string" use="required"/>
 </xs:complexType>
</xs:element>
</xs:schema>
```

http://www.w3schools.com/Schema/schema_example.asp

# Parsing XML

2020年7月17日　　13:50



Web Services – Part 4　　　　　　　　　　PYTHON FOR EVERYBODY

```python
import xml.etree.ElementTree as ET          xml1.py
data = '''<person>
  <name>Chuck</name>
  <phone type="intl">
     +1 734 303 4456
  </phone>
  <email hide="yes"/>
</person>'''

tree = ET.fromstring(data)
print('Name:',tree.find('name').text)
print('Attr:',tree.find('email').get('hide'))
```

Web Services – Part 4　　　　　　　　　　PYTHON FOR EVERYBODY

```python
import xml.etree.ElementTree as ET          xml2.py
input = '''<stuff>
  <users>
    <user x="2">
      <id>001</id>
      <name>Chuck</name>
    </user>
    <user x="7">
      <id>009</id>
      <name>Brent</name>
    </user>
  </users>
</stuff>'''

stuff = ET.fromstring(input)
lst = stuff.findall('users/user')
print('User count:', len(lst))
for item in lst:
    print('Name', item.find('name').text)
    print('Id', item.find('id').text)
    print('Attribute', item.get("x"))
```

# JSON and the REST Architecture

| XML | JSON |
|-----------|------|
| Tag | |
| Attribute | |
| text | |

# JavaScript Object Notation (JSON)

2020年8月15日　　19:44

```python
import json                                    json1.py
data = '''{
  "name" : "Chuck",
  "phone" : {
     "type" : "intl",
     "number" : "+1 734 303 4456"
  },
  "email" : {
     "hide" : "yes"
  }
}'''

info = json.loads(data)
print('Name:',info["name"])
print('Hide:',info["email"]["hide"])
```

JSON represents data
as nested "lists" and
"dictionaries"

```python
import json                                    json2.py
input = '''[
  { "id" : "001",
    "x" : "2",
    "name" : "Chuck"
  } ,
  { "id" : "009",
    "x" : "7",
    "name" : "Chuck"
  }
]'''

info = json.loads(input)
print('User count:', len(info))
for item in info:
    print('Name', item['name'])
    print('Id', item['id'])
    print('Attribute', item['x'])
```

JSON represents data
as nested "lists" and
"dictionaries"

## Web Services – Part 7

PYTHON FOR

```python
import urllib.request, urllib.parse, urllib.error
import json

serviceurl = 'http://maps.googleapis.com/maps/api/geocode/json?'

while True:
    address = input('Enter location: ')
    if len(address) < 1: break

    url = serviceurl + urllib.parse.urlencode({'address': address})

    print('Retrieving', url)
    uh = urllib.request.urlopen(url)
    data = uh.read().decode()
    print('Retrieved', len(data), 'characters')

    try:
        js = json.loads(data)
    except:
        js = None

    if not js or 'status' not in js or js['status'] != 'OK':
        print('==== Failure To Retrieve ====')
        print(data)
        continue

    lat = js["results"][0]["geometry"]["location"]["lat"]
    lng = js["results"][0]["geometry"]["location"]["lng"]
    print('lat', lat, 'lng', lng)
    location = js['results'][0]['formatted_address']
    print(location)
```

Enter location: Ann Arbor, MI
Retrieving http://maps.googleapis.com/...
Retrieved 1669 characters
lat 42.2808256 lng -83.7430378
Ann Arbor, MI, USA
Enter location:

geojson.py

# Object Oriented Python

2020年8月18日　　18:22

# Unicode characters and Strings

2020年8月18日      16:41

# An HTTP Request in Python

```python
import socket

mysock = socket.socket(socket.AF_INET, socket.SOCK_STREAM)
mysock.connect(('data.pr4e.org', 80))
cmd = 'GET http://data.pr4e.org/romeo.txt HTTP/1.0\n\n'.encode()
mysock.send(cmd)

while True:
    data = mysock.recv(512)
    if (len(data) < 1):
        break
    print(data.decode())
mysock.close()
```

Decode :Bytes------Unicode(ASCII/UTF-8)

bytes.decode(encoding="utf-8", errors="strict")
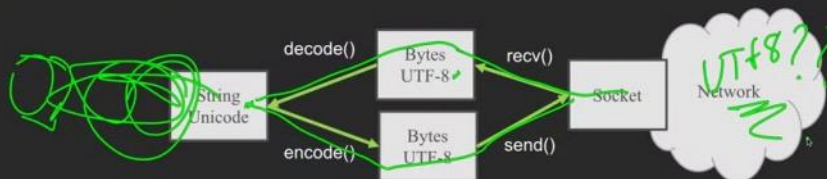bytearray.decode(encoding="utf-8", errors="strict")
Return a string decoded from the given bytes. Default encoding is 'utf-8'. errors may be given to set a different error handling scheme. The default for errors is 'strict', meaning that encoding errors raise a UnicodeError. Other possible values are 'ignore', 'replace' and any other name registered via codecs.register_error(), see section Error Handlers. For a list of possible encodings, see section Standard Encodings.

str.encode(encoding="utf-8", errors="strict")
Return an encoded version of the string as a bytes object. Default encoding is 'utf-8'. errors may be given to set a different error handling scheme. The default for errors is 'strict', meaning that encoding errors raise a UnicodeError. Other possible values are 'ignore', 'replace', 'xmlcharrefreplace', 'backslashreplace' and any other name registered via codecs.register_error(), see section Error Handlers. For a list of possible encodings, see section Standard Encodings.

https://docs.python.org/3/library/stdtypes.html#bytes.decode
https://docs.python.org/3/library/stdtypes.html#str.encode

```python
import socket

mysock = socket.socket(socket.AF_INET, socket.SOCK_STREAM)
mysock.connect(('data.pr4e.org', 80))
cmd = 'GET http://data.pr4e.org/romeo.txt HTTP/1.0\n\n'.encode()
mysock.send(cmd)

while True:
    data = mysock.recv(512)
    if (len(data) < 1):
        break
    print(data.decode())
mysock.close()
```

# Class and Object

2020年8月18日　　17:29



Objects - Part 2

PYTHON FOR EVERYBODY

```
class PartyAnimal:
    x = 0

    def party(self) :
        self.x = self.x + 1
        print("So far",self.x)

an = PartyAnimal()

print("Type", type(an))
print("Dir ", dir(an))
```

We can use dir() to find the "capabilities" of our newly created class.

```
$ python party3.py
Type <class '__main__.PartyAnimal'>
Dir  ['__class__', ... 'party', 'x']
```

# Object Life Cycle

2020年8月18日　　18:06

## Delete/Deconstruct

# Object Inheritance

2020年8月18日　　18:11

# Basic Structured Query Language

2020年8月18日　18:21