# Opening URLs

- The module used for opening URLs is `urllib2`
- The method used is similar to the file `open` in syntax
- Returns a handler to the URL, which could be used as a handle to a file (`readlines`, `read` etc.)

```
1 >>> import urllib2
2 >>> r = urllib2.urlopen('http://python.org/')
3 >>> html = r.read(300)
4 >>> print(html)
5 <!DOCTYPE html PUBLIC "-//W3C//DTD XHTML 1.0
6 Transitional//EN" "http://www.w3.org/TR/xhtml1/
7 DTD/xhtml1-transitional.dtd">
8
9
10 <html xmlns="http://www.w3.org/1999/xhtml" xml:
11 lang="en" lang="en">
12
13 <head>
14   <meta http-equiv="content-type"
15 content="text/html; charset=utf-8" />
16   <title >Python Programming Language --
17 Official Website</title >
```

# General Way

- Not all urls can be opened this way.
- There could be complicated operations such as communicating with he cgi-bin of the server; or some ftp server; etc.
- For that purpose, there are Requests and Opener objects
  - Requests can send along extra data to the server
  - Opener can be used for complicated operations.

```
1 >>> from urllib2 import Request
2 >>> req = Request('http://www.google.com/')
3 >>> brwser = 'Mozilla/4.0 (compatible; MSIE 6.0;
4                   Windows NT 5.0)'
5 >>> req.add_header('User-Agent', brwser)
6 >>> opener = urllib2.build_opener()
7 >>> opened = opener.open(req)
8 >>> print(opened.read(150))
9 <!doctype html><head><meta http-equiv=content-typ
10 content="text/html; charset=UTF-8"><title >hallo -
11 Google Search</title ><script >window.google={kEl:"
12 >>>
```

# On Error?

- In case of errors, one can use the exception to show the error messages
- Two Exceptions which come handy are `HTTPError` and `URLError`
- They have to be used in the same order when you write the code. Because HTTPError is a subclass of URLError
- See the example below.

```
1
2 from urllib2 import Request, urlopen, URLError, H
3 req = Request(someurl)
4 try:
5     response = urlopen(req)
6 except HTTPError, e:
7     print('The server didn't fulfill the req.')
8     print('Error code: ', e.code)
9 except URLError, e:
10    print('We failed to reach a server.')
11    print('Reason: ', e.reason)
12 else:
13    print('everything is fine')
```

# Regular Expressions - A recap

- What are they?
  A means to find out string patters, To match strings, To find substrings and so forth
- When not to use them?
  When they are unavoidable. In normal cases where one needs to check whether a string is a substring of another, then is could be easier and more understandable and perhaps more efficient to use the normal string methods.
- When to use them?
  When you know they must be.

# Regular Expressions in Theory

- Finite Automata - NFA and DFA, Alphabets
- Books on Compilers give a good account of these
- Limitations : ($a^n b^n$), palindromes

# Meta Characters

- If you want to search for `''test''`, then easy.
- What if you don't know what you want to search for. For example a telephone number? (Which you don't know)
- There are some classes of characters which are dedicated to make the using of regular expressions possible.
- Normal characters match for themselves. E.g. `t` matches *t*.
- Some special characters don't match themselves.

```
. ^ $ * + ? { [ ] \ | ( )
```

[ and ] : These can be used to specify a class of characters.

[a-z] : stands for all the lowercase characters. The literal '-' has special meaning inside the square brackets.

[abc$] stands for the characters 'a', 'b', 'c' and the dollar sign.

Even though $ has special meaning in RE context, but inside [ and ]

^ : For negation/complementing a set

[^a-z] means everything which is
not lowercase.

\ is perhaps the most important
metacharacter.

It is used when a meta-character
is to be matched.

```
\d : Every decimal digit. [0-9]

\D : Everything non-digit; [^0-9]

\s : Any whitespace; [ \t\n\r\f\b]

\S : Any nonwhitespace character

\w : Any alpha-numeric; [a-zA-Z0-9_]

\W : Any non-alpha-numeric-character
```

# Importance of DOT

The character ".." matches everything but a newline.
Even that can be done using a different mode of the RE module, using `re.DOTALL`

# Repeating Things

- `*` : `ca*t` would match ct, cat, caat, caaaat, ...

- `+` : `ca+t` would match all of them except for ct

- `?` : `ca?t` would match only ct or cat

- $\{m, n\}$ : Minimum m times, maximum n times.
  $ca\{2,4\}t$ would match caat, caaat and caaaat. But not anything else.

# Repeating Things

- It is easy to see that * is nothing but $\{0,\}$
- Similarly, + is nothing but $\{1,\}$ and
- ? is $\{0,1\}$

```
a|b matches a or b.

^, $ match the beginning
      and ending of a line.

\A, \Z match the beginning
        and end of a string

'\A[abc]*\Z'  matches all strings
  which are combinations of a, b and c

\b matches word boundaries:
  'class\b' match 'class next Thursday'
```

'class\b' doesn't match 'classified'

a[bcd]*b against 'abcbd'

a       The a in the RE matches.

abcbd   The engine matches [bcd]*,
         going as far as it can,
         which is to the end of the
         string.

Failure The engine tries to match b,
         but the current position is
         at the end of the string,

```
          so it fails.

abcb      Back up, so that [bcd]*
          matches one less character.

Failure   Try b again, but the current
          position is at the last
          character, which is a "d".

abc       Back up again, so that [bcd]*
          is only matching "bc".

abcb      Try b again. This time but the
          character at the current
```

position is "b", so it succeeds.

# Using Them

- Compile them
- Match them
    - `match()`: Determine if the re matches the string
    - `search()`: Scan and find the matches
    - `findall()`: Find all the matches
    - `finditer()`: Return and iterator
- Use them

```
1 >>> import re
2 >>> p = re.compile('[a-z]+')
3 >>> p
4 <_sre.SRE_Pattern object at 80c3c28>
5 >>> p.match("")
6 >>> print(p.match(""))
7 None
8 >>> m = p.match( 'tempo')
9 >>> print(m)
10 <_sre.SRE_Match object at 80c4f68>
```

# Using Them

- `group()` : The string matched
- `start()` : Start of the string
- `end()` : The End of the string
- `span()` : A tuple with (start, end)

```
1 >>> m.group()
2 'tempo'
3 >>> m.start(), m.end()
4 (0, 5)
5 >>> m.span()
6 (0, 5)
7 >>> print(p.match('::: message'))
8 None
9 >>> m = p.search('::: message') ; print(m)
10 <re.MatchObject instance at 80c9650>
11 >>> m.group()
12 'message'
13 >>> m.span()
14 (4, 11)
```

```
1 p = re.compile( ... )
2 m = p.match( 'string goes here' )
3 if m:
4     print('Match found: ', m.group())
5 else:
6     print('No match')
7 ─────────────────────────────────────────
8 >>> p = re.compile('\d+')
9 >>> p.findall('12 drummers drumming,
10                11 pipers piping,
11                10 lords a-leaping')
12 ('12', '11', '10')
13 ─────────────────────────────────────────
14 >>> iterator = p.finditer('12 drummers drumming,
15                     11 ... 10 ...')
16 >>> iterator
17 <callable-iterator object at 0x401833ac>
```

```
18 >>> for match in iterator:
19 ...      print(match.span())
20 ...
21 (0, 2)
22 (22, 24)
23 (29, 31)
```

# Generators

# Generators

- Iterator creators (So to speak)
- Regular functions which return without returning.
- Uses `yield` statement
- Each call of `next` resumes from where it left off.
- State/Data values stored

```
1
2 def reverse(data):
3     for index in range(len(data)−1, −1, −1):
4         yield data(index)
5
6 >>> for char in reverse('golf'):
7 ...     print(char)
8 ...
9 f
10 l
11 o
12 g
```

# Generators ...

- Generators are the equivalent of class based Iterators
- `__iter__` and `next` are created automatically
- Saving the vales makes it easier. No need to separate initialization/storage of index.
- Automatic raising of Exception on termination.

# Simulating Generators

- Can be simulated with normal functions.
  1. Start with an empty list.
  2. Fill in the list instead of the `yield` statement
  3. Then return an iterator of the list
  4. Same result

```python
1
2  def r(data):
3      for index in range(len(data)-1, -1, -1):
4          yield data(index)
5
6
7  def rS(data):
8      _list = ()
9      for index in range(len(data)-1, -1, -1):
10         _list.append(data(index))
11     return iter(_list)
```

```
 1
 2 >>> import gs
 3 >>> for x in gs.r('this is cool'):
 4 ...      print(x)
 5 ...
 6 l o o c   s i   s i h t
 7 >>> for x in gs.rS('this is cool'):
 8 ...      print(x)
 9 ...
10 l o o c   s i   s i h t
11 >>>
12 >>>
```

# Parsers in Python

- XML
- HTML

# XML Parser

- SAX (Simple API for XML)
  - Reads the file as required
  - Special methods are called when tags are opened/closed

- DOM
  - Reads the whole file in a go
  - The whole structure is readily accessible for use.

# SAX Parser

- `xml.sax.make_parser()` gives a generic parser object.
- The parser object is an instance of XMLReader. (It can read and output structured XML)
- A content handler has to be implemented for the XMLReader (example)
- Contenthandler is a class which is implemented for the specific needs

# ContentHandler

- `startDocument()/endDocument()` are called from reading and processing the XML-Codes

- `startElement(name, attrs)` is called whenever a new tag is opened
  - `name` is the name of the tag
  - `attrs` contains the attributes part of the tag. It is an attribute object.

# Contenthandler

- `endElement(name)` is called when a tag is closed.
- `characters(str)` gives the CDATA in the parameter to be used.
- There is no guarantee that all the data inside would be given in a single instance. One has to collect data if needed. (Example)

```
1
2 from xml.sax.handler import ContentHandler
3 class CDATAPrinter(ContentHandler):
4   def startElement(self, name, attrs):
5     self.cdata=''
6   def endElement(self, name):
7     if len(self.cdata.strip()) > 0:
8       print(name, ':', self.cdata.strip())
9   def characters(self, str):
10     self.cdata += str
```

```
1 <something>
2          <string>HA HA HA </string>
3          <number>12 34 43 </number>
4          <nothing> nothing </nothing>
5 </something>
6 ——————————————————————————————————————————
7 >>> import boo
8 >>> import xml.sax
9 >>> parser = xml.sax.make_parser()
10 >>> parser.setContentHandler(boo.CDATAPrinter())
11 >>> parser.parse('cal.xml')
12 string : HA HA HA
13 number : 12 34 43
14 nothing : nothing
15 something : nothing
16 >>>
```

# HTML Parsing

- HTML is sometimes XML
- HTML tags need not be closed always
- HTML tags can have attributes and some have always

# HTML Parsing

- Similar to XML parsing
- There is an abstract class `HTMLParser` which needs to be implemented for own purposes
- It contains the following methods
    - `handle_starttag(tag, attrs)`
    - `handle_endttag(tag)`
    - `handle_startendtag(tag,attrs)`
    - `handle_data(data) (for characters(str))`

# HTML Parsing

- The HTMLParser has its own ContentHandler. Just calling `HTMLParser()` gives an instance of the class.
- For parsing, one has to feed the html-text to the parser. `parser.feed(hstring)`
- As far as it can, it would ignore the errors in the string. Sometimes EOF reaches before the error-limit is reached.
- To read a URL, the following code would be useful.
  `parser.feed(urllib2.open(URL).read())`

```
1
2 from HTMLParser import HTMLParser
3 class MyHTMLParser(HTMLParser):
4         def handle_starttag(self, tag, attrs):
5                 print("Breaking In: ", tag)
6         def handle_endtag(self, tag):
7                 print("Getting Out: ", tag)
8         def handle_startendtag(self, tag, attrs):
9                 print("Empty Tag??: ", tag)
```

```
1 >>> import myhtmlparser
2 >>> import urllib2
3 >>> parser = myhtmlparser.MyHTMLParser()
4 >>> parser.feed(urllib2.urlopen("http://www.bing.
5 Breaking In:    html
6 Breaking In:    head
7 Empty Tag??:    meta
8 Breaking In:    script
9 Getting Out:    script
10 Breaking In:    script
11 Getting Out:    script
12 ...
13 ...
14 ...
15 Breaking In:    script
16 Getting Out:    script
17 Getting Out:    body
```

18 Getting Out: html

# Decorator Pattern

- In object-oriented programming, the decorator pattern is a design pattern that allows new/additional behaviour to be added to an existing class dynamically.
- In Python one cannot say that to be the same with the Decorator; even though one can achieve the same functionality with decorators in python.
- So, what are decorators IN Python?

# Functions taking Functions

- Functions can have pointers to other functions as parameters.

- A function which can take another function as its parameter and can achieve something there by could be mainly classified as a decorator.
  See example.

```
1 >>> def ourdecorator2(foo):
2 ...     def newfoo():
3 ...         print("We are to call", foo.__name__)
4 ...         return foo()
5 ...     return newfoo
6 ...
7 >>>
8 >>> foo = ourdecorator2(foo1)
9 >>>
10 >>> foo()
11 We are to call foo1
12 Hello World
13 >>>
```

```
 1 >>>
 2 >>> def ourdecorator(foo):
 3 ...     print("We are to call", foo.__name__)
 4 ...     return foo()
 5 ...
 6 >>> def foo1():
 7 ...     print("Hello World")
 8 ...
 9 >>>
10 >>> foo  = ourdecorator(foo1)
11 We are to call foo1
12 Hello World
13 >>>
14 >>>
```

# Similar to Macros

- Decorators are similar to MACROS in other programming languages
- They are usually used to make a wrapper around functions
- And of course, classes too.

# Usage, then Creation

- Function decorators are placed above the function with the key-character '@'
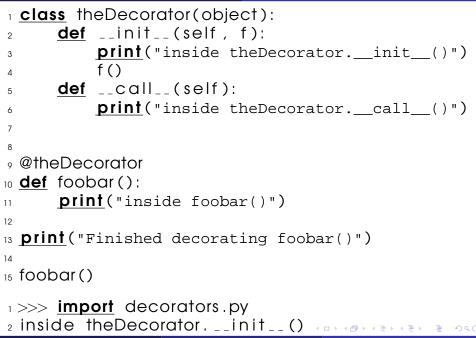
  ```
  @thedecorator
  def foo():
  ....
  ```

- The interpreter compiles `foo` and calls the decorator with that as argument.

- The result of that replaces the code for `foo`

# How to implement decorator

- Could be Functions or Classes.
- The condition is that whatever the decorator returns, that should be callable.
- An object is callable, if the method `__call__` is implemented.

```python
1 class theDecorator(object):
2     def __init__(self, f):
3         print("inside theDecorator.__init__()")
4         f()
5     def __call__(self):
6         print("inside theDecorator.__call__()")
7
8
9 @theDecorator
10 def foobar():
11     print("inside foobar()")
12
13 print("Finished decorating foobar()")
14
15 foobar()
```

```python
1 >>> import decorators.py
2 inside theDecorator.__init__()
```

```
3 inside foobar()
4 Finished decorating foobar()
5 inside theDecorator.__call__()
```

# Observation

- From the output, it is clear that the `init` is called when the decorator is used.
- So, usually, the call to the function is done only in the `call` function.
- Once a function is decorated, the behaviour totally changes. The call goes only to the decorated code. (line number 4 of the output)

```python
 1 class loggerdeco(object):
 2
 3     def __init__(self, f):
 4         self.f = f
 5
 6     def __call__(self):
 7         print("Entering", self.f.__name__)
 8         self.f()
 9         print("Exited", self.f.__name__)
10
11 @loggerdeco
12 def func1():
13     print("inside func1()")
14
15 @loggerdeco
16 def func2():
17     print("inside func2()")
```

```
 1
 2 func1()
 3 func2()
 4 Entering func1
 5
 6
 7 inside func1()
 8 Exited func1
 9 Entering func2
10 inside func2()
11 Exited func2
```

# Using Functions

- The same can be achieved using functions, instead of classes.

- The decorator functions usually enclose the decorated function in between the decoration.

- This is done inside a subfunction (equivalent of `__call__` and the pointer to the subfunction is returned.

```
1 def loggerdeco(f):
2     def new_f():
3         print("Entering", f.__name__)
4         f()
5         print("Exited", f.__name__)
6     return new_f
7
8 @loggerdeco
9 def func1():
10     print("inside func1()")
11
12 @loggerdeco
13 def func2():
14     print("inside func2()")
```

```
1 func1()
2 func2()
3 print(func1.__name__)
4
5
6 Entering func1
7 inside func1()
8 Exited func1
9 Entering func2
10 inside func2()
11 Exited func2
12 new_f
```

# Observation

- The name of the functions have been changed to `new_f`.
- This can be changed by reassigning `new_f.__name__ = f.__name__`

There are many cool uses of decorators. You can see more examples at
`http://wiki.python.org/moin/PythonDecorat`

```python
1  class memoized(object):
2    def __init__(self, func):
3      self.func = func
4      self.cache = {}
5    def __call__(self, *args):
6      try:
7        return self.cache(args)
8      except KeyError:
9        self.cache(args) = value = self.func(*args)
10        return value
11      except TypeError:
12        return self.func(*args)
13    def __repr__(self):
14      return self.func.__doc__
15
16  @memoized
17  def fibonacci(n):
```

```
18      "Return the nth fibonacci number."
19      if n in (0, 1):
20          return n
21      return fibonacci(n−1) + fibonacci(n−2)
22
23  for i in xrange(1, 100, 9):
24      print(fibonacci(i))
```

```
1 (sadanand@lxmayr10 @ ~) time python memorized.py
2 1
3 55
4 4181
5 317811
6 24157817
7 1836311903
8 139583862445
9 10610209857723
10 806515533049393
11 61305790721611591
12 4660046610375530309
13
14 real      0m0.014s
15 user      0m0.008s
16 sys       0m0.000s
17 (sadanand@lxmayr10 @ ~)
```

◄ □ ► ◄ 🗗 ► ◄ 📃 ► ◄ 📃 ►      ≣      ♡ ९ ℃

# Static Variables and Methods

- A static variable in a class has always the same value, independent of the instances.
- Static variables are class variables, they belong to the class than to the instances
- They are accessed by the name of the Class, rather than the instance.

```
1 class myStatic:
2     instances = 0
3     def __init__(self):
4         myStatic.instances += 1
5
6     def howmany(self):
7         return myStatic.instances
8
9 x = myStatic()
10 print(x.howmany())
11 y = myStatic()
12 print(y.howmany())
13 print(x.howmany())
14 ——————————
15 print(1, 2, 2)
```

# Static Methods

- They have the same return value independent of the class instance
- They don't have the `self` parameter
- For the same reason, they cannot access any of the `self.*` objects.
- The keyword is a decorator named `@staticmethod`

```python
1  class myStatic:
2      instances = 0
3
4      def __init__(self):
5          myStatic.instances += 1
6
7      @staticmethod
8      def howmany():
9          return myStatic.instances
10
11
12 x = myStatic()
13 print(myStatic.howmany())
14 y = myStatic()
15 print(myStatic.howmany())
16 =======================
17 prints 1, 2 as expected.
```

# Classmethod

- `@classmethod` is perhaps a special thing for python.
- The methods decorated with this gets as the initial variable a class which is the original class (not the instance)
- That helps the function to act like a normal method of the class, by accepting all the attributes and treat them as static as well.

```python
1  class myStatic:
2      instances = 0
3      def __init__(self):
4          self.addinstance()
5
6      @classmethod
7      def howmany(cls):
8          return cls.instances
9      @classmethod
10     def addinstance(cls):
11         cls.instances += 1
12
13 x = myStatic()
14 print(myStatic.howmany())
15 y = myStatic()
16 print(myStatic.howmany())
17 ============================
```

18 prints 1, 2 as expected.

```python
 1 class myStatic :
 2     instances = 0
 3
 4     def __init__(self):
 5         self.addinstance()
 6
 7     @classmethod
 8     def howmany(cls):
 9         return cls.instances
10
11     @classmethod
12     def addinstance(cls):
13         cls.instances += 1
14
15     def nastything(self):
16         print("trying to be nasty")
17         self.instances = -1
```

```
1 x = myStatic()
2 x.nastything()
3 print(myStatic.howmany())
4 print(x.howmany())
5 y = myStatic()
6 x.nastything()
7 print(myStatic.howmany())
8 ====================
9 trying to be nasty
10 1
11 1
12 trying to be nasty
13 2
```
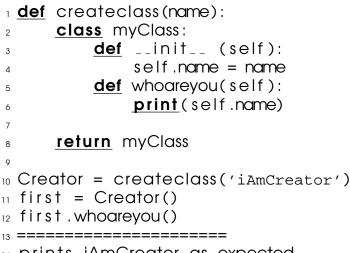
# Single Instance

- The static methods can be used to create a singleton object/pattern
- They are classes for which there is only one instance at any given time.
- They could be implemented using
  1. The class instance could lie in a static variable
  2. The method which gets the instance can be made static.

# Anonymous Classes

- New classes could be defined inside functions and returned.
- Such are called anonymous classes
- Anonymous classes can also be created using `classobj`

```python
1  def createclass(name):
2      class myClass:
3          def __init__(self):
4              self.name = name
5          def whoareyou(self):
6              print(self.name)
7
8      return myClass
9
10 Creator = createclass('iAmCreator')
11 first = Creator()
12 first.whoareyou()
13 ======================
14 prints iAmCreator as expected.
```

```python
1 from new import classobj
2 class foo:
3     def __init__ (self):
4         self.x = 'x'
5
6 foo2 = classobj('foo2', (foo,),
7   {'bar': lambda self, x: 'got ' + str(x)})
8
9 print(foo2().bar(3))
10 print(foo2().x)
11 ===============
12 prints got2, x as expected.
```

# Meta Classe

- Not in the scope of our course.
- Creating tailormade classes / customized ones.
  ```
  __metaclass__
  ```

# Some Philosophy?

If you'd like to know some python philosophy,
then you may import the module `this`

The Zen of Python, by Tim Peters

Beautiful is better than ugly.
Explicit is better than implicit.
Simple is better than complex.
Complex is better than complicated.
Flat is better than nested.
Sparse is better than dense.
Readability counts.
Special cases aren't special enough to break the rules.
Although practicality beats purity.
Errors should never pass silently.
Unless explicitly silenced.
In the face of ambiguity, refuse the temptation to guess.
There should be one– and preferably only one –obvious way to do it.
Although that way may not be obvious at first unless

you're Dutch.
Now is better than never.
Although never is often better than *right* now.
If the implementation is hard to explain, it's a bad idea.
If the implementation is easy to explain, it may be a good idea.
Namespaces are one honking great idea – let's do more of those!

# Tab Complete

- Getting Tab Complete (like Bash) in python prompt.

```
1
2 try :
3     import readline
4 except ImportError:
5     print("Unable to load readline module.")
6 else :
7     import rlcompleter
8     readline.parse_and_bind("tab: complete")
```

# Problems

- Open an URL (from commandline)
- Use HTML parser to parse and filter the links / Write a re.query
- NeverEnding Iterator
- Single instance class
- Recursive function tracer.