

- 1 Useful; Unnecessary
- 2 Parsing
- 3 Shelves
- 4 The searching Details
- 5 Problems

# Functions

We have already seen functions. But only the simplest forms. We can have functions

- With arguments having default values
- With keywords as arguments
- With multiple arguments.

```
1 >>> def myfoo(bar, foobar=True):
2 ...     print bar,
3 ...     if foobar:
4 ...         print "ha ha ha!"
5 ...
6 >>> myfoo("hello")
7 hello ha ha ha!
8 >>> myfoo("hello", foobar=False)
9 hello
10 >>>
```

# Default values taken only once

- The default value of the parameter is initialised only once and it stays the same if not specifically called. Look at the following example.

```
1 >>> def add(this , tothat=()):
2 ...     for e in this:
3 ...         tothat.append(e+1)
4 ...     return tothat
5 ...
6 >>> add((23 , 34))
7 (24 , 35)
8 >>> add((23 , 34))
9 (24 , 35 , 24 , 35)
10 >>> add((23 , 34))
11 (24 , 35 , 24 , 35 , 24 , 35)
12 >>> add((23 , 34))
13 (24 , 35 , 24 , 35 , 24 , 35 , 24 , 35)
14 >>> add((23 , 34) , (1 , 2))
15 (1 , 2 , 24 , 35)
16 >>>
```

# Multiple Arguments

- Functions with a \*-ed argument can have multiple arguments.
- The arguments would be packed in a tuple
- The \*-ed argument must follow the other typed of arguments.

```
1 >>> def mularg(i, j, *rest):
2     ...     print i+j
3     ...     for k in rest:
4     ...         print k
5     ...
6 >>> mularg(1, 2)
7 3
8 >>> mularg(1, 2, 4)
9 3
10 4
11 >>> mularg('hello', 'world',
12           'this', 'is', 'cool!')
13 helloworld
14 this
15 is
16 cool!
17 >>>
```

# Docstrings

- Strings surrounded by three quotes at the beginning of functions could be used for documentation purposes.
- These strings contain newlines in them.



```
1
2 >>> def simpledoc():
3     ...     """This is a simple hello
4     ...         world program - just to reveal
5     ...         the beauty of docstrings"""
6     ...     print "Hello World"
7     ...
8 >>> simpledoc.__doc__
9 'This is a simple hello\n                world program - j
10 the beauty of docstrings'
11 >>> print simpledoc.__doc__
12 This is a simple hello
13     world program - just to reveal
14     the beauty of docstrings
15 >>> help(simpledoc)
16 ...
```

# With expression

- Files are to be always closed after use.
- A keyword named `with`
- Using `with` helps automatic closing of files after use.
- The object which is used with `with` must have the methods - `__enter__` and `__exit__` implemented

1  
2  
3  
4  
5  
6

```
with open(filename) as f:  
    for line in f:  
        print line
```

# 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>
```

---

```
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_endtag(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  
3  
4  
5  
6  
7  
8  
9

```
from HTMLParser import HTMLParser
class MyHTMLParser(HTMLParser):
    def handle_starttag(self , tag , attrs):
        print "Breaking In: " , tag
    def handle_endtag(self , tag):
        print "Getting Out: " , tag
    def handle_startendtag(self , tag , attrs):
        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
```

# Shelves

- A shelf is a persistent dictionary object in python
- A dictionary in the secondary storage
- Could be opened and used as needed.
- `open` and `close` are the usual methods needed.



```
1 >>> import shelve
2 >>> d = shelve.open("myfile.shelf")
3 >>> d('lala') = 'booboo'
4 >>> d('kiki') = 'myamya'
5 >>> d
6 {'lala': 'booboo', 'kiki': 'myamya'}
7 >>> d('xx') = range(4)
8 >>> d
9 {'lala': 'booboo', 'xx': (0, 1, 2, 3), 'kiki': 'm
10 >>> d.close()
11 >>>
12 (sadanand@lxmayr10
13 myfile.shelf.bak  myfile.shelf.dat  myfile.shelf.
14 (sadanand@lxmayr10
15 >>> import shelve
16 >>> d = shelve.open("myfile.shelf")
17 >>> d
18 {'kiki': 'myamya', 'xx': (0, 1, 2, 3), 'lala': 'k
```

# Relevance of a Word

- Searching and indexing is done based on the relevance of the word.
- The simplest method could be the frequency of occurrence.
  - ▶ That would lead to a problem that 'the', 'a', etc. would get more relevance.
- A better method:  $tf-idf$

# tf-idf

- `tf-idf` is a measure or a benchmark to find the relevance of each word on the basis of its occurrence and frequency in each file.
- It can be calculated as follows.
  - ▶  $n_{i,j}$  is the number of occurrences of word  $w_i$  in the document  $d_j$
  - ▶  $D$  is the number of documents
  - ▶  $D_i$  is the number of documents in which the word  $w_i$  occurs.

$$tf_{i,j} = n_{i,j} / \sum_k n_{k,j} \quad (1)$$

$$idf_i = \ln \frac{D}{D_i} \quad (2)$$

$$tfidf_{i,j} = tf_{i,j} \times idf_i \quad (3)$$

# What with it?

- A higher value of tf-idf implies that the word has a higher frequency of occurrence in the less number of files where it appear.
- The common words 'the' or 'a' would occur in every file and that makes the denominator in  $idf_i$  larger - thereby making the tf-idf value smaller.
- Every word has a tfidf value for each file.
- <http://en.wikipedia.org/wiki/Tf-idf>

# Problems

- Implement an HTMLParser
- Use the parser to filter the text from the documents pointed by the nodes of the graph
- Create tf-idf values.

# Looking Forward

- How many more lectures?
- What more to be done?