

Introduction to programming

Lecture 6

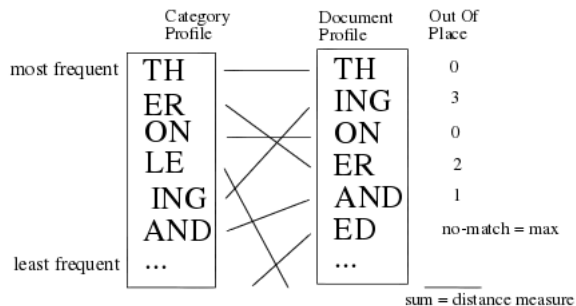


**UNIVERSITY OF
GOTHENBURG**

Richard Johansson

October 6, 2015

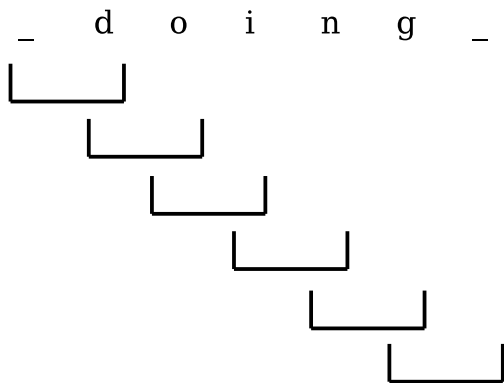
comparing a document profile to a language profile



collecting n -gram statistics

```
def collect_ngram_statistics(words, dictionary, n):  
    pad = ' '*(n-1)  
    for word in words:  
        padded_word = '%s%s%s' % (pad,word,pad)  
        index = 0  
        while index+n <= len(padded_word):  
            ngram = padded_word[index:index+n]  
            if ngram in dictionary:  
                dictionary[ngram] += 1  
            else:  
                dictionary[ngram] = 1  
            index += 1
```

example: collecting bigrams



command-line arguments

- ▶ when running a program, you can get its **command-line arguments** from the list `sys.argv`:

```
import sys
if __name__ == '__main__':
    inputfile = sys.argv[1]
    outputfile = sys.argv[2]
    ... do something with inputfile and outputfile ...
```


data persistence

- ▶ when our program has carried out some work, we might want to **save** it so that we can reuse it later
- ▶ we have already seen how to write to a text file:

```
with open('output.txt', 'w') as f:
    print('this is the output to the text file', file=f)
```

- ▶ **pickling** (in other languages called **serializing**): converting a Python object to raw data (a string) so that it can be written to a file and later reloaded
- ▶ we can save our data without having to define a file format

```
import pickle
with open('output.data', 'wb') as f:
    pickle.dump(some_object, f)
```

```
...
with open('output.data', 'rb') as f:
    reloaded = pickle.load(f)
```

example: saving and loading a frequency table

```
import nltk
import sys
import pickle
```

```
def compute_frequencies(filename):
    ...
    return frequencies
```

```
if __name__ == '__main__':
    inputfile = sys.argv[1]
    outputfile = sys.argv[2]
    freqs = compute_frequencies(inputfile)
    with open(outputfile, 'wb') as f:
        pickle.dump(freqs, f)
```

```
import sys
import pickle
```

```
if __name__ == '__main__':
    inputfile = sys.argv[1]
    testword = sys.argv[2]
    with open(inputfile, 'rb') as f:
        freqs = pickle.load(f)
    print(freqs[testword])
```


using a Counter

- ▶ a **Counter** (note the capital C) is a dictionary specialized for frequency counting

```
from collections import Counter
```

```
document = ["this", "is", "a", "collection", "of",
            "words", "and", "it", "is", "extracted",
            "from", "a", "text"]
```

```
# almost like before
freqs = Counter()
for word in document:
    freqs[word] += 1
```

```
# to get the most frequent:
print(freqs.most_common(3))
```

using a Counter (even simpler)

```
from collections import Counter
document = ["this", "is", "a", "collection", "of",
            "words", "and", "it", "is", "extracted",
            "from", "a", "text"]

freqs = Counter(document)
```

example: counting the part-of-speech tags for each word

The	DT
last	JJ
thing	NN
they	PRP
needed	VBD
was	VBD
another	DT
drag-down	JJ
blow	NN
.	.
,,	,,
That	DT
measure	NN
could	MD
compel	VB
...	...

example: counting the part-of-speech tags for each word
(code)

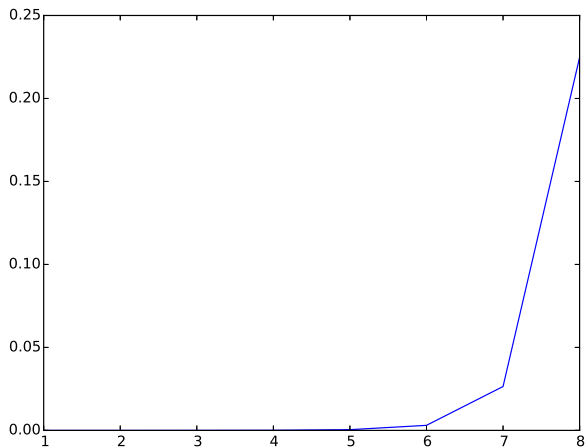
```
from collections import defaultdict, Counter

stats = defaultdict(Counter)

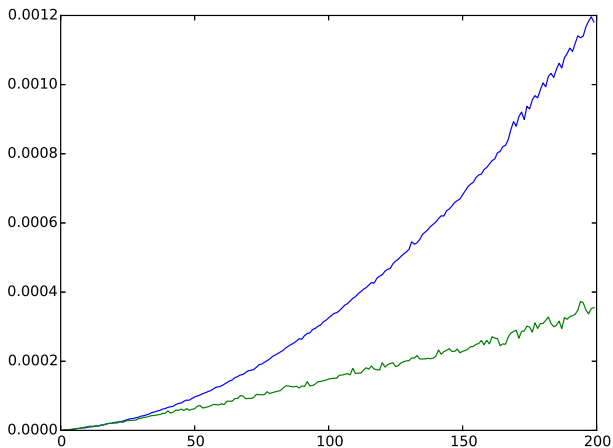
with open('tagged_corpus.txt') as f:
    for l in f:
        word, tag = l.split()
        stats[word][tag] += 1

print(stats['measure'])
```


measuring the time of the idiot sort algorithm



execution time of selection sort (blue) and quicksort (green)



reasoning about time complexity

- ▶ when determining the time complexity, we try to reason about how many steps the algorithm will take, depending on the input size N
- ▶ in general
 - ▶ a single loop over the whole input gives $O(N)$
 - ▶ (assuming each step takes constant time!)
 - ▶ a loop inside a loop gives $O(N) \cdot O(N) = O(N^2)$
 - ▶ or equivalently, calling an $O(N)$ function inside a loop
 - ▶ **but** one loop after another gives $O(N) + O(N) = O(N)$

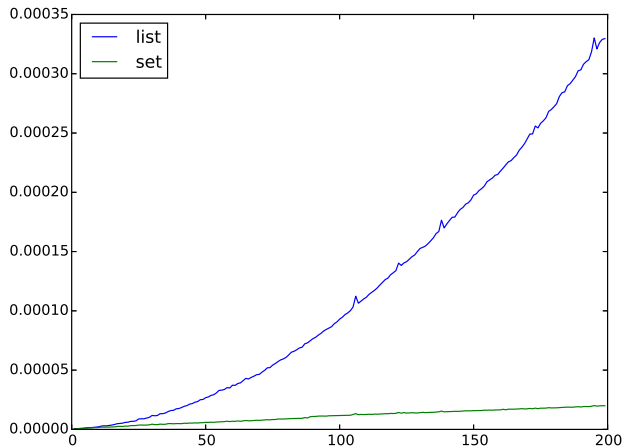
data structures

- ▶ we use **data structures** to store the data that our program processes
 - ▶ lists, sets, dictionaries, ...
- ▶ the selection of a data structure is a tradeoff
 - ▶ list: we remember order; fast access; quite fast to add elements at the end but slow elsewhere; slow to test membership
 - ▶ set: we don't remember order; fast to add elements; fast to test membership
 - ▶ dictionary: key-value mapping; we don't remember insertion order; fast lookup by key; slow lookup by value
- ▶ in some cases, we may need to develop our own data structures
 - ▶ see last part of this lecture

example: counting the number of unique elements in a list

```
def count_unique(lst):  
    seen_before = []  
    for x in lst:  
        if x not in seen_before:  
            seen_before.append(x)  
    return len(seen_before)
```

counting unique elements



recap: declaring a class

The name of the class

The name of its superclass

```
class MyClassName (object):
```

Constructor definition

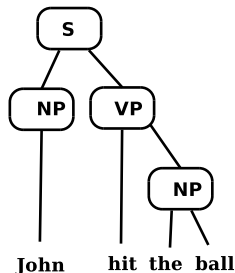
```
    def __init__(self, inputs):
        ...
        (do something)
        ...
```

Method definitions

```
    def some_method(self, inputs):
        ...
        (do something)
        ...
    def some_other_method(self, inputs):
        ...
        return something
```

example: phrase structure trees

- ▶ a **phrase structure tree** is a tree commonly used to represent syntactic structure
- ▶ it consists of **phrases** and **words**
- ▶ a phrase consists of a phrase label (e.g. NP, VP, ...) and a list of children (words or other phrases)



example: phrase structure trees

```
w1 = Word("John")
w2 = Word("hit")
w3 = Word("the")
w4 = Word("ball")
```

```
p1 = Phrase("NP", [w1])
p2 = Phrase("NP", [w3, w4])
p3 = Phrase("VP", [w2, p2])
p4 = Phrase("S", [p1, p3])
```

```
p4.printout()
```

```
print(w4.depth())
```

S:

NP:

John

VP:

hit

NP:

the

ball

3

phrase structure trees: the code

```
class Word(object):
    def __init__(self, word):
        self.word = word
        self.parent = None

    def printout(self, ind):
        print(" "*ind + self.word)

    def depth(self):
        if not self.parent:
            return 0
        else:
            return 1 + self.parent.depth()

class Phrase(object):
    def __init__(self, label, children):
        self.parent = None
        for c in children:
            c.parent = self
        self.children = children
        self.label = label

    def printout(self, ind = 0):
        print(" "*ind + self.label+":")
        for c in self.children:
            c.printout(ind + 4)

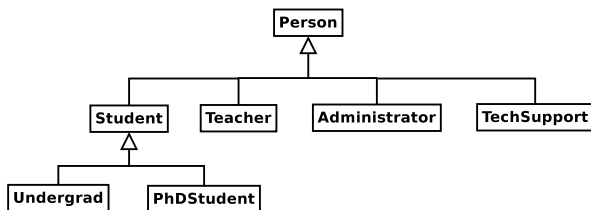
    def depth(self):
        if not self.parent:
            return 0
        else:
            return 1 + self.parent.depth()
```


example of design with inheritance

- ▶ a hierarchy of inheritance can be quite deep:

```
class Person(object):  
    ...  
class Student(Person):  
    ...  
class Teacher(Person):  
    ...  
class Undergrad(Student):  
    ...
```

- ▶ note: if `isinstance(x, Undergrad)`, then we also have `isinstance(x, Person)` and `isinstance(x, Student)`



phrase structure trees with inheritance

```
class Node(object):
    def depth(self):
        if not self.parent:
            return 0
        else:
            return 1 + self.parent.depth()

class Word(Node):
    def __init__(self, word):
        self.word = word
        self.parent = None

    def printout(self, ind):
        print(" "*ind + self.word)

class Phrase(Node):
    def __init__(self, label, children):
        self.parent = None
        for c in children:
            c.parent = self
        self.children = children
        self.label = label

    def printout(self, ind = 0):
        print(" "*ind + self.label+":")
        for c in self.children:
            c.printout(ind + 4)
```

overview

introduction to the theory of algorithms

more object-oriented programming

making your own data structures

developing our own data structures: a linked list

- ▶ a **linked list** is a data structure consisting of a chain of links, where each link contains a piece of data
- ▶ advantages compared to a normal Python list: easy and fast to insert, especially at the start of the list
- ▶ disadvantages: complicated and slow to get the n -th item

```
class LinkedList(object):
```

```
...
```

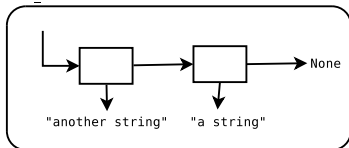
```
the_list = LinkedList()
```

```
the_list.add_first("a string")
```

```
the_list.add_first("another string")
```

```
print(the_list.get_first())
```

the_list



linked list implementation

```
class Link(object):
    def __init__(self, data, next):
        self.data = data
        self.next = next

class LinkedList(object):
    def __init__(self):
        self.first = None
    def add_first(self, data):
        self.first = Link(data, self.first)
    def get_first(self):
        if self.first:
            return self.first.data
    def remove_first(self, data):
        if self.first:
            self.first = self.first.next

the_list = LinkedList()
the_list.add_first("a string")
the_list.add_first("another string")
print(the_list.get_first())
```

iterators and iterables

- ▶ an **iterator** is an object that has a method called `__next__`
 - ▶ `__next__` will generate a new item each time it is called
- ▶ an **iterable** is an object that has a method called `__iter__` that will return an iterator
- ▶ if an object is iterable, then we can use it in a `for`
 - ▶ all data structures such as lists, sets, dictionaries are iterable
 - ▶ if we are going through a list, the iterator will remember the position where we are currently looking

```
for x in some_iterable:
    ... do something with x ...
```

example

```
class NumberGenerator(object):
    def __init__(self):
        self.current = 0
    def __next__(self):
        self.current += 1
        if self.current > 10:
            raise StopIteration
        return self.current

class NumberSequence(object):
    def __iter__(self):
        return NumberGenerator()

numbers = NumberSequence()
for n in numbers:
    print(n)
```

making the linked list iterable

```
class Link(object):
    def __init__(self, data, next):
        self.data = data
        self.next = next

class LinkedListIterator(object):
    def __init__(self, start):
        self.current = start
    def __next__(self):
        if not self.current:
            raise StopIteration
        else:
            out = self.current.data
            self.current = self.current.next
            return out

class LinkedList(object):
    def __init__(self):
        self.first = None
    ...
    def __iter__(self):
        return LinkedListIterator(self.first)
```

```
the_list = LinkedList()

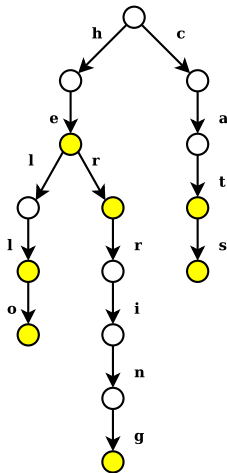
the_list.add_first("test1")
the_list.add_first("test2")
the_list.add_first("test3")

for x in the_list:
    print(x)

print(list(the_list))
```

developing our own data structures: letter tree (trie)

- ▶ dictionaries are efficient for storing words, but sometimes we need to do more complex things:
 - ▶ finding all words starting with *h*, or alphabetically between *abc* and *abx*
 - ▶ finding the words most similar to the misspelled word *hering*
 - ▶ finding anagrams
 - ▶ ...
- ▶ a **trie** is a data structure for strings where all strings sharing a prefix are represented as a tree node



implementing the trie

```
class TrieNode(object):
    def __init__(self):
        self.children = {}
        self.end = False

    def insert(self, s, position):
        if position == len(s):
            self.end = True
        else:
            letter = s[position]
            if not self.children.has_key(letter):
                child = TrieNode()
                self.children[letter] = child
            self.children[letter].insert(s, position+1)

    def contains(self, s, position):
        if position == len(s):
            return self.end
        letter = s[position]
        if not self.children.has_key(letter):
            return False
        return self.children[letter].contains(s, position+1)
```

```
class Trie(object):
    def __init__(self):
        self.root = TrieNode()

    def insert(self, s):
        self.root.insert(s, 0)

    def contains(self, s):
        return self.root.contains(s,0)

t = Trie()

t.insert("hello")
t.insert("he-man")
t.insert("herring")
t.insert("horrible")

print(t.contains("herring"))
print(t.contains("hell"))
```