

Introduction to programming

Lecture 4: processing files and counting words



**UNIVERSITY OF
GOTHENBURG**

Richard Johansson

September 22, 2015

example: most frequent word in GP

```
with open('gp.txt', encoding='utf-8') as f:
    table = {}
    for line in f:
        for word in line.split():
            if word in table:
                table[word] += 1
            else:
                table[word] = 1
print(max(table, key=table.get))
```


opening a file for reading

- ▶ before Python can access the contents of a file, the file needs to be **opened**
- ▶ use builtin function `open` to open a file for reading

```
with open("textfile.txt") as f:
```

```
    ...
```

- ▶ `f` is a **file object**

exception handling

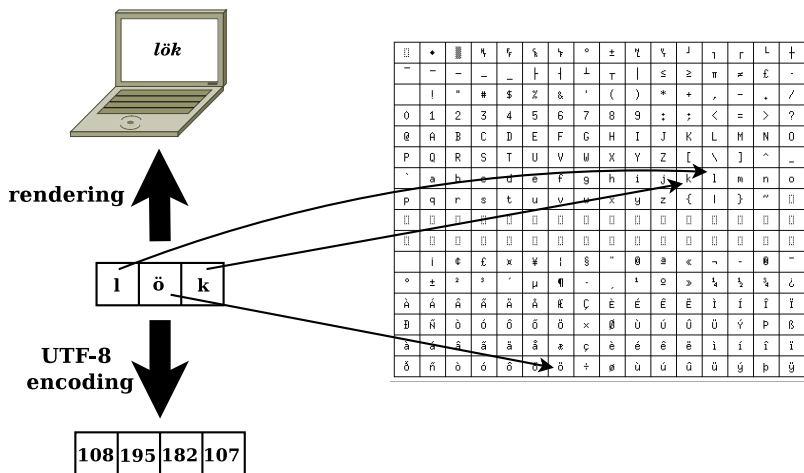
- ▶ what happens if we try to read a file that does not exist?

```
with open("doesnotexist.txt") as f:
    content = f.read()
    print(content)
```

- ▶ an **exception** will be **raised** when something goes wrong
- ▶ we will exit whatever we were doing, and if the exception is not **handled**, the program will stop

```
try:
    with open("doesnotexist.txt") as f:
        content = f.read()
        print(content)
except IOError:
    print("I couldn't open the file!")
```


three levels of character processing ...



example: counting words

```
from nltk.tokenize import sent_tokenize, word_tokenize

def compute_word_frequencies(filename):
    frequencies = {}
    with open(filename) as f:
        content = f.read()
        for sen in sent_tokenize(content):
            for word in word_tokenize(sen):
                if word in frequencies:
                    frequencies[word] += 1
                else:
                    frequencies[word] = 1
    return frequencies

freqs = compute_word_frequencies("test.txt")
print(freqs["the"])
```

example: counting bigrams

```
import nltk

def compute_bigram_frequencies(filename):
    ...

bfreqs = compute_bigram_frequencies("test.txt")
print(bfreqs["New York"])
```

example: what's the probability of the next word?

$$P(\text{next word is York} | \text{current word is New}) = \frac{\text{count}(\text{New York})}{\text{count}(\text{New})}$$

```
def transition_probability(w1, w2):  
    ...  
  
print(transition_probability("New", "York"))
```


sorting

- ▶ sometimes we need to **sort** elements of a list (or other collection) into some order:
 - ▶ `some_list.sort()` sorts a list in place
 - ▶ `sorted(some collection)` creates a new list and sorts it

```
the_list = [ 8, 7, 3, 6, 11 ]
```

```
print(sorted(the_list))
```

```
the_list.sort()
```

```
print(the_list)
```

detour: default input values

- ▶ we can define a **default value** for a function input:

```
def count_words(sentence, separator=" "):
    return len(sentence.split(separator))

print(count_words("this is a test sentence"))

print(count_words("this_is_another_sentence"))

print(count_words("this_is_another_sentence", "_"))
```

detour: calling a function with named inputs

- ▶ the inputs can be specified by name instead of order
- ▶ this is particularly useful when there are many inputs

```
def count_words(sentence, separator=" "):
    return len(sentence.split(separator))

print(count_words(sentence="this_is_another_sentence",
                  separator="_"))

print(count_words(separator="_",
                  sentence="this_is_another_sentence"))

print(count_words("this_is_another_sentence",
                  separator="_"))

# illegal!
#print(count_words(separator="_",
                  "this_is_another_sentence"))
```


sorting example

```
def number_of_vowels(w):  
    count = 0  
    for c in w:  
        if c in ['a', 'e', 'i', 'o', 'u']:  
            count += 1  
    return count  
  
word_list = [ "This", "is", "a", "list", "of", "words" ]  
  
print(sorted(word_list))  
print(sorted(word_list, key=len))  
print(sorted(word_list, key=number_of_vowels))  
print(sorted(word_list, key=len, reverse=True))
```

this program will print:

```
['This', 'a', 'is', 'list', 'of', 'words']  
['a', 'is', 'of', 'This', 'list', 'words']  
['This', 'is', 'a', 'list', 'of', 'words']  
['words', 'This', 'list', 'is', 'of', 'a']
```


the most frequent word in a frequency dictionary

```
frequencies = compute_word_frequencies('some_file.txt')  
  
print(max(frequencies, key=frequencies.get))
```

one more data type: tuples

- ▶ **tuples** are fixed-size lists that cannot be changed
 - ▶ a tuple with 2 items is called a *pair*
 - ▶ a tuple with 3 items is called a *triple*
 - ▶ a tuple with n items is called an *n-tuple*
- ▶ tuples are more efficient than normal lists
- ▶ they are written with round brackets: `t = (3, "xyz")`
- ▶ useful fact about tuples: they can be compared and sorted
 - ▶ will sort by first item, then by second item, ...

```
pairs1 = [ (6, "xyz"), (3, "ghi"), (5, "abc") ]
pairs2 = [ ("xyz", 6), ("ghi", 3), ("abc", 5) ]

print(sorted(pairs1))
print(sorted(pairs2))
```


example: sorting alphabetically and by frequency

```
import nltk
def compute_word_frequencies(filename):
    ...
    return frequencies

def get_frequency(word_freq_pair):
    return word_freq_pair[1]

freqs = compute_word_frequencies("test.txt")
word_freq_pairs = freqs.items()
for word_freq_pair in sorted(word_freq_pairs):
    print(word_freq_pair)

for word_freq_pair in sorted(word_freq_pairs,
                             key=get_frequency,
                             reverse=True):
    print(word_freq_pair)
```