

# String Operations



MARCH 19<sup>TH</sup> 2012

# Python has lots of string operations...



- You can find a bunch of these in Section 5.6.1. “String Methods” of the Python documentation (v.2.7.2).
- These are in addition to the operations we studied that are common to lists *and* strings.
  - indexing, slicing
  - membership testing (`in` and `not in`) and concatenation (`+`).
  - `index`, `count`

# String operations



- Here is a categorization (of some of these methods) that might help you navigate the long list of available string operations:
  - **Boolean methods:** isalpha, isalnum, isdigit, islower, isupper, isspace, startswith.
  - **Reformatting methods:** lower, upper, swapcase, capitalize, center, strip, lstrip, rstrip, ljust, rjust.
  - **Split methods:** split, lsplit, rsplit, splitlines.
  - **Join methods:** join.
  - **Replace methods:** replace

# Examples: boolean methods



```
>>> "hello".isalpha()
True
>>> "hello".isalnum()
True
>>> "1234".isdigit()
True
>>> "39.78".isalnum()
False
>>> "hello?".islower()
True
>>> "Hello?".islower()
False
>>> "hello?".startswith("he")
True
>>> "hello?".startswith("He")
False
```

# Examples: Reformatting methods



```
>>> "Hello, how are you?".lower()
'hello, how are you?'
>>> "Hello, how are you?".swapcase()
'hELLO, HOW ARE YOU?'
>>> "jack".capitalize()
'Jack'
>>> " this string has spaces.. ".strip()
'this string has spaces..'
>>> " this string has spaces.. ".lstrip()
'this string has spaces..'
>>> " this string has spaces.. ".rstrip()
'this string has spaces..'
>>> "test".center(20)
'      test      '
>>> "hello??" .rjust(20)
'          hello??'
```

# Split and Join



```
>>> "hello, how are you?".split()
['hello,', 'how', 'are', 'you?']
>>> "Other, characters, can, be, used,,to,split?".split(",")
['Other', ' characters', ' can', ' be', ' used', '', 'to',
'split?']
>>> '''This string
... spans a
... few lines'''.splitlines()
['This string', 'spans a ', 'few lines']
>>> "".join(["hello", "are", "you", "ok?"])
'hello are you ok?'
>>> "??".join(["hello", "are", "you", "ok?"])
'hello??are??you??ok?'
```

# Replace



```
>>> "hello how are you?".replace(" ", "!")
```

```
'hello!how!are!you?'
```

```
>>> "hello, how are you?".replace("h", "")
```

```
'ello, ow are you?'
```

# Problem: Creating a dictionary



Read a text file (e.g., a large novel such as “War and Peace”) that is guaranteed to consist of words that are correctly spelled. Extract “words” from this file and write these out in alphabetical order in a file called “dictionary.txt”.

A word is a contiguous sequence of letters, preceded by a non-letter and followed by a non-letter.

Words in “dictionary.txt” should be unique and should be in lower case.

**Extra credit:** An attempt should be made to avoid proper nouns.



# Version 1



```
#Programmer: Sriram Pemmaraju
#Date: March 13, 2012
#Version 1

# Open the input file "war.txt"
fin = open("smallWar.txt", "r")

wordList = []

# Loop that processes each line of the file
for line in fin:
    wordList = wordList + line.split()

#Close the input file
fin.close()

#Block of code that produces output
fout = open("dictionary1.txt", "w")
for word in wordList:
    fout.write(word+"\n")
fout.close()
```

# Comments on Version 1



- I tried to run this on the full text of “War and Peace” but it was taking too long to complete. So I decided to run it on a smaller version of file consisting of just the first 1000 lines of the original.
- We will revisit the reason for this inefficiency later.
- The output file contains one word per line.
- But, the words contain non-letters and upper case letters.
- They are not in sorted order and surely contain duplicates.

# Version 2



- I made three changes to Version 1 in order to deal with the issue of non-letters.
- First I created a list of all non-letter characters that might be in the text file.
- Notice the use of `ord`, `chr`, and `map` in this code.

```
# List of all non-letter characters
```

```
punctuationMarks = map(chr, range(0, ord("A")) + range(ord("Z")+1, ord("a")) +  
range(ord("z")+1, 127))
```

## Version 2



- I defined a function that takes a line (i.e., a string) and replaces every non-letter in this string by a blank.
- Notice the use of the `replace` method in this code.

```
# Replaces each non-letter character by a blank
def filterOutPunctuation(punctuationMarks, s):
    for mark in punctuationMarks:
        s = s.replace(mark, " ")
    return s
```

## Version 2



- Finally, I process each line by first replacing non-letters by blanks and then splitting at blanks.

# Loop that processes each line of the file

for line in fin:

    newLine = filterOutPunctuation(punctuationMarks, line)

    wordList = wordList + newLine.split()

# Comments on Version 2



- The list of words produced as output no longer contains non-letters.
- However, it does contain upper case letters, is not sorted, and contains duplicates!

# Version 3



- I made three changes to Version 2 in order to deal with the upper case letter issue.
- First I defined a function that takes a list of words and turns this list into words in lower case.

```
def toLower(s):  
    return s.lower()
```

```
def makeListLower(wordList):  
    return map(toLower, wordList)
```

# Version 3



- When each line is processed , the words are turned into lower case words.

# Loop that processes each line of the file

for line in fin:

```
newLine = filterOutPunctuation(punctuationMarks, line)
```

```
wordList = wordList + makeListLower(newLine.split())
```



# Version 3



- The word list is sorted before being printed.

```
#Block of code that produces output
fout = open("dictionary3.txt", "w")
wordList.sort()
for word in wordList:
    fout.write(word+"\n")
fout.close()
```

# Version 4



- Now I created a new version that dealt with the issue of duplicates.
- The only change I made is to the code that produces output.

```
fout = open("dictionary4.txt", "w")
wordList.sort()
previousWord = "" # keeps track of the word most recently output
for word in wordList:
    # only print out new words
    if word != previousWord:
        fout.write(word+"\n")
        previousWord = word
fout.close()
```

# File I/O



- More often than not programs read from files, rather than from input typed at the keyboard.
- Often one program reads what another program outputs.
- More and more, programs are reading data produced by other hardware, e.g., sensors, telescopes, microarrays, etc.
- In these instances very little, if any, input is provided at the keyboard.

# File objects



- Simplest Python statement for opening a file:  

```
f = open("war.txt")
```
- Assuming that there is a file called “war.txt” in the same directory as your Python program, this statement *opens* the file for reading.
- Subsequently, the file can be accessed via the variable **f**.
- Since **f** is a variable, it has a type. Try `type(f)`.

# File objects



- The variable `f` is often called a *file object*.
- If the file is missing from the directory, an error message is issued.

```
>>> g = open("hello.txt")
```

```
Traceback (most recent call last):
```

```
File "<string>", line 1, in <fragment>
```

```
IOError: [Errno 2] No such file or directory: 'hello.txt'
```

- Once a file object is successfully connected to a file residing on your machine, we can use the file object to read from the file in a variety of ways.

# Reading from a file



- `s = f.read()`  
Reads everything from the file into the string `s`
- `s = f.readline()`  
Reads the next line from the files into `s`
- `for line in f:`  
    `print line.split()`  
Allows us to read and process the file line by line

# Let us solve these problems on “War and Peace”



1. Build a dictionary of words extracted from the text that we might be able to use later, maybe in a spellchecker.
2. Compute the number of sentences in the text.
3. Compute the frequencies of letters in the text.

Two useful built-in Python functions that can help in solving Problem 3 are `ord` and `chr`.

# Two useful functions



- `ord(ch)`  
if `ch` is a single character string, this function returns the ASCII code for `ch`
- `chr(i)`  
returns a string of one character whose ASCII code is the integer `i`

What is ASCII?

It stands for the *American Standard Code for Information Interchange*. It assigns a number in the range 0..255 to every character that can be entered at the keyboard.



# More on ASCII



- The numbers 0..31 are reserved for unprintable characters, e.g., the tab character (“\t”), the end of line character (“\n”), etc.
- 32 is the ASCII value of the space character (“ ”)
- 33..47 is used for some punctuation characters
- 48..57 is used for digits “0” through “9”
- 65..90 is used for upper case letters
- 97..122 is used for lower case letters

# ASCII Table



| Dec | Hx | Oct | Char                               | Dec | Hx | Oct | Html  | Chr          | Dec | Hx | Oct | Html  | Chr      | Dec | Hx | Oct | Html   | Chr        |
|-----|----|-----|------------------------------------|-----|----|-----|-------|--------------|-----|----|-----|-------|----------|-----|----|-----|--------|------------|
| 0   | 0  | 000 | <b>NUL</b> (null)                  | 32  | 20 | 040 | &#32; | <b>Space</b> | 64  | 40 | 100 | &#64; | <b>@</b> | 96  | 60 | 140 | &#96;  | <b>`</b>   |
| 1   | 1  | 001 | <b>SOH</b> (start of heading)      | 33  | 21 | 041 | &#33; | <b>!</b>     | 65  | 41 | 101 | &#65; | <b>A</b> | 97  | 61 | 141 | &#97;  | <b>a</b>   |
| 2   | 2  | 002 | <b>STX</b> (start of text)         | 34  | 22 | 042 | &#34; | <b>"</b>     | 66  | 42 | 102 | &#66; | <b>B</b> | 98  | 62 | 142 | &#98;  | <b>b</b>   |
| 3   | 3  | 003 | <b>ETX</b> (end of text)           | 35  | 23 | 043 | &#35; | <b>#</b>     | 67  | 43 | 103 | &#67; | <b>C</b> | 99  | 63 | 143 | &#99;  | <b>c</b>   |
| 4   | 4  | 004 | <b>EOT</b> (end of transmission)   | 36  | 24 | 044 | &#36; | <b>\$</b>    | 68  | 44 | 104 | &#68; | <b>D</b> | 100 | 64 | 144 | &#100; | <b>d</b>   |
| 5   | 5  | 005 | <b>ENQ</b> (enquiry)               | 37  | 25 | 045 | &#37; | <b>%</b>     | 69  | 45 | 105 | &#69; | <b>E</b> | 101 | 65 | 145 | &#101; | <b>e</b>   |
| 6   | 6  | 006 | <b>ACK</b> (acknowledge)           | 38  | 26 | 046 | &#38; | <b>&amp;</b> | 70  | 46 | 106 | &#70; | <b>F</b> | 102 | 66 | 146 | &#102; | <b>f</b>   |
| 7   | 7  | 007 | <b>BEL</b> (bell)                  | 39  | 27 | 047 | &#39; | <b>'</b>     | 71  | 47 | 107 | &#71; | <b>G</b> | 103 | 67 | 147 | &#103; | <b>g</b>   |
| 8   | 8  | 010 | <b>BS</b> (backspace)              | 40  | 28 | 050 | &#40; | <b>(</b>     | 72  | 48 | 110 | &#72; | <b>H</b> | 104 | 68 | 150 | &#104; | <b>h</b>   |
| 9   | 9  | 011 | <b>TAB</b> (horizontal tab)        | 41  | 29 | 051 | &#41; | <b>)</b>     | 73  | 49 | 111 | &#73; | <b>I</b> | 105 | 69 | 151 | &#105; | <b>i</b>   |
| 10  | A  | 012 | <b>LF</b> (NL line feed, new line) | 42  | 2A | 052 | &#42; | <b>*</b>     | 74  | 4A | 112 | &#74; | <b>J</b> | 106 | 6A | 152 | &#106; | <b>j</b>   |
| 11  | B  | 013 | <b>VT</b> (vertical tab)           | 43  | 2B | 053 | &#43; | <b>+</b>     | 75  | 4B | 113 | &#75; | <b>K</b> | 107 | 6B | 153 | &#107; | <b>k</b>   |
| 12  | C  | 014 | <b>FF</b> (NP form feed, new page) | 44  | 2C | 054 | &#44; | <b>,</b>     | 76  | 4C | 114 | &#76; | <b>L</b> | 108 | 6C | 154 | &#108; | <b>l</b>   |
| 13  | D  | 015 | <b>CR</b> (carriage return)        | 45  | 2D | 055 | &#45; | <b>-</b>     | 77  | 4D | 115 | &#77; | <b>M</b> | 109 | 6D | 155 | &#109; | <b>m</b>   |
| 14  | E  | 016 | <b>SO</b> (shift out)              | 46  | 2E | 056 | &#46; | <b>.</b>     | 78  | 4E | 116 | &#78; | <b>N</b> | 110 | 6E | 156 | &#110; | <b>n</b>   |
| 15  | F  | 017 | <b>SI</b> (shift in)               | 47  | 2F | 057 | &#47; | <b>/</b>     | 79  | 4F | 117 | &#79; | <b>O</b> | 111 | 6F | 157 | &#111; | <b>o</b>   |
| 16  | 10 | 020 | <b>DLE</b> (data link escape)      | 48  | 30 | 060 | &#48; | <b>0</b>     | 80  | 50 | 120 | &#80; | <b>P</b> | 112 | 70 | 160 | &#112; | <b>p</b>   |
| 17  | 11 | 021 | <b>DC1</b> (device control 1)      | 49  | 31 | 061 | &#49; | <b>1</b>     | 81  | 51 | 121 | &#81; | <b>Q</b> | 113 | 71 | 161 | &#113; | <b>q</b>   |
| 18  | 12 | 022 | <b>DC2</b> (device control 2)      | 50  | 32 | 062 | &#50; | <b>2</b>     | 82  | 52 | 122 | &#82; | <b>R</b> | 114 | 72 | 162 | &#114; | <b>r</b>   |
| 19  | 13 | 023 | <b>DC3</b> (device control 3)      | 51  | 33 | 063 | &#51; | <b>3</b>     | 83  | 53 | 123 | &#83; | <b>S</b> | 115 | 73 | 163 | &#115; | <b>s</b>   |
| 20  | 14 | 024 | <b>DC4</b> (device control 4)      | 52  | 34 | 064 | &#52; | <b>4</b>     | 84  | 54 | 124 | &#84; | <b>T</b> | 116 | 74 | 164 | &#116; | <b>t</b>   |
| 21  | 15 | 025 | <b>NAK</b> (negative acknowledge)  | 53  | 35 | 065 | &#53; | <b>5</b>     | 85  | 55 | 125 | &#85; | <b>U</b> | 117 | 75 | 165 | &#117; | <b>u</b>   |
| 22  | 16 | 026 | <b>SYN</b> (synchronous idle)      | 54  | 36 | 066 | &#54; | <b>6</b>     | 86  | 56 | 126 | &#86; | <b>V</b> | 118 | 76 | 166 | &#118; | <b>v</b>   |
| 23  | 17 | 027 | <b>ETB</b> (end of trans. block)   | 55  | 37 | 067 | &#55; | <b>7</b>     | 87  | 57 | 127 | &#87; | <b>W</b> | 119 | 77 | 167 | &#119; | <b>w</b>   |
| 24  | 18 | 030 | <b>CAN</b> (cancel)                | 56  | 38 | 070 | &#56; | <b>8</b>     | 88  | 58 | 130 | &#88; | <b>X</b> | 120 | 78 | 170 | &#120; | <b>x</b>   |
| 25  | 19 | 031 | <b>EM</b> (end of medium)          | 57  | 39 | 071 | &#57; | <b>9</b>     | 89  | 59 | 131 | &#89; | <b>Y</b> | 121 | 79 | 171 | &#121; | <b>y</b>   |
| 26  | 1A | 032 | <b>SUB</b> (substitute)            | 58  | 3A | 072 | &#58; | <b>:</b>     | 90  | 5A | 132 | &#90; | <b>Z</b> | 122 | 7A | 172 | &#122; | <b>z</b>   |
| 27  | 1B | 033 | <b>ESC</b> (escape)                | 59  | 3B | 073 | &#59; | <b>;</b>     | 91  | 5B | 133 | &#91; | <b>[</b> | 123 | 7B | 173 | &#123; | <b>{</b>   |
| 28  | 1C | 034 | <b>FS</b> (file separator)         | 60  | 3C | 074 | &#60; | <b>&lt;</b>  | 92  | 5C | 134 | &#92; | <b>\</b> | 124 | 7C | 174 | &#124; | <b> </b>   |
| 29  | 1D | 035 | <b>GS</b> (group separator)        | 61  | 3D | 075 | &#61; | <b>=</b>     | 93  | 5D | 135 | &#93; | <b>]</b> | 125 | 7D | 175 | &#125; | <b>}</b>   |
| 30  | 1E | 036 | <b>RS</b> (record separator)       | 62  | 3E | 076 | &#62; | <b>&gt;</b>  | 94  | 5E | 136 | &#94; | <b>^</b> | 126 | 7E | 176 | &#126; | <b>~</b>   |
| 31  | 1F | 037 | <b>US</b> (unit separator)         | 63  | 3F | 077 | &#63; | <b>?</b>     | 95  | 5F | 137 | &#95; | <b>_</b> | 127 | 7F | 177 | &#127; | <b>DEL</b> |

# Some examples of `chr` and `ord` in action



```
>>> ord("a")
```

```
97
```

```
>>> chr(97)
```

```
'a'
```

```
>>> ord(" ")
```

```
32
```

```
>>> ord("o")
```

```
48
```

```
>>> chr(48)
```

```
'0'
```

```
>>> chr(49)
```

```
'1'
```

```
>>> ord("A")
```

```
65
```

```
>>> ord("B")
```

```
66
```

# How are these functions useful?



- Because of the the fact that all the upper case letters occur consecutively in the ASCII table, the expression  $\text{ord}(ch) - \text{ord}("A")$  has value 0 for  $ch = "A"$ , value 1 for  $ch = "B"$ , has value 2 for  $ch = "C"$ , etc.
- Similarly,  $\text{ord}(ch) - \text{ord}("a")$  has value 0 for  $ch = "a"$ , has value 1 for  $ch = "b"$ , has value 2 for  $ch = "c"$ , etc.

# A program to count letter frequencies



```
f = open("war.txt")
L = [0]*26
s = f.read()
for ch in s:
    if ch.isupper():
        L[ord(ch)-ord("A")] = L[ord(ch)-ord("A")] + 1
    elif ch.islower():
        L[ord(ch)-ord("a")] = L[ord(ch)-ord("a")] + 1
print L
```

Notice how `ord(ch)-ord("A")` and `ord(ch)-ord("a")` are used to index into the list `L`.

# Another example



- The `ord` and `chr` functions can be used to perform Caesar's Cipher (Problem 3, HW 7).
- Try this: `chr(ord("a") + 4)`
- What does this expression evaluate to?