# Programming Principles in Python (CSCI 503/490)

### Strings & Files

Dr. David Koop

(some slides adapted from Dr. Reva Freedman)





# Memoization

- memo dict =  $\{\}$ def memoized slow function(s, t): if (s, t) not in memo dict: return memo dict[(s, t)] • for s, t in [(12, 10), (4, 5), (5, 4), (12, 10)]: pass else: c = compute fast function(s, t)
- Second time executing for s=12, t=10, we don't need to compute!
- Tradeoff memory for compute time

### D. Koop, CSCI 503/490, Fall 2024

memo dict[(s, t)] = compute slow function(s, t)if s > t and (c := memoized slow function(s, t) > 50):





2

# Functional Programming

- Programming without imperative statements like assignment
- In addition to comprehensions & iterators, have functions:
  - map: iterable of n values to an iterable of n transformed values
  - filter: iterable of n values to an iterable of m (m  $\leq$  n) values
- Eliminates need for concrete looping constructs









# Lambda Functions

- def is even(x): return (x % 2) == 0
- filter(is even, range(10) # generator
- Lots of code to write a simple check
- Lambda functions allow inline function definition
- Usually used for "one-liners": a simple data transform/expression
- filter(lambda x: x & 2 == 0, range(10))
- Parameters follow lambda, no parentheses
- No return keyword as this is implicit in the syntax
- JavaScript has similar functionality (arrow functions): (d = 0)





# Unicode and ASCII

- Conceptual systems
- ASCII:
  - old, English-centric, 7-bit system (only 128 characters)
- Unicode:
  - Can represent over 1 million characters from all languages + emoji 🎉
  - Characters have hexadecimal representation:  $\acute{e} = U+00E9$  and name (LATIN SMALL LETTER E WITH ACUTE)
  - Python allows you to type "é" or represent via code "\u00e9"
- Codes: ord  $\rightarrow$  character to integer, chr  $\rightarrow$  integer to character





![](_page_4_Picture_14.jpeg)

![](_page_4_Picture_15.jpeg)

# Strings

- Objects with methods
- Finding and counting substrings: count, find, startswith
- Removing leading & trailing substrings/whitespace: strip, removeprefix
- Transforming Text: replace, upper, lower, title
- Checking String Composition: isalnum, isnumeric, isupper

![](_page_5_Picture_7.jpeg)

![](_page_5_Picture_9.jpeg)

![](_page_5_Picture_10.jpeg)

![](_page_5_Picture_11.jpeg)

# <u>Assignment 4</u>

- Assignment will cover strings and files
- Reading & writing data to files
- Deals with characters and formatting

![](_page_6_Picture_7.jpeg)

![](_page_6_Picture_9.jpeg)

# Splitting

- s = "Venkata, Ranjit, Pankaj, Ali, Karthika"
- names = s.split(',') # names is a list
- names = s.split(',', 3) # split by commas, split <= 3 times
- separator may be multiple characters
- if no separator is supplied (sep=None), runs of consecutive whitespace delimit elements
- rsplit works in reverse, from the right of the string
- partition and rpartition for a single split with before, sep, and after
- splitlines splits at line boundaries, optional parameter to keep endings

### D. Koop, CSCI 503/490, Fall 2024

![](_page_7_Picture_10.jpeg)

![](_page_7_Picture_12.jpeg)

8

# Joining

- join is a method on the **separator** used to join a list of strings
- ', '.join(names)
  - names is a list of strings, ', ' is the separator used to join them
- Example:
  - def orbit(n): # ... return orbit as list print(','.join(orbit as list))

![](_page_8_Picture_10.jpeg)

![](_page_8_Picture_12.jpeg)

![](_page_8_Picture_14.jpeg)

# Formatting

- with specified width
- s.zfill: ljust with zeroes
- s.format: templating function
  - Replace fields indicated by curly braces with corresponding values
  - "My name is {} {}".format(first name, last name)
  - "My name is {1} {0}".format(last name, first name)
  - "My name is {first name} {last name}.format(
  - Braces can contain number or name of keyword argument
  - Whole format mini-language to control formatting

### D. Koop, CSCI 503/490, Fall 2024

### • s.ljust, s.rjust: justify strings by adding fill characters to obtain a string

# first name=name[0], last name=name[1])

![](_page_9_Picture_13.jpeg)

![](_page_9_Picture_15.jpeg)

![](_page_9_Picture_16.jpeg)

# Format Strings

- Formatted string literals (f-strings) prefix the starting delimiter with f
- Reference variables **directly**!
  - f"My name is {first name} {last name}"
- Can include expressions, too:
  - f"My name is {name[0].capitalize()} {name[1].capitalize()}"
- Same format mini-language is available

![](_page_10_Picture_10.jpeg)

![](_page_10_Picture_12.jpeg)

11

# Format Mini-Language Presentation Types

- Not usually required for obvious types
- :d for integers
- : c for characters
- :s for strings
- :e or :f for floating point
  - e: scientific notation (all but one digit after decimal point)
  - f: fixed-point notation (decimal number)

![](_page_11_Picture_11.jpeg)

![](_page_11_Picture_13.jpeg)

# Field Widths and Alignments

- After : but before presentation type
  - f'[{27:10d}]' # '[ 27]'
  - f'[{"hello":10}]' # '[hello
- Shift alignment using < or >:
  - f'[{"hello":>15}]' # '[
- Center align using ^:
  - f'[{"hello":^7}]' # '[ hello ]'

![](_page_12_Figure_9.jpeg)

![](_page_12_Picture_10.jpeg)

![](_page_12_Picture_11.jpeg)

![](_page_12_Picture_13.jpeg)

# Numeric Formatting

- Add positive sign: - f'[{27:+10d}]' # '[ +27]'
- Add space but only show negative numbers:
- Separators:
  - f' { 12345678:, d } ' # '12, 345, 678'

### D. Koop, CSCI 503/490, Fall 2024

### - print(f'{27: d}\n{-27: d}') # note the space in front of 27

![](_page_13_Picture_11.jpeg)

![](_page_13_Picture_13.jpeg)

![](_page_13_Picture_14.jpeg)

# Raw Strings

- Raw strings prefix the starting delimiter with r
- Disallow escaped characters
- '\\n is the way you write a newline, \\\\ for \\.'
- r"\n is the way you write a newline, \\ for  $\$ ."
- Useful for regular expressions

![](_page_14_Picture_9.jpeg)

![](_page_14_Picture_11.jpeg)

## Regular Expressions

- AKA regex
- A syntax to better specify how to decompose strings
- Look for patterns rather than specific characters
- "31" in "The last day of December is 12/31/2016."
- May work for some questions but now suppose I have other lines like: "The last day of September is 9/30/2016."
- ...and I want to find dates that look like:
- {digits}/{digits}/{digits}
- Cannot search for every combination!
- \d+/\d+/\d+ # \d is a character class

![](_page_15_Picture_12.jpeg)

![](_page_15_Picture_14.jpeg)

### Metacharacters

- this is optional.
- . ^ \$ \* + ? { } [ ] \ | ( )
- []: define character class
- ^: complement (opposite)
- \: escape, but now escapes metacharacters and references classes
- \*: repeat zero or more times
- +: repeat one or more times
- ?: zero or one time
- {m,n}: at least m and at most n

### D. Koop, CSCI 503/490, Fall 2024

### Need to have some syntax to indicate things like repeat or one-of-these or

![](_page_16_Picture_13.jpeg)

![](_page_16_Picture_15.jpeg)

# Predefined Character Classes

	Character class
	\d
Any c	$\setminus D$
Any whitespace chara	\ S
Any character	\ S
Any word character	\ W
Any charac	$\setminus W$

D. Koop, CSCI 503/490, Fall 2024

### Matches

Any digit (0–9).

haracter that is *not* a digit.

acter (such as spaces, tabs and newlines).

that is *not* a whitespace character.

(also called an **alphanumeric character**)

cter that is *not* a word character.

![](_page_17_Picture_11.jpeg)

![](_page_17_Picture_13.jpeg)

![](_page_17_Picture_14.jpeg)

# Performing Matches

Method/Attribute	
match()	Determine
search()	Scan throu
findall()	Find all su
finditer()	Find all su

D. Koop, CSCI 503/490, Fall 2024

### Purpose

if the RE matches at the beginning of the string.

ugh a string, looking for any location where this RE matches.

bstrings where the RE matches, and returns them as a list.

bstrings where the RE matches, and returns them as an iterator.

![](_page_18_Picture_8.jpeg)

![](_page_18_Picture_10.jpeg)

# Regular Expressions in Python

- import re
- re.match(<pattern>, <str to check>)
  - Returns None if no match, information about the match otherwise
  - Starts at the **beginning** of the string
- re.search(<pattern>, <str to check>)
  - Finds single match anywhere in the string
- re.findall(<pattern>, <str to check>)
  - Finds **all** matches in the string, search only finds the first match
- Can pass in flags to alter methods: e.g. re.IGNORECASE

![](_page_19_Picture_13.jpeg)

![](_page_19_Picture_15.jpeg)

![](_page_19_Picture_17.jpeg)

### Examples

- s0 = "No full dates here, just 02/15"s1 = "02/14/2021 is a date"
  - s2 = "Another date is 12/25/2020"
- re.match(r'\d+/\d+/\d+',s1) # returns match object
- re.match(r'\d+/\d+/\d+',s0) # None
- re.match(r'\d+/\d+/\d+',s2) # None!
- re.search(r'\d+/\d+/\d+',s2) # returns 1 match object
- re.search(r'\d+/\d+/\d+',s3) # returns 1! match object
- re.findall(r'\d+/\d+/\gammas3) # returns list of strings
- re.finditer(r'\d+/\d+/\d+',s3) # returns iterable of matches

![](_page_20_Picture_14.jpeg)

![](_page_20_Picture_16.jpeg)

![](_page_20_Picture_17.jpeg)

# Grouping

- Parentheses capture a group that can be accessed or used later • Access via groups () or group (n) where n is the number of the group, but
- numbering starts at 1
- Note: group (0) is the full matched string
- for match in re.finditer(r'(d+)/(d+)/(d+)', s3): print(match.groups())
- for match in re.finditer(r'(d+)/(d+)/(d+)', s3): print (  $\{2\} - \{0:02d\} - \{1:02d\}$  '.format ( \*[int(x) for x in match.groups()])) operator expands a list into individual elements

![](_page_21_Picture_9.jpeg)

![](_page_21_Picture_11.jpeg)

![](_page_21_Picture_12.jpeg)

![](_page_21_Picture_13.jpeg)

# Modifying Strings

Method/Attribute	Purpose
split()	Split the strir RE matches
sub()	Find all substructions for the second
subn()	Does the sar string and the

### D. Koop, CSCI 503/490, Fall 2024

ng into a list, splitting it wherever the

trings where the RE matches, and n with a different string

me thing as sub(), but returns the new e number of replacements

![](_page_22_Picture_6.jpeg)

![](_page_22_Picture_8.jpeg)

![](_page_22_Picture_9.jpeg)

![](_page_22_Picture_10.jpeg)

# Substitution

- Do substitution in the middle of a string: • re.sub(r'(\d+)/(\d+)/(\d+)',r'\3-\1-\2',s3)
- All matches are substituted
- First argument is the regular expression to match
- Second argument is the substitution
- $-1, 2, \dots$  match up to the **captured groups** in the first argument Third argument is the string to perform substitution on
- Can also use a **function**:
- to date = lambda m: f'{m.group(3)}-{int(m.group(1)):02d}-{int(m.group(2)):02d}' re.sub(r'(\d+)/(\d+)/(\d+)', to date, s3)

![](_page_23_Picture_10.jpeg)

![](_page_23_Picture_12.jpeg)

![](_page_23_Picture_13.jpeg)

![](_page_23_Picture_14.jpeg)

### D. Koop, CSCI 503/490, Fall 2024

### Files

![](_page_24_Picture_2.jpeg)

![](_page_24_Picture_4.jpeg)

![](_page_24_Picture_5.jpeg)

![](_page_24_Picture_6.jpeg)

### Files

- A file is a sequence of data stored on disk.
- Python uses the standard Unix newline character (n) to mark line breaks.
  - On Windows, end of line is marked by  $\r\n$ , i.e., carriage return + newline.
  - On old Macs, it was carriage return  $\r$  only.
  - Python **converts** these to n when reading.

![](_page_25_Picture_8.jpeg)

![](_page_25_Picture_10.jpeg)

![](_page_25_Picture_11.jpeg)

![](_page_25_Picture_12.jpeg)

# Opening a File

- handle).
- We access the file via the file object.
- <filevar> = open(<name>, <mode>)
- Mode 'r' = read or 'w' = write, 'a' = append
- read is default

### • Opening associates a file on disk with an object in memory (file object or file)

### • Also add 'b' to indicate the file should be opened in binary mode: 'rb','wb'

![](_page_26_Picture_11.jpeg)

![](_page_26_Picture_13.jpeg)

27

# Standard File Objects

- When Python begins, it associates three standard file objects:
  - sys.stdin: for input
  - sys.stdout: for output
  - sys.stderr: for errors
- In the notebook
  - sys.stdin isn't really used, get input can be used if necessary
  - sys.stdout is the output shown after the code
  - sys.stderr is shown with a red background

![](_page_27_Picture_12.jpeg)

![](_page_27_Picture_14.jpeg)

![](_page_27_Picture_15.jpeg)

![](_page_27_Picture_16.jpeg)

# Files and Jupyter

- You can **double-click** a file to see its contents (and edit it manually) • To see one as text, may need to right-click
- Shell commands also help show files in the notebook
- The ! character indicates a shell command is being called
- These will work for Linux and macos but not necessarily for Windows
- !cat <fname>: print the entire contents of <fname>
- !head -n <num> <fname>: print the first <num> lines of <fname>
- !tail -n <num> <fname>: print the last <num> lines of <fname>

![](_page_28_Picture_10.jpeg)

![](_page_28_Picture_12.jpeg)

![](_page_28_Picture_13.jpeg)

![](_page_28_Picture_14.jpeg)

# Reading Files

• Use the open () method to open a file for reading

- f = open('huck-finn.txt')

- f = open('huck-finn.txt', 'r')
- Usually, add an 'r' as the second parameter to indicate read (default) • Can iterate through the file (think of the file as a collection of lines):
  - for line in f:

if 'Huckleberry' in line: print(line.strip())

- Using line.strip() because the read includes the newline, and print writes a newline so we would have double-spaced text
- Closing the file: f.close()

![](_page_29_Picture_12.jpeg)

![](_page_29_Picture_14.jpeg)

![](_page_29_Picture_15.jpeg)