Coding Bootcamp Code in Python

# FILES:
# I/O AND DATA FORMATS

# Reading from files

- Reading from text files, line by line
  - E.g., read file line by line, convert to uppercase, and print

```
1   with open(file_name, 'r') as text_file:
2       for line in text_file:
3           print(line.upper())
```

with …: context manager

- Reading from a binary file, value by value
  - E.g., read doubles (8 bytes) and print

```
1   from struct import unpack
2   with open(file_name, 'rb') as bin_file:
3       double_bytes = bin_file.read(8)
4       while double_bytes:
5           print(unpack('d', double_bytes)[0])
6           double_bytes = bin_file.read(8)
```

Not portable!!!:
data type size?
Encoding?
little /big endian?

# Libraries & data formats

- Standard library (Python 3.x)
  - Comma separated value files: `csv`
  - Configuration files: `ConfigParser`
  - Semi-structured data: `json`, `htmllib`, `sgmllib`, `xml`
- Non-standard libraries
  - Images: scikit-image
  - HDF5: pytables
  - pandas
  - Bioinformatics: Biopython

Use the "batteries" that are included!

# Data formats: CSV

Let Sniffer figure out
CSV dialect (e.g., Excel)

```
0   from csv import Sniffer, DictReader
1   with open(file_name, 'rb') as csv_file:
2       dialect = Sniffer().sniff(csv_file.read(1024))
3       csv_file.seek(0)
4       sum = 0.0
5       csv_reader = DictReader(csv_file, fieldnames=None,
6                                  restkey='rest', restval=None,
7                                  dialect=dialect)
8       for row in csv_reader:
9           print('{name} --- {weight}'.format(name=row['name'],
10                                 weight=row['weight']))
11          sum += float(row['weight'])
12      print('sum = {0}'.format(sum))
```

DictReader uses first
row to deduce field names

Access fields by name,
thanks to DictReader

Drawback: you still need to know field types

# Writing to files

- Writing to text files

  Note: `'w'` *replaces* existing file, use `'x'` to avoid

  – E.g., compute and write to file

```
1  with open(file_name, 'w') as text_file:
2      for i in range(0, 10):
3          text_file.write('{0}: {1}\n'.format(i, i*i))
```

- Append to text files
  – E.g., add some more squares to same file

```
1  with open(file_name, 'a') as text_file:
2      for i in range(10, 20):
3          text_file.write('{0}: {1}\n'.format(i, i*i))
```
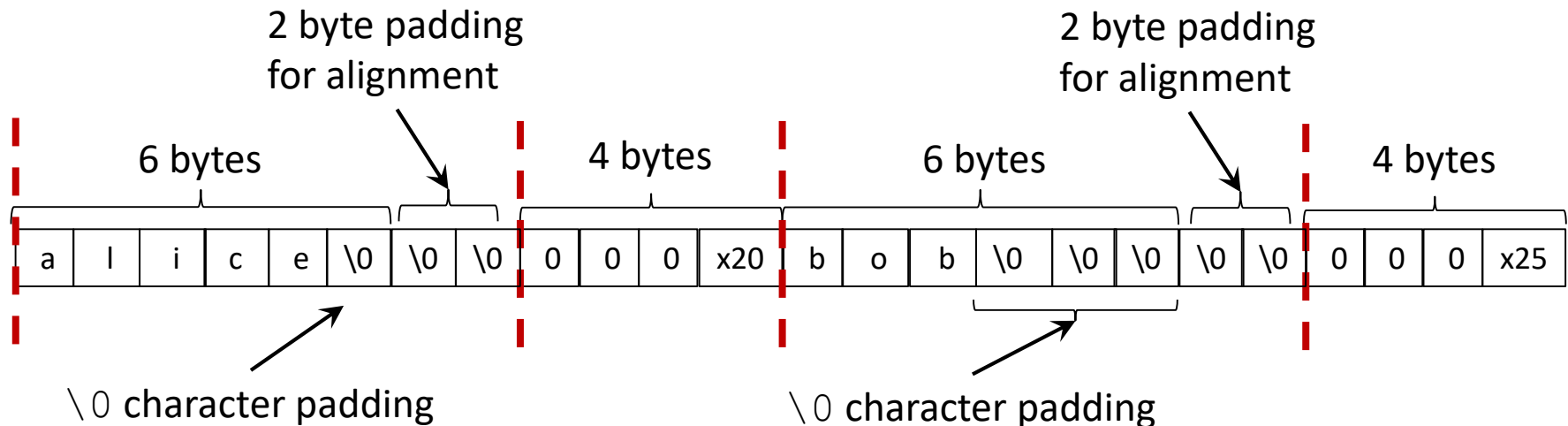
- Writing binary files: don't go there…

# … unless you have to

- ## Writing to binary files

byte representation of name truncated to 6 characters

```
1    from struct import pack
2    with open(file_name, 'wb') as bin_file:
3        for name, age in people:
4            bin_file.write(pack('6si',
5                               bytes(name, 'ascii'),
6                               age))
```

2 byte padding for alignment

2 byte padding for alignment

6 bytes    4 bytes    6 bytes    4 bytes

| a | l | i | c | e | \0 | \0 | \0 | 0 | 0 | 0 | x20 | b | o | b | \0 | \0 | \0 | \0 | \0 | 0 | 0 | 0 | x25 |

\0 character padding

\0 character padding

# Data formats: XML output

```xml
<?xml version="1.0" ?>
<blocks>
  <block name="block_01">
    <item>
      0.1
    </item>
    <item>
      1.1
    </item>
  </block>
  <block name="block_02">
    <item>
      0.2
    </item>
    <item>
      1.2
    </item>
  </block>
</blocks>
```

# Data formats: creating XML

```python
1    from xml.dom.minidom import Document
2    nr_blocks = 2
3    nr_items = 2
4    doc = Document()
5    blocks = doc.createElement('blocks')
6    doc.appendChild(blocks)
7    for block_nr in range(1, nr_blocks + 1):
8        block = doc.createElement('block')
9        block_name = 'block_{0:02d}'.format(block_nr)
10       block.setAttribute('name', block_name)
11       blocks.appendChild(block)
12       for item_nr in range(0, nr_items):
13           item = doc.createElement('item')
14           text = '{0}.{1}'.format(item_nr, block_nr)
15           text_node = doc.createTextNode(text)
16           item.appendChild(text_node)
17           block.appendChild(item)
18   print(doc.toprettyxml(indent='  '))
```

# Code Pack 09

A. Create a database with Text Files

B. Working with CSV

Coding Bootcamp Code in Python

# WEB SCRAPING:
# GATHERING DATA FROM THE WEB

# Introduction

- Caveat: web scraping code is brittle, typically not robust against
  - page layout changes (unless proper use of CSS)
  - page content changes
  - site redesign

Use site APIs (e.g., REST interface) whenever available!

- Many frameworks available, here Beautiful Soup
- However, for tables only, consider pandas

# Beautiful Soup

- ## Open web page using `urllib`

```
import urllib
…
page = urllib.request.urlopen(page_url)
```

Note: urllib2 for Python 2.x

- ## Cook soup out of opened page

```
from bs4 import BeautifulSoup
…
soup = BeautifulSoup(page, "html5lib")
```

- ## Eat soup

```
print('looking at {0}'.format(soup.title.string))
```

Assumes page has a `title` element

# Finding stuff

- First element with tag, e.g., `a`

```
print('looking at {0}'.format(soup.a))
```

- All element with tag, e.g., `a`

```
for a in soup.find_all('a'):
    print('a element: {0}'.format(a))
```

- Element content, e.g., `a`

```
for a in soup.find_all('a'):
    print('link text: {0}'.format(a.string))
```

- Element attribute, e.g., `href` in `a`

```
for a in soup.find_all('a'):
    print('link url: {0}'.format(a.get('href')))
```

# Code Pack 10

A. Try the Web Scraping Code

Coding Bootcamp Code in Python

# ERRORS:
# DEALING WITH EXCEPTIONS

# Errors

```
…
def main():
    file_name = sys.argv[1]
    with open(file_name) as in_file:
        for line in in_file:
            print('|{0}|'.format(line.rstrip('\r\n')))
    return 0
…
```

```
$ python quote.py
Traceback (most recent call last):
  File "./quote.0.py", line 13, in <module>
    status = main()
  File "./quote.py", line 6, in main
    file_name = sys.argv[1]
IndexError: list index out of range
```

exception thrown

Either check length of sys.argv, or deal with error!

# Playing catch

```
…
def main():
    try:
        file_name = sys.argv[1]
    except IndexError as e:
        sys.stderr.write('### error: no input file\n')
        return 1
    with open(file_name) as in_file:
        for line in in_file:
            print('|{0}|'.format(line.rstrip('\r\n')))
    return 0
…
```

```
$ python quote.py
### error: no input file
```

# More trouble

```
$ python quote.py  bla
Traceback (most recent call last):
  File "./quote.py", line 17, in <module>
    status = main()
  File "./quote.py", line 11, in main
    with open(file_name) as in_file:
IOError: [Errno 2] No such file or directory: 'bla'
```

exception thrown

# Catching more

```
…
def main():
    try:
        file_name = sys.argv[1]
        in_file = open(file_name)
        with in_file:
            for line in in_file:
                print('|{0}|'.format(line.rstrip('\r\n')))
    except IndexError as e:
        sys.stderr.write('### error: no input file\n')
        return 1
    except IOError as e:
        msg = "### I/O error on '{0}': {1}".format(e.filename,
                                                   e.strerror)
        sys.stderr.write(msg)
        return 2
    return 0
…
```

# All handled!

- Now all exceptions are handled

```
$ python quote.py  bla
### I/O error on 'bla': No such file or directory
```

- Note that code size increased from 5 to 16 lines
  - Handling errors takes effort
  - Worthwhile if others are using your software!
- One can create own exceptions, derive class from `Exception`

# Code Pack 11

- A. Python fundamentals:
- ~~1. Primitive Datatypes and Operators~~
- ~~2. Variables and Collections~~
- ~~3. Control Flow and Iterables~~
- ~~4. Functions~~
- ~~5. Modules~~
- ~~6. Classes~~
- 7. Advanced