

# Programming Data Management & Analysis

Daniel Dörr

Faculty of Technology, Bielefeld University

```
332
333
334     if extrapolate is None:
335         extrapolate = self.extrapolate
336     x = np.asarray(x)
337     x_shape, x_ndim = x.shape, x.ndim
338     x = np.ascontiguousarray(x.ravel(), dtype=np
339
340     # With periodic extrapolation we map x to the
341     # [self.t[k], self.t[n]].
342     if extrapolate == 'periodic':
343         n = self.t.size - self.k - 1
344         x = self.t[self.k] + (x - self.t[self.k]) *
345         extrapolate = False
346
347     out = np.empty((len(x), prod(self.c.shape[1:])),
348                   dtype=self._ensure_c_contiguous())
349     self._evaluate(x, nu, extrapolate, out)
350     out = out.reshape(x_shape + self.c.shape[1:])
351     if self.axis != 0:
352         # transpose to move the calculated values to t
353         l = list(range(out.ndim))
354         l = l[x_ndim:x_ndim+self.axis] + l[:x_ndim] +
355         out = out.transpose(l)
356     return out
357
358 def _evaluate(self, xp, nu, extrapolate, out):
359     _bspl.evaluate_spline(self.t, self.c.reshape(self.c
360     self.k, xp, nu, extrapolate, out)
361
362 def _ensure_c_contiguous(self):
363     """
364     c and t may be modified by the user. The Cython code
365     c and t are C contiguous.
366     """
367     if not self.c.flags.c_contiguous:
368         self.c = np.ascontiguousarray(self.c)
369     if not self.t.flags.c_contiguous:
370         self.t = np.ascontiguousarray(self.t)
```

# Recap

# for-Loop: Iteration over ordered collections

## Loop over elements

```
1 # tuple filled with arbitrary elements  
2 my_tuple = (1, 2.0, 'text', list(), dict())  
3  
4 # for-loop over my_tuple with control  
   variable 'el'  
5 for el in my_tuple:  
6     msg = 'element: {}'.format(el)  
7     print(msg)
```

# while loop: conditional iteration

Loops until condition becomes True

```
1 x = 5
2 while x > 0:
3     print(x)
4     x -= 1 # shorthand for x = x - 1
```

# Functions and classes—examples of code reuse

```
1 class Library:
2     description = 'This is a Library'
3
4     def __init__(self, name):
5         # name the library
6         self.name = name
7         # create empty book storage on initialization
8         self.storage = list()
9
10    def addBook(self, book):
11        self.storage.append(book)
12
13    def getAllBooks(self):
14        return tuple(self.storage)
15
16 myLib = Library('Bodleian Library')
17 myLib.addBook('The Art of Computer Programming (D. Knuth)')
```

# Modules—examples of code reuse

## mystringutils.py

```
1 #
2 # A module for all kinds of string utils
3 #
4
5 def findSubstringInStrings(stringCollection,
6                             pattern):
7     occ = list()
8     for i, s in enumerate(stringCollection):
9         j = s.find(pattern)
10        while j != -1:
11            occ.append((i, j))
12            j = s.find(pattern, j+1)
13    return occ
```

## mymain.py

```
1 #!/usr/bin/env python3
2
3 from mystringutils import
4     findSubstringInStrings
5
6 if __name__ == '__main__':
7     myStringList = ['the rain in spain',
8                    'ain\t no sunshine',
9                    'she was greeted with disdain']
10
11     occOfAin = findSubstringInStrings(
12         myStringList, 'ain')
13     print(occOfAin)
```

# Python debugger—example

```
1 # dictionary filled with arbitrary elements
2 my_dict = {'key': 'value', 1: 'text', (1, 2)
3           : 'text'}
4
5 # invoke Python debugger
6 breakpoint()
7
8 # for-loop over keys of my_dict with control
9 variable 'key'
10 for key in my_dict:
11     my_dict[(key, 1, 2, 3)] = 'new element'
```

**Input & Out-  
put**

**File Formats**

**Jupyter  
Notebook**

**Text mining**



# Interactive reading from console

- ❖ Reading a string from console:

```
my_string = input()
```

- ❖ Specify prompt:

```
year_str = input('When did the Lakers win'  
+ ' their last championship? ')
```

# Reading from command line

---

## example\_input\_argument.py

---

```
1 #!/usr/bin/env python3
2 from sys import argv
3
4 if __name__ == '__main__':
5     my_arg1 = argv[1]
6     my_arg2 = argv[2]
7     print('1st input argument:', my_arg1)
8     print('2nd input argument:', my_arg2)
```

# Reading from file

... like a book: open & read!

```
1 f = open('Frankenstein.txt')  
2 my_text = f.read()
```

# Reading from file

alternatively, use “with” statement:

```
1 with open('Frankenstein.txt') as f:  
2     my_text = f.read()
```

# Reading from file

read file line-by-line:

```
1 lines = list()
2 f = open('Frankenstein.txt',
3         newline='\n')
4 for line in f:
5     lines.append(line)
```

# Reading from file

Dynamic: read from file with name requested by prompt

```
1 fName = input('Input file: ')
2 lines = list()
3 f = open(fName)
4 for line in f:
5     lines.append(line)
```

# Writing to file

... just as simple as reading!

```
1 f = open('letter1.txt', 'w')  
2 f.write('TO Mrs. Saville, England')  
3 f.flush()
```

# Writing to file

`close()` flushes, then closes the file:

```
1 f = open('letter1.txt', 'w')  
2 f.write('TO Mrs. Saville, England')  
3 f.close()
```



# Writing to file

with automatically closes the file:

```
1 with open('letter1.txt', 'w') as f:  
2     f.write('TO Mrs. Saville, England')
```

# Writing to file

Direct printout to file::

```
1 with open('output.txt') as f:  
2     print('TO Mrs. Saville, England',  
          file = f)
```

# Quiz

- Command line arguments are received through the

`input()` function      `argv` list

- Complete the code for reading a file by filling in the blanks:

```
1 _____ open('myfile.txt') as _____:  
2     contents = f._____()
```

- Which function(s) invoke(s) the writing of file buffer data to the file system?

`clear`      `close`      `write`      `buffer`      `flush`

# Quiz

- Command line arguments are received through the

`input()` function      argv list ✓

- Complete the code for reading a file by filling in the blanks:

```
1 with open('myfile.txt') as f:  
2     contents = f.read()
```

- Which function(s) invoke(s) the writing of file buffer data to the file system?

clear      close ✓      write      **buffer**      flush ✓

**File I/O**

**File Formats**

**Jupyter  
Notebook**

**Text mining**

# Unstructured data: plain text

... like a book: open & read!

```
1 f = open('Frankenstein.txt')  
2 my_text = f.read()
```

# Structured data: XML

EXtensible Markup Language: *a hierarchical data structure*

```
1 <book category="Python">
2   <title lang="en">The Quick Python Book</title>
3   <isbn>1884777740</isbn>
4   <pageCount>444</pageCount>
5   <publishedDate>
6     <date>1999-10-01T00:00:00.000-0700</date>
7   </publishedDate>,
8   <authors>
9     <author>Daryl Harms</author>
10    <author>Kenneth McDonald</author>
11  </author>
12 </book>
```

# Structured data: JSON

JavaScript Object Notation: *similar to XML, but more compact*

```
1 {  
2   "title" : "The Quick Python Book",  
3   "isbn"  : "1884777740",  
4   "pageCount" : 444,  
5   "publishedDate" : { "date" : "1999-10-01T00:00:00.000-0700" },  
6   "authors" : [ "Daryl Harms", "Kenneth McDonald" ],  
7   "categories" : [ "Python" ]  
8 }
```



# Structured Data: tables

## Extract from file “books.tsv”

title	isbn	pageCount	publishedDate	authors	categories
Unlocking Android	1933988673	416	2009-04-01	W. Frank Ableson, Charlie Collins, Robi Sen	Open Source, Mobile
Specification by Example	1617290084	-	2011-06-03	Gojko Adzic	Software Engineering
Flex 4 in Action	1935182420	600	2010-11-15	Tariq Ahmed, Dan Orlando, John C. Bland II, Joel Hooks	Internet
Zend Framework in Action	1933988320	432	2008-12-01	Rob Allen, Nick Lo, Steven Brown	Web Development
Flex on Java	1933988797	265	2010-10-15	Bernerd Allmon, Jeremy Anderson	Internet
Griffon in Action	1935182234	375	2012-06-04	Andres Almiray, Danno Ferrin, , James Shingler	Java
OSGi in Depth	193518217X	325	2011-12-12	Alexandre de Castro Alves	Java
Flexible Rails	1933988509	592	2008-01-01	Peter Armstrong	Web Development
Hello! Flex 4	1933988762	258	2009-11-01	Peter Armstrong	Internet
Coffeehouse	1884777384	316	1997-07-01	Levi Asher, Christian Crumlish	Miscellaneous
MongoDB in Action	1935182870	-	2011-12-12	Kyle Banker	Next Generation Databases
Taming Jaguar	1884777686	362	2000-07-01	Michael J. Barlotta, Jason R. Weiss	PowerBuilder
Hibernate in Action	193239415X	400	2004-08-01	Christian Bauer, Gavin King	Java
Java Persistence with Hibernate	1932394885	880	2006-11-01	Christian Bauer, Gavin King	Java
JSTL in Action	1930110529	480	2002-07-01	Shawn Bayern	Internet
iBATIS in Action	1932394826	384	2007-01-01	Clinton Begin, Brandon Goodin, Larry Meadors	Web Development
Designing Hard Software	133046192	350	1997-02-01	Douglas W. Bennett	Object-Oriented Programming
Hibernate Search in Action	1933988649	488	2008-12-21	Emmanuel Bernard, John Griffin	Java
...					

# Structured data: tables

## Reading tables using the csv module

```
1 import csv
2
3 f = open('books.tsv')
4 table = list()
5
6 for row in csv.reader(f, delimiter = '\t'):
7
8     # ignore rows that are empty or start with '#'
9     if not row or row[0].startswith('#'):
10         continue
11
12     table.append(row)
13
14 # print first row of table
15 print(table[0])
```

# Structured data: Matrices

$$A = \begin{pmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1} & a_{m2} & \cdots & a_{mn} \end{pmatrix}$$

# Quiz

*True or false?*

- ❖ XML tags have opening and closing elements
- ❖ XML and JSON are archaic data formats
- ❖ The *delimiter* parameter of `csv` reader specifies the the character that separates rows
- ❖ Each column of a table represents a single data point

# Quiz

*True or false?*

- ❖ XML tags have opening and closing elements true
- ❖ XML and JSON are archaic data formats false
- ❖ The *delimiter* parameter of `csv` reader specifies the character that separates rows false
- ❖ Each column of a table represents a single data point false

**File I/O**

**File Formats**

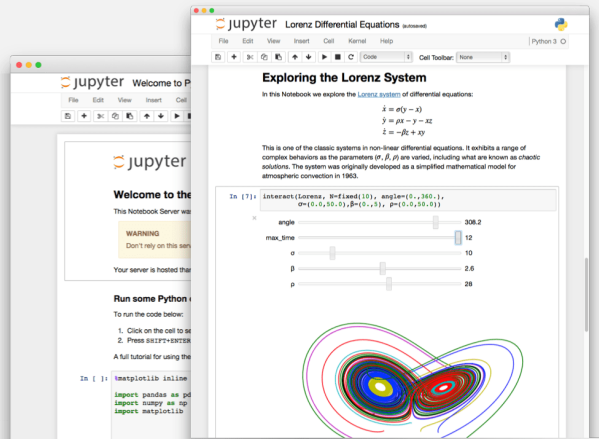
**Jupyter  
Notebook**

**Text mining**

# Jupyter Notebook

## Why use Jupyter Notebook in Data Science?

- Simultaneous documentation & analysis
- Step-by-step processing
- Ensures reproducibility



# Quiz

If you haven't done already, now is a good time to launch Jupyter and familiarize yourself with the tool.

- ❖ Create your own Jupyter notebook
- ❖ Run this chapter's notebook—you can find it in the course material
- ❖ Familiarize yourself with the markdown formatting language
- ❖ Have a look at the shortcuts table. What are the shortcuts for:
  - ❖ Run the current cell, select next
  - ❖ Run selected cells
  - ❖ Save and checkpoint



# Quiz

If you haven't done already, now is a good time to launch Jupyter and familiarize yourself with the tool.

- ❖ Create your own Jupyter notebook
- ❖ Run this chapter's notebook—you can find it in the course material
- ❖ Familiarize yourself with the markdown formatting language
- ❖ Have a look at the shortcuts table. What are the shortcuts for:
  - ❖ Run the current cell, select next `↑ + ↵ / Shift + Enter`
  - ❖ Run selected cells `ctrl + ↵ / Ctrl + Enter`
  - ❖ Save and checkpoint `⌘ + S / Ctrl + S`

**File I/O**

**File Formats**

**Jupyter  
Notebook**

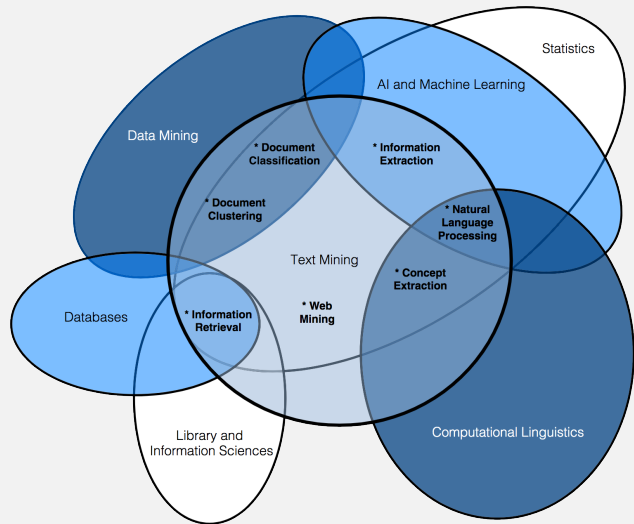
**Text mining**

# Text mining

Relies on *Natural Language Processing* (NLP)

Main (constitutive) tasks:

- ❖ Document summarization, clustering & classification
- ❖ Information extraction
- ❖ Information discovery



source: Miner, Gary. *Practical Text Mining and Statistical Analysis for Non-structured Text Data Applications*. 1st ed. Amsterdam: Academic Press, 2012.

# Document summarization, clustering & classification

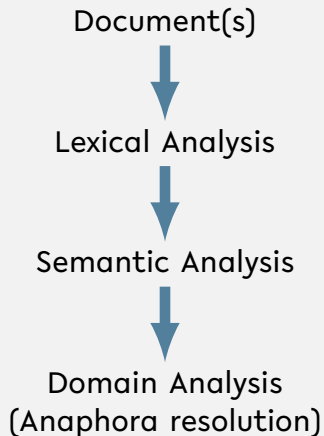
## ❖ Document summarization

- ❖ Goal: Extract essence of a text
- ❖ TextRank
  - Method for ranking sentences
  - Similar to Google's PageRank

## ❖ Document clustering & classification

- ❖ Uses classic data mining techniques
- ❖ Popular: Supervised Learning methods
- ❖ Applied to terms, documents or parts thereof

# Information extraction



source: Miner, Gary. Practical Text Mining and Statistical Analysis for Non-structured Text Data Applications. 1st ed. Amsterdam: Academic Press, 2012.

# Lexical Analysis

- **Tokenization:** decomposition into sentences or words
- **Stemming:** reduction of words to their roots
- **Lemmatization:** inflection & reduction of words to roots

# Semantic & Domain Analysis

## ❖ Semantic Analysis

- ❖ Infers relationships of words
- ❖ Often relies on *parse trees*

## ❖ Domain Analysis

- ❖ Establishes references between parts of text

# Natural Language Toolkit – NLTK

A comprehensive library for natural language processing

NLTK supports

- ❖ Text corpora and lexical resources
- ❖ Tools for
  - ❖ Document summarization & classification,
  - ❖ Information extraction

Read the free book to learn more about NLTK at

<https://www.nltk.org/book/>



# Stemming

- ❖ Process of reducing a word to its *root* (*stem*)
- ❖ Porter Stemmer
  - ❖ Proposed by Martin Porter in 1979
  - ❖ Idea: Each word can be represented by the form  $[C](VC)^m[V]$  where
    - $C$  := consecutive consonants and
    - $V$  := consecutive vowels
    - $m \geq 0$
  - ❖ Simple set of suffix reduction rules, e.g.
    - $ses \rightarrow s$
    - $ies \rightarrow i$
    - $y \rightarrow i$
  - ❖ **Outcome is not true root of the word**, but works well in practice to find words with same root of the English language

# Lemmatization

- ❖ reduces the *inflected* words e.g.: runs, running, ran → run
- ❖ Requires additional information of the language
- ❖ WordNet Lemmatizer:
  - ❖ Uses WordNet database to inflect words
  - ❖ Works best if *part-of-speech* (POS) information is provided: e.g. is word a verb or noun?

```
1 from nltk.stem.wordnet import WordNetLemmatizer
2 wnl = WordNetLemmatizer()
3
4 # calling lemmatizer without POS information
5 wnl.lemmatize('ran')
6 # returns 'ran'
7
8 wnl.lemmatize('ran', 'v') # 'v' for 'verb'
9 # returns 'run'
```

# Quiz

## *True or false?*

- ❖ Stemming is the process of decomposing text into smaller units
- ❖ Inflection is the change of a word's form
- ❖ The Porter Stemmer requires no adaptations to work well on any language
- ❖ The quality of lemmatization depends on the utilized data base
- ❖ Semantic analysis often relies on parse trees

# Quiz

## *True or false?*

- ❖ Stemming is the process of decomposing text into smaller units false
- ❖ Inflection is the change of a word's form true
- ❖ The Porter Stemmer requires no adaptations to work well on any language false
- ❖ The quality of lemmatization depends on the utilized data base true
- ❖ Semantic analysis often relies on parse trees false

# Recap

# Summary

- ❖ reading and writing files, command line arguments, prompt
- ❖ structured and unstructured data formats
  - ❖ Text
  - ❖ JSON & XML
  - ❖ Tables & matrices
- ❖ Jupyter Notebook
- ❖ Text mining, lexical analysis

# What comes next?

- ❖ Familiarize yourself with Jupyter Notebook
- ❖ Play with NLTK
- ❖ Due date for this week's exercises is **Sunday, May 17, 2020.**

*Next lecture:* Numerical Data Analysis, NumPy, ... ..