

# TEXT MINING INTRO TO PYTHON

Johan Falkenjack (based on slides by Mattias Villani)

**NLPLAB**  
**Dept. of Computer and Information Science**  
**Linköping University**

# OVERVIEW

- ▶ What is Python? How is it special?
- ▶ Python's objects
- ▶ If-else, loops and list comprehensions
- ▶ Functions
- ▶ Classes
- ▶ Modules

# WHAT IS PYTHON?

- ▶ First version in 1991
- ▶ **High-level language**
- ▶ Emphasizes **readability**
- ▶ **Interpreted** (bytecode .py and .pyc) [can be compiled via C/Java]
- ▶ Automatic memory management
- ▶ Strongly dynamically typed
- ▶ **Functional and/or object-oriented**
- ▶ **Glue** to other programs (interface to C/C++ or Java etc)
- ▶ Popular in data science (“Prototype in R, implement in Python”)
- ▶ Two currently developed versions, 2.x and 3.x
  - ▶ This course uses Python 3

# THE BENEVOLENT(?) DICTATOR FOR LIFE (BDFL) GUIDO VAN ROSSUM



# PYTHON PECULIARITIES (COMPARED TO R/MATLAB)

- ▶ Not primarily a numerical language.
- ▶ **Indexing begins at 0**, as indexes refer to breakpoints between elements.
- ▶ It follows that `myVector[0:2]` returns the first and second element, but not the third.
- ▶ **Indentation matters!**
- ▶ Can import specific functions from a module.
- ▶ Assignment **by object**, **NOT by copy** or **by reference**.
  - ▶ Approximately, assignment **by copy of reference**.
- ▶ `a = b = 1` assigns 1 to both a and b.

# PYTHON'S OBJECTS

- ▶ Built-in types: **numbers**, **strings**, **lists**, **dictionaries**, **tuples** and **files**.
- ▶ **Vectors**, **arrays** and functions to manipulate them are available in the **numpy/scipy** modules.
- ▶ Python is a **strongly typed** language. `'johan' + 3` gives an error.
- ▶ Python is a **dynamically typed** language. No need to declare a variables type before it is used. Python figures out the object's type.
- ▶ Implication: Polymorphism by default:
  - ▶ “In other words, don't check whether it IS-a duck: check whether it QUACKS-like-a duck, WALKS-like-a duck, etc, etc, depending on exactly what subset of duck-like behaviour you need to play your language-games with.” - Alex Martelli

# STRINGS

- ▶ No character type, strings are strings, even if they are length 1.
- ▶ `s = 'Spam'`
- ▶ `s[0]` returns first character, `s[-2]` return next to last character.  
`s[0:2]` returns first **two** characters (slicing).
- ▶ `len(s)` returns the number of letters.
- ▶ `s.lower()`, `s.upper()`, `s.count('m')`, `s.endswith('am')`,  
...
- ▶ **Which methods are available for my object?** Try in Jupyter: type `s.` followed by TAB.
- ▶ `+` operator **concatenates strings** by creating a new string (computationally expensive if you do it often).
- ▶ `sentence = 'Guido is the benevolent dictator for life'.sentence.split()`
- ▶ `s*3` returns `'SpamSpamSpam'`

# INDEXING

- ▶ Strings are always read from left to right
- ▶ Indexes are the positions between characters and start at 0, i.e. the position before the first character.
- ▶ `s[0]` starts at index 0 and reads once character forward, i.e. returns the first character in `s`.
- ▶ Slicing are when you take a “slice” from a string based on indexes, for instance `s[1:2]` returns the substring from position 1 to 3, i.e. the second and third characters.
- ▶ `s[:3]` returns the string upto position 3, `s[3:]` returns the string from position 3 to the end.
- ▶ `s[2]` is equivalent to `s[2:3]` (with strings).
- ▶ Indexes can be negative and are then counted from the end of the sequence (they are NOT excluders as in R).
- ▶ `s[-2]` returns next to last character, `s[-3:]` returns the substring containing the last three characters.



# THE LIST OBJECT

- ▶ A list is an ordered **container of several values**, possibly of different types.
- ▶ `myList = ['spam','spam','bacon',2]`
- ▶ The list object has several associated **methods**
  - ▶ `myList.append('egg')`
  - ▶ `myList.count('spam')`
  - ▶ `myList.sort()`
- ▶ `+` operator concatenates lists: `myList + myOtherList` merges the two lists as one list.
- ▶ Elements are **not named** and lists are **not vectors** in a mathematical sense.

# THE LIST OBJECT

- ▶ Extract elements from a list: `myList[1]`
- ▶ Lists inside lists (nested lists):
  - ▶ `myOtherList = ['monty', 'Python']`
  - ▶ `myList[1] = myOtherList`
  - ▶ `myList[1]` returns the list `['monty', 'Python']`
  - ▶ `myList[1][1]` returns the string `'Python'`
- ▶ Indexing and slicing works as with strings...
  - ▶ **EXCEPT** a bit more generally, in a string, all elements are strings, but not all elements of lists are necessarily lists.
  - ▶ `myOtherList[1]` returns an element (the string `'Python'`) while `myOtherList[1:2]` returns a list (`['Python']`).

# STRINGS AGAIN

- ▶ Strings are **immutable**, i.e. can't be changed after creation.
- ▶ Every “change” creates a new string (remember concatenation).
- ▶ Try to avoid creating more strings than necessary:
  - ▶ Avoid: `my_string = 'Python' + ' ' + 'is' + ' ' + 'fun!'`
  - ▶ Instead: `my_string = ' '.join(['Python', 'is', 'fun!'])`
- ▶ In loops where you construct strings, add constituents to a list and join after loop finishes.
- ▶ Remember that you have to explicitly change the type of, for instance, numbers, when concatenating: `str(1)`

# DICTIONARIES

- ▶ **Unordered** collection of pairs, often names and some value.
- ▶ `myDict = {'Sarah':29, 'Erik':28, 'Evelina':26}`
- ▶ Elements are **accessed by keyword not by index**:  
`myDict['Evelina']` returns 26.
- ▶ **Values can contain any object**: `myDict = {'Marcus':[23,14], 'Cassandra':17, 'Per':[12,29]}`. `myDict['Marcus'][1]` returns 14.
- ▶ **Keys must be immutable**: `myDict = {2:'contents of box2', (3, 'a'):'content of box 4', 'blackbox':10}`
- ▶ `myDict.keys()`
- ▶ `myDict.values()`
- ▶ `myDict.items()`

# TUPLES

- ▶ `myTuple = (3,4,'johan')`
- ▶ **Like lists, but immutable**
- ▶ Why?
  - ▶ Faster than lists
  - ▶ Protected from change
  - ▶ Can be used as keys in dictionaries
  - ▶ Multiple return object from function
  - ▶ Swapping variable content `(a, b) = (b, a)`
  - ▶ Sequence unpacking `a, b, c = myTuple`
- ▶ `list(myTuple)` returns `myTuple` as a list. `tuple(myList)` does the opposite.

# SETS

- ▶ **Set.** Contains objects in **no order** with **no identification**.
  - ▶ With a **list**, elements are ordered and identified by position. `myList[2]`
  - ▶ With a **dictionary**, elements are unordered but identified by some key. `myDict['myKey']`
  - ▶ With a **set**, elements stand for themselves. No indexing, no key-reference.
- ▶ Declaration: `fib=set( (1,1,2,3,5,8,13) )` returns the set `{1, 2, 3, 5, 8, 13}`.
- ▶ Supported operations: `len(s)`, `x in s`, `set1 < set2`, `union`, `intersection`, `add`, `remove`, `pop ...`

# VECTORS AND ARRAYS (AND MATRICES)

- ▶ `from scipy import *`
- ▶ `x = array([1,7,3])`
- ▶ 2-dimensional **array** (matrix): `X = array([[2,3],[4,5]])`
- ▶ **Indexing arrays**
  - ▶ First row: `X[0,]`
  - ▶ Second column: `X[:,1]`
  - ▶ Element in position 1,2: `X[0,1]`
- ▶ Array **multiplication** (`*`) is element-wise, for matrix multiplication use `dot()`.
- ▶ There is also a **matrix object**: `X = matrix([[2,3],[4,5]])`
  - ▶ **Arrays are preferred** (not matrices).
- ▶ Submodule **scipy.linalg** contains a lot of **matrix-functions** applicable to arrays (`det()`, `inv()`, `eig()` etc). I recommend: `from scipy.linalg import *`

# BOOLEAN OPERATORS

- ▶ True/False
- ▶ and
- ▶ or
- ▶ not
- ▶ `a = True; b = False; a and b` [returns False].



# IF-ELSE CONSTRUCTS

## IF-ELSE STATEMENT

```
a = 1
if a==1:
    print('a is one')
elif a==2:
    print('a is two')
else:
    print('a is not one or two')
```

- ▶ **Switch statements** via dictionaries (see Jackson's Python book) if absolutely necessary.

# WHILE LOOPS

## WHILE LOOP

```
a =10
while a>1:
    print('bigger than one')
    a = a - 1
else:
    print('smaller than one')
```

# FOR LOOPS

- ▶ **for loops can iterate over any iterable.**
- ▶ **iterables:** strings, lists, tuples, ranges

## FOR LOOP

```
word = 'johan'
for letter in word:
    print(letter)
myList = ['']*10
for i in range(10):
    myList = 'johan' + str(i)
```

# LIST COMPREHENSIONS

- ▶ As in R, loops can be slow. For small loops executed many times to generate many results, try comprehensions:
- ▶ Set definition in mathematics

$$\{x \text{ for } x \in \mathcal{X}\}$$

where  $\mathcal{X}$  is some a finite set.

$$\{f(x) \text{ for } x \in \mathcal{X}\}$$

- ▶ Comprehensions in Python:
  - ▶ `myList = [x for x in range(10)]`
  - ▶ `mySet = {sqrt(x) for x in range(10)}` (don't forget from math import sqrt)
  - ▶ `myDict = {x : sqrt(x) for x in sample if x >= 0}`

# DEFINING FUNCTIONS AND CLASSES

```
def mySquare(x):  
    return x**2
```

```
class myClass(object):  
    def __init__():  
        self.count = 0  
    def method(self, x):  
        self.count += x
```

- ▶ If `a` is instance of `myClass`, `a.method(b)` can be thought of as `myClass.method(a, b)` where `self` becomes the object `a` itself
- ▶ Make you own module by putting several classes or functions in a `.py` file. Then import what you need from that file.

# STRINGS, A THIRD TIME

- ▶ String formatting is often used to create recurrent string injections:
  - ▶ "My name is {0}".format('Fred')
  - ▶ "The story of {0}, {1}, and {c}".format(a, b, c=d)
- ▶ Regular expressions.
  - ▶ Used to find and manipulate patterns in strings, extremely powerful stuff.
  - ▶ Available in Python through the module `re`.
  - ▶ Introduction: [https://regexone.com/lesson/introduction\\_abcs](https://regexone.com/lesson/introduction_abcs)
  - ▶ Practice: <https://alf.nu/RegexGolf>

- ▶ Comments on individual lines starts with #
- ▶ Doc-strings can be used as comments spanning over multiple lines but this should be avoided `"""This is a looooong comment"""`