

COM6115: Text Processing

OO Programming: Python basics
Configuring Program Behaviour
Programming Tips

Mark Hepple

Department of Computer Science
University of Sheffield

Object Oriented Programming

- So far, we have used a *procedural programming* paradigm
 - ◇ focus is on writing *functions* or *procedures* to operate on data
- Alternative paradigm: **Object Oriented Programming (OOP)**
 - ◇ focus is on creating *classes* and *objects*
 - ◇ *objects* contain both *data* and *functionality*
- OOP has become the *dominant* programming paradigm
 - ◇ developed to make it easier to create and/or modify large, complex software systems
- These slides introduce *basics* of OOP in Python (*without inheritance*)
- See the '*extended presentation*' slides (on module homepage) for:
 - ◇ more on background and motivation for OOP
 - ◇ basics of using *inheritance* in Python

Let's talk about meaning

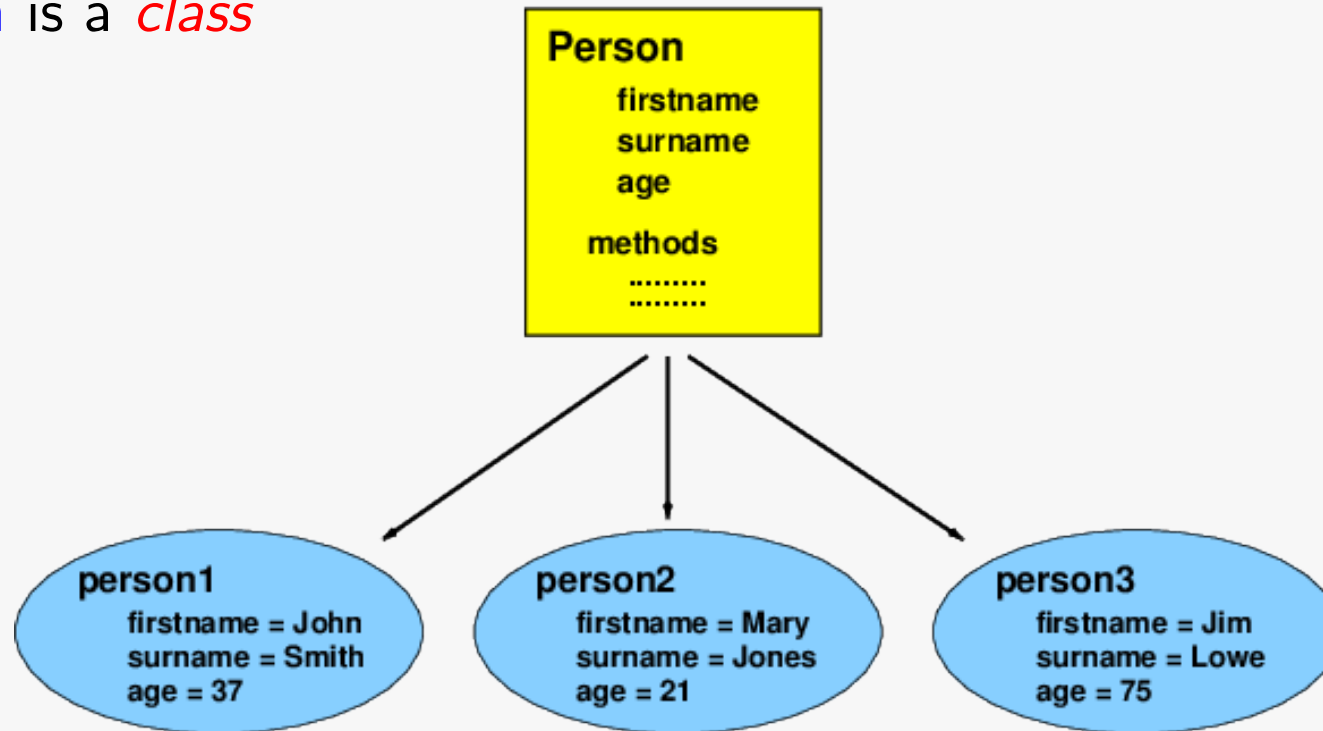
- Key notion: *CONCEPT*
 - ◇ general idea of a *class of things* with particular properties in common
e.g. concepts: *person, bird, animal, vehicle, chair, etc.*
- A *concept* has *INSTANCES*
 - ◇ actual occurrences in the world
e.g. concept *person* has *instances* such as: *Me! You! Beyoncé!*
- For a given concept, expect certain *attributes*
 - ◇ a specific *actual* person will *instantiate* these attributes
e.g. for *person*, expect: *age, gender, height, etc.*
- Concept may also have associated *expected behaviours*
e.g. for *person* — *walk, talk, read, hoover, give birth*
- These ideas approximate key ideas of OOP, especially:
 - ◇ *concept* \approx *CLASS*
 - ◇ *instance* \approx *OBJECT*

Objects and Classes — *an example*

- A **Person** class might:
 - ◇ have *attributes* (variables) for:
 - name, age, height, address, tel.no., job, etc
 - ◇ have *methods* (functions) to:
 - update address
 - update job status
 - work out if they are adult or child
 - work out if they pay full fare on the bus
 - *etc.*
- There might be many *objects* of the **Person** class
 - ◇ each representing a *different person*
 - *with different specific data*
 - ◇ but all store *similar information* and *behave similarly*

Objects and Classes — *an example*

- Person is a *class*



- ◇ person1, person2 & person3 are *objects*

Defining Classes in Python

- Definition opens with keyword `class` + class name
- Class needs an *initialisation* method
 - ◇ called when an instance is created
 - ◇ has 'special' name: `__init__`
 - ◇ establishes the *attributes* (i.e. vars) belonging to objects

```
class Person:
    def __init__(self):
        self.firstname = None
        self.surname = None
        self.age = None
        self.species = 'homo sapiens'
```

- ◇ note use of *special variable* `self` here
- ◇ it is the instance's way of *referring to itself*
e.g. `self.species` above means "the `species` attribute of *this instance*"

Defining Classes in Python (ctd)

- **Person** class with its *initialisation* method, again:

```
class Person:
    def __init__(self):
        self.firstname = None
        self.surname = None
        self.age = None
        self.species = 'homo sapiens'
```

- Can create an **object** (i.e. *instance*) of this class as follows:

```
>>> p1 = Person()
>>> p1.species
'homo sapiens'
```

- ◇ here, call to **Person()** creates a new instance of the **Person** class
 - the `__init__` method is called automatically, to initialise the object
 - the object is assigned to **p1**

Defining Classes in Python (ctd)

- Last example again:

```
>>> p1 = Person()
>>> p1.species
'homo sapiens'
```

- ◇ statement `p1.species` accesses `p1`'s species *attribute* directly
i.e. that value is accessed in the e.g. above, and printed by the interpreter

- Can think of objects as being like “*bundles of data*”

- ◇ each object is a different *bundle of data*, storing info about a *different instance of the class*

- ◇ Note extra *self* arg in:

```
class Person:
    def __init__(self):
        self.firstname = None
        ...
```

- is object's way of talking about *itself*, i.e. *the bundle that I am*
- info stored with a `self.` attribute becomes *part of the bundle*
— is carried around with it, and is *always available*

Defining Classes in Python (ctd)

- More generally, **initialisation** method can have *parameters*
 - ◇ can be used to set *initial values of attributes*

```
def __init__(self, firstname, surname, age):  
    self.firstname = firstname  
    self.surname = surname  
    self.age = age  
    self.species = 'homo sapiens'
```

- ◇ example of creating an instance:

```
>>> p1 = Person('John', 'Smith', 37)  
>>> p1.firstname  
'John'  
>>> p1.age  
37
```

- ◇ note `__init__` has 4 args, but 3 given when object created – *Why?*
 - first `self` is left *implicit* – stands for *this object* (i.e. *bundle of data*)
 - that object stored as `p1`, can access bundle data directly, e.g. `p1.age`

Defining Classes — *adding functionality*

- Can define (more) functions — in OOP, are known as *methods*

```
class Person:
    def __init__(self):
        ...

    def greeting_informal(self):
        print('Hi', self.firstname)

    def greeting_formal(self):
        print('Welcome, Citizen', self.surname)
```

- ◇ as before, *self* appears as 1st arg of every method
 - shows that this is an *object method*, i.e. will be *called from an object*
- ◇ *self* again refers to *this instance*, allowing access to *its own data*
 - thus, *self.firstname* above *means* value of *my firstname attribute*
 - that value, stored with this bundle of data, is accessed and used

Defining Classes — *adding functionality* (ctd)

- Example: here create two instances:

```
>>> p1 = Person('Harry', 'Potter', 12)
>>> p2 = Person('Hermione', 'Grainger', 12)
```

- Call newly defined methods from instances:

```
>>> p1.greeting_informal()
Hi Harry
>>> p1.greeting_formal()
Welcome, Citizen Potter
>>> p2.greeting_formal()
Welcome, Citizen Grainger
```

- ◇ note that 1st `self` arg from definition again absent – *i.e. is left implicit*
- ◇ when `p1.greeting_informal()` is called, `p1` stores an instance, and `self` aspects of definition are *about that instance*
- ◇ thus, method calls access data (e.g. `surname`) from given instance (`p1` or `p2`), and output depends on that

Defining Classes — *adding functionality* (ctd)

- Another method ...

```
class Person:
    ...

    def greeting_age_based(self):
        if self.age < 18:
            print('Welcome, Young', self.firstname)
        elif self.age > 60:
            print('Welcome - oh Venerable', self.firstname)
        else:
            self.greeting_formal()
```

- ◇ here see behaviour that *uses* instance data (*firstname*) and that *is conditioned on* instance data (*age*)
- ◇ note: 'else' case *calls* another method of the instance
 - does so in form: *self.greeting_formal()*
 - uses *self*, as it is *this object's* method being used
 - but *self* is *prefixed*, not supplied as arg

Defining Classes — *adding functionality* (ctd)

- Example:

```
>>> p1 = Person('Harry', 'Potter', 12)
>>> p2 = Person('Sirius', 'Black', 38)
>>> p3 = Person('Minerva', 'McGonagall', 66)
```

- ◇ call methods — observe behaviour is conditioned on person's age

```
>>> p1.greeting_age_based()
Welcome, Young Harry
>>> p2.greeting_age_based()
Welcome, Citizen Black
>>> p3.greeting_age_based()
Welcome - oh Venerable Minerva
```

- Have introduced *basics* of OOP in Python (*without inheritance*)
- See the '*extended presentation*' slides (on module homepage) for:
 - ◇ more on background and motivation for OOP
 - ◇ basics of using *inheritance* in Python

Configuring Program Behaviour

- Often want to *configure* the behaviour of a program, e.g. to:
 - ◇ specify files from which to take *input*
 - ◇ name of files to which to write *output/results*
 - ◇ *set* various *parameters*:
 - e.g. weight/threshold values, number of results to print, etc
- For *scientific computing*, often want to run program under a wide *range* of different *settings*:
 - ◇ i.e. so alternative results can be *compared*, *plotted*, etc.
- Might configure via a *GUI*, *but*
 - ◇ time-consuming to develop
 - ◇ time-consuming to use, if each configuration must be entered separately
- Alternative: configure via the *command line*
 - ◇ use '*flag*' symbols (e.g. '*-s*') to *name* specific command line options

Command Line Options

- Using command line options — e.g. might have call:

```
python myCode.py -w -t 0.5 -d data1.txt -r results1.txt
```

- ◇ with options to specify the input data file (**-d**), the results file (**-r**), and a threshold value (**-t**) affecting the process
- ◇ and a *boolean* option **-w** to direct some aspect of behaviour
e.g. whether terms are *weighted* or not

- **Help option:** good practice to include a boolean *help* option **-h**:

- ◇ if present, code *just prints help message* and *then quits*
- ◇ help message says how to call program, lists options, etc

- Allows for use of *batch files*:

- i.e. text file containing commands to invoke program under a range of parameter settings
- ◇ easy way to generate a range of experimental results

The getopt Module

- The `getopt` module helps with parsing command line options
 - ◇ allows both *short* options (`-s`) and *long* ones (`--long-option`)
 - ◇ here consider only short options
- Specify allowed options via a string, e.g. `'hi:o:I'`
 - ◇ each letter in string accepted as an option
 - ◇ letters followed by `:"` require an arg string, e.g. `-i` here
 - ◇ otherwise flag is boolean, e.g. `-h` here
- Parsing usually applied to `sys.argv[1:]`

e.g. `opts, args = getopt.getopt(sys.argv[1:], 'hi:o:I')`

- ◇ here, `opts` is the options found — given as a *list of pairs*
 - convert to *convert* to a *dictionary*, e.g. with `opts = dict(opts)`
- ◇ `args` is any remaining 'bare' arguments – as a list
 - flag options should *precede* bare args on command line
- ◇ note that `sys.argv[0]` is name of your *code file* – don't pass this

The getopt Module (ctd)

- See 'demo' code file on using **getOpts** module (on module homepage)
 - ◇ run in a CMD window (or linux/mac terminal)
 - ◇ invoke python on code directly, as follows:

```
> python getOptsDemo.py -h
```

```
-----  
USE: python getOptsDemo.py (options)
```

```
OPTIONS:
```

```
    -h : print this help message
```

```
    -s FILE : use stoplist file FILE (required)
```

```
    -b : use binary weighting (default is off)
```

```
-----  
> python getOptsDemo.py -s stops.txt -b file1.txt file2.txt
```

```
SUMMARY
```

```
Command line strings: ['getOptsDemo.py', '-s', 'stops.txt', '-b', 'file1.txt',
```

```
Arguments: ['file1.txt', 'file2.txt']
```

```
Options:
```

```
    Stopwords file: stops.txt
```

```
    Binary weighting: 1
```

```
>
```

Python Tips — *the Good, the Bad, and the Ugly*

- *Elegance* is important:
 - ◇ clear, readable coding helps rapid/effective code development
- Learn to use the clean constructs Python provides
 - e.g. use `k in dict` rather than `dict.has_key(k)`
- Know the *default iteration* behaviour of your data structure
 - ◇ so can usually address content via a simple *for*-loop
- Understand the importance of *hash-based* data structures
 - ◇ allow *constant time* look-up / update
 - ◇ usually much more efficient than *sequence-based* data structures
 - ◇ *beware* of doing *sequence-based* look-up in *hash-based* structures

Python Tips — *know the default iteration behaviour*

- Simple *for*-loop provides clean, readable way to address content of an iterable data structure:

```
for item in Iterable:  
    do_something(item)
```

- ◊ so, useful to know *default iteration behaviour* for *common cases*
- Iterating over *X* gives items *Y* ...
 - ◊ a *string* gives *chars* in their given (left-to-right) order
 - ◊ a *list* gives its *elements*, in their given order
 - ◊ a *tuple* gives its *elements*, in their given order
 - ◊ a *set* gives its *elements*, in no particular order
 - ◊ a *dictionary* gives its *keys*, in no particular order
 - ◊ a *file-stream* gives its *lines of text*, in file order

Python Tips — *hash-based data structures*

- In *text processing*, often want to handle info about *very many items*
e.g. counts for 100K words, or *millions* of ngrams
- *Hash-based data structures* are very suitable for this
i.e. Python *dictionary* and *set* data structures
- Why? — allow (roughly) *constant time* access to info for a key/item
i.e. in a *fixed* (small) amount of time *irrespective of how many items stored*
- Using *sequential* data structs (e.g. list) for similar tasks is a *bad idea*
 - ◇ gives (typically) *linear time* access (i.e. \propto num items stored)
- Test “*item in D*” uses look-up method appropriate to *D*
 - e.g. if it's a *list*, look-up is by *left-to-right sequential comparison*
 - e.g. if it's a *set*, look-up uses *hash-based* method
 - e.g. if it's a *dictionary*, look-up uses *hash-based* method

Python Tips — *hash-based data structures* (ctd)

- Avoid changing hash look-up to sequential one — *common error*
- If `D` is a dictionary, `D.keys()` gives a 'smart iterator' over `D`'s keys
 - ◇ so `x in D.keys()` as efficient as `x in D` (but *less elegant!*)
- BUT all of `list(D)`, `list(D.keys())`, `sorted(D)` return a *list*
 - ◇ so (e.g.) `x in sorted(D)` is *sequential* and *v.inefficient*
- Also v.inefficient is following attempt to check for `x` in `D`:

```
for k in D.keys():  
    if k == x:  
        ...
```

- ◇ recreates sequential character of look-up
- ◇ surprisingly commonly seen!

Python Tips — *avoid piecemeal coding solutions*

- Desire to break task into manageable ‘chunks’ sometimes leads to *inelegant ‘piecemeal’ solutions*
 - ◇ avoid this, *unless the task really requires it*
- *Example*: task = count the non-stoplist words in a file
 - ◇ might be tempted to handle as follows (assume stoplist loaded):
 - read the lines of text into a list
 - iterate over list to split each line into a list of tokens
 - iterate again, to delete stop list words
 - iterate again, counting tokens (into a dictionary)

— this is a poor solution !!
 - ◇ better solution — more efficient, and simpler to code:
 - read the text line by line (i.e. using a `for`-loop)
 - for each line read, access tokens
 - e.g. using `.split()` string method, or using a `regex`+`findall`
 - for each token: if it’s a stopword, skip it, otherwise count it