# Computerlinguistische Anwendungen
# Python/Git

Thang Vu

CIS, LMU
thangvu@cis.uni-muenchen.de

April 15, 2015

## Introduction

- Review of Python
- Introduction to NumPy
- Introduction to Git

# Core Data Types

| Object type | Example creation |
| --- | --- |
| Numbers | 123, 3.14 |
| Strings | 'this class is cool' |
| Lists | [1, 2, [1, 2]] |
| Dictionaries | {'1': 'abc', '2': 'def'} |
| Tuples | (1, 'Test', 2) |
| Files | open('file.txt'), open('file.bin', 'wb') |
| Sets | set('a', 'b', 'c') |
| Others | boolean, None |
| Program unit types | Functions, modules, classes |

## Variables

- store data, e.g., numbers
- content can be changed (is variable)
- have a data type
- assignment: var_name = value, e.g., `num = 17`

## Dynamic Typing

- dynamic typing model
- types are determined automatically at runtime
- type of a variable can change
- check type with `type(var)`

## Number Data Types

- integers, floating-point numbers, complex numbers, decimals, rationals
- Numbers support the basic mathematical operations, e.g.:
    - $+$ addition
    - $*$, $/$ multiplication, division
    - $**$ exponentiation
    - $<$, $>$, $<=$, $>=$ comparison
    - $!=$, $==$ (in)equality

```
1  >>> 1/4
2  >>> float(1/4)
3  >>> float(1)/4
```

## String Data Types

- immutable sequence of single characters
- **ASCII**: 256 characters: `'tree'`,`'2.1'`,`'two tokens'`
- **Unicode**: > 110,000 characters: $u'tree'$,$u'\sigma'$, $u'\u2B0000'$

## Operations

```
s1 = 'the'
```

| Operation | Description | Output |
| --- | --- | --- |
| len(s1) | length of the string | 3 |
| s1[0] | indexing, 0-based | 't' |
| s1[-1] | backwards indexing | 'e' |
| s1[0:3] | slicing, extracts a substring | 'the' |
| s1[:2] | slicing, extracts a substring | 'th' |
| s1 + ' sun' | concatenation | 'the sun' |
| s1 * 3 | repetition | 'thethethe' |
| != , == | (in)equality | True, False |

8

## String-Specific Methods

```
s1 = 'these'
```

| Operation | Description | Output |
|---|---|---|
| `'-'.join(s1)` | concatenate with delimiter '-' | `'t-h-e-s-e'` |
| `s1.find('se')` | finds offsets of substrings | 3 |
| `s1.replace('ese', 'at')` | replace substrings, s1 is still the initial string | `'that'` |
| `s1.split('s')` | splits string at delimiter | `['the', 'e']` |
| `s1.upper()` | upper case conversions | `'THESE'` |
| `s1.lower()` | lower case conversions | `'these'` |

## Lists

- collection of arbitrarily typed objects
- mutable
- positionally ordered
- no fixed size
- initialization: `L = [123, 'spam', 1.23]`
- empty list: `L = []`

## Operations

```
L = [123, 'spam', 1.23]
```

| Operation | Description | Output |
|---|---|---|
| len(L) | length of the list | 3 |
| L[1] | indexing, 0-based | 'spam' |
| L[0:2] | slicing, extracts a sublist | [123, 'spam', 1.23] |
| L + [4, 5, 6] | concatenation | [123, 'spam', 1.23, 4, 5, 6] |
| L * 2 | repetition | [123, 'spam', 1.23, 123, 'spam', 1.23] |

## List-Specific Methods

```
L = [123, 'spam', 1.23]
```

| Operation | Description | Output |
|---|---|---|
| L.append('NI') | append to the end | [123, 'spam', 1.23, 'NI'] |
| L.pop(2) | delete item | [123, 'spam'] |
| L.insert(0, 'aa') | insert item at index | ['aa', 123, 'spam', 1.23] |
| L.remove(123) | remove given item | ['spam', 1.23] |
| L.sort() | sort list | [1.23, 123, 'spam'] |
| L.reverse() | reverse list | [1.23, 'spam', 123] |

## Nested Lists

Let us consider the 3*x*3 matrix of numbers M = [[1, 2, 3], [4, 5, 6], [7, 8, 9]]. *M* is a list of 3 objects, which are in turn lists as well and can be referred to as rows.

- M[1] – returns the second row in the main list: [4, 5, 6]
- M[1][2] – returns the third object situated in the in the second row of the main list: 6

## Dictionaries

- Dictionaries are not sequences, they are known as mappings
- They are mutable like lists
- They represent a collection of key-value pairs
- e.g.

```
1  >>> D = {'food':'Spam', 'quantity':4, 'color':'pink'}
```

# Dictionary Operations

```
1  >>> D = {'food':'Spam', 'quantity':4, 'color':'pink'}
2  >>> D['food']              #Fetch value of key 'food'
3  'Spam'
4  >>> D['quantity'] += 1   #Add 1 to the value of 'quantity'
5  >>> D
6  D = {'food':'Spam', 'quantity':5, 'color':'pink'}
```

# Dictionary Operations (cont.)

```
1  >>> D = {}
2  >>> D['name'] = 'Bob'      #Create keys by assignment
3  >>> D['job'] = 'researcher'
4  >>> D['age'] = 40
5
6  >>> D
7  D = {'name':'Bob', 'job':'researcher', 'age':40}
8
9  >>> print D['name']
10 Bob
```

## Dictionary Operations (cont.)

```
1  >>> #Alternative construction techniques:
2  >>> D = dict(name='Bob', age=40)
3  >>> D = dict([('name', 'Bob'), ('age', 40)])
4  >>> D = dict(zip(['name', 'age'], ['Bob', 40]))
5  >>> D
6  {'age': 40, 'name': 'Bob'}
7  >>> #Check membership of a key
8  >>> 'age' in D
9  True
10 >>> D.keys()      #Get keys
11 ['age', 'name']
12 >>> D.values()    #Get values
13 [40, 'Bob']
14 >>> D.items()     #Get all keys and values
15 [('age', 40), 'name', 'Bob']
16 >>> len(D)        #Number of entries
17 2
```

# Dictionary Operations (cont.)

```
1  >>> D = {'name': 'Bob'}
2  >>> D2 = {'age': 40, 'job': 'researcher'}
3  >>> D.update(D2)
4  >>> D
5  {'job': 'researcher', 'age': 40, 'name': 'Bob'}
6  >>> D.get('job')
7  'researcher'
8  >>> D.pop('age')
9  40
10 >>> D
11 {'job': 'researcher', 'name': 'Bob'}
```

# Tuples

- Sequences like lists but immutable like strings
- Used to represent fixed collections of items

```
1  >>> T = (1, 2, 3, 4)      #A 4-item tuple
2  >>> len(T)                #Length
3  4
4  >>> T + (5, 6)            #Concatenation
5  (1, 2, 3, 4, 5, 6)
6  >>> T[0]                  #Indexing, slicing and more
7  1
8  >>> len(T)
9  ???
```

## Sets

- Mutable

- Unordered collections of **unique** and **immutable objects**

```
1  >>> set([1, 2, 3, 4, 3])
2  set([1, 2, 3, 4])
3  >>> set('spaam')
4  set(['a', 'p', 's', 'm'])
5  >>> {1, 2, 3, 4}
6  set([1, 2, 3, 4])
7  >>> S = {'s', 'p', 'a', 'm'}
8  >>> S
9  set(['a', 'p', 's', 'm'])
10 >>> S.add('element')
11 >>> S
12 set(['a', 'p', 's', 'm', 'element'])
```

# Files

- The main interface to access files on your computer
- Can be used to read and write text

```
1  >>> f = open('data.txt','w')  #Make a new file in output
       mode 'w'
2  >>> f.write('Hello\n')  #Write a string to it
3  4
4  >>> f.write('World\n')
5  (1, 2, 3, 4, 5, 6)
6  >>> f.close()              #Close to flush output puffers to
       disk
7  1
8  >>> #Continue writing at the end of an existing file
9  >>> f.write('data.txt', 'a')
10 >>> f.write('Cont.'\n)
11 >>> f.close()
```

# Files

```
1  >>> f = open('data.txt')    #'r' is default processing mode
2  >>> text = f.read()         #Read entire file in a string
3  >>> text
4  Hello\nWorld\nCont.
5  >>> print text              #print interprets control
       characters
6  Hello
7  World
8  Cont.
9  >>> text.split()            #File content is always a string
10 ['Hello', 'World', 'Cont.']
```

**Larger files are always better read line by line!!!**

```
1  >>> for line in open('data.txt','r'): print line
```

## Immutable vs Mutable

Immutable:

- numbers
- strings
- tuples

Mutable:

- lists
- dictionaries
- sets
- newly coded objects

# Testing: if statements

```
1  >>> x = 'killer rabbit'
2  >>> if x == 'roger':
3  ...     print 'shave and a haircut'
4  ... elif x == 'bugs':
5  ...     print 'whats up?'
6  ... else:
7  ...     print 'run away!'
8  ...
9  run away!
```

### Note!

The `elif` statement is the equivalent of `else if` in Java or `elsif` in Perl.

# Looping: while loops

```
1  >>> while True:
2  ...        print 'Type Ctrl-C to stop me!'
3  ...
4  >>> x == 'spam'
5  ... while x:              #while x is not empty
6  ...      print x
7  ...      x = x[1:]
8  ...
9  spam
10 pam
11 am
12 m
```

## Looping: for loops

The for loop is a generic iterator in Python: it can step through the
items in any ordered sequence or other iterable objects (strings, lists,
tuples, and other built-in iterables, as well as new user-defined
objects).

```python
L = [1, 2, 3, 4]
for i in L:
    print i
```

```python
for i in range(0, 5):
    print i
```

# Looping: for loops

The most efficient file scanner in Python:

```
1  #use iterator: best for text input
2  >>> for line in open('data.txt'):
3  ...        print line
```

## Note!

This is not only the shortest but as well the most efficient coding for reading a file in Python. It relies on file iterators to automatically read one line on each loop iteration, which allows it to work with arbitrarily large files – often needed in NLP!

## Function

- A function is a device that groups a set of statements so they can be run more than once in a program
- Why use functions?
    - Maximizing code reuse and minimizing redundancy
    - Procedural decomposition

## Function: def statements

```
1  def name(arg1, arg2, ..., argN):
2      statements
```

```
1  def name(arg1, arg2, ..., argN):
2      ...
3      return value
```

# Function: def statements

```
1  def func(): ..      #Create function object
2  func()              #Call object
3  func.attr = value   #Attach attributes
```

```
1  >>> def times(x, y):   #Create and assign function
2  ...       return x*y    #Body executed when called
3  ...
4  >>> times(2, 5)
5  10
```

# Module

- Packaging of program code and data for reuse
- Provides self contained namespaces that minimize variable name clashes across programs
- The names that live in a module are called its attributes
- Typically correspond to Python program
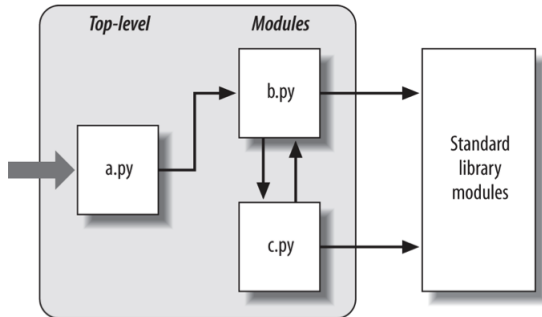- Module might be extensions coded in external languages such C++ or Java

## Module

- import – Lets a client (importer) fetch a module as a whole
- from – Allows clients to fetch particular names from a module

# Imports and Attributes

```
1   # save in b.py
2   def spam(text):
3       print text + ' spam'
```

```
1   # File a.py
2   import b           #Import module b
3   b.spam('hallo')   #Print 'hallo spam'
```

# Imports and Attributes

## Regular Expressions

- Used to generate patterns that can be used to search for strings
- Is an algebraic formula whose value is a pattern consisting of a set of strings

| **regex** | | **string** |
|-----------|---------------|-----------------------|
| a | $\rightarrow$ | a |
| ab | $\rightarrow$ | ab |
| a$\star$ | $\rightarrow$ | a, aa, aaa, aaa ... |
| a$\star$b$\star$ | $\rightarrow$ | ab, abb, aabb, aab ... |

## Regular Expressions

- What can you match with the following regular expressions?

```
1  1. ^[Tt]the\b.*
2  2. [:;]-?[\|opPD\)\(]
3  3. <.*?>
4  4. \d+\-year\-old
```

## Regular Expressions in Python

- To use Regular Expressions in Python, import the module `re`
- Then, there are two basic ways that you can use to match patterns:
    - `re.match()`
    - `re.search()`
- Both return `None` when there is no match

```python
import re

wordlist = ['farmhouse', 'greenhouse', 'guesthouse']

for w in wordlist:
    if re.match('(g.*?)(?=house)', w):
        print w
```

```python
match = re.search(pattern, string)
if match:
    process(match)
```

# Regular Expressions in Python

Another way of using regular expressions is to compile them and
reuse them as objects in your code.

```python
import re

wordlist = ['farmhouse', 'greenhouse', 'guesthouse']
regex = re.compile('(g.*?)(?=house)')

for w in wordlist:
    if regex.match(w):
        print w
```

## Python classes

```
1   class Classifier:
2       def __init__(self, lambda1, lambda2):
3           self.l1 = lambda1
4           self.l2 = lambda2
5       def train(self, data):
6           ....
7       def test(self, data):
8           ....
9   if __name__ = '__main__':
10      data = 'This is training data'
11      testdata = 'This is test data'
12      lambda1 = 0.002
13              lambda2 = 0.0005
14              model = Classifier(lambda1, lambda2)
15      model.train(data)
16      model.test(testdata)
```

- Access the data and the methods of each objects using
  `objectName.attributes` and `objectName.methods`

## Storing objects

- **Objects** save data which we might want to reuse in the future
- Use `pickle`, you can save them and load them for reuse

# import `pickle`

```python
import pickle

class Classifier:
        def __init__(self, params):
            self.params = params
        def setParams(self, params): ...
        def train(self, data): ....
        def test(self, testdata): ...

if __name__ = '__main__':
        params = [param1, param2, param3]
        data = 'This is training data'

        model = Classifier(params)
        model.train(data)
        #Store the model somewhere to reuse
        pickle.dump( model, open( 'model.p', 'wb' ) )
```

# import `pickle`

```python
import pickle

class Classifier:
    def __init__(self, params):
        self.params = params
    def setParams(self, params): ...
    def train(self, data): ....
    def test(self, testdata): ....

if __name__ = '__main__':
    testdata = 'This is test data'

    model = pickle.load(open('model.p', 'rb'))
    model.test(testdata)
```

## NumPy

- NumPy is a package supporting for **large, multi-dimensional arrays and matrices**, along with a large library of **high-level mathematical functions to operate on these arrays**
    - **ndarray** object is the core of the NumPy package
    - **ndarray** = n-dimensional arrays of homogeneous data
    - The standard mathematical and scientific packages in Python uses NumPy arrays
    - More information in http://www.numpy.org/
- NumPy will be helpful since machine learning works with high dimensional arrays

## ndarray vs. list

- An ndarray is like a list
- However, there are several differences:
    - All the elements in an ndarray should have the same type. In a list, you can have different types
    - The number of elements in an ndarray is fixed, i.e. the number of elements cannot be changed
    - ndarray in NumPy is more efficient and faster than list

## How to use NumPy

- Install NumPy (see HOWTO in
  http://www.scipy.org/scipylib/download.html)
- Take a look on the NumPy tutorial in
  www.scipy.org/Tentative_NumPy_Tutorial

# NumPy: Should know

```python
>>> import numpy as np
# Several ways to create a numpy array
>>> arr = np.array([[1, 2, 3], [4, 5, 6], [7, 8, 9]])
>>> a = np.ones((3,3), dtype=float)
>>> b = np.zeros((3,3), dtype=float)
>>> c = np.ones_like(arr)
>>> d = np.zeros_like(arr)
>>> e = np.identity(3, dtype = float)
array([[1, 0, 0],
       [0, 1, 0],
       [0, 0, 1]])
```

# NumPy: Should know

```
1  >>> import numpy as np
2  # create a numpy array from a list
3  >>> arr = np.array([[1, 2, 3], [4, 5, 6], [7, 8, 9]])
4  # returns the array dimension
5  >>> arr.ndim
6  2
7  # returns a tuple with the size of each array dimension
8  >>> arr.shape
9  (3, 3)
10 # returns the number of all the elements
11 >>> arr.size
12 9
13 # returns the transposed matrix
14 >>> arr.T
15 array([[1, 4, 7],
16        [2, 5, 8],
17        [3, 6, 9]])
18 # returns the type of all the elements
19 >>> arr.dtype
20 dtype('int64')
```

## NumPy: Should know

```
1  >>> import numpy as np
2  # create a numpy array from a list
3  >>> arr = np.array([[1, 2, 3], [4, 5, 6], [7, 8, 9]])
4  # array slicing and indexing
5  >>> arr[1]
6  array([4, 5, 6])
7  >>> arr[1][2]
8  6
9  >>> arr[:2]
10 array([1, 2, 3],
11        [4, 5, 6])
12 >>> arr[1:]
13 array([4, 5, 6],
14        [7, 8, 9])
15 >>> arr[1][:2]
16 array[4, 5]
```
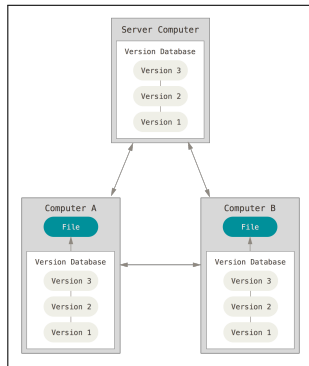
# NumPy: Should know

```python
>>> import numpy as np
>>> arr = np.array([[1, 2, 3], [4, 5, 6], [7, 8, 9]])
array([[1, 2, 3],
       [4, 5, 6],
       [7, 8, 9]])
# array methods
>>> arr.flatten()
array([1, 2, 3, 4, 5, 6, 7, 8, 9])
>>> b = np.array([[1], [4], [7])
array([[1],
       [4],
       [7]])
>>> c = np.concatenate((arr, b), axis=1)
array([[1, 2, 3, 1],
       [4, 5, 6, 4],
       [7, 8, 9, 7]])
```
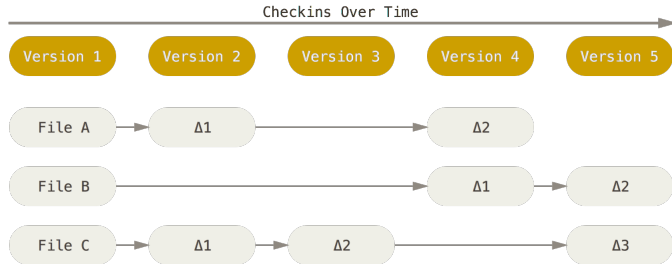
## Version control system

- >1 people work on the same documents (often code)
- It helps even if you are the only person who works on the code
  - Sometimes, you wish to go back to the past and undo your changes
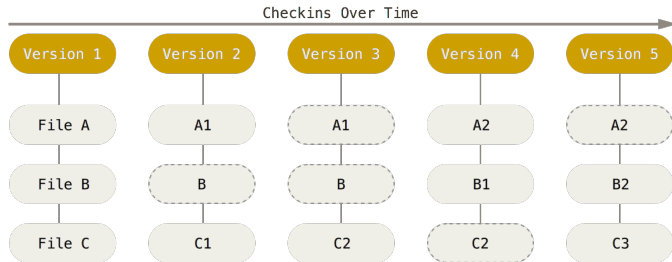- There are several VCS softwares such as svn, **git**, ...

## Git - a distributed VCS

+ A free and open
source distributed
version control system
+ Is easy to use
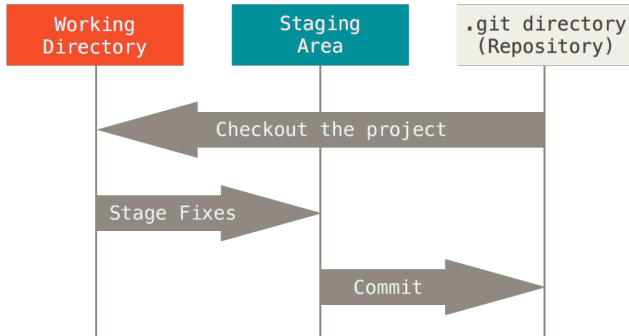+ Very helpful in many
contexts, especially in
software development

# Other subversion control systems

# Git snapshot

# A Local Git Project

## Git workflow

- Modify the files in local directory
- Stage the files and adding their snapshots in the staging area
- Do a commit, which takes the files as they are in the staging area and stores that snapshot permanently to YOUR Git directory
- At the end, if everything is done, submit the changes to the repository on the server

## How to start

- Get access to a Git server, e.g. CIP-Pool
    - Every student has a CIP-Pool access to the git server from IFI
    - You can request an account at: https://tools.rz.ifi.lmu.de/cipconf
- Install git on your machine
    - Generally you can find git at: http://git-scm.com/
    - For Unix Systems (Linux and MacOS): Most package managers have a git already as a package
    - For Windows: Windows users can get an easy installer at the above mentioned site

# Starting a project

- There are two ways: **create** or **clone**
- Create a new project

```
1  $ git init [project_name]
2  # create new repository with specified name
3  # [project_name] can be omitted in make the working
       directory the git controlled folder
```

- Clone an existing project

```
1  $ git clone /path/to/repository
2  # clone a local repository
3  $ git clone [URL]
4  # clone remote repository
```

# Make a change in your local

- Get status

```
1  $ git status
2  # Lists all new/modified files to be commited
```

- You can propose changes using:

```
1  $ git add [file]
2  # Snapshots the file in preparation for versioning
```

- Or sometime you want to remove them:

```
1  $ git reset [file]
2  # Unstages the file, but preserve its contents
```

- Then commit these changes

```
1  $ git commit -m "[descriptive message]"
2  # Records file snapshots permanently in version history
```

## Submit the local change to the server

- With the command line:

```
1 $ git push [alias] [branch]
2 # Uploads all local branch commits
```

- Alias?
  - So that you don't have to use the full URL of a remote repository every time → Git stores an alias or nickname for each remote repository URL
  - By default, if you cloned the project (as opposed to creating a new one locally), Git will automatically add the URL of the repository that you cloned from under the name 'origin'

- Branch? At that point, you have only the 'master' branch

- Therefore, you will often use

```
1 $ git push origin master
```

## Branch

- Branches are used to develop features isolated from each other
- The master branch is the "default" branch when you create a repository
- Use other branches for development and merge them back to the master branch upon completion

```
1  $ git branch
2  # Lists all local branches in the current repository
3  $ git branch [branch-name]
4  # Creates a new branch
5  $ git checkout [branch-name]
6  # Switches to the specified branch and updates the
      working directory
7  $ git merge [branch]
8  # Combines the specified branchs history into the
      current branch
9  $ git branch -d [branch-name]
10 # Deletes the specified branch
```

## Branch

- However, if you changed the same part of the same file differently in the two branches, Git will not be able to merge them

- Git adds standard conflict-resolution markers to the files that have **conflicts**, so you can open them **manually and resolve those conflicts**

```
<<<<<<< HEAD:myfile
This line is written without newlines
=======
This line is written with a
newline
>>>>>>> [branch_name]:myfile
```

- To resolve this conflict you have to remove the annotation and keep one/rewrite the lines with the conflicts.

# Update

- To get the latest updates you can use:

```
1  $ git pull
2  # Downloads bookmark history and incorporates changes
3  $ git diff
4  # Shows file differences not yet staged
5  $ git diff --staged
6  # Shows file differences between staging and the last
       file version
7  $ git fetch [bookmark]
8  # Downloads all history from the repository bookmark
```

# Redo changes & Undo files

- Undo a commit

```
1  $ git reset [commit]
2  # Undoes all commits after [commit], preserving changes
       locally
```

- You should almost never do this

```
1  $ git reset --hard [commit]
2  # Discards all history and changes back to the
       specified commit
```

## More about Git

- A compact list of all the important command lines on our website
- A good tutorial: http://gitimmersion.com/index.html
- Take the time and practice yourself

# Q&A