

lab1__python_intro

February 6, 2024

1 [COM6513] Introduction to Python for NLP

1.1 Instructor: Nikos Aletras

The goal of this lab session (**not assessed**) is to introduce you to [Python 3](#), Jupyter notebooks and main “[data science](#)” packages that we will use throughout the course. Specifically, you will be presented to NumPy, SciPy, Pandas and Matplotlib libraries which are the backbone for data manipulation, visualisation, and scientific computing in Python. You will also be presented with essential/basic approaches to text pre-processing.

1.2 Learning objectives

By the end of this session, you will be able to:

- Setup, configure and run Python Jupyter notebooks.
- understand Python basic syntax and know where to find further help (e.g. inline help in notebooks).
- remember basic text processing tricks
- use the basic Numpy, SciPy, Pandas and Matplotlib functionalities.
- have a good overview of popular packages useful for scientific computing.

1.3 Practicalities

It is strongly recommended to use [Anaconda](#). In general, for developing in Python, you could use standard IDEs like [PyCharm](#). For all the assignments we will be using [IPython](#) and [Jupyter](#) notebooks.

1.4 Python for Natural Language Processing/Machine Learning/Data Science

1.4.1 Pros

- Open source - free to install
- Large scientific community (all new cool machine learning libraries are mostly introduced in Python)
- Easy to learn
- Widely used in industry for building data science products
- Allows interfacing with C/C++ via the Cython library (<http://cython.org/>)
- Easy GPU computing via CUDA + ML libraries
- Parallelisation with OpenMP

1.4.2 Cons

- Interpreted (not-compiled) language.
- Python code might be slower compared to C/C++
- Not ideal for multithread applications; the interpreter prevents from executing one Python bytecode at a time, Global Interpreter Lock (GIL)

1.5 Essential Scientific Python Libraries that we will use in the course

1.5.1 NumPy

Numerical Python (<http://www.numpy.org/>) is the foundational package for scientific computing in Python. Most of the other scientific computing libraries are built on top of NumPy.

- Provides a fast and efficient multidimensional array object *ndarray*
- Computations and operations between arrays
- Serialisation of array objects to disk
- Linear algebra operations and other mathematical operations
- Integrating C, C++ and Fortran code to Python (e.g. BLAS and Lapack libraries)

1.5.2 SciPy

SciPy (<https://www.scipy.org/>) is scientific computing and technical computing library. SciPy contains modules for optimization, linear algebra, integration, interpolation.

1.5.3 Pandas

Pandas (<https://pandas.pydata.org/>) is a library that provides richer data structures compared to NumPy called *DataFrames*. DataFrames are similar to the ones used in R (*data.frame*) and allow sophisticated indexing functionality for reshaping, slicing and dicing, aggregating data sets.

1.5.4 Matplotlib

Matplotlib (<https://matplotlib.org/>) is the basic plotting library in Python.

1.5.5 Seaborn

Seaborn (<https://seaborn.pydata.org/>) is a visualization library based on matplotlib. It makes it easier for drawing graphs by providing a high-level interface to matplotlib. It can also directly plot Pandas dataframes.

1.5.6 IPython

IPython (Interactive Python <https://ipython.org/>) provides a platform for interactive computing (shells) that offers introspection, rich media, shell syntax, tab completion and history.

1.5.7 Jupyter Notebook

Jupyter (<https://jupyter.org/>) is a web-based application that allows to create documents (i.e. notebooks) that contain executable code, rich media and markdown inter alia. Jupyter interacts with Python via an IPython kernel.

1.5.8 Anaconda

Anaconda (<https://www.anaconda.com/>) is a free and open-source Python (and R) distribution that comes bundled with all the essential (and many more) packages for data science and machine learning. Anaconda also offers management of packages, dependencies and environments.

1.5.9 Other popular NLP and ML Libraries

- SpaCy (<https://spacy.io/>) is an open-source library for Natural Language Processing.
- NLTK (<https://www.nltk.org/>) is a package that provides text processing libraries for classification, tokenization, stemming, tagging, parsing etc..
- Scikit-learn (<https://scikit-learn.org/>) is an open-source machine learning library.
- Tensorflow (<https://www.tensorflow.org/>) and PyTorch (<https://pytorch.org/>) are popular open-source libraries for implementing and training neural network architectures.
- Keras (<https://keras.io/>) is a library that provides high level abstractions for implementing neural network architectures (built on top of Tensorflow)

Note that you are not allowed (unless it is explicitly specified) to use any of these six libraries in the assignments

1.6 Installation and Setup

Thanks to Anaconda, all above packages come in one bundle so if you use your own machine, you just need to install Anaconda for Python 3 following the instructions here: <http://docs.anaconda.com/anaconda/install/> (already installed on University's machines)

Note that Anaconda supports Windows, MacOS and Linux. Choose your OS and follow the instructions!

You can load this notebook by running:

- `$jupyter notebook lab1_python_intro.ipynb`

1.7 The Basics

You can skip this section if you are already familiar with basic Python functionality.

1.7.1 Help in Jupyter

```
[ ]: ?len

#get help for a property or a method on Jupyter/Ipypthon:
# ? followed by the method or property
```

```
[ ]: str.<TAB>

#do not run this cell, if you hit tab after typing `str.` you will get a
#list with all available properties and methods of an object
```

1.7.2 Basic Arithmetic

```
[ ]: 1+1
```

```
[ ]: 1/3
```

1.7.3 Variables

```
[ ]: x = 1+1
```

```
[ ]: x
```

Note that Python 3 supports dynamic typing:

```
[ ]: y = 5.0
     y = True
     y = 'data'
     y
```

1.7.4 Whitespace Formatting

Python uses indentation to delimit blocks of code.

```
[ ]: for i in range(3):
     print(i+10)
```

Indentation could be either 5 spaces or a tab, however it should be consistent throughout the code:

```
[ ]: for i in range(3):
     print(i+10)
     print(i-10)
```

Whitespace is ignored inside parentheses and brackets:

```
[ ]: x = [
     [1,2,3],
     [4,5,6]
     ]
     x
```

You can use a backslash to indicate that a statement continues to next line:

```
[ ]: 1 + \
     1
```

1.7.5 Modules

Not all available functionality is loaded by default, however we can load build-in or third-party modules (packages).

```
[ ]: import random
      random.gauss(0,10)
```

```
[ ]: import numpy as np
      np.random.normal([0, 0], [1,10], size=[5,2])
```

1.7.6 Functions

Functions take zero or more inputs and return an output

```
[ ]: def x_squared(x):
      return x**2 #x power of 2

      # using positional arguments
      def x_squared_new(x=3):
          return x**2 #x power of 2

      a = x_squared(2)
      b = x_squared_new()
      a, b
```

1.7.7 Strings

```
[ ]: a = "natural"
      b = 'language'
      c = a+' '+b # concatenate strings
      c, len(c) # string length
```

1.7.8 Lists

```
[ ]: l1 = [1,2,3,4,5] # list with 3 integers
      l2 = [1,'a'] # list with one int and one char
      l2[1] = 2 # update the value of the second element of the list
      l3 = [l1, l2, []] #list of lists
```

```
[ ]: # use : to slice the list
      l1[:] # [1,2,3,4,5] - copy of l1
      l1[1:4] # [2,3,4]
      l1[:2] # [1,2]
      l1[2:] # [3,4,5]
```

```
[ ]: l1[-1] # choose the last element
      l1[-3:] # last three elements
```

```
[ ]: # check list membership
      1 in [1,2] #True
      1 in [2,3] #False
```

```
[ ]: # list concatenation
x = [1,2,3]
y = [4]
x+y
```

```
[ ]: x.append(5) # append an item at the end of the list
x
```

```
[ ]: len(x) # length of the list, that's 4
```

List comprehensions

```
[ ]: [x for x in range(5)]
```

```
[ ]: [x+1 for x in range(50,100,10)]
```

1.7.9 Tuples

Similar to list but no element modifications are allowed

```
[ ]: a = (1,2,3) +(4,5)
```

```
[ ]: a
```

```
[ ]: a[0] = 6
```

1.7.10 Dictionaries

Data structures that associate keys to values

```
[ ]: d = {} #empty dictionary
d = dict() #empty dictionary
```

```
[ ]: d['Joe'] = 10
d['Mary'] = 30
d
```

```
[ ]: d.keys()
```

```
[ ]: d.values()
```

```
[ ]: d.items()
```

```
[ ]: 'Mary' in d
```

```
[ ]: len(d) # size of the dictionary
```

```
[ ]: del d['Mary'] # delete key
```

```
[ ]: d
```

1.7.11 Sets

```
[ ]: s = set()
s.add(1)
s.add(2)
s
```

```
[ ]: v = set([4,5,5,5,5,4,4,4,4,2])
v
```

```
[ ]: s & v # intersection
```

```
[ ]: s | v # union
```

1.7.12 Control Flow

```
[ ]: if 1>2:
    print("Yes!")
else:
    print('No')
```

```
[ ]: i=5
if i<0:
    i=10
elif i>0 and i<=5:
    i=20
else:
    i=30
i
```

```
[ ]: x = 0
while x < 2:
    print(x, "is less than 2")
    x+=1
```

```
[ ]: for i in range(100):
    print('Sure!')
    break
```

1.7.13 Randomness

```
[ ]: import random
```

```
[ ]: #random.random() produces numbers uniformly between 0 and 1  
v = [random.random() for _ in range(10)]
```

```
[ ]: v
```

```
[ ]: random.seed(123) # set the random seed to get reproducible results!  
random.random()
```

```
[ ]: random.seed(123)  
random.random()
```

```
[ ]: random.uniform(-1, 1) # set the boundaries to sample uniformly
```

1.7.14 Object-Oriented Programming

```
[ ]: class MLmodel:  
  
    # member functions for a linear regression model  
    # y = input * w (input: vector, w: weights)  
  
    def __init__(self, w=None, params_size=3):  
        # constructor to initialise a model  
        # self.w is the weight vector (parameters)  
        # if w is not set, we initialise it randomly  
        # note that self.w and w are different!  
        if w == None:  
            self.w = [random.uniform(-0.1,0.1) for _ in range(params_size)]  
        else:  
            self.w = w  
  
    def predict(self,X):  
        return sum([X[i]*self.w[i] for i in range(len(self.w))])  
  
    def train(self, X, Y):  
        #not implemented  
        return None
```

```
[ ]: clf = MLmodel()
```

```
[ ]: x = [2., 5., 0]  
clf.predict(x)
```



```
[ ]: clf = MLmodel(params_size=5)
      x = [2., 5., 0, 18, 9]
      clf.predict(x)
```

1.8 Text Processing Basics

```
[ ]: # raw text
      d = """
          the cat sat on the mat
      """
```

```
[ ]: d_tok = d.split() # simple whitespace tokenisation
      d_tok
```

```
[ ]: vocab = set(d_tok) # obtain a vocabulary using a set
      vocab
```

```
[ ]: #create a vocab_id to word dictionary
      id2word = enumerate(vocab)
      id2word = dict(id2word)
      id2word
```

Can you generate a word2id dictionary? E.g. {'on':0, 'the':2 ...}

```
[ ]:
```

1.8.1 Regular expressions

```
[ ]: import re

      numRE = re.compile('[0-9]+')

      numRE.findall('45 09 dfs 56352 tta& 1')
```

1.8.2 Counters

```
[ ]: from collections import Counter

      a = Counter(['a', 'foo', 'foo', 'a', 'foo'])
      a
```

1.9 Advanced list handling

```
[ ]: # create a new list by aligning elements from two lists
a = list(zip([1,2,3,4], [2,3,4,5,6]))
print(a)
```

```
[ ]: list(zip([1,2,3,4], [2,3,4,5,6], [3,4,5,6,7]))
```

```
[ ]: list(zip(*a)) # unzip a list
```

```
[ ]: # enumerate elements of a list
list(enumerate(['a','b','c']))
```

1.10 NumPy

```
[ ]: import numpy as np
```

1.10.1 Arrays

NumPy array are similar to Python lists, except that every element of an array must be of the same type.

```
[ ]: a = np.array([1, 2, 3], np.float32)
type(a)
```

```
[ ]: # equivalent to range(3) or np.arange(0,3,1)
# last argument is the step
np.arange(3)
```

```
[ ]: a[:2]
```

```
[ ]: a[0]=80
a
```

Multidimensional arrays:

```
[ ]: b = np.array([[1,2,3,4,5],[6,7,8,9,10]])
b
```

Slicing across dimensions:

```
[ ]: b[1,:]
```

```
[ ]: b[:,3]
```

```
[ ]: b[1:,2:4]
```

The *shape* property returns a tuple with the size of each dimension

```
[ ]: b.shape
```

The *dtype* property returns the data type of the array

```
[ ]: b.dtype
```

```
[ ]: # create a copy of the array with a different data type
c = b.astype(np.float)
c
```

Manipulate arrays:

```
[ ]: c.reshape((5,2)) # change shape, keep the same number of elements
```

```
[ ]: c.transpose() # transpose an array -- same as c.T
```

```
[ ]: c.flatten() # flatten an array, resulting to 1-D array
```

```
[ ]: a = np.array([1,2])
b = np.array([10,20])
c = np.array([100, 200])
# concatenate two or more arrays
np.concatenate((a, b, c))
```

```
[ ]: np.vstack((a,b)) # stack vertically
```

```
[ ]: np.hstack((b,a)) # stack horizontally
```

EXERCISE Form the 2-D array (without typing it in explicitly):

```
[[1, 6, 11],
 [2, 7, 12],
 [3, 8, 13],
 [4, 9, 14],
 [5, 10, 15]]
```

```
[ ]: # Type your answer here
```

1.10.2 Array mathematics

```
[ ]: a = np.array([1,2,3,4,5])
```

```
[ ]: a*2
```

```
[ ]: a+4
```

```
[ ]: a/2
```

```
[ ]: a**2
[ ]: np.sqrt(a)
[ ]: np.exp(a)
[ ]: a.dot(np.array([1,2,3,4,5])) # dot product
[ ]: a = np.array([[1,2],[3,4]])
      b = np.array([[0,1],[2,3]])
[ ]: a*b #elementwise multiplication
[ ]: a+b
```

1.10.3 Basic Array Operations

```
[ ]: a = np.array([1,2,3,4,5])
[ ]: a.mean()
[ ]: a.min()
[ ]: a.argmax()
[ ]: b = np.array([[0,1],[2,3]])
[ ]: b.mean(axis=0)
[ ]: b.mean(axis=1)
[ ]: c = np.array([100,-1,20,5.6])
      c.sort() # sort an array (inplace)
      c
[ ]: a
```

1.10.4 Comparison operators, value testing, item selection

```
[ ]: a = np.array([1,2,3,4,9])
      b = np.array([4,2,8,5,7])
[ ]: a == b
[ ]: a >= b
```

```
[ ]: idx = np.nonzero(a>=b) # indices of the nonzero elements e.g. True  
idx
```

```
[ ]: np.where(a>=b) # check where a>= b
```

```
[ ]: a[idx] #select elements in a where a >= b
```

```
[ ]: a.nonzero() # non-zero elements
```

```
[ ]: np.isnan(a) # check for NaN values
```

```
[ ]: idx = np.array([3,1,2]) # array of indices
```

```
[ ]: a[idx] # subset of a containing the elements in idx
```

1.10.5 Vector and matrix mathematics - Basic Linear Algebra

```
[ ]: a = np.array([[1, 2, 3],[2,3,4],[5,6,7]], np.float)  
b = np.array([0, 1, 1], np.float)
```

```
[ ]: np.dot(a,b) # dot product
```

```
[ ]: np.dot(a.T,b)
```

```
[ ]: np.inner(a, b) # inner product
```

```
[ ]: np.outer(a,b) # outer product
```

```
[ ]: np.cross(a, b) #cross product
```

```
[ ]: np.linalg.det(a) # determinant of a
```

```
[ ]: vals, vecs = np.linalg.eig(a) # eigenvalues and eigenvectors  
vals, vecs
```

```
[ ]: b = np.linalg.inv(a) # invert a matrix  
b
```

```
[ ]: U, s, Vh = np.linalg.svd(a) # Singular Value Decomposition  
U
```

1.11 Scipy

1.11.1 Statistics

```
[ ]: import numpy as np
     from scipy import stats

[ ]: x1 = np.random.uniform(-1,1, size=5)
     x2 = np.random.uniform(-1,1, size=5)

     # t-test to test whether the mean of two samples are statistical significant.
     stats.ttest_ind(x1, x2)
```

1.11.2 Sparse Matrices

Sometimes our data might be so large that cannot fit memory but contain many zeros that we do not need to store. In that case SciPy provides memory efficient sparse matrix data structures.

- csc_matrix: Compressed Sparse Column format
- csr_matrix: Compressed Sparse Row format
- bsr_matrix: Block Sparse Row format
- lil_matrix: List of Lists format
- dok_matrix: Dictionary of Keys format
- coo_matrix: COOrdinate format (aka IJV, triplet format)
- dia_matrix: DIAGONal format

See <https://docs.scipy.org/doc/scipy/reference/sparse.html> for more details.

```
[ ]: import numpy as np
     from scipy.sparse import *

[ ]: A = csr_matrix([[1, 2, 0], [0, 0, 3], [4, 0, 5]])
     # 5 non-zero elements
     A

[ ]: v = np.array([1, 0, -1])
     A.dot(v)
```

1.12 Pandas

```
[ ]: import pandas as pd
```

1.12.1 Object Creation

Series is a one-dimensional *ndarray* with axis labels.

```
[ ]: s = pd.Series([1,2,3,np.nan,6,100])
```

```
[ ]: s
```

```
[ ]: # series can contain any data type casted to an object
s = pd.Series([10,True,5,'test',6,8])
```

```
[ ]: s
```

DataFrame is a Two-dimensional size-mutable, potentially heterogeneous tabular data structure with labeled axes.

```
[ ]: df = pd.DataFrame(np.random.randn(6,4), index=None, columns=list('ABCD'))
```

```
[ ]: df
```

```
[ ]: # Specify an index
df = pd.DataFrame(np.random.randn(6,4),
                  index=pd.date_range('20180101', periods=6),
                  columns=list('ABCD'))
```

```
[ ]: df
```

1.12.2 Viewing Data

```
[ ]: df.head()
```

```
[ ]: df.tail(2)
```

```
[ ]: df.describe()
```

```
[ ]: df.T
```

```
[ ]: df.sort_index(axis=1, ascending=False)
```

```
[ ]: df.sort_values(by='D')
```

1.12.3 Data Selection

```
[ ]: df['A'] # selecting a single column
```

```
[ ]: df[0:3] # slicing rows by index
```

```
[ ]: df['20180101':'20180102'] # slicing with a labelled index

[ ]: df.loc['20180101'] # selecting by label

[ ]: df.loc[:,['A','D']] # selecting multiple columns

[ ]: df.loc['20180102':'20180105',['C','A']] # row and column slicing

[ ]: df.loc['20180102','A'] # access to scalar values

[ ]: df.iloc[3] # selecting rows using numeric indices

[ ]: # selecting rows and columns using numeric indices
df.iloc[2:4,0:2]

[ ]: df.iloc[1,0] # getting a value

[ ]: df[df['A'] > 0] # boolean indexing

[ ]: df[df > 0]

[ ]: df.iat[0,1] = 0 # setting one value

# setting the values of an entire column using a numpy array
df.loc[:, 'C'] = np.ones(len(df))

[ ]: df
```

1.12.4 Handling Missing Data

```
[ ]: df1 = df[df > -0.5]
df1

[ ]: # drop any rows that have missing data
df1.dropna(how='any')

[ ]: df1.fillna(value=0) # filling missing data
```

1.12.5 Operations

```
[ ]: df.mean() # Performing a descriptive statistic

[ ]: # Performing a descriptive statistic on the axis 1
df.mean(1)

[ ]: df.apply(np.cumsum) # Applying functions to the data
```



```
[ ]: df['D'].apply(np.exp) # Applying a function to a single column
```

```
[ ]: df.values # get DataFrame values as a numpy array
```

1.12.6 Grouping

- Splitting the data into groups based on some criteria
- Applying a function to each group independently
- Combining the results into a data structure

```
[ ]: df = pd.DataFrame({'A' : ['red', 'blue', 'blue', 'red',  
                             'red', 'red', 'blue', 'blue'],  
                      'B' : ['red', 'blue', 'green', 'green',  
                             'blue', 'blue', 'red', 'green'],  
                      'C' : np.random.randn(8),  
                      'D' : np.random.randn(8)})  
  
df
```

```
[ ]: # group by A then apply sum (think it as a SQL operation)  
df.groupby('A').sum()
```

```
[ ]: df.groupby(['A', 'B']).sum() # hierachical index
```

```
[ ]: pd.pivot_table(df, values='D', index=['A'], columns=['B'])
```

```
[ ]:
```

1.12.7 Data I/O

CSV

```
[ ]: df.to_csv('foo.csv') # write to CSV
```

```
[ ]: df = pd.read_csv('foo.csv')
```

```
[ ]: df
```

```
[ ]: df = pd.read_csv('foo.csv', index_col=0)
```

```
[ ]: df
```

HDF5

```
[ ]: df.to_hdf('foo.h5', 'df')  
pd.read_hdf('foo.h5', 'df')
```

1.13 Plotting with Matplotlib

```
[ ]: #allow plotting inline
    %matplotlib inline

import matplotlib.pyplot as plt

[ ]: loss = [4.2, 3.8, 3.7, 3.6, 3.55]
    epochs = [1,2,3,4,5]

[ ]: plt.plot(epochs,loss);

[ ]: plt.plot(epochs,loss,color='g',
            linestyle='dashdot',
            label='validation loss');

plt.title("loss monitoring")
plt.xlabel("epochs")
plt.ylabel("loss");
plt.legend();
```

Excercise: Plot a figure containing two lines

```
[ ]:
```

1.14 Wrap Up

In this session, you saw:

- How to setup, configure and run Python Jupyter notebooks.
- How Python differs from other languages; its basic syntax and know where to find further help (e.g. inline help in notebooks).
- Remember basic text processing tricks.
- Use of basic Numpy, SciPy, Pandas and Matplotlib functionalities.

More Practice:

- Go through this [Python tutorial](#) and this [one](#) (“1. Getting started with Python for science” section) on NumPy, SciPy and matplotlib libraries.

Extras:

- [Jupyter tutorial](#)
- [Introduction to Unix/Linux](#)

```
[ ]:
```