# Chapter 1: Python introduction

Werner Bauer

- Python is an open source (freely available) programming language. There are many ways to install and use Python. One possible way is:
  - **Anaconda**: a Python distribution for scientific computing,
  - **Spyder**: a Python IDE, particularly suited for Matlab users.
- Using Anaconda, you can install Spyder as follows:

```
conda install spyder # execute this in a shell
```
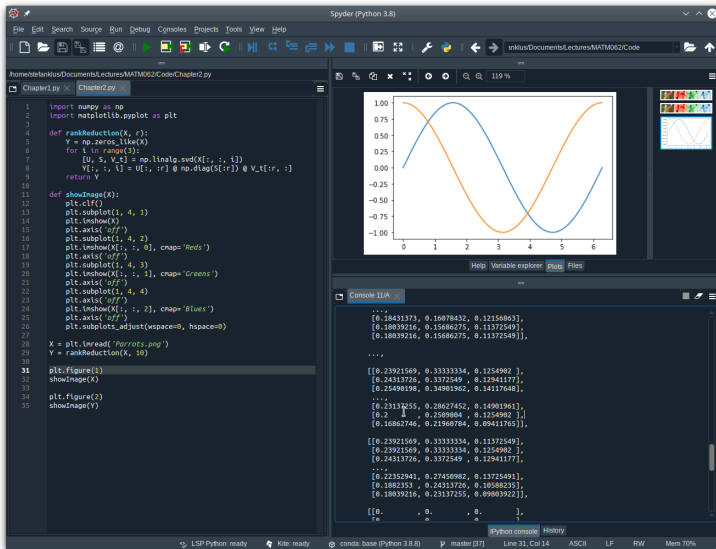
- Install also already NumPy, SciPy, and Matplotlib:
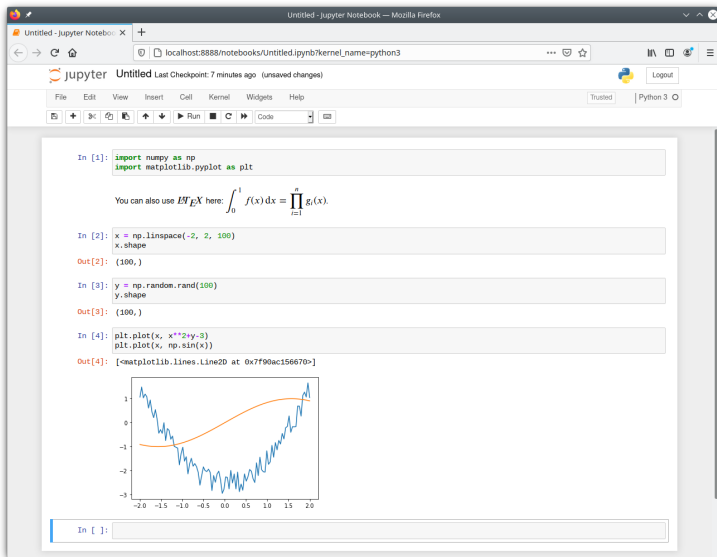
```
conda install numpy scipy matplotlib
```

- If you want to use Jupyter notebooks:

```
conda install jupyter
```

- The obligatory "Hello, World!"program:

```
print("Hello, World!")
```

# Spyder



3

Untitled - Jupyter Notebook — Mozilla Firefox

localhost:8888/notebooks/Untitled.ipynb?kernel_name=python3

Jupyter  **Untitled** Last Checkpoint: 7 minutes ago  (unsaved changes)  Logout

File  Edit  View  Insert  Cell  Kernel  Widgets  Help

Trusted | Python 3 ○

Code

```
In [1]: import numpy as np
        import matplotlib.pyplot as plt
```

You can also use $\LaTeX$ here: $\int_0^1 f(x)\,dx = \prod_{i=1}^n g_i(x)$.

```
In [2]: x = np.linspace(-2, 2, 100)
        x.shape
```
Out[2]: (100,)

```
In [3]: y = np.random.rand(100)
        y.shape
```
Out[3]: (100,)

```
In [4]: plt.plot(x, x**2+y-3)
        plt.plot(x, np.sin(x))
```
Out[4]: [<matplotlib.lines.Line2D at 0x7f90ac156670>]



```
In [ ]:
```

- Variable names must start with a letter or an underscore:

```
1  b = True      # boolean
2  n = 3         # integer
3  x = 1.41      # floating point number
4  z = 1 + 2j    # complex number
5  s = 'text'    # string
6  t = None      # empty value
```

- Print the variable:

```
1  print(b) # True
2  print(n) # 3
3  print(x) # 1.41
4  print(z) # (1+2j)
5  print(s) # text
6  print(t) # None
```

- Python automatically assigns a data type:

```
1  print(type(b)) # <class 'bool'>
2  print(type(n)) # <class 'int'>
3  print(type(x)) # <class 'float'>
4  print(type(z)) # <class 'complex'>
5  print(type(s)) # <class 'str'>
6  print(type(t)) # <class 'NoneType'>
```

```
1   b = [False, True]
2
3   for i in b:
4       print("!{} = {}".format(i, not i))
5
6   # !False = True
7   # !True = False
8
9   for i in b:
10      for j in b:
11          print("{} & {} = {}".format(i, j, i and j))
12
13  # False & False = False
14  # False & True = False
15  # True & False = False
16  # True & True = True
17
18  for i in b:
19      for j in b:
20          print("{} | {} = {}".format(i, j, i or j))
21
22  # False | False = False
23  # False | True = True
24  # True | False = True
25  # True | True = True
```

| $A$ | $\neg A$ |
|-----|----------|
| 0   | 1        |
| 1   | 0        |

| $A$ | $B$ | $A \wedge B$ |
|-----|-----|--------------|
| 0   | 0   | 0            |
| 0   | 1   | 0            |
| 1   | 0   | 0            |
| 1   | 1   | 1            |

| $A$ | $B$ | $A \vee B$ |
|-----|-----|------------|
| 0   | 0   | 0          |
| 0   | 1   | 1          |
| 1   | 0   | 1          |
| 1   | 1   | 1          |

- Note that Python uses 0-based indexing (unlike Matlab).
- Tuples are ordered and cannot be modified:

```
1  t = (2, 'a', 11.2, (3, 4), True)
2  print(t[2]) # 11.2
```

- Lists are ordered and can be modified:

```
1   l = [5, 4, 3]
2   l.append(2)     # add 2 to the list
3   print(l)        # [5, 4, 3, 2]
4   print(len(l))   # 4
5
6   l = l[::-1]     # revert list
7   print(l)        # [2, 3, 4, 5]
8
9   l.insert(1, 7)  # insert 7 at the first index
10  print(l)        # [2, 7, 3, 4, 5]
11
12  l.pop()         # remove last item
13  print(l)        # [2, 7, 3, 4]
```

- Access the last elements of a list using negative indices:

```
1  print(l[-1]) # 4
2  print(l[-2]) # 3
3  print(l[-3]) # 7
```

- Sets are unordered collections of variables:

```
1  s = {'a', 'b', 'c'}
2
3  print(s == {'b', 'a', 'c', 'a'}) # True
4
5  s.add('d') # add an element to the set
6  print(s)   # {'c', 'b', 'a', 'd'}
7
8  print('a' in s) # True
9  print('z' in s) # False
```

- Dictionaries are unordered collections of key-value pairs:

```
1  d = {'red':  [1, 0, 0],
2       'green':[0, 1, 0],
3       'blue': [0, 0, 1]}
4
5  print(d['red'])   # [1, 0, 0]
6  print(d.keys())   # dict_keys(['red', 'green', 'blue'])
7  print(d.values()) # dict_values([[1, 0, 0], [0, 1, 0], [0, 0, 1]])
8
9  print(d['yellow']) # KeyError: 'yellow'
```

- Use `defaultdict` if you want to provide a default value.

- Additional modules can be imported as follows:

```
1  import math
2  print(math.factorial(5)) # 120
3  dir(math) # prints a list of functions contained in math
```

- You can also write your own modules and import them.
- The imported modules need to be in the same directory or contained in the search path for module files.
- In Spyder, open the path manager 🐍 or use:

```
1  sys.path.insert(0, "/home/xyz/FolderName")
```

- Import only specific functions:

```
1  from math import factorial
2  factorial(5) # note that we can now write factorial instead of math.factorial
```

- It is also possible to give imported modules a new name:

```
1  import numpy as np            # these are common aliases
2  import scipy as sp            # which will often be used
3  import matplotlib.pyplot as plt  # in what follows
```

# For loops, while loops, and range

Python:

```python
for i in range(10):
    print(i)

for i in range(10, 20):
    print(i)

for i in range(20, 30, 2):
    print(i)
```

Matlab:

```matlab
for i = 0:9
    fprintf('%d\n', i);
end
for i = 10:19
    fprintf('%d\n', i);
end
for i = 20:2:28
    fprintf('%d\n', i);
end
```

Loop 1: $0, 1, 2, 3, 4, 5, 6, 7, 8, 9$
Loop 2: $10, 11, 12, 13, 14, 15, 16, 17, 18, 19$
Loop 3: $20, 22, 24, 26, 28$

Python:

```python
d = 256
while d > 1:
    d = d//2
    print(d)
```

Matlab:

```matlab
d = 256;
while d > 1
    d = d/2;
    fprintf('%d\n', d);
end
```

Output: $128, 64, 32, 16, 8, 4, 2, 1$

Use / for floating point division and // for integer division.
That is, $5/2 = 2.5$ and $5//2 = 2$.

Python:

```python
1  s = 0
2  for i in range(10):
3      s = s + i**2
4      if s > 50:
5          break
6  print(s) # 55
```

Matlab:

```matlab
1  s = 0;
2  for i = 0:9
3      s = s + i^2;
4      if s > 50
5          break
6      end
7  end
8  disp(s); % 55
```

The break statement terminates the current loop.

Python:

```python
1  for i in range(10):
2      if i % 2 == 0:
3          continue
4      print(i) # 1, 3, 5, 7, 9
```

Matlab:

```matlab
1  for i = 0:9
2      if mod(i, 2) == 0
3          continue;
4      end
5      disp(i); % 1, 3, 5, 7, 9
6  end
```

The continue statement skips all the remaining statements and
returns the control to the beginning of the loop.

Python:

```python
def fibonacci():
    x, y = 0, 1
    while True:
        x, y = y, x+y
        yield y

g = fibonacci()

print(next(g)) # 1
# ...
print(next(g)) # 2
# ...
l = [next(g) for i in range(10)] # [3, 5, 8, 13, 21, 34, 55, 89, 144, 233]
```

- If a function contains at least one yield statement, it becomes a generator function.
- Yield pauses the function, saves its current state, and continues when the function is called again.
- There is no comparable Matlab feature.

Python:

```python
1  for i in range(1, 5):
2      if i % 2 == 0:
3          print('even')
4      else:
5          print('odd')
```

Matlab:

```matlab
1  for i = 1:4
2      if mod(i, 2) == 0
3          fprintf('even\n')
4      else
5          fprintf('odd\n')
6      end
7  end
```

Output: odd, even, odd, even

Python:

```python
1  if expression:
2      statements
3  elif expression:
4      statements
5  else:
6      statements
```

Matlab:

```matlab
1  if expression
2      statements
3  elseif expression
4      statements
5  else
6      statements
7  end
```

## Indentation

Note that Python relies on the indentation to define a block of code. Other programming languages often use `begin` and `end` or { and } for this purpose.
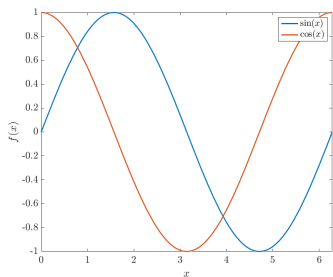
Python:

```
1   f = lambda x: math.sin(x)
2   g = lambda x: math.cos(x)
3   h = lambda x, y: x**2 + y**2
4   f(math.pi)
5   h(2, 3)
6   x = np.linspace(0, 2*math.pi, 100)
7   plt.plot(x, np.vectorize(f)(x))
8   plt.plot(x, np.vectorize(g)(x))
```

Matlab:

```
1   f = @(x) sin(x);
2   g = @(x) cos(x);
3   h = @(x, y) x^2 + y^2;
4   f(pi)
5   h(2, 3)
6   figure; hold on;
7   fplot(f, [0, 2*pi]);
8   fplot(g, [0, 2*pi]);
```

- Apply `len` to all items of the list:

```
1  tuple_lengths = list(map(len, [(1,), (1, 2), (1, 2, 3)])) # tuple_lengths = [1, 2, 3]
```

- Apply lambda function sequentially to pairs of values:

```
1  from functools import reduce
2  r = reduce((lambda x, y: x * y), [1, 2, 3, 4]) # r = 24
```

- This corresponds to:

```
1  r = 1
2  for i in [1, 2, 3, 4]:
3      r = r * i
```

- Filter values for which the lambda function returns `True`:

```
1  x = [7, -1, 3, 2, -3, -2]
2  y = [i for i in filter(lambda x: x>0, x)] # y = [7, 3, 2]
```

- You could also write:

```
1  y = [i for i in x if i > 0]
```

- Generate a list of all permutations of $[1,2,3]$:

```
1  x = [1, 2, 3]
2  y = list(itertools.permutations(x)) # the generator is converted to a list
3
4  print(y) # [(1, 2, 3), (1, 3, 2), (2, 1, 3), (2, 3, 1), (3, 1, 2), (3, 2, 1)]
```

- Generate a list of all combinations of $[1,2,3,4]$:

```
1  import itertools
2
3  x = [1, 2, 3, 4]
4  y = list(itertools.combinations(x, 3))
5
6  print(y) # [(1, 2, 3), (1, 2, 4), (1, 3, 4), (2, 3, 4)]
```

- Compute accumulated sums of $[1,2,3,4,5]$:

```
1  import itertools
2
3  x = [1, 2, 3, 4, 5]
4  y = list(itertools.accumulate(x))
5
6  print(y) # [1, 3, 6, 10, 15]
```

- Passing single argument:

```
1  def double(n):
2      return 2*n
3  double(2) # 4
```

- Passing multiple arguments to a function:

```
1  def func(*args): # args is a tuple that contains all parameters
2      for i in args:
3          print(i)
4
5  func(1, 'a', —2)
6  # 1
7  # a
8  # —2
9  l = [1,'a',—2]
10 func(*l) # * unpacks the list, same output as above
11
12 func(l) # [1, 'a', —2]
```

- Passing multiple named arguments to a function:

```
1  def func(**kwargs): # kwargs is a dictionary that contains keyword/value pairs
2      for keyword, value in kwargs.items():
3          print(f'{keyword} = {value}')
4
5  func(kw1=1, kw2='a', kw3=—2)
6  # kw1 = 1
7  # kw2 = a
8  # kw3 = —2
```

- Can of course be combined.

```
 1   class Fraction(object):
 2       def __init__(self, a, b):
 3           '''Initializer or constructor, called when you create a new object.'''
 4           d = math.gcd(a, b) # greatest common divisor of a and b
 5           self.a = a // d
 6           self.b = b // d
 7
 8       def __repr__(self):
 9           '''Prints fraction.'''
10           return "Fraction {}/{}".format(self.a, self.b)
11
12       def __add__(self, rhs):
13           '''Adds two fractions.'''
14           p = self.a * rhs.b + self.b * rhs.a
15           q = self.b * rhs.b
16           return Fraction(p, q)
17
18       def __eq__(self, rhs):
19           '''Checks if two fractions are equal.'''
20           return self.a * rhs.b == rhs.a * self.b
21
22   x1 = Fraction(1, 3) # x1 = 1/3
23   x2 = Fraction(1, 6) # x2 = 1/6
24   x3 = x1 + x2
25
26   print(x3)                    # Fraction 1/2
27   print(x3 == Fraction(2, 4)) # True
28
29   x4 = x1 — x2 # not implemented yet, define __sub__ function
```

```
 1  class Person(object):
 2      type = 'Person' # this is a class attribute, everything is by default public
 3
 4      def __init__(self, name, age):
 5          self.name = name
 6          self.age = age
 7
 8      def __str__(self):
 9          return f'Name: {self.name}\nAge: {self.age}'
10
11      def increaseAge(self):
12          self.age += 1
13
14  class Student(Person): # inherits properties and functions from Person
15      type = 'Student'
16
17      def __init__(self, name, age, urn):
18          super().__init__(name, age)
19          self.urn = urn
20
21      def __str__(self):
22          return super().__str__() + f'\nURN: {self.urn}'
23
24  c1 = Person('Alice', 21)
25  c2 = Student('Bob', 22, 720283)
26  print(c1)          # Name: Alice, Age: 21
27  print(c1.type)     # Person
28  print(c2)          # Name: Bob, Age: 22, URN: 720283
29  print(c2.type)     # Student
30  c2.increaseAge()   # calls the inherited function
31  print(c2.age)      # 23
```

- Different string formatting options:

```
1  x = 5
2  print("%d^2 = %d" % (x, x**2))      # 5^2 = 25 (old C-style formatting)
3  print("{}^2 = {}".format(x, x**2)) # 5^2 = 25 (newer)
4  print(f"{x}^2 = {x**2}")           # 5^2 = 25 (even newer)
```

- Don't add a newline to the string:

```
1  print("Don't end the line ", end="")
2  print("here!")
```

- Conversion of data types:

```
1  a = 123    # integer
2  b = '456'  # string
3
4  a = str(a) # '123'
5  b = int(b) # 456
```

- In-place operations:

```
1  a, b, c = 1, 2, 3
2  a += 2  # 3
3  b *= 3  # 6
4  c **= 4 # 81
```

- Conditional expressions:

```python
def maxVal(a, b):
    return a if a > b else b

print(maxVal(2, 5)) # 5
print(maxVal(7, 5)) # 7
```

- This is similar to the ternary operator in C++:

```cpp
#include <iostream>

template<typename T>
T maxVal(T a, T b)
{
    return a > b ? a : b;
}

int main()
{
    std::cout << maxVal(2, 5) << "\n" << maxVal(7, 5) << std::endl;
}
```

- List comprehension:

```python
l1 = [1, 3, 2, 5, 4, 6]
l2 = [i**2 for i in l1]             # [1, 9, 4, 25, 16, 36]
l3 = [2*i for i in l2 if i % 2 == 0] # [8, 32, 72]
```

# Other useful features

- Use the keyword `pass` as a placeholder for future code:

```
1  def toBeImplemented():
2      pass
```

- Enumerate items:

```
1  for index, item in enumerate(['th', 'st', 'nd', 'rd', 'th']):
2      print('{}{}'.format(index, item)) # 0th, 1st, 2nd, 3rd, 4th
```

- Create dictionaries:

```
1  d = dict(x=11, y=22, z=33)
2  print(d['y']) # 22
```

- Document your functions and classes with docstrings:

```
1  def sqr(x):
2      '''This function returns x squared.'''
3      return x**2 # note that in Matlab we would write x^2
4
5  print(sqr.__doc__) # This function returns x squared.
```

- Check whether objects are identical:

```
1  a = [1, 2, 3]
2  b = a.copy()
3  print(a == b) # True  (compares the value of a and b)
4  print(a is b) # False (compares the identities of a and b)
5
6  s = 'text'
7  t = s
8  print(s == t) # True (compares the value of s and t)
9  print(s is t) # True (compares the identities of s and t)
```

- Note that Python uses reference counting:

```
1  from sys import getrefcount
2
3  s = 'some random text'
4  print(getrefcount(s)) # 3
5  t = s
6  print(getrefcount(s)) # 4
7  t = 'new text'
8  print(getrefcount(s)) # 3
```

- You can also compare the addresses:

```
1  s = 'some random text'
2  t = s
3  print(f"Address of s: {id(s)}") # 140512670168288
4  print(f"Address of t: {id(t)}") # 140512670168288
```

# Scalars, vectors, and matrices

Python:

```
1   n = 3
2   a = 1.5
3   x = np.array([1, 0, 2])
4   A = np.array([[1, 2, 3], [4, 5, 6], [7, 8, 9]])
5   B = np.zeros((n, n))
6   C = np.ones((n, n))
7   D = np.eye(n)
8   y = A @ x
9   M = A @ C
10  N = A.T
```

Matlab:

```
1   n = 3
2   a = 1.5
3   x = [1; 0; 2]
4   A = [1, 2, 3; 4, 5, 6; 7, 8, 9]
5   B = zeros(n)
6   C = ones(n)
7   D = eye(n)
8   y = A*x
9   M = A*C
10  N = A'
```

$$n = 3,\ a = 1.5,\ x = \begin{bmatrix} 1 \\ 0 \\ 2 \end{bmatrix}$$

$$A = \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \\ 7 & 8 & 9 \end{bmatrix},\ B = \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix},\ C = \begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix},\ D = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

$$y = \begin{bmatrix} 7 \\ 16 \\ 25 \end{bmatrix},\ M = \begin{bmatrix} 6 & 6 & 6 \\ 15 & 15 & 15 \\ 24 & 24 & 24 \end{bmatrix},\ N = \begin{bmatrix} 1 & 4 & 7 \\ 2 & 5 & 8 \\ 3 & 6 & 9 \end{bmatrix}$$

Python:

```
1  n = 4
2  A = sp.linalg.hilbert(n)  # Hilbert matrix of size 4
3  c = A[:, 1]               # second column
4  r = A[0, :]               # first row
5  s = A[::2, :]             # rows one and three, all columns
```

Matlab:

```
1  n = 4
2  A = hilb(n)
3  c = A(:, 2)
4  r = A(1, :)
5  s = A(1:2:end, :)
```

$n = 4$

$$A = \begin{bmatrix} 1 & \frac{1}{2} & \frac{1}{3} & \frac{1}{4} \\ \frac{1}{2} & \frac{1}{3} & \frac{1}{4} & \frac{1}{5} \\ \frac{1}{3} & \frac{1}{4} & \frac{1}{5} & \frac{1}{6} \\ \frac{1}{4} & \frac{1}{5} & \frac{1}{6} & \frac{1}{7} \end{bmatrix}, c = \begin{bmatrix} \frac{1}{2} \\ \frac{1}{3} \\ \frac{1}{4} \\ \frac{1}{5} \end{bmatrix}$$

$$r = \begin{bmatrix} 1 & \frac{1}{2} & \frac{1}{3} & \frac{1}{4} \end{bmatrix}$$

$$s = \begin{bmatrix} 1 & \frac{1}{2} & \frac{1}{3} & \frac{1}{4} \\ \frac{1}{3} & \frac{1}{4} & \frac{1}{5} & \frac{1}{6} \end{bmatrix}$$

Python:

```
1   def f(x):
2       for i in range(len(x)):
3           x[i] = i+1
4
5   x = np.zeros(4) # [0, 0, 0, 0]
6   f(x)
7   print(x) # [1, 2, 3, 4]
```

Matlab:

```
1   x = zeros(4, 1);
2   f(x)
3   disp(x) % [0, 0, 0, 0]
4
5   function f(x)
6       for i = 1:length(x)
7           x(i) = i;
8       end
9   end
```

Python:

```
1   v = np.array(range(6)) # [0 1 2 3 4 5]
2   w = v[3:]              # [3 4 5]
3   w[0] = 99              # [99  4  5]
4   print(v)              # [0  1  2 99  4  5]
```

Matlab:

```
1   v = [0:5]     % [0 1 2 3 4 5]
2   w = v(4:end)  % [3 4 5]
3   w(1) = 99     % [99  4  5]
4   disp(v)       % [0 1 2 3 4 5]
```

Matlab creates copies of vectors and matrices, Python does not.
Use the `.copy()` function to create copies of NumPy arrays.

- Save variables into a Matlab-style MAT file:

```python
1   import numpy as np
2   import scipy as sp
3   import scipy.io
4
5   n = 100
6   x = np.linspace(-1, 1, n)
7   y = x**2
8   A = np.eye(n)
9
10  sp.io.savemat('Results.mat', {'n':n, 'x':x, 'y':y, 'A':A})
```

- Load file in Matlab:
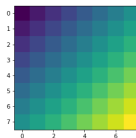
```matlab
1   load Results
2   plot(x, y)
```
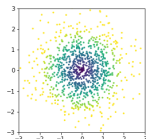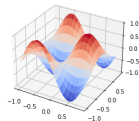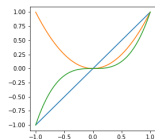
- Load variables from MAT file into main scope:

```python
1   data = scipy.io.loadmat('Results.mat', squeeze_me=True)
2   for s in data.keys():
3       if s[:2] == '__' and s[-2:] == '__': continue
4       exec('%s = data["%s"]' % (s, s))
```

```
1   import numpy as np
2   import scipy as sp
3   import scipy.linalg
4   import matplotlib.pyplot as plt
5   from matplotlib import cm
6
7   fig = plt.figure(1); plt.clf()
8
9   plt.subplot(1, 4, 1)
10  x = np.linspace(—1, 1, 100)
11  plt.plot(x, x)
12  plt.plot(x, x**2)
13  plt.plot(x, x**3)
14
15  ax = fig.add_subplot(1, 4, 2, projection='3d')
16  x = np.arange(—1, 1, 0.1)
17  y = np.arange(—1, 1, 0.1)
18  X, Y = np.meshgrid(x, y)
19  Z = np.sin(np.pi*X) * np.cos(np.pi*Y)
20  ax.plot_surface(X, Y, Z, cmap=cm.coolwarm)
21
22  plt.subplot(1, 4, 3)
23  x = np.random.randn(2, 1000)
24  plt.scatter(x[0, :], x[1, :], s=4, \
25              c=np.sqrt(x[0, :]**2 + x[1, :]**2), vmin=0, vmax=2)
26  plt.xlim(—3, 3); plt.ylim(—3, 3);
27
28  plt.subplot(1, 4, 4)
29  plt.imshow(1/sp.linalg.hilbert(8))
```

# Scalar-vector multiplication and matrix-vector multiplication

## Compute $y = \alpha x$.

Python:

```
1  import numpy as np
2  n = 5
3  alpha = 2
4  x = np.linspace(1,5,5) # [1., 2., 3., 4., 5.]
5
6  y = np.zeros(n)
7  for i in range(n):
8      y[i] = alpha*x[i]
```

Matlab:

```
1  n = 5;
2  alpha = 2;
3  x = [1:n]';
4
5  y = zeros(n, 1);
6  for i = 1:n
7      y(i) = alpha*x(i);
8  end
```

## Compute $y = Ax$.

Python:

```
1  import numpy as np
2  A = np.random.rand(n, n)
3  x = np.random.rand(n)
4
5  y = np.zeros(n)
6  for i in range(n):
7      for j in range(n):
8          y[i] = y[i] + A[i, j]*x[j]
```

Matlab:

```
1  n = 5;
2  A = rand(n);
3  x = rand(n, 1);
4
5  y = zeros(n, 1);
6  for i = 1:n
7      for j = 1:n
8          y(i) = y(i) + A(i, j)*x(j);
9      end
10 end
```

Compute $C = AB$.

Python:

```python
import numpy as np
n = 5
A = np.random.rand(n, n)
B = np.random.rand(n, n)

C = np.zeros((n, n))
for i in range(n):
    for j in range(n):
        for k in range(n):
            C[i, j] = C[i, j] + A[i, k]*B[k, j]
```

Matlab:

```matlab
n = 5;
A = rand(n);
B = rand(n);

C = zeros(n);
for i = 1:n
    for j = 1:n
        for k = 1:n
            C(i, j) = C(i, j) + A(i, k)*B(k, j);
        end
    end
end
```

- Importing Pandas (docu: `https://pandas.pydata.org/`)

```
1   import pandas as pd
2   # Creating series object 'data':
3   data = pd.Series([0.25, 0.5, 0.75, 1.0], index=['a', 'b', 'c', 'd'])
4   data
5   # Output:
6   a    0.25
7   b    0.50
8   c    0.75
9   d    1.00
10  dtype: float64 # NOTE: in the following we skip this line of output
```

- Accessing data (using dictionary-like expressions)

```
1   data['b'] # 0.5
2   'a' in data # True
3   data.keys() # Index(['a', 'b', 'c', 'd'], dtype='object')
4   list(data.items()) # [('a', 0.25), ('b', 0.5), ('c', 0.75), ('d', 1.0)]
```

- Modify series objects (similar to dict-syntax):

```
1   data['e'] = 1.25 # extend Series by assigning to a new index a value:
2   data # Output:
3   a    0.25
4   b    0.50
5   c    0.75
6   d    1.00
7   e    1.25
```

- slicing by explicit and implicit index

```
1  data['a':'c'] # explicit indexing
2  # Output:
3  a    0.25
4  b    0.50
5  data[0:2] # implicit indexing
6  # Output:
7  a    0.25
8  b    0.50
```

- Masking

```
1   data[(data > 0.3) & (data < 0.8)]
2  # Output:
3  b    0.50
4  c    0.75
```

- Fancy indexing (pass a list of indices):

```
1  data[['a', 'e']]
2  # Output:
3  a    0.25
4  e    1.25
```

Danger of confusion: if a series has explicit integer index an indexing operation such as `data[1]` will use the explicit indices, while a slicing operation like `data[1:3]` will use the implicit Python-style index

Alternative indexing:

```
1    data = pd.Series(['a', 'b', 'c'], index=[1, 3, 5])
```

- The loc attribute allows indexing and slicing that always references the explicit index:

```
1    data.loc[1]
2    # Output:
3    'a'
4
5    data.loc[1:3] # this includes key—value pair with external index 3
6    # Output:
7    1    a
8    3    b
```

- The iloc attribute allows indexing and slicing that always references the implicit Python-style index:

```
1    data.iloc[1]
2    # Output:
3    'b'
4
5    data.iloc[1:2] # using implicit indexing, last value for last index not included
6    # Output:
7    3    b
8
9    data.iloc[1:3]
10   # Output:
11   3    b
12   5    c
```

- DataFrame as a dictionary of related Series objects:

```
area = pd.Series({'California': 423967, 'Texas': 695662, 'New York': 141297,
                  'Florida': 170312, 'Illinois': 149995})
pop = pd.Series({'California': 38332521, 'Texas': 26448193, 'New York': 19651127,
                 'Florida': 19552860, 'Illinois': 12882135})
data = pd.DataFrame({'area':area, 'pop':pop})

print(data) # Output:
              area        pop
California  423967   38332521
Texas       695662   26448193
New York    141297   19651127
Florida     170312   19552860
Illinois    149995   12882135

print(data['area']) # Output:
California    423967
Texas         695662
New York      141297
Florida       170312
Illinois      149995
```

- Dictonary-style syntax similar to Series objects, e.g.:

```
data['density'] = data['pop'] / data['area']

print(data) # Output:
              area        pop     density
California  423967   38332521   90.413926
Texas       695662   26448193   38.018740
New York    141297   19651127  139.076746
Florida     170312   19552860  114.806121
Illinois    149995   12882135   85.883763
```

- We can view DataFrame as enhanced two-dimensional array:

```
1   data.T # calculate transpose
2   # Output:
3              California        Texas       New York        Florida       Illinois
4   area       4.239670e+05  6.956620e+05  1.412970e+05  1.703120e+05  1.499950e+05
5   pop        3.833252e+07  2.644819e+07  1.965113e+07  1.955286e+07  1.288214e+07
6   density    9.041393e+01  3.801874e+01  1.390767e+02  1.148061e+02  8.588376e+01
7
8   data.values
9   # Output:
10  array([[4.23967000e+05, 3.83325210e+07, 9.04139261e+01],
11         [6.95662000e+05, 2.64481930e+07, 3.80187404e+01],
12         [1.41297000e+05, 1.96511270e+07, 1.39076746e+02],
13         [1.70312000e+05, 1.95528600e+07, 1.14806121e+02],
14         [1.49995000e+05, 1.28821350e+07, 8.58837628e+01]])
```

- DataFrame indexing follows style of dictionary:

```
1   data.values[0]
2   # output
3   array([4.23967000e+05, 3.83325210e+07, 9.04139261e+01])
4
5   data['area']
6   # Output:
7   California    423967
8   Texas         695662
9   New York      141297
10  Florida       170312
11  Illinois      149995
12
13  data['area']['California'] # Output:
14  423967
```

# DataFrame (DF): Access and modify data in numpy like fashion:

- `iloc` (implicit indexing): index and column labels maintained

```
1  data.iloc[:3, :2] # Output:
2             area        pop
3  California  423967  38332521
4  Texas       695662  26448193
5  New York    141297  19651127
```

- `loc` (explicit indexing): index and column labels maintained

```
1  data.loc[:'Illinois', :'pop'] # Output:
2             area        pop
3  California  423967  38332521
4  Texas       695662  26448193
5  New York    141297  19651127
6  Florida     170312  19552860
7  Illinois    149995  12882135
```

- Numpy masking, filtering, indexing works:

```
1   data.loc[data.density > 100, ['pop', 'density']] # Output:
2                pop      density
3   New York  19651127  139.076746
4   Florida   19552860  114.806121
5
6   data[0:1] # Output:
7                area        pop      density
8   California  423967  38332521   90.413926
9
10  data.iloc[0, 2] = 90 # modify value similar to numpy
11  data[0:1] # Output:
12               area        pop  density
13  California  423967  38332521     90.0
```

40

- Concatenating:

```
df1 = pd.DataFrame({'A': ['A0', 'A1', 'A2', 'A3'], 'B': ['B0', 'B1', 'B2', 'B3'],
                    'C': ['C0', 'C1', 'C2', 'C3'], 'D': ['D0', 'D1', 'D2', 'D3']},
                    index=[0, 1, 2, 3])

df2 = pd.DataFrame({'A': ['A4', 'A5', 'A6', 'A7'], 'B': ['B4', 'B5', 'B6', 'B7'],
                    'C': ['C4', 'C5', 'C6', 'C7'], 'D': ['D4', 'D5', 'D6', 'D7']},
                    index=[4, 5, 6, 7])

df3 = pd.DataFrame({'A': ['A8', 'A9', 'A10', 'A11'], 'B': ['B8', 'B9', 'B10', 'B11'],
                    'C': ['C8', 'C9', 'C10', 'C11'], 'D': ['D8', 'D9', 'D10', 'D11']},
                    index=[8, 9, 10, 11])


pd.concat([df1,df2,df3])

# Output:
      A    B    C    D
0    A0   B0   C0   D0
1    A1   B1   C1   D1
2    A2   B2   C2   D2
3    A3   B3   C3   D3
4    A4   B4   C4   D4
5    A5   B5   C5   D5
6    A6   B6   C6   D6
7    A7   B7   C7   D7
8    A8   B8   C8   D8
9    A9   B9   C9   D9
10   A10  B10  C10  D10
11   A11  B11  C11  D11
```

- Concatenating:

```
pd.concat([df1,df2,df3],axis=1)

# Output:

       A    B    C    D    A    B    C    D    A    B    C    D
0     A0   B0   C0   D0  NaN  NaN  NaN  NaN  NaN  NaN  NaN  NaN
1     A1   B1   C1   D1  NaN  NaN  NaN  NaN  NaN  NaN  NaN  NaN
2     A2   B2   C2   D2  NaN  NaN  NaN  NaN  NaN  NaN  NaN  NaN
3     A3   B3   C3   D3  NaN  NaN  NaN  NaN  NaN  NaN  NaN  NaN
4    NaN  NaN  NaN  NaN   A4   B4   C4   D4  NaN  NaN  NaN  NaN
5    NaN  NaN  NaN  NaN   A5   B5   C5   D5  NaN  NaN  NaN  NaN
6    NaN  NaN  NaN  NaN   A6   B6   C6   D6  NaN  NaN  NaN  NaN
7    NaN  NaN  NaN  NaN   A7   B7   C7   D7  NaN  NaN  NaN  NaN
8    NaN  NaN  NaN  NaN  NaN  NaN  NaN  NaN   A8   B8   C8   D8
9    NaN  NaN  NaN  NaN  NaN  NaN  NaN  NaN   A9   B9   C9   D9
10   NaN  NaN  NaN  NaN  NaN  NaN  NaN  NaN  A10  B10  C10  D10
11   NaN  NaN  NaN  NaN  NaN  NaN  NaN  NaN  A11  B11  C11  D11
```

- Merging DF: merge via a common Series (same column name by default)

```
left = pd.DataFrame({'key': ['K0', 'K1', 'K2', 'K3'],
                     'A': ['A0', 'A1', 'A2', 'A3'],
                     'B': ['B0', 'B1', 'B2', 'B3']})

right = pd.DataFrame({'key': ['K0', 'K1', 'K2', 'K3'],
                      'C': ['C0', 'C1', 'C2', 'C3'],
                      'D': ['D0', 'D1', 'D2', 'D3']})


pd.merge(left,right,how='inner',on='key')

# Output:
  key   A    B    C    D
0  K0  A0   B0   C0   D0
1  K1  A1   B1   C1   D1
2  K2  A2   B2   C2   D2
3  K3  A3   B3   C3   D3
```

- Merging dataframes: inner, outer, left, right

```
left = pd.DataFrame({'key1': ['K0', 'K0', 'K1', 'K2'],
                     'key2': ['K0', 'K1', 'K0', 'K1'],
                      'A': ['A0', 'A1', 'A2', 'A3'],
                      'B': ['B0', 'B1', 'B2', 'B3']})
right = pd.DataFrame({'key1': ['K0', 'K1', 'K1', 'K2'],
                      'key2': ['K0', 'K0', 'K0', 'K0'],
                       'C': ['C0', 'C1', 'C2', 'C3'],
                       'D': ['D0', 'D1', 'D2', 'D3']})

print(pd.merge(left, right, how='inner', on=['key1', 'key2'])) #Output:
  key1 key2   A   B   C   D
0  K0   K0  A0  B0  C0  D0
1  K1   K0  A2  B2  C1  D1
2  K1   K0  A2  B2  C2  D2

print(pd.merge(left, right, how='outer', on=['key1', 'key2'])) #Output:
  key1 key2    A    B    C    D
0  K0   K0   A0   B0   C0   D0
1  K0   K1   A1   B1  NaN  NaN
2  K1   K0   A2   B2   C1   D1
3  K1   K0   A2   B2   C2   D2
4  K2   K1   A3   B3  NaN  NaN
5  K2   K0  NaN  NaN   C3   D3

print(pd.merge(left, right, how='left',on=['key1', 'key2'])) #Output:
  key1 key2   A   B    C    D
0  K0   K0  A0  B0   C0   D0
1  K0   K1  A1  B1  NaN  NaN
2  K1   K0  A2  B2   C1   D1
3  K1   K0  A2  B2   C2   D2
4  K2   K1  A3  B3  NaN  NaN
# ..., how = 'right', .. works similarly
```

# Groupby

To group rows together and call aggregate functions:

```python
data = {'Company':['GOOG','GOOG','MSFT','MSFT','FB','FB'],
        'Person':['Sam','Charlie','Amy','Vanessa','Carl','Sarah'],
        'Sales':[200,120,340,124,243,350]}
df = pd.DataFrame(data)
```

Based on column name, .groupby() groups rows together, e.g. by 'Company'. This will create a DataFrameGroupBy object.

```python
print(df.groupby('Company').describe()) # Output:
          Sales
        count   mean         std    min     25%    50%     75%    max
Company
FB        2.0  296.5   75.660426  243.0  269.75  296.5  323.25  350.0
GOOG      2.0  160.0   56.568542  120.0  140.00  160.0  180.00  200.0
MSFT      2.0  232.0  152.735065  124.0  178.00  232.0  286.00  340.0

print(df.groupby('Company').count()) # Output:
         Person  Sales
Company
FB            2      2
GOOG         2      2
MSFT         2      2

print(df.groupby('Company').mean()) # Output
           Sales
Company
FB         296.5
GOOG       160.0
MSFT       232.0

print(df.groupby('Company').std()) # Output similar to .mean()
```

- isnull() and notnull() return Boolean mask over data

```python
import numpy as np
data = pd.Series([1, np.nan, 'hello', None])
data
# Output:
0        1
1      NaN
2    hello
3     None
data.isnull()
# Output:
0    False
1     True
2    False
3     True
data[data.notnull()] # similarly: data.dropna()
# Output:
0        1
2    hello
```

- .dropna() to drop null values

```
1   data = pd.Series([1, np.nan, 'hello', None])
2
3   data.dropna()
4   # Output:
5   0          1
6   2      hello
7
8   df = pd.DataFrame({'A':[1,2,np.nan],
9                      'B':[5,np.nan,np.nan],
10                     'C':[1,2,3]})
11  df
12  # Output:
13          A       B       C
14  0      1.0     5.0     1
15  1      2.0     NaN     2
16  2      NaN     NaN     3
17
18  df.dropna(axis=0)
19  # Output:
20       A    B  C
21  0  1.0  5.0  1
22
23  df.dropna(axis=1)
24  # Output:
25          C
26  0       1
27  1       2
28  2       3
```

# Filling null values with .fillna()

- Sometimes we want to fill rather then drop null values

```
df = pd.DataFrame({'A':[1,2,np.nan], 'B':[5,np.nan,np.nan], 'C':[1,2,3]})
print(df) # Output:
     A    B  C
0  1.0  5.0  1
1  2.0  NaN  2
2  NaN  NaN  3

df.fillna(0) # fill NA with single values such as 0
     A    B  C
0  1.0  5.0  1
1  2.0  0.0  2
2  0.0  0.0  3
df.fillna(method='ffill') # forward—fill (within columns); also backward—fill 'bfill' possible
     A    B  C
0  1.0  5.0  1
1  2.0  5.0  2
2  2.0  5.0  3

df['A'].fillna(value=df['A'].mean()) # fill missing values per columns with e.g. mean:
0    1.0
1    2.0
2    1.5

print(df) # NOTE: df is not changed !!!
     A    B  C
0  1.0  5.0  1
1  2.0  NaN  2
2  NaN  NaN  3
```

- NOTE: use `inplace=True` to overwrite values in df itself (default, as above: `inplace=False`)

48