

Data Science Cheat Sheet

Python Basics

BASICS, PRINTING AND GETTING HELP

x = 3 - Assign 3 to the variable **x**
print(x) - Print the value of **x**
type(x) - Return the type of the variable **x** (in this case, **int** for integer)

help(x) - Show documentation for the **str** data type
help(print) - Show documentation for the **print()** function

READING FILES

```
f = open("my_file.txt", "r")
file_as_string = f.read()
- Open the file my_file.txt and assign its contents to s
```

```
import csv
f = open("my_dataset.csv", "r")
csvreader = csv.reader(f)
csv_as_list = list(csvreader)
- Open the CSV file my_dataset.csv and assign its data to the list of lists csv_as_list
```

STRINGS

s = "hello" - Assign the string "hello" to the variable **s**

```
s = """She said,
there's a good idea.
"""
```

- Assign a multi-line string to the variable **s**. Also used to create strings that contain both " and ' characters

len(s) - Return the number of characters in **s**

s.startswith("hel") - Test whether **s** starts with the substring "hel"

s.endswith("lo") - Test whether **s** ends with the substring "lo"

"{} plus {} is {}".format(3,1,4) - Return the string with the values 3, 1, and 4 inserted

s.replace("e", "z") - Return a new string based on **s** with all occurrences of "e" replaced with "z"

s.split(" ") - Split the string **s** into a list of strings, separating on the character " " and return that list

NUMERIC TYPES AND

MATHEMATICAL OPERATIONS

i = int("5") - Convert the string "5" to the integer 5 and assign the result to **i**

f = float("2.5") - Convert the string "2.5" to the float value 2.5 and assign the result to **f**

5 + 5 - Addition

5 - 5 - Subtraction

10 / 2 - Division

5 * 2 - Multiplication

3 ** 2 - Raise 3 to the power of 2 (or 3²)

27 ** (1/3) - The 3rd root of 27 (or ³√27)

x += 1 - Assign the value of **x + 1** to **x**

x -= 1 - Assign the value of **x - 1** to **x**

LISTS

l = [100, 21, 88, 3] - Assign a list containing the integers 100, 21, 88, and 3 to the variable **l**

l = list() - Create an empty list and assign the result to **l**

l[0] - Return the first value in the list **l**

l[-1] - Return the last value in the list **l**

l[1:3] - Return a slice (list) containing the second and third values of **l**

len(l) - Return the number of elements in **l**

sum(l) - Return the sum of the values of **l**

min(l) - Return the minimum value from **l**

max(l) - Return the maximum value from **l**

l.append(16) - Append the value 16 to the end of **l**

l.sort() - Sort the items in **l** in ascending order

" ".join(["A", "B", "C", "D"]) - Converts the list ["A", "B", "C", "D"] into the string "A B C D"

DICTIONARIES

d = {"CA": "Canada", "GB": "Great Britain", "IN": "India"} - Create a dictionary with keys of "CA", "GB", and "IN" and corresponding values of "Canada", "Great Britain", and "India"

d["GB"] - Return the value from the dictionary **d** that has the key "GB"

d.get("AU", "Sorry") - Return the value from the dictionary **d** that has the key "AU", or the string "Sorry" if the key "AU" is not found in **d**

d.keys() - Return a list of the keys from **d**

d.values() - Return a list of the values from **d**

d.items() - Return a list of (key, value) pairs from **d**

MODULES AND FUNCTIONS

The body of a function is defined through indentation.

import random - Import the module **random**

from math import sqrt - Import the function **sqrt** from the module **math**

```
def calculate(addition_one, addition_two,
              exponent=1, factor=1):
    result = (value_one + value_two) ** exponent * factor
    return result
```

- Define a new function **calculate** with two required and two optional named arguments which calculates and returns a result.

addition(3, 5, factor=10) - Run the **addition** function with the values 3 and 5 and the named argument **10**

BOOLEAN COMPARISONS

x == 5 - Test whether **x** is equal to 5

x != 5 - Test whether **x** is not equal to 5

x > 5 - Test whether **x** is greater than 5

x < 5 - Test whether **x** is less than 5

x >= 5 - Test whether **x** is greater than or equal to 5

x <= 5 - Test whether **x** is less than or equal to 5

x == 5 or name == "alfred" - Test whether **x** is equal to 5 or **name** is equal to "alfred"

x == 5 and name == "alfred" - Test whether **x** is equal to 5 and **name** is equal to "alfred"

5 in l - Checks whether the value 5 exists in the list **l**

"GB" in d - Checks whether the value "GB" exists in the keys for **d**

IF STATEMENTS AND LOOPS

The body of if statements and loops are defined through indentation.

```
if x > 5:
    print("{} is greater than five".format(x))
elif x < 0:
    print("{} is negative".format(x))
else:
    print("{} is between zero and five".format(x))
```

- Test the value of the variable **x** and run the code body based on the value

for value in l:

```
    print(value)
```

- Iterate over each value in **l**, running the code in the body of the loop with each iteration

while x < 10:

```
    x += 1
```

- Run the code in the body of the loop until the value of **x** is no longer less than **10**

Data Science Cheat Sheet

Python - Intermediate

KEY BASICS, PRINTING AND GETTING HELP

This cheat sheet assumes you are familiar with the content of our [Python Basics Cheat Sheet](#)

s - A Python string variable
i - A Python integer variable
f - A Python float variable

l - A Python list variable
d - A Python dictionary variable

LISTS

l.pop(3) - Returns the fourth item from **l** and deletes it from the list

l.remove(x) - Removes the first item in **l** that is equal to **x**

l.reverse() - Reverses the order of the items in **l**

l[1::2] - Returns every second item from **l**, commencing from the 1st item

l[-5:] - Returns the last 5 items from **l** specific axis

STRINGS

s.lower() - Returns a lowercase version of **s**

s.title() - Returns **s** with the first letter of every word capitalized

"23".zfill(4) - Returns **"0023"** by left-filling the string with **0**'s to make it's length **4**.

s.splitlines() - Returns a list by splitting the string on any newline characters.

Python strings share some common methods with lists

s[:5] - Returns the first **5** characters of **s**

"fri" + "end" - Returns **"friend"**

"end" in s - Returns **True** if the substring **"end"** is found in **s**

RANGE

Range objects are useful for creating sequences of integers for looping.

range(5) - Returns a sequence from **0** to **4**

range(2000, 2018) - Returns a sequence from **2000** to **2017**

range(0, 11, 2) - Returns a sequence from **0** to **10**, with each item incrementing by **2**

range(0, -10, -1) - Returns a sequence from **0** to **-9**

list(range(5)) - Returns a list from **0** to **4**

DICTIONARIES

max(d, key=d.get) - Return the key that corresponds to the largest value in **d**

min(d, key=d.get) - Return the key that corresponds to the smallest value in **d**

SETS

my_set = set(l) - Return a **set** object containing the unique values from **l**

len(my_set) - Returns the number of objects in **my_set** (or, the number of unique values from **l**)

a in my_set - Returns **True** if the value **a** exists in **my_set**

REGULAR EXPRESSIONS

import re - Import the Regular Expressions module

re.search("abc", s) - Returns a **match** object if the regex **"abc"** is found in **s**, otherwise **None**

re.sub("abc", "xyz", s) - Returns a string where all instances matching regex **"abc"** are replaced by **"xyz"**

LIST COMPREHENSION

A one-line expression of a for loop

[i ** 2 for i in range(10)] - Returns a list of the squares of values from **0** to **9**

[s.lower() for s in l_strings] - Returns the list **l_strings**, with each item having had the **.lower()** method applied

[i for i in l_floats if i < 0.5] - Returns the items from **l_floats** that are less than **0.5**

FUNCTIONS FOR LOOPING

```
for i, value in enumerate(l):
    print("The value of item {} is {}".format(i, value))
```

- Iterate over the list **l**, printing the index location of each item and its value

```
for one, two in zip(l_one, l_two):
    print("one: {}, two: {}".format(one, two))
```

- Iterate over two lists, **l_one** and **l_two** and print each value

```
while x < 10:
    x += 1
```

- Run the code in the body of the loop until the value of **x** is no longer less than **10**

DATETIME

import datetime as dt - Import the **datetime** module

now = dt.datetime.now() - Assign **datetime** object representing the current time to **now**

wks4 = dt.datetime.timedelta(weeks=4) - Assign a **timedelta** object representing a timespan of 4 weeks to **wks4**

now - wks4 - Return a **datetime** object representing the time 4 weeks prior to **now**

newyear_2020 = dt.datetime(year=2020, month=12, day=31) - Assign a **datetime** object representing December 25, 2020 to **newyear_2020**

newyear_2020.strftime("%A, %b %d, %Y") - Returns **"Thursday, Dec 31, 2020"**

dt.datetime.strptime('Dec 31, 2020', "%b %d, %Y") - Return a **datetime** object representing December 31, 2020

RANDOM

import random - Import the **random** module

random.random() - Returns a random float between **0.0** and **1.0**

random.randint(0, 10) - Returns a random integer between **0** and **10**

random.choice(l) - Returns a random item from the list **l**

COUNTER

from collections import Counter - Import the **Counter** class

c = Counter(l) - Assign a **Counter** (dict-like) object with the counts of each unique item from **l**, to **c**

c.most_common(3) - Return the 3 most common items from **l**

TRY/EXCEPT

Catch and deal with Errors

l_ints = [1, 2, 3, "", 5] - Assign a list of integers with one missing value to **l_ints**

```
l_floats = []
for i in l_ints:
    try:
        l_floats.append(float(i))
    except:
        l_floats.append(i)
```

- Convert each value of **l_ints** to a float, catching and handling **ValueError: could not convert string to float:** where values are missing.

Data Science Cheat Sheet

NumPy

KEY

We'll use shorthand in this cheat sheet

arr - A numpy Array object

IMPORTS

Import these to start

`import numpy as np`

IMPORTING/EXPORTING

`np.loadtxt('file.txt')` - From a text file

`np.genfromtxt('file.csv', delimiter=',')`
- From a CSV file

`np.savetxt('file.txt', arr, delimiter='')`
- Writes to a text file

`np.savetxt('file.csv', arr, delimiter=',')`
- Writes to a CSV file

CREATING ARRAYS

`np.array([1,2,3])` - One dimensional array

`np.array([(1,2,3), (4,5,6)])` - Two dimensional array

`np.zeros(3)` - 1D array of length 3 all values 0

`np.ones((3,4))` - 3x4 array with all values 1

`np.eye(5)` - 5x5 array of 0 with 1 on diagonal (Identity matrix)

`np.linspace(0,100,6)` - Array of 6 evenly divided values from 0 to 100

`np.arange(0,10,3)` - Array of values from 0 to less than 10 with step 3 (eg [0,3,6,9])

`np.full((2,3),8)` - 2x3 array with all values 8

`np.random.rand(4,5)` - 4x5 array of random floats between 0-1

`np.random.rand(6,7)*100` - 6x7 array of random floats between 0-100

`np.random.randint(5,size=(2,3))` - 2x3 array with random ints between 0-4

INSPECTING PROPERTIES

`arr.size` - Returns number of elements in **arr**

`arr.shape` - Returns dimensions of **arr** (rows, columns)

`arr.dtype` - Returns type of elements in **arr**

`arr.astype(dtype)` - Convert **arr** elements to type **dtype**

`arr.tolist()` - Convert **arr** to a Python list

`np.info(np.eye)` - View documentation for **np.eye**

COPYING/SORTING/RESHAPING

`np.copy(arr)` - Copies **arr** to new memory

`arr.view(dtype)` - Creates view of **arr** elements with type **dtype**

`arr.sort()` - Sorts **arr**

`arr.sort(axis=0)` - Sorts specific axis of **arr**

`two_d_arr.flatten()` - Flattens 2D array **two_d_arr** to 1D

`arr.T` - Transposes **arr** (rows become columns and vice versa)

`arr.reshape(3,4)` - Reshapes **arr** to 3 rows, 4 columns without changing data

`arr.resize((5,6))` - Changes **arr** shape to 5x6 and fills new values with 0

ADDING/REMOVING ELEMENTS

`np.append(arr, values)` - Appends **values** to end of **arr**

`np.insert(arr,2, values)` - Inserts **values** into **arr** before index 2

`np.delete(arr,3,axis=0)` - Deletes row on index 3 of **arr**

`np.delete(arr,4,axis=1)` - Deletes column on index 4 of **arr**

COMBINING/SPLITTING

`np.concatenate((arr1,arr2),axis=0)` - Adds **arr2** as rows to the end of **arr1**

`np.concatenate((arr1,arr2),axis=1)` - Adds **arr2** as columns to end of **arr1**

`np.split(arr,3)` - Splits **arr** into 3 sub-arrays

`np.hsplit(arr,5)` - Splits **arr** horizontally on the 5th index

INDEXING/SLICING/SUBSETTING

`arr[5]` - Returns the element at index 5

`arr[2,5]` - Returns the 2D array element on index [2][5]

`arr[1]=4` - Assigns array element on index 1 the value 4

`arr[1,3]=10` - Assigns array element on index [1][3] the value 10

`arr[0:3]` - Returns the elements at indices 0,1,2 (On a 2D array: returns rows 0,1,2)

`arr[0:3,4]` - Returns the elements on rows 0,1,2 at column 4

`arr[:2]` - Returns the elements at indices 0,1 (On a 2D array: returns rows 0,1)

`arr[:,1]` - Returns the elements at index 1 on all rows

`arr<5` - Returns an array with boolean values

`(arr<3) & (arr>5)` - Returns an array with boolean values

`~arr` - Inverts a boolean array

`arr[arr<5]` - Returns array elements smaller than 5

SCALAR MATH

`np.add(arr,1)` - Add 1 to each array element

`np.subtract(arr,2)` - Subtract 2 from each array element

`np.multiply(arr,3)` - Multiply each array element by 3

`np.divide(arr,4)` - Divide each array element by 4 (returns **np.nan** for division by zero)

`np.power(arr,5)` - Raise each array element to the 5th power

VECTOR MATH

`np.add(arr1,arr2)` - Elementwise add **arr2** to **arr1**

`np.subtract(arr1,arr2)` - Elementwise subtract **arr2** from **arr1**

`np.multiply(arr1,arr2)` - Elementwise multiply **arr1** by **arr2**

`np.divide(arr1,arr2)` - Elementwise divide **arr1** by **arr2**

`np.power(arr1,arr2)` - Elementwise raise **arr1** raised to the power of **arr2**

`np.array_equal(arr1,arr2)` - Returns True if the arrays have the same elements and shape

`np.sqrt(arr)` - Square root of each element in the array

`np.sin(arr)` - Sine of each element in the array

`np.log(arr)` - Natural log of each element in the array

`np.abs(arr)` - Absolute value of each element in the array

`np.ceil(arr)` - Rounds up to the nearest int

`np.floor(arr)` - Rounds down to the nearest int

`np.round(arr)` - Rounds to the nearest int

STATISTICS

`np.mean(arr,axis=0)` - Returns mean along specific axis

`arr.sum()` - Returns sum of **arr**

`arr.min()` - Returns minimum value of **arr**

`arr.max(axis=0)` - Returns maximum value of specific axis

`np.var(arr)` - Returns the variance of array

`np.std(arr,axis=1)` - Returns the standard deviation of specific axis

`arr.corrcoef()` - Returns correlation coefficient of array

Data Science Cheat Sheet

Pandas

KEY

We'll use shorthand in this cheat sheet

df - A pandas DataFrame object
s - A pandas Series object

IMPORTS

Import these to start
import pandas as pd
import numpy as np

IMPORTING DATA

pd.read_csv(filename) - From a CSV file
pd.read_table(filename) - From a delimited text file (like TSV)
pd.read_excel(filename) - From an Excel file
pd.read_sql(query, connection_object) - Reads from a SQL table/database
pd.read_json(json_string) - Reads from a JSON formatted string, URL or file.
pd.read_html(url) - Parses an html URL, string or file and extracts tables to a list of dataframes
pd.read_clipboard() - Takes the contents of your clipboard and passes it to **read_table()**
pd.DataFrame(dict) - From a dict, keys for columns names, values for data as lists

EXPORTING DATA

df.to_csv(filename) - Writes to a CSV file
df.to_excel(filename) - Writes to an Excel file
df.to_sql(table_name, connection_object) - Writes to a SQL table
df.to_json(filename) - Writes to a file in JSON format
df.to_html(filename) - Saves as an HTML table
df.to_clipboard() - Writes to the clipboard

CREATE TEST OBJECTS

Useful for testing

pd.DataFrame(np.random.rand(20,5)) - 5 columns and 20 rows of random floats
pd.Series(my_list) - Creates a series from an iterable **my_list**
df.index = pd.date_range('1900/1/30', periods=df.shape[0]) - Adds a date index

VIEWING/INSPECTING DATA

df.head(n) - First n rows of the DataFrame
df.tail(n) - Last n rows of the DataFrame
df.shape - Number of rows and columns
df.info() - Index, Datatype and Memory information
df.describe() - Summary statistics for numerical columns
s.value_counts(dropna=False) - Views unique values and counts
df.apply(pd.Series.value_counts) - Unique values and counts for all columns

SELECTION

df[col] - Returns column with label **col** as Series
df[[col1, col2]] - Returns Columns as a new DataFrame
s.iloc[0] - Selection by position
s.loc[0] - Selection by index
df.iloc[0, :] - First row
df.iloc[0,0] - First element of first column

DATA CLEANING

df.columns = ['a', 'b', 'c'] - Renames columns
pd.isnull() - Checks for null Values, Returns Boolean Array
pd.notnull() - Opposite of **s.isnull()**
df.dropna() - Drops all rows that contain null values
df.dropna(axis=1) - Drops all columns that contain null values
df.dropna(axis=1, thresh=n) - Drops all rows have have less than n non null values
df.fillna(x) - Replaces all null values with **x**
s.fillna(s.mean()) - Replaces all null values with the mean (mean can be replaced with almost any function from the statistics section)
s.astype(float) - Converts the datatype of the series to float
s.replace(1, 'one') - Replaces all values equal to 1 with 'one'
s.replace([1,3], ['one', 'three']) - Replaces all 1 with 'one' and 3 with 'three'
df.rename(columns=lambda x: x + 1) - Mass renaming of columns
df.rename(columns={'old_name': 'new_name'}) - Selective renaming
df.set_index('column_one') - Changes the index
df.rename(index=lambda x: x + 1) - Mass renaming of index

FILTER, SORT, & GROUPBY

df[df[col] > 0.5] - Rows where the **col** column is greater than 0.5
df[(df[col] > 0.5) & (df[col] < 0.7)] - Rows where 0.7 > col > 0.5
df.sort_values(col1) - Sorts values by **col1** in ascending order
df.sort_values(col2, ascending=False) - Sorts values by **col2** in descending order
df.sort_values([col1, col2], ascending=[True, False]) - Sorts values by

col1 in ascending order then **col2** in descending order

df.groupby(col1) - Returns a groupby object for values from one column
df.groupby([col1, col2]) - Returns a groupby object values from multiple columns
df.groupby(col1)[col2].mean() - Returns the mean of the values in **col2**, grouped by the values in **col1** (mean can be replaced with almost any function from the statistics section)
df.pivot_table(index=col1, values=[col2, col3], aggfunc=mean) - Creates a pivot table that groups by **col1** and calculates the mean of **col2** and **col3**
df.groupby(col1).agg(np.mean) - Finds the average across all columns for every unique column 1 group
df.apply(np.mean) - Applies a function across each column
df.apply(np.max, axis=1) - Applies a function across each row

JOIN/COMBINE

df1.append(df2) - Adds the rows in **df1** to the end of **df2** (columns should be identical)
pd.concat([df1, df2], axis=1) - Adds the columns in **df1** to the end of **df2** (rows should be identical)
df1.join(df2, on=col1, how='inner') - SQL-style joins the columns in **df1** with the columns on **df2** where the rows for **col1** have identical values. **how** can be one of 'left', 'right', 'outer', 'inner'

STATISTICS

These can all be applied to a series as well.
df.describe() - Summary statistics for numerical columns
df.mean() - Returns the mean of all columns
df.corr() - Returns the correlation between columns in a DataFrame
df.count() - Returns the number of non-null values in each DataFrame column
df.max() - Returns the highest value in each column
df.min() - Returns the lowest value in each column
df.median() - Returns the median of each column
df.std() - Returns the standard deviation of each column



Python For Data Science NumPy Cheat Sheet

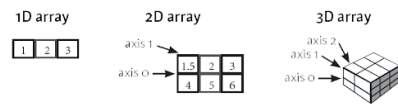
Learn NumPy online at www.DataCamp.com

NumPy

The NumPy library is the core library for scientific computing in Python. It provides a high-performance multidimensional array object, and tools for working with these arrays

Use the following import convention:
`>>> import numpy as np`

NumPy Arrays



> Creating Arrays

```
>>> a = np.array([1,2,3])
>>> b = np.array([(1,5,2,3), (4,5,6)], dtype = float)
>>> c = np.array([(1,5,2,3), (4,5,6)], [(3,2,1), (4,5,6)], dtype = float)
```

Initial Placeholders

```
>>> np.zeros(3,4) #Create an array of zeros
>>> np.ones(2,3,4, dtype=np.int16) #Create an array of ones
>>> d = np.arange(10,25,5) #Create an array of evenly spaced values (step value)
>>> np.linspace(0,2,9) #Create an array of evenly spaced values (number of samples)
>>> e = np.full((2,2),7) #Create a constant array
>>> f = np.eye(2) #Create a 2x2 identity matrix
>>> np.random.random((2,2)) #Create an array with random values
>>> np.empty((3,2)) #Create an empty array
```

> I/O

Saving & Loading On Disk

```
>>> np.save('my_array', a)
>>> np savez('array.npy', a, b)
>>> np.load('my_array.npy')
```

Saving & Loading Text Files

```
>>> np.loadtxt('myfile.txt')
>>> np.genfromtxt('my_file.csv', delimiter=',')
>>> np.savetxt('myarray.txt', a, delimiter='')
```

> Asking For Help

```
>>> np.info(np.ndarray.dtype)
```

> Inspecting Your Array

```
>>> a.shape #Array dimensions
>>> len(a) #length of array
>>> b.ndim #Number of array dimensions
>>> e.size #Number of array elements
>>> b.dtype #Data type of array elements
>>> b.dtype.name #Name of data type
>>> b.astype(int) #Convert on array to a different type
```

> Data Types

```
>>> np.int64 #Signed 64-bit integer types
>>> np.float32 #Standard double-precision floating point
>>> np.complex #Complex numbers represented by 128 floats
>>> np.bool #Boolean type storing TRUE and FALSE values
>>> np.object #Python object type
>>> np.string_ #Fixed-length string type
>>> np.unicode_ #Fixed-length unicode type
```

> Array Mathematics

Arithmetic Operations

```
>>> g = a - b #Subtraction
array([[ -0.5,  0. ,  0. ],
       [ -3. , -3. ,  3. ]])
>>> np.subtract(a,b) #Subtraction
>>> b + a #Addition
array([[ 2.5,  4. ,  6. ],
       [ 5. ,  7. ,  9. ]])
>>> np.add(b,a) #Addition
>>> a / b #Division
array([[ 0.66666667,  1. ,  1. ],
       [ 0.25 ,  0.4 ,  0.5 ]])
>>> np.divide(a,b) #Division
>>> a * b #Multiplication
array([[ 1.5,  4. ,  9. ],
       [  4. , 10. , 18. ]])
>>> np.multiply(a,b) #Multiplication
>>> np.exp(b) #Exponentiation
>>> np.sqrt(b) #Square root
>>> np.sin(a) #Print sines of an array
>>> np.cos(b) #Element-wise cosine
>>> np.log(a) #Element-wise natural logarithm
>>> e.dot(f) #Dot product
array([[ 7. ,  7.]])
```

Comparison

```
>>> a == b #Element-wise comparison
array([[False,  True,  True],
       [False, False, False]], dtype=bool)
>>> a < 2 #Element-wise comparison
array([True, False, False], dtype=bool)
>>> np.array_equal(a, b) #Array-wise comparison
```

Aggregate Functions

```
>>> a.sum() #Array-wise sum
>>> a.min() #Array-wise minimum value
>>> b.max(axis=0) #Maximum value of an array row
>>> b.cumsum(axis=1) #Cumulative sum of the elements
>>> a.mean() #Mean
>>> np.median(b) #Median
>>> np.corrcoef(a) #Correlation coefficient
>>> np.std(b) #Standard deviation
```

> Copying Arrays

```
>>> h = a.view() #Create a view of the array with the same data
>>> np.copy(a) #Create a copy of the array
>>> h = a.copy() #Create a deep copy of the array
```

> Sorting Arrays

```
>>> a.sort() #Sort on array
>>> c.sort(axis=0) #Sort the elements of an array's axis
```

> Subsetting, Slicing, Indexing

```
Subsetting
>>> a[2] #Select the element at the 2nd index
>>> b[1,2] #Select the element at row 1 column 2 (equivalent to b[[1][2]])
6.9

Slicing
>>> a[0:2] #Select items at index 0 and 1
array([1, 2])
>>> a[0:1] #Select items at rows 0 and 1 in column 1
array([ 2. ,  5. ])
>>> b[:,1] #Select all items at row 0 (equivalent to b[0:1, :])
array([[1.5,  2. ,  3. ]])
>>> c[1,...] #Same as [1, :, :]
array([[ 3. ,  2. ,  1. ],
       [ 4. ,  8. ,  6. ]])
>>> a1 = a[::-1] #Reversed array a array([3, 2, 1])

Boolean Indexing
>>> a[a<2] #Select elements from a less than 2
array([1])

Fancy Indexing
>>> b[[1, 0, 1, 0],[8, 1, 2, 0]] #Select elements (1,0),(0,1),(1,2) and (0,0)
array([[ 4. ,  2. ,  8. ,  1.5]])
>>> a[[0, 8, 1, 0]]#[0,1:0,2:0] #Select a subset of the matrix's rows and columns
array([[ 4. ,  5. ,  6. ,  4. ],
       [ 1.5,  2. ,  3. ,  1.5],
       [ 4. ,  8. ,  6. ,  4. ],
       [ 1.5,  2. ,  3. ,  1.5]])
```

> Array Manipulation

```
Transposing Array
>>> t = np.transpose(b) #Permute array dimensions
>>> t.T #Permute array dimensions

Changing Array Shape
>>> b.ravel() #Flatten the array
>>> a.reshape(b,2) #Reshape, but don't change data

Adding/Removing Elements
>>> h.resize((2,6)) #Return a new array with shape (2,6)
>>> np.append(a,g) #Append items to an array
>>> np.insert(a, 1, 5) #Insert items in an array
>>> np.delete(a,[1]) #Delete items from an array

Combining Arrays
>>> np.concatenate((a,b),axis=0) #Concatenate arrays
array([ 1. ,  2. ,  3. ,  10. ,  20])
>>> np.vstack((a,b)) #Stack arrays vertically (row-wise)
array([[ 1. ,  2. ,  3. ],
       [ 1.5,  2. ,  3. ],
       [ 4. ,  5. ,  6. ]])
>>> np.r_[a,f] #Stack arrays vertically (row-wise)
array([[ 7. ,  7. ,  1. ,  0. ],
       [ 7. ,  7. ,  0. ,  1.]])
>>> np.column_stack((a,b)) #Create stacked column-wise arrays
array([[ 1, 10],
       [ 2, 15],
       [ 3, 20]])
>>> np.c_[a,d] #Create stacked column-wise arrays

Splitting Arrays
>>> np.split(a,3) #Split the array horizontally at the 3rd index
[array([1]),array([2]),array([3])]
>>> np.vsplit(a,2) #Split the array vertically at the 2nd index
[array([[1,5,2, , 3, ],
       [ 4. ,  5. ,  6. ]]),
array([[ 3. ,  2. ,  1. ],
       [ 4. ,  8. ,  6.]])]
```

Learn Data Skills Online at www.DataCamp.com



Python For Data Science Matplotlib Cheat Sheet

Learn Matplotlib online at www.DataCamp.com

Matplotlib

Matplotlib is a Python 2D plotting library which produces publication-quality figures in a variety of hardcopy formats and interactive environments across platforms.

Prepare The Data

1D Data

```
>>> import numpy as np
>>> x = np.linspace(0, 10, 100)
>>> y = np.cos(x)
>>> z = np.sin(x)
```

2D Data or Images

```
>>> data = 2 * np.random.random((10, 10))
>>> data2 = 3 * np.random.random((10, 10))
>>> Y, X = np.mgrid[0:10, 0:10]
>>> U = -1 - X**2 + Y
>>> V = 1 + X - Y**2
>>> from matplotlib.cbook import get_sample_data
>>> img = np.load(get_sample_data("axes_grid/obliviate_normal.npy"))
```

Create Plot

```
>>> import matplotlib.pyplot as plt
```

Figure

```
>>> fig = plt.figure()
>>> fig2 = plt.figure(figsize=plt.figaspect(2.0))
```

Axes

All plotting is done with respect to an Axes. In most cases, a subplot will fit your needs. A subplot is an axes on a grid system.

```
>>> fig.add_axes()
>>> ax1 = fig.add_subplot(221) #row-col-num
>>> ax2 = fig.add_subplot(212)
>>> fig3, axes = plt.subplots(nrows=2, ncols=2)
>>> fig4, axes2 = plt.subplots(ncols=3)
```

Save Plot

```
>>> plt.savefig("fig.png") #Save figure
>>> plt.savefig("fig.png", transparent=True) #Save transparent figure
```

Show Plot

```
>>> plt.show()
```

Plotting Routines

1D Data

```
>>> fig, ax = plt.subplots()
>>> lines = ax.plot(x,y) #Draw points with lines or markers connecting them
>>> ax.scatter(x,y) #Draw unconnected points, colored or colored
>>> axes[0,0].barh([1,2,3],[5,4,5]) #Plot vertical rectangles (constant width)
>>> axes[1,0].barh([8.5,1,2.5],[0,1,2]) #Plot horizontal rectangles (constant height)
>>> axes[1,1].axhline(0.45) #Draw a horizontal line across axes
>>> axes[0,1].axvline(0.65) #Draw a vertical line across axes
>>> ax.fill(x,y,color='blue') #Draw filled polygons
>>> ax.fill_between(x,y,color='yellow') #Fill between y-values and 0
```

2D Data

```
>>> fig, ax = plt.subplots()
>>> im = ax.imshow(img, #Colormapped or RGB arrays
                cmap='gist_earth',
                interpolation='nearest',
                vmin=2,
                vmax=2)
>>> axes[0].pcolormesh(data2) #Pseudocolor plot of 2D array
>>> axes[0].pcolormesh(data) #Pseudocolor plot of 2D array
>>> CS = plt.contour(X,Y,U) #Plot contours
>>> axes[2].contourf(data) #Plot filled contours
>>> axes[2].ax.contour(CS) #Label a contour plot
```

Vector Fields

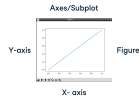
```
>>> axes[0,1].arrow(0,0,0.5,0.5) #Add an arrow to the axes
>>> axes[1,1].quiver(Y,X,U) #Plot a 2D field of arrows
>>> axes[0,1].streamplot(X,Y,U,V) #Plot a 2D field of arrows
```

Data Distributions

```
>>> ax1.hist(y) #Plot a histogram
>>> ax3.boxplot(y) #Make a box and whisker plot
>>> ax3.violinplot(x) #Make a violin plot
```

Plot Anatomy & Workflow

Plot Anatomy



Workflow

The basic steps to creating plots with matplotlib are:

- 1. Prepare Data
- 2. Create Plot
- 3. Plot Customized Plot
- 4. Save Plot
- 5. Show Plot

```
>>> import matplotlib.pyplot as plt
>>> x = [1,2,3,4] #Step 1
>>> y = [10,20,25,30]
>>> fig = plt.figure() #Step 2
>>> ax = fig.add_subplot(111) #Step 3
>>> ax.plot(x, y, color='lightblue', linewidth=3) #Step 3, 4
>>> ax.scatter([2,4,6],
             [15,15,25],
             color='darkgreen',
             marker='^')
>>> ax.set_xlim(1, 6.5)
>>> plt.savefig("fig.png") #Step 5
>>> plt.show() #Step 6
```

Close and Clear

```
>>> plt.cla() #Clear an axis
>>> plt.clf() #Clear the entire figure
>>> plt.close() #Close a window
```

Plotting Customize Plot

Colors, Color Bars & Color Maps

```
>>> plt.plot(x, x, x, x**2, x, x**3)
>>> ax.plot(x, y, alpha=0.4)
>>> ax.plot(x, y, c='r')
>>> fig.colorbar(m, orientation='horizontal')
>>> im = ax.imshow(img,
                 cmap='seismic')
```

Markers

```
>>> fig, ax = plt.subplots()
>>> ax.scatter(x,y,marker='^')
>>> ax.plot(x,y,marker='o')
```

Linestyles

```
>>> plt.plot(x,y,linewidth=4.0)
>>> plt.plot(x,y,ls='solid')
>>> plt.plot(x,y,ls='--')
>>> plt.plot(x,y,'--',x**2,y**2,...)
>>> plt.setp(lines,color='r',linewidth=4.0)
```

Text & Annotations

```
>>> ax.text(1,
         -2.1,
         'Example Graph',
         style='italic')
>>> ax.annotate("Sine",
              xy=(6, 0),
              xycoords='data',
              xytext=(10.5, 0),
              textcoords='data',
              arrowprops=dict(arrowstyle="->",
                              connectionstyle="arc3"),)
```

Mathtext

```
>>> plt.title("$\sigma_1=15\%$", fontsize=20)
```

Limits, Legends and Layouts

Limits & Autoscaling

```
>>> ax.margins(x=0,y=0.1) #Add padding to a plot
>>> ax.axis('equal') #Set the aspect ratio of the plot to 1
>>> ax.set_xlim(0,10.5), ylim[-1.5,1.5] #Set limits for x-and y-axis
>>> ax.set_xlim(0,10.5) #Set limits for x-axis
```

Legends

```
>>> ax.set(title="An Example Axes", #Set a title and x-and y-axis labels
         ylabel="Y-Axis",
         xlabel="X-Axis")
>>> ax.legend(loc='best') #No overlapping plot elements
```

Ticks

```
>>> ax.xaxis.set(ticks=range(1,5), #Manually set x-ticks
              ticklabels=[5,10,12,"foo"])
>>> ax.tick_params(axis='x', #Make y-ticks longer and go in and out
                 direction='inout',
                 length=8)
```

Subplot Spacing

```
>>> fig.subplots_adjust(wspace=0.5, #Adjust the spacing between subplots
                    hspace=1,
                    left=0.125,
                    right=0.9,
                    top=0,
                    bottom=0.1)
```

Axis Spines

```
>>> ax1.spines['top'].set_visible(False) #Make the top axis line for a plot invisible
>>> ax1.spines['bottom'].set_position(("outward",10)) #Move the bottom axis line outward
```



Learn Data Skills Online at www.DataCamp.com



Python For Data Science Pandas Basics Cheat Sheet

Learn Pandas Basics online at www.DataCamp.com

Pandas

The **Pandas** library is built on NumPy and provides easy-to-use **data structures** and **data analysis** tools for the Python programming language.

Use the following import convention:

```
>>> import pandas as pd
```

Pandas Data Structures

Series

A **one-dimensional** labeled array capable of holding any data type

a	3
b	-5
c	7
d	4

Index →

```
>>> s = pd.Series([3, -5, 7, 4], index=['a', 'b', 'c', 'd'])
```

Dataframe

A **two-dimensional** labeled data structure with columns of potentially different types

Country	Capital	Population
Belgium	Brussels	1190845
India	New Delhi	130379335
Brazil	Brasilia	207847528

```
>>> data = {'Country': ['Belgium', 'India', 'Brazil'],  
          'Capital': ['Brussels', 'New Delhi', 'Brasilia'],  
          'Population': [1190845, 130379335, 207847528]}  
>>> df = pd.DataFrame(data,  
                    columns=['Country', 'Capital', 'Population'])
```

Dropping

```
>>> s.drop(['a', 'c']) #Drop values from rows (axis=0)  
>>> df.drop('Country', axis=1) #Drop values from columns (axis=1)
```

Asking For Help

```
>>> help(pd.Series.loc)
```

Sort & Rank

```
>>> df.sort_index() #Sort by labels along an axis  
>>> df.sort_values(by='Country') #Sort by the values along an axis  
>>> df.rank() #Assign ranks to entries
```

I/O

Read and Write to CSV

```
>>> pd.read_csv('file.csv', header=None, nrows=5)  
>>> df.to_csv('myDataFrame.csv')
```

Read and Write to Excel

```
>>> pd.read_excel('file.xlsx')  
>>> df.to_excel('dir/myDataFrame.xlsx', sheet_name='Sheet1')
```

Read multiple sheets from the same file

```
>>> xlsx = pd.ExcelFile('file.xls')  
>>> df = pd.read_excel(xlsx, 'Sheet1')
```

Read and Write to SQL Query or Database Table

```
>>> from sqlalchemy import create_engine  
engine = create_engine('sqlite://:memory:')  
>>> pd.read_sql("SELECT * FROM my_table", engine)  
>>> pd.read_sql_table('my_table', engine)  
>>> pd.read_sql_query("SELECT * FROM my_table", engine)  
  
read_sql() is a convenience wrapper around read_sql_table() and read_sql_query()  
>>> df.to_sql('myDF', engine)
```

Selection

Also see NumPy Arrays

Getting

```
>>> s['b'] #Get one element  
-5  
>>> df[1] #Get subset of a DataFrame  
Country Capital Population  
1 India New Delhi 1383171835  
2 Brazil Brasilia 207847528
```

Selecting, Boolean Indexing & Setting

```
By Position  
>>> df.iloc[[0],[0]] #Select single value by row & column  
'Belgium'  
>>> df.iat[[0],[0]]  
'Belgium'  
  
By Label  
>>> df.loc[[0], ['Country']] #Select single value by row & column labels  
'Belgium'  
>>> df.at[[0], ['Country']]  
'Belgium'  
  
By Label/Position  
>>> df.ix[[2]] #Select single row of subset of rows  
Country Brazil  
Capital Brasilia  
Population 207847528  
>>> df.ix[:, 'Capital'] #Select a single column of subset of columns  
0 Brussels  
1 New Delhi  
2 Brasilia  
>>> df.ix[:, 'Capital'] #Select rows and columns  
'New Delhi'  
  
Boolean Indexing  
>>> s[s > 5] #Series s where value is not > 1  
>>> s[s < -2] | (s > 2) #Series s where value is < -2 or > 2  
>>> df[df['Population'] > 120000000] #Use filter to adjust DataFrame  
  
Setting  
>>> s['a'] = 6 #Set index a of Series s to 6
```

Retrieving Series/DataFrame Information

Basic Information

```
>>> df.shape #rows, columns  
>>> df.index #describe index  
>>> df.columns #describe DataFrame columns  
>>> df.info() #Info on DataFrame  
>>> df.count() #Number of non-NA values
```

Summary

```
>>> df.sum() #Sum of values  
>>> df.mean() #Cumulative sum of values  
>>> df.min()/df.max() #Minimum/Maximum values  
>>> df.idxmin()/df.idxmax() #Minimum/Maximum index value  
>>> df.describe() #Summary statistics  
>>> df.mean() #Mean of values  
>>> df.median() #Median of values
```

Applying Functions

```
>>> f = lambda x: x*2  
>>> df.apply(f) #Apply function  
>>> df.applymap(f) #Apply function element-wise
```

Data Alignment

Internal Data Alignment

NA values are introduced in the indices that don't overlap

```
>>> s3 = pd.Series([7, -2, 3], index=['a', 'c', 'd'])  
>>> s = s3  
a 10.0  
b NaN  
c 5.0  
d 7.0
```

Arithmetic Operations with Fill Methods

You can also do the internal data alignment yourself with the help of the fill methods:

```
>>> s.add(s3, fill_value=8)  
a 10.0  
b -5.0  
c 5.0  
d 7.0  
>>> s.sub(s3, fill_value=2)  
>>> s.div(s3, fill_value=4)  
>>> s.mul(s3, fill_value=3)
```

Learn Data Skills Online at
www.DataCamp.com

